



**8.3 Appendix 3: Reference 4:** Persoonsgegevens **“Belastingdienst IH Risk Model: Code and Data Storage,” /IH selectie/Project Managment/Quality/Project Initiation Quality Documentation/.**

**IH Risk Model  
Code and Data Storage Structure  
Version: 1.0**



**Table of Contents**

- 1. AWS+ PROJECT SPACE ..... 88**
- 1.1 MAIN STRUCTURE ..... 88
- 1.1.1 *Utilities* ..... 89
- 1.2 /CODE ..... 90
- 1.2.1 *General Guidelines* ..... 90
- 1.2.2 *Breakdown by Type* ..... 90
- 1.2.3 */CODE Subject Areas* ..... 90
- 1.2.4 */CODE Structure of each Subject Area* ..... 91
- 1.3 /DATA ..... 92
- 1.3.1 *Breakdown by Type* ..... 92
- 1.3.2 */DATA Subject Areas* ..... 92
- 1.3.3 */DATA Structure of each Subject Area* ..... 93
- 2. AWS+ RAW DATA ..... 94**
- 3. TERADATA ..... 94**
- 3.1 READ-ONLY RAW DATA ..... 94
- 3.2 WORKSPACE ..... 94
- 3.3 TEMPORARY WORKSPACE ..... 94



## 1 AWS+ Project Space

Our shared project workspace is at:

/home/<your user id>/AD005/shared/analyse/PMI/IHriskmodel

<b>../AD005/shared/analyse/PMI/IHriskmodel</b>	
../CODE	Contains all the project's SAS code and Enterprise Guide projects
../DATA	Contains the project's data excluding the raw data and the data stored on Teradata
../METADATA	Contains the project's metadata i.e. look-up tables
../MODEL	Contains the Enterprise Miner projects that build the model(s)
../PERSONAL	Personal 'sandbox' areas

### 1.1 Main Structure

<b>../CODE</b>	
../CODE/aangifte	Code that builds features from the aangifte
../CODE/aanslag	Code that builds features/target from the aanslag
../CODE/toetsing	Code that analyses the audit data
../CODE/ABT	Code to build the main Analytics Base Table (ABT) by bringing all the features together
../CODE/utilities	Useful utility or general code

<b>../DATA</b>	
../DATA/aangifte	Data from the aangifte
../DATA/aanslag	Data from the aanslag
../DATA/toetsing	Data from the audit
../DATA/ABT	Data to do with the main ABT

<b>../METADATA</b>	
../METADATA	LIBNAME Metadata "home/&user_id./AD005/shared/analyse/PMI/IHriskmodel/METADATA";

Where

%LET user\_id =  /\* use your id \*/



../MODEL	
../LabPhase	The Enterprise Miner project(s) used to build the lab phase model.
../PilotPhase1	The Enterprise Miner project(s) used to build the pilot phase model.
../PilotPhase2	The Enterprise Miner project(s) used to build the second pilot phase model.

../PERSONAL	
Persoonsgegevens	

Use these personal areas for whatever you like. They're a 'sandbox' area for experimentation.

Note that no code or data needed to build your subject area datasets should be in here.

1.1.1 Utilities

/CODE/utilities will contain general bits of code that everyone needs access to.

What's in here will grow over time but initially I'll put in the macros that we'll use to do the initial data exploration.

As we get to know the data we'll find that there are tasks that we're doing again and again on different subject areas. When that happens we'll wrap the code for that up in a macro and store it here so that we can leverage each other's work.



## 1.2 /CODE

### 1.2.1 *General Guidelines*

You can use the 'point-and-click' interface in Enterprise Guide for initial data exploration and so on.

However, everything used to build the ABT **must be written in SAS code**, and the **code must be saved as a .sas file**.

Under no circumstances save code in the project as embedded code, as these are impossible to restore in the event of corruption and even moving/renaming the project can cause code to be lost.

In general, if you have code files that have to be run in a certain order then include that order in the naming.

Example:

```
1_NW_income_set_up.sas
2_NW_income_extract.sas
3_NW_income_transpose.sas
4_NW_income_clean.sas
etc.
```

### 1.2.2 *Breakdown by Type*

This is currently defined just for the aangifte NW data. The other source areas will follow a similar pattern.

../CODE/aangifte	
../CODE/aangifte/NW	Neit Winst – not profit aka individuals
../CODE/aangifte/W	Winst – profit aka companies

## 1.3 /CODE Subject Areas

../CODE/aangifte/NW	
../CODE/aangifte/NW/ABT	Code that brings all the aangifte features together
../CODE/aangifte/NW/assets	Code that builds features from the assets, savings and investments subject area
../CODE/aangifte/NW/benefits	Code that builds features from benefit (pension, sickness etc.) subject area
../CODE/aangifte/NW/debt	Code that builds features from debt subject area
../CODE/aangifte/NW/expenses	Code that builds features from expenses and tax allowances subject area
../CODE/aangifte/NW/foreign	Code that builds features from the foreign interests



	subject area
../CODE/aangifte/NW/partner	Code that builds features from the partner subject area
../CODE/aangifte/NW/income	Code that builds features from the personal income subject area
../CODE/aangifte/NW/property	Code that builds features from the property subject area
../CODE/aangifte/NW/credits	Code that builds features from the tax credits subject area
../CODE/aangifte/NW/person	Code that builds features from the person (i.e. unique ID, age etc.) subject area

### 1.3.1 CODE Structure of each Subject Area

../CODE/aangifte/NW/<area>	
../CODE/aangifte/NW/<area>/build	Code that's ready to be used to build the assets features
../CODE/aangifte/NW/<area>/dev	Code that's still in development
../CODE/aangifte/NW/<area>/projects	Enterprise Guide projects
../CODE/aangifte/NW/<area>/archive	Code not currently in use



## 1.4 /DATA

As we figure out how to use the Teradata server some/all of this may move from AWS+ to Teradata.

### 1.4.1 Breakdown by Type

This is currently defined just for the aangifte NW data. The other source areas will follow a similar pattern.

../DATA/aangifte	
../DATA/aangifte/NW	Neit Winst – not profit aka individuals
../DATA/aangifte/W	Winst – profit aka companies

### 1.4.2 /DATA Subject Areas

These LIBNAMEs will be used by other parts of the project that need to reference the subject area data.

**Note:**

`&LET user_id = ; /* use your own id */`

../DATA/aangifte/NW	
../DATA/aangifte/NW/aangifte_ABT	LIBNAME <b>NWaangif</b> "home/&user_id./AD005/shared/analyse/PMI/IHriskmodel/DATA/aangifte/NW/aangifte_ABT/final";
../DATA/aangifte/NW/assets	LIBNAME <b>NWassets</b> "home/&user_id./AD005/shared/analyse/PMI/IHriskmodel/DATA/aangifte/NW/assets/final";
../DATA/aangifte/NW/benefits	LIBNAME <b>NWben</b> "home/&user_id./AD005/shared/analyse/PMI/IHriskmodel/DATA/aangifte/NW/benefits/final";
../DATA/aangifte/NW/debt	LIBNAME <b>NWdebt</b> "home/&user_id./AD005/shared/analyse/PMI/IHriskmodel/DATA/aangifte/NW/debt/final";
../DATA/aangifte/NW/expenses	LIBNAME <b>NWexpen</b> "home/&user_id./AD005/shared/analyse/PMI/IHriskmodel/DATA/aangifte/NW/expenses/final";
../DATA/aangifte/NW/foreign	LIBNAME <b>NWforeig</b> "home/&user_id./AD005/shared/analyse/PMI/IHriskmodel/DATA/aangifte/NW/foreign/final";
../DATA/aangifte/NW/partner	LIBNAME <b>NWpart</b> "home/&user_id./AD005/shared/analyse/PMI/IHriskmodel/DATA/aangifte/NW/partner/final";
../DATA/aangifte/NW/income	LIBNAME <b>NWincome</b> "home/&user_id./AD005/shared/analyse/PMI/IHriskmodel/DATA/aangifte/NW/income/final";
../DATA/aangifte	LIBNAME <b>NWprop</b> "home/&user_id./AD005/shared/analyse/PMI/IHriskmodel/DATA/aangifte/NW/prop/final";



/NW/property	TA/aangifte/NW/property/final";
../DATA/aangifte/NW/credits	LIBNAME <b>NWcredit</b> "home/&user_id./AD005/shared/analyse/PMI/IHriskmodel/DATA/aangifte/NW/credits/final";
../DATA/aangifte/NW/person	LIBNAME <b>NWperson</b> "home/&user_id./AD005/shared/analyse/PMI/IHriskmodel/DATA/aangifte/NW/person/final";

### 1.4.3/DATA Structure of each Subject Area

../DATA/aangifte/NW/<area>	
../DATA/aangifte/NW/<area>/input	The input data you're using, extracted from the raw data.
../DATA/aangifte/NW/<area>/tmp	Temporary tables, in theory anything in here could be deleted
../DATA/aangifte/NW/<area>/intermediate	Intermediate tables that should be kept
../DATA/aangifte/NW/<area>/final	Final table(s), one row per id
../DATA/aangifte/NW/<area>/EDA	Data used/generated in the exploratory data analysis (EDA) stage

#### LIBNAMES internal to each area

You use these when you are working within a subject area. Make sure these definitions are at the top of your code.

```
%LET user_id =  /* use your own id */
%LET area = income; /* for example */
```

```
LIBNAME input
"home/&user_id./AD005/shared/analyse/PMI/IHriskmodel/DATA/aangifte/NW/&area
./input";
```

```
LIBNAME tmp
"home/&user_id./AD005/shared/analyse/PMI/IHriskmodel/DATA/aangifte/NW/&area
./tmp";
```

```
LIBNAME inter
"home/&user_id./AD005/shared/analyse/PMI/IHriskmodel/DATA/aangifte/NW/&area
./intermediate";
```

```
LIBNAME final
"home/&user_id./AD005/shared/analyse/PMI/IHriskmodel/DATA/aangifte/NW/&area
./final";
```





## 2 AWS+ Raw Data

Read-only raw data. The aanslag and the toetsing (audit) raw data are stored here.

Access them using:

```
%LET user_id =  /* use your own id */  
  
LIBNAME raw  
"home/&user_id./AD005/shared/datasets/BI-A_data/IH/IHriskmodel";
```

## 3 Teradata

These are the LIBNAMEs to access the Teradata space.

```
%LET user_id =  /* use your own id */
```

### 3.1 Read-Only Raw Data

The aangifte data is stored here.

```
LIBNAME TD_raw TERADATA  
bulkload=yes server='tdata01' schema='DL_INCASSO_PILOT_IN'  
user="&user_id" password=
```

### 3.2 Workspace

The area we can write to. More detail will be added once I've developed the schemas.

```
LIBNAME TD_write TERADATA  
bulkload=yes server='tdata01' schema='DL_INCASSO_PILOT_OUT'  
user="&user_id" password=
```

### 3.3 Temporary Workspace

This will be wiped at the end of each work day. More detail will be added once I've developed the schemas.

```
LIBNAME TD_tmp TERADATA  
bulkload=yes server='tdata01' schema='DL_INCASSO_PILOT_WRK'  
user="&user_id" password=
```



**Revision History**

Date	Version	Description	Author
21/03/2014	1.0	Initial version	Persoonsgegevens



**8.4 Appendix 4: Reference 5:** Persoonsgegevens "Belastingdienst IH Risk Model: Standard Methods ver 1," /IH selectie/Documentatie/Software process/.

**Belastingdienst IH Risk Model  
Standard Methods**

**NOTE:**

**HANDLING NEW / RETIRING OLD Metacodes:**

**NEW – TWO Year Change management process:**

Year 1: Collect, analyse, use metric from "this year"

Year 2: Compare to previous year (i.e. year 1)

**OLD – One year change process – retire code**

i.e. don't do if no-one had it in previous years and don't do if no-one has it this year

**1 Analysing a numeric metacode with volgnummers (sequence numbers)**

*example: total loon (salary) income.*

o *Current Year:*

**COUNT:** number of non-zero, non-missing volgnummers (NB: numeric fields only)

**SUM:** contents over volgnummers

**IF:** there is a precalculated total field then COMPARE SUM to COUNT

**THEN SET:** flag to 1 if different OR 0 if totals are equal

o *V-A comparison*

*We want to handle both the count of volgnummers and also the values in the volgnummers. Example: We want to know how many salaries an individual has (3 salaries) and also the total monetary value of those salaries (35.000).*

*Create/Return the following:*

Ratio: (V/A)

- 1 if V sum = A sum
- < 1 if V < A sum
- > 1 V sum > A sum
- -1 if Divide by 0 (example Salary 1000 last year, 0 this year)



## Absolute Change (V-A)

- 0 if V count = A count
- +ve if V count > A count
- -ve if V count < A count

## V Value

## Flags:

- Compare count V to count A
  - i. Binary flag: 0 if counts are the same, 1 otherwise
- Compare sum V to sum A
  - i. Binary flag: 0 if counts are the same, 1 otherwise

○ *Prior year comparison:*

Do not do any prior year comparisons if this is a brand new metacode this year or if prior year was the final year for this metacode – see notes at start of document.

Create the following per Metacode:

Ratio: (Current Year /Previous Year)

- 1 if this year sum = prev year sum
- < 1 if this year sum < prev year sum
- > 1 if this year sum > prev year sum

## Absolute Change (Current Year – Previous Year)

- 0 if this year count = prev year count
- +ve if this year count > prev year count
- -ve if this year count < prev year count

## Flags:

- **New** (to this sofi) metacode
  - 0 if sofi had this metacode this year and last year
  - 1 if sofi has it this year, but didn't last year
- **Dropped** (from this sofi) metacode
  - 0 if sofi had this metacode this year and last year
  - 1 if sofi had it last year, but doesn't this year
- **Compare** count of metacodes used this year to count prior year



- **Binary flag: 0 if counts are the same, 1 otherwise**

## 2 Character metacode

**Example:** *list of houses*

**NOTE: In our data all character metacodes are 0 if present as we don't have the actual character data.**

*Current year:*

- Count how many non-missing volnummers there are
- V-A comparison:
  - Compare count V to count A
    - Binary flag: 0 if counts are the same, 1 otherwise
- Compare sum V to sum A
  - Binary flag: 0 if counts are the same, 1 otherwise

○ *Prior year comparison:*

Do not do any prior year comparisons if this is a brand new metacode this year or if prior year was the last year for this metacode

- New (to this sofi) metacode
  - 0 if sofi had this metacode this year and last year
  - 1 if sofi has it this year, but didn't last year
- Dropped (from this sofi) metacode
  - 0 if sofi had this metacode this year and last year
  - 1 if sofi had it last year, but doesn't this year
- Compare count this year to count prior year
  - Binary flag: 0 if counts are the same, 1 otherwise
- Absolute difference in the counts
  - this year – prev year
  - 0 if this year count = prev year count
  - +ve if this year count > prev year count
  - -ve if this year count < prev year count

## 3 Date Metacode

**Example:** *date of birth of dependants*

**DoB - Available for all children.**

- Current year:



- Happened this year
- > 1 if event happened this tax year, 0 otherwise
  - Valid whole tax year
- > 1 if event happened before the start of the tax year (and has not ended)
  
- V-A comparison:
  - Compare V date to A date
- > Binary flag: 0 if dates are the same, 1 otherwise
  
- Prior year comparison:
  - Do not do any prior year comparisons if this is a brand new metacode this year or if prior year was the last year for this metacode
  - New (to this sofi) metacode
    - 0 if sofi has metacode this year and had metacode last year
    - 1 if sofi has it this year, but didn't last year
  - Dropped (from this sofi) metacode
    - 0 if sofi had this metacode this year and last year
    - 1 if sofi had it last year, but doesn't this year
  - Compare date this year to date prior year
    - ??

#### 4 Analysing related metacodes with sequence numbers

Tao has a method to compare related metacodes – such as a list of employments and the salaries from those employments.

Need to take a look at that – Ideas of metrics

#### 5 Complexity measures

For each subject area

- Number of metacodes used
  - i.e. number of questions answered
- % of the subject area covered
  - i.e. number of questions answered / number of questions in the subject area
- Max sequence number used
- Number of metacodes new to this aangifte when compared to last year
- Number of metacodes dropped in this aangifte when compared to last year
- Fiscal partnerships!!!! Important
- New Sofinum - more chance of mistake



## 6 Consistency measures

This should be calculated per subject area

- Number of times the V value differs from the A value
  - Common sets of metacodes used together in different sub topics– i.e. 25 and 107 always appear together – if they do not then flag – Needs business understanding – To be addressed.
  - Count of flags for current & prior year comparison. Increase in # of new/drop flags then less consistent
  - Fiscal partnerships – Need to explore how this will be integrated
  - New Sofinum - need to make an assumption regarding their handling
- **Other**
- VIA Data – How do we integrate this?

### *Document Revision History*

Date	Version	Description	Author
19/05/14	1.0	Initial document created	Persoonsgegevens



**8.5 Appendix 5: Reference 6:** Persoonsgegevens **“Belastingdienst IH Risk Model: Business Requirements ver 0.7,” /IH selectie/Project Managment/Quality/.**

**Belastingdienst IH Risk Model  
Business Requirements  
Version: 0.7**

**DRAFT**





**Table of Contents**

- 1. OVERVIEW ..... 105**
  - 1.1 PURPOSE OF THIS DOCUMENT ..... 6
  - 1.2 CURRENT STATE..... 105
  - 1.3 BASELINES FOR LAB PHASE 1 ..... 106
  - 1.4 BASELINES FOR LAB PHASE 2 ..... 107
  - 1.5 SOLUTION DESIGN ..... 108
    - 1.5.1 *Types of Risk* ..... 108
    - 1.5.2 *Types of Taxpayer* ..... 108
    - 1.5.3 *Refining the Model* ..... 108
    - 1.5.4 *Value Constraint*..... 109
  - 1.6 EXPECTED OUTPUT..... 110
  - 1.7 DESCRIPTION OF THE DATA TO BE USED ..... 110
  - 1.8 KEY LINKS ..... 110
- 2. LAB PHASE 1..... 111**
  - 2.1 DESCRIPTION ..... 111
  - 2.2 OUTCOME ..... 111
  - 2.3 DESCRIPTION OF THE DATA TO BE USED ..... 111
  - 2.4 DATE RANGE ..... 111
  - 2.5 TIMELINE ..... 111
- 3. LAB PHASE 2..... 112**
  - 3.1 DESCRIPTION ..... 112
  - 3.2 OUTCOME ..... 113
  - 3.3 DESCRIPTION OF THE DATA TO BE USED ..... 113
    - 3.3.1 *Whole Aangifte Model* ..... 113
    - 3.3.2 *Whole Aangifte Target*..... 113
    - 3.3.3 *Subject Models*..... 114
    - 3.3.4 *Subject Model Targets* ..... 114
  - 3.4 DATE RANGE ..... 114
  - 3.5 TIMELINE ..... 114
- 4. PILOT PHASE..... 115**
  - 4.1 DESCRIPTION ..... 115
  - 4.2 VALIDATION ..... 115
  - 4.3 OUTCOME ..... 115
  - 4.4 DESCRIPTION OF THE DATA TO BE USED ..... 115
  - 4.5 DATE RANGE ..... 116



4.6 TIMELINE ..... 116

**5. APPENDIX A: DETAILED BASELINE RESULTS, LAB PHASE 1 ..... 117**

5.1 ALL AUDITED AANGIFTEN FROM 2010 AND 2011 ..... 118

    5.1.1 *Counts* ..... 118

    5.1.2 *Percentages* ..... 119

5.2 AUDITED AANGIFTEN FROM 2010 AND 2011 SELECTED BY RISK RULES ..... 120

    5.2.1 *Counts* ..... 120

    5.2.2 *Percentages* ..... 121

5.3 AUDITED AANGIFTEN FROM 2010 AND 2011 SELECTED USING RANDOM SELECTION.. 123

    5.3.1 *Counts* ..... 123

    5.3.2 *Percentages* ..... 123

    5.3.3 *Addendum* ..... 123

5.4 AUDIT COUNTS BY TAX YEAR ..... 124

**6. APPENDIX B: DETAIL FOR THE BASELINE, LAB PHASE 2 ..... 124**

**Revision History**

Date	Version	Description	Author
13/03/14	0.1	High level Business Understanding	Persoonsgegevens
	0.2	Comments from <input type="text" value="Persoonsgegevens"/>	
27/3/14	0.3	Added baseline counts and notes on a value constraint	
15/4/14	0.4	Recalculated baselines. Some edits to text.	
23/4/14	0.5	Baselines recalculated again	
27/5/14	0.6	Fiscal partner added	
28/8/14	0.7	Updated after August programme board	

# 1 Overview

## 1.1 Purpose of this document

This is the technical requirements document for the 'IH Risk Model' analytics project and forms part of the Accenture Analytics quality control process. There are related business needs and objectives that are stated elsewhere.

This is a 'live' document with later stages being more closely defined as the project progresses. As such it will remain a work in progress until the start of the last phase of the project.

## 1.2 Current State

About 12 million aangiften are returned yearly. The IH selection module selects approximately 1,1 million aangiften for audit. Of these about 1,1 million are randomly selected and the rest are selected on the basis of business rules.

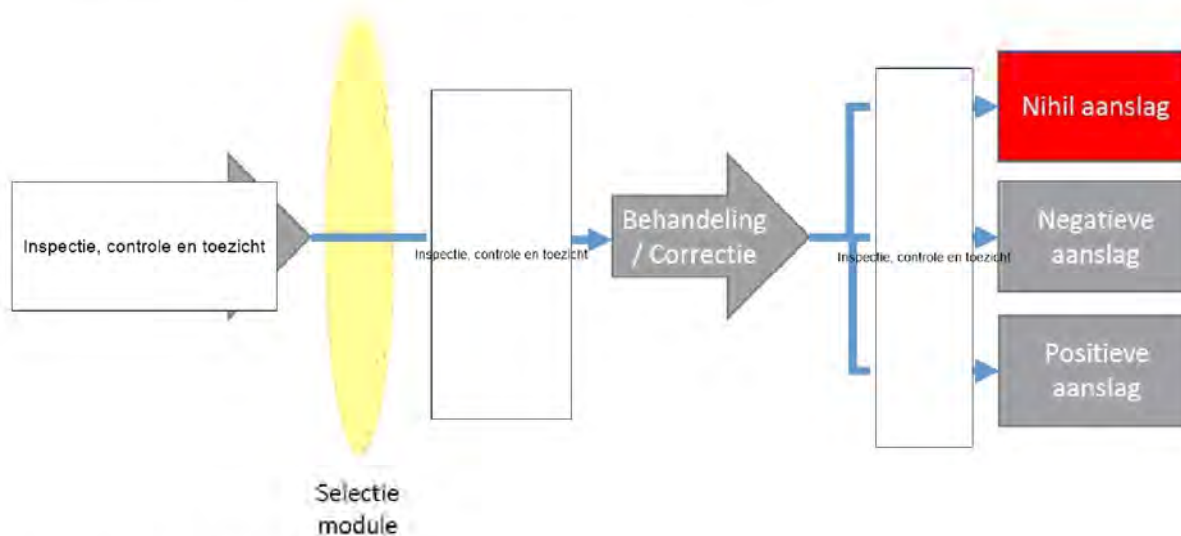


Figure 14: IH Selectie Module



### 1.3 Baselines for Lab Phase 1

Counts representing baseline figures for counts and percentages of positive, negative and nihil aanslag are shown in the following two tables. These baselines are currently just given for 'Niet Winst'<sup>28</sup> aangiften as the first phases of this project will concentrate on these.

**Table 26: Niet Winst Selection and Audit Results - Part 1**

Results from 2010 and 2011	Number Selected	Number With Completed Audits	Number with Positive Aanslag	Number with Negative Aanslag	Number with Nihil Aanslag
Risk rated	1.529.326	1.457.916	Inspectie, controle en toezicht		
Random sample	37.666	37.666			

**Table 27: Niet Winst Selection and Audit Results - Part 2**

Results from 2010 and 2011	% of Completed Audits with Positive Aanslag	% of Completed Audits with Negative Aanslag	% of Completed Audits with Nihil Aanslag
Risk rated	Inspectie, controle en toezicht		
Random sample			

The results of the audit come from the SAS dataset

`aang_vord_aans_totaal.sas7bdat`.

The indication that a sofi number is in the random sample comes from the SAS datasets

`sp2010.sas7bdat` and `sp2011.sas7bdat`.

These files are stored in the AWS+ directory `/AD005/shared/datasets/BI-A_data/IH/IHriskmodel`.

Counts are as of 27th March 2014. Note that here:

- Positive Aanslag means that after audit the taxpayer is found to owe more tax than originally declared in the aangifte.
- Negative Aanslag means that after audit the taxpayer is found to owe less tax than originally declared in the aangifte.

<sup>28</sup> aka 'Non Profit' aka individuals

<sup>29</sup> Apparent failure to add up to 100% is due to rounding. Results to a finer level of detail are available on request.



- Nihil Aanslag means that after audit there is no change in the amount of tax the taxpayer owes. This means that either the declaration of amount of tax owed was exactly correct, or that any change is within limits and is therefore not actioned.

The full breakdown of audit results is given in Appendix A: Detailed Baseline Results.

### 1.4 Baselines for Lab Phase 2

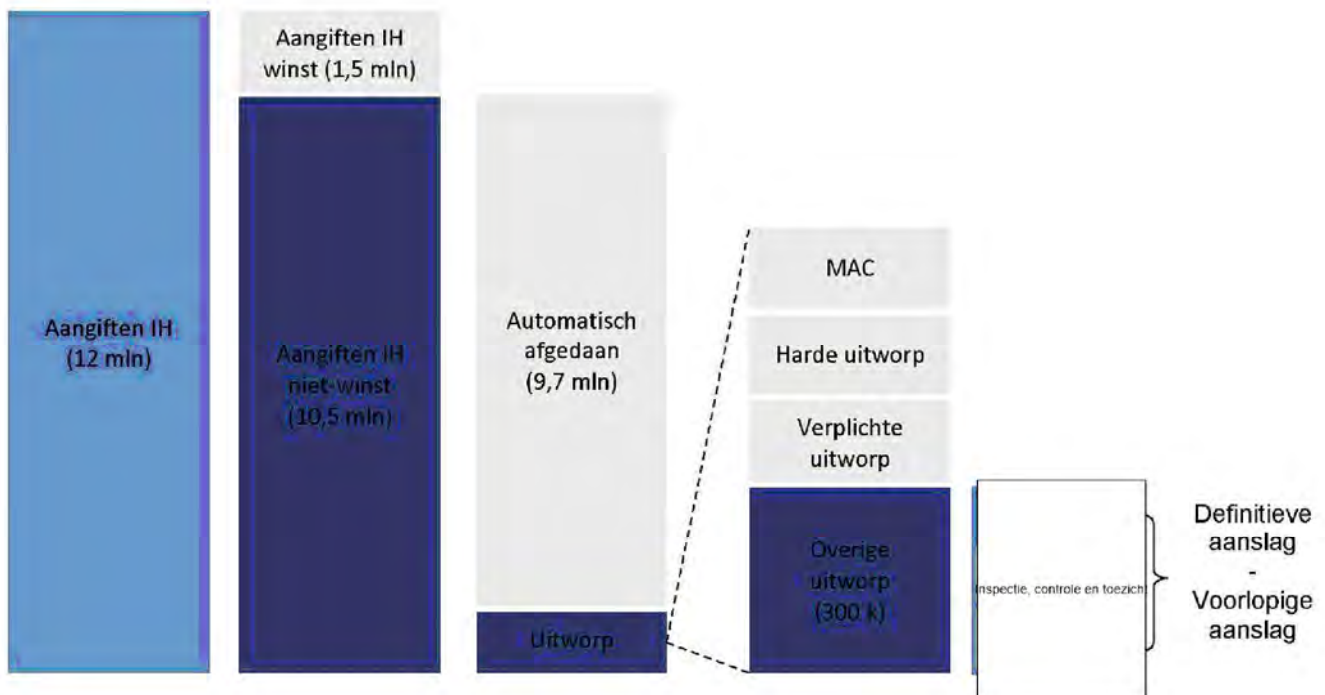
After Lab Phase 1 was completed the baseline was refined.

The UWB (risk rules) break into several parts:

- MAC – an automatic correction process where data in the aangifte is replaced with the equivalent from the contra data.
- Harde uitworp - these are manually audited break down into
  - rules that are being developed to become MAC in future years
  - political rules such as 'select everyone who claimed young person's disability allowance'
- Verplichte uitworp - the rules that select the random sample and similar categories
- Risk UWB

MAC is an efficient process, in only 4% of the aangifte selected by MAC rules is there no correction. Since MAC is so efficient type (2b) rules, those rules that are used to develop and test new MAC rules, need to be retained. The 'political rules' need to be retained to allow Belastingdienst to focus on specific categories of aangifte. Similarly the verplichte uitworp also need to be retained.

Hence it is only the last of these categories, type (4) that we are benchmarking ourselves against. For these the positive correction rate is 30% and this is the baseline that we aim to improve on.





## 1.5 Solution Design

This is the overall solution design. The individual phases that make up the project are defined later on in this document.

The project will build a predictive model to assess aangiften for risk.

We aim to increase the percentage of positive aanslagen by  percentage points over the current IH selection module baseline.

In addition, the project will build a modelling methodology and team and support structure that allows for a cyclical model development cycle.

This cycle will enable the model to be enhanced by the inclusion of more data sources as they become available and will allow the model to be maintained, refreshed and rebuilt over time.

The 'fab' stages are fully defined at this point but later stages may altered as discussions and planning continues.

### 1.5.1 Types of Risk

Initially we will be modelling the risk of a positive aangifte. That is, the taxpayer owes more tax than was originally declared.

The cycle will allow the exact definition of risk to be refined or for the creation of multiple models assessing for these different types of risk. Examples include:

- risk of a negative aanslag
- risk of incomplete payment
- risk of late payment
- risk of disputes.

### 1.5.2 Types of Taxpayer

Initially we will concentrate on 'niet winst' aangiften, that is aangiften from individuals rather than businesses (winst).

However the modelling cycle will allow the creation of multiple models looking at different types of taxpayer, such as:

- individuals or private persons ('niet winst')
- the self-employed
- small businesses
- small to medium enterprises
- larger enterprises.

It is not envisaged that we will model large multinational corporations as there are usually too few of them, and their accountancy and tax processes are too complicated, for them to be efficiently modeled by a predictive model of the type discussed here.

The exact timing of the development of the 'Winst' models is still under discussion.

### 1.5.3 Refining the Model



An aangifte that has been selected by the risk rules is not audited as a whole<sup>30</sup>, rather only the sections indicated by the rule are checked. This means that it is necessary to generate a related set of risk models that run in parallel, each aimed at different parts of the aangifte.

#### 1.5.4 Value Constraint

The plan as it stands does not include a 'value' constraint. This means that it is possible that the cases selected by the risk model(s) may be predominately very low value.

A value constraint could be created by the building of one or more secondary value predictive models to be run after the risk model(s). These value models would predict the chance that the value of the positive aanslag is greater than a cut-off threshold.

The cases forwarded for audit by the models would then be those that:

- were predicted to have a high chance of a positive correction
- and
- were predicted to have a high chance of the values of the correction to be greater than a threshold.

---

<sup>30</sup> Aangiften selected as part of the random sample are fully checked.



## 1.6 Expected Output

There are interim stages each with their own end point (see sections 2, 3, 4 and **Fout! Verwijzingsbron niet gevonden.**) but the desired result of this process is a suite of SAS processes that will be used to asses aangiften for risk and return the risk scores to the 'live' Belastingdienst system.

These processes will run in the SAS environment and so will need:

- a data feed from ABS to bring the data from the aangifte into SAS
- a data feed from SAS back into ABS to return the risk rating(s).

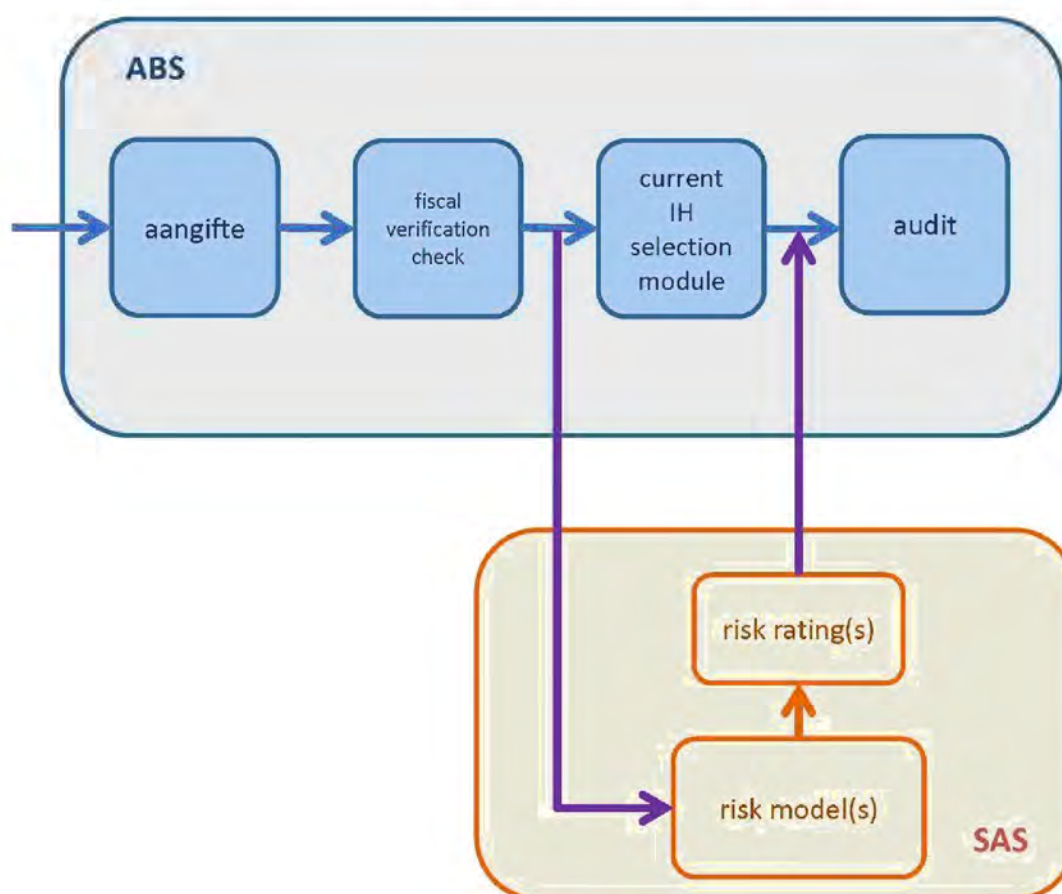


Figure 15: the 'live' system

## 1.7 Description of the Data to be Used

This will depend on the phase of the project. For the initial phases the data will come from ABS.

## 1.8 Key Links

Data will be linked on the basis of the Dutch social security number, known as Sofi number, to give a per person view of the data.



## 2 Lab Phase 1

### 2.1 Description

The lab phase will build a preliminary risk model, looking at a limited selection of audited aangiften and using a limited set of data sources.

This process will be used to mobilise the team and ensure that all data and other processes are functioning as needed.

### 2.2 Outcome

We will model the risk of a positive aanslag. That is, the risk that the taxpayer owes more tax than declared on the aangifte.

### 2.3 Description of the Data to be Used

We will use aangiften that were part of the random selection process in 2010 and 2011. We need two years as there are not enough cases in one year.

The actual data we will build the model on is:

1. the data in the aangifte
2. the result of the fiscal verification check (FVC), which is an automatic consistency check performed on the aangifte.
3. the results of the audit as shown in the definitive aanslag. This will be used to form the target.

In addition, **we will only model private persons** (aka 'non-profit'), and so will exclude all aangifte from businesses.

### 2.4 Date Range

Random sample aangiften from tax years 2010 and 2011. Auditing the random sample from tax year 2012 is still in progress so these years represent the mode recent available.

### 2.5 Timeline

Lab phase 1 is scheduled to take 6 to 8 weeks.



## 3 Lab Phase 2

### 3.1 Description

The second lab phase will build a predictive risk model that builds on and extends the modelling work that was done in lab phase 1.

We will build a model that looks at the whole aangifte, as well as 10 subject level models that are built up from the sections in the aangifte. These will provide specificity. The subjects are:

- **Assets and Debts**
  - Section 19 – result of providing assets
  - Section 20 – balance sheet
  - Section 28 – substantial interest
  - Section 29 – assets
  - Section 30 – debts
  - Section 31 – summary of sections 29-30
  - Section 47 – separated private assets
- **Expenses**
  - Section 25 – (back) amount previously deducted
  - Section 32 – alimony paid out
  - Section 33 – expenses for a child under 30
  - Section 34 – medical expenses
  - Section 35 – disabled children
  - Section 36 – educational costs
  - Section 37 – listed building expenses
  - Section 38 – lost venture capital
  - Section 49 – donations
  - Section 40 – previous year's expenses
- **Fiscal Partner**
  - Section 57 - selection of fields from the partner's aangifte
- **Foreign Income**
  - Section 15 – salary from working abroad
  - Section 16 – pension from abroad
  - Section 48 – foreign dividends and interest
  - Section 52 – Belgium
- **Healthcare**
  - Section 34 – medical expenses
  - Section 53 – Health Insurance Act
- **Other Income**
  - Section 18 – income from other work
  - Section 22 – alimony received
  - Section 23 – periodical income received
  - Section 24 – other income received
  - Section 48 – foreign dividends and interest
  - Section 50 – income to preserve
  - Section 51 – income taxed elsewhere
- **Own Home**
  - Section 21 – own home
- **Salary and Pensions**
  - Section 12 – salary from normal job
  - Section 13 – income from pension
  - Section 14 – salary from an international organization



- Section 15 – salary from working abroad
- Section 16 – pension from abroad
- Section 26 – insurance contributions
- Section 27 – negative expenses
- Section 49 – revisionary interest
- Section 52 – Belgium
- Tax Credits
  - Section 41 – low earner
  - Section 42 – life discount
  - Section 43 – tax credits parents with children under 27 years
  - Section 44 – single elderly discount
  - Section 45 – young disabled discount
  - Section 46 – tax i.v.m. investments
- Travel Deductions
  - Section 17 – travel deduction public transport

In addition every subject also includes section 1, personal details, as most tax affairs need to be considered the context of the taxpayer's marital status, number of children etc.

The mapping of sections to subject was built in conjunction with advice from auditors and LTO.

### 3.2 Outcome

We will model the risk of a positive correction. That is, the risk that the taxpayer owes more tax than declared on the aangifte.

### 3.3 Description of the Data to be Used

#### 3.3.1 Whole Aangifte Model

We will build the model on niet winst aangiften from 2011 that have been manually audited, excluding aangiften:

- that are automatically corrected by MAC UWBs
- that are selected as part of the 'harde uitwerp' process, and hence manually audited, and have a single UWB on them, where that UWB is expected to become a MAC UWB in future tax years.

This excludes cases that were, or could have been, corrected by an automatic process that checks the aangifte values against the contra data.

The actual data we will build the model on is:

4. the data in the aangifte
5. the result of the fiscal verification check (FVC)
6. history from the previous year's aangifte
7. a selection of external or 'contra' information
8. a selection of 'additional data' such as:
  - a. the number of aangiften previously submitted for this sofi number and this year
  - b. was VIA data used
  - c. was the aangifte submitted electronically or on paper
9. the results of the audit. This will be used to form the targets.

#### 3.3.2 Whole Aangifte Target



For the whole aangifte model the target is the positive correction. Specifically that the tax owed on the DA (definitive aanslag, the result of an audit or an automatic correction) is greater the tax owed on the VA (preliminary aanslag, this contains the tax as calculated from the values on the aangifte). Out of process aangiften that do not have a VA followed by a DA are not included in the model build data set as we cannot construct an adequate target.

Note that in Lab Phase1 the target we used included fines and interest and so was corrupted by aangifte that were late but did not otherwise have a positive correction.

### 3.3.3 Subject Models

An aangifte will be included in the model build data for a subject if that aangifte is part of the overall model build (i.e. it fulfills the data requirements in section 3.3.1) and it has a UWB code on it that indicates that one or more of the metacodes in that subject should have been inspected<sup>31</sup>.

If an aangifte was part of the random sample or VIP set then we assume the whole aangifte was looked at and it is included in the model build data set for all subjects.

### 3.3.4 Subject Model Targets

The target for the subject models is set thus:

If the data indicates that at least one metacode in the sections that make up that subject was altered during the audit<sup>32</sup>

and the overall target is 1 (i.e. there was a positive correction),

then the subject target is set to 1 (a hit).

Otherwise it the subject target is 0 (a miss).

## 3.4 Date Range

Tax year 2011.

In addition we will use historical data from the previous tax year.

## 3.5 Timeline

Lab phase 2 is scheduled to take 10 to 12 weeks.

---

<sup>31</sup> In theory an auditor only audits the parts of the aangifte indicated by the UWB, in practice they look at and edit other parts of the aangifte. Since we have no way of knowing exactly which fields were inspected, assuming that the auditor followed the work instructions for the UWB is the best we can do.

<sup>32</sup> We calculate this directly from the aangifte data, so here we do not have to rely on assumptions about work practices.



## 4 Pilot Phase

### 4.1 Description

The pilot phase will refine and validate the models from the second lab phase and prepare them for roll-out.

We will select 2,000 cases to be audited in a controlled setting. The bulk of these, 1800, will be 'live' unaudited cases and these will be used to validate the model. The remaining 200 will be historical cases from the dataset that was used to build the model. These will be used to help us understand any anomalies that may occur in the model.

Additionally we will gather data from the results of the audit for the entire random sample.

The results of the audits, whether there was a positive correction or not, and which sections of the aangiften were updated, will then be compared with the model predictions.

### 4.2 Validation

We will score the niet winst random sample for tax year 2013, using both the whole aangifte model and the 10 subject models.

We will then

- rank them by overall risk where the overall level of risk is the maximum risk score across all models
- divide this ordered set into two strata
  - the top 5% by volume – these are the aangiften that the model would select as risky
  - the lower 95% by volume - these are the aangiften the model would not select.

From the top 5% of cases we will then randomly select 1,400 cases.

From the bottom 95% of cases we will randomly select 400 cases.

These 1,800 cases will then be audited and the results compared with the model's predictions.

### 4.3 Outcome

We will model the risk of a positive aanslag. That is, the risk that the taxpayer owes more tax than declared on the aangifte.

We will validate the Niet Winst model.

### 4.4 Description of the Data to be Used

We will use aangiften from random sample for tax year 2013.

The actual data we will build the model on is:

10. the data in the aangifte
11. the result of the fiscal verification check (FVC), which is an automatic consistency check performed on the aangifte.
12. historical aangifte data
13. external or 'contra' information.



In addition we will audit 200 historical cases from the 2011 data used to build the model, using the data as described in section 3.3

#### 4.5 Date Range

For the 1,800 'live' cases we will use:

1. Un-audited aangiften from tax year 2013 for the sample for validation.
2. The previous year's data (2012) for those aangiften.

For the 200 'historical' cases we will use:

1. Data from tax year 2011
2. The previous year's data (2010) for those aangiften.

#### 4.6 Timeline

Pilot phase is scheduled to take 10 weeks



## 5 Appendix A: Detailed Baseline Results, Lab Phase 1

Counts come from the SAS dataset `aang_vord_aans_totaal.sas7bdat`, stored in the AWS+ directory `/AD005/shared/datasets/BI-A_data/IH/IHriskmodel`.

Counts as of 27th March 2014.

Tax year is defined by

```
belastingjaar = 2010 or 2011
```

Individuals aka Not Profit aka Niet Winst is defined by

```
aangiftesoort = Niet Winst
```

Note that this means that

```
aangiftesoort = Niet Winst btl
```

```
aangiftesoort = Winst
```

are excluded from our baseline counts.

Random Sample defined by

```
RISCIOCAT = '201000' or '201001'
```

Which, for 2010 and 2011, is equivalent to

```
risciobenaming = "Steekproef" or "Steekproef met indicatie VIA"
```

As of 23<sup>rd</sup> April we were given the files `sp2010` and `sp2011` which contain the definitive list of sofi numbers in the random sample. This reduced the number of sofi numbers we had to work with by about 16%.

The Positive Aanslag audit result is defined by

```
POS_NEG = P
```

The Negative Aanslag audit result is defined by

```
POS_NEG = N
```

The Nihil Aanslag audit result is defined by

```
POS_NEG = Z
```





## 5.1 All Audited Aangiften from 2010 and 2011

This section gives the breakdown by tax year and aangiften type for all aangiften audited in 2010 and 2011 for which the analytics team had results.

### 5.1.1 Counts

# Audit Result	Tax Year		TOTAL
	2010	2011	
Missing	Inspectie, controle en toezicht		
Negative			
Nihil			
Positive			
TOTAL			
TOTAL WITH RESULTS			

Table 28: all audited 2010 and 2011 aangiften

# Audit Result	Tax Year		TOTAL
	2010	2011	
Missing	Inspectie, controle en toezicht		
Negative			
Nihil			
Positive			
TOTAL			
TOTAL WITH RESULTS			

Table 29: all audited 2010 and 2011 'Winst' aangiften

# Audit Result	Tax Year		TOTAL
	2010	2011	
Missing	Inspectie, controle en toezicht		
Negative			
Nihil			
Positive			
TOTAL			
TOTAL WITH RESULTS			

Table 30: all audited 2010 and 2011 'Niet Winst btl' aangiften



# Audit Result	Tax Year		TOTAL
	2010	2011	
Missing	Inspectie, controle en toezicht		ht
Negative			
Nihil			
Positive			
TOTAL			
TOTAL WITH RESULTS			

Table 31: all audited 2010 and 2011 'Niet Winst' aangiften

5.1.2 Percentages

These tables represents the general baselines for audit with results in 2010 and 2011.

% Audit Result	Tax Year		TOTAL
	2010	2011	
Negative	Inspectie, controle en toezicht		
Nihil			
Positive			

Table 32: percentage of all aangiften audited with results for tax years 2010 and 2011

% Audit Result	Tax Year		TOTAL
	2010	2011	
Negative	Inspectie, controle en toezicht		
Nihil			
Positive			

Table 33: percentages for 'Winst' aangiften audited with results for tax years 2010 and 2011

% Audit Result	Tax Year		TOTAL
	2010	2011	
Negative	Inspectie, controle en toezicht		icht
Nihil			
Positive			

Table 34: percentages for 'Niet Winst btl' aangiften audited with results for tax years 2010 and 2011

% Audit Result	Tax Year		TOTAL
	2010	2011	
Negative	Inspectie, controle en toezicht		
Nihil			
Positive			

Table 35: percentages for 'Niet Winst' aangiften audited with results for tax years 2010 and 2011



## 5.2 Audited Aangiften from 2010 and 2011 Selected By Risk Rules

This section gives the breakdown by tax year and aangiften type for all 2010 and 2011 aangiften that were selected for auditing based on risk rules in the IH Selection Module.

### 5.2.1 Counts

# Audit Result	Tax Year		TOTAL
	2010	2011	
Missing	Inspectie, controle en toezicht.		
Negative			
Nihil			
Positive			
TOTAL			
TOTAL WITH RESULTS			

Table 36: 2010 and 2011 aangiften selected for audit by risk rules

# Audit Result	Tax Year		TOTAL
	2010	2011	
Missing	Inspectie, controle en toezicht.		
Negative			
Nihil			
Positive			
TOTAL			
TOTAL WITH RESULTS			

Table 37: 'Winst' 2010 and 2011 aangiften selected for audit by risk rules

# Audit Result	Tax Year		TOTAL
	2010	2011	
Missing	Inspectie, controle en toezicht.		
Negative			
Nihil			
Positive			
TOTAL			
TOTAL WITH RESULTS			

Table 38: 'Niet Winst bt!' 2010 and 2011 aangiften selected for audit by risk rules



# Audit Result	Tax Year		TOTAL
	2010	2011	
Missing	Inspectie, controle en toezicht		
Negative			
Nihil			
Positive			
TOTAL			
TOTAL WITH RESULTS			

Table 39: 'Niet Winst' 2010 and 2011 aangiften selected for audit by risk rules

### 5.2.2 Percentages

These tables represents the baselines of audit results aangiften selected by risk for tax years 2010 and 2011.

% Audit Result	Tax Year		TOTAL
	2010	2011	
Negative	Inspectie, controle en toezicht		cht
Nihil			
Positive			

Table 40: percentages for all aangiften selected by risk rules, audited with results for tax years 2010 and 2011

% Audit Result	Tax Year		TOTAL
	2010	2011	
Negative	Inspectie, controle en toezicht		
Nihil			
Positive			

Table 41: percentages for 'Winst' aangiften selected by risk rules, audited with results for tax years 2010 and 2011

% Audit Result	Tax Year		TOTAL
	2010	2011	
Negative	Inspectie, controle en toezicht		
Nihil			
Positive			

Table 42: percentages for 'Niet Winst btl' aangiften selected by risk rules, audited with results for tax years 2010 and 2011

%	Tax Year	TOTAL
---	----------	-------



Audit Result	2010	2011	
Negative	Inspectie, controle en toezicht		
Nihil			
Positive			

Table 43: percentages for 'Niet Winst' aangiften selected by risk rules, audited with results for tax years 2010 and 2011



### 5.3 Audited Aangiften from 2010 and 2011 Selected Using Random Selection

This section gives the breakdown by tax year and aangiften type for all 2010 and 2011 aangiften that were selected for auditing based on a random selection.

In theory only 'Niet Winst' aangiften are selected for the random sample. In the dataset available at the time of writing there are 6 'Winst' or 'Niet Winst bt' aangiften that are marked as being randomly selected in 2011 and 2012. Since this is a tiny number they will be ignored.

#### 5.3.1 Counts

# Audit Result	Tax Year		TOTAL
	2010	2011	
Missing	Inspectie, controle en toezicht		
Negative			
Nihil			
Positive			
TOTAL			
TOTAL WITH RESULTS			

Table 44: 'Niet Winst' 2010 and 2011 aangiften randomly selected for audit

#### 5.3.2 Percentages

This table represents the baselines or randomly selected aangiften with audit results for tax years 2010 and 2011

% Audit Result	Tax Year		TOTAL
	2010	2011	
Negative	Inspectie, controle en toezicht		
Nihil			
Positive			

Table 45: Percentages of Completed 'Niet Winst' 2010 and 2011 aangiften randomly selected for audit

#### 5.3.3 Addendum

As of 23<sup>rd</sup> April we were given the files `sp2010` and `sp2011` which contain the definitive list of sofi numbers in the random sample as the risk category can change over time. This has reduced the number of sofi numbers we had to work with by about 16%.



# Audit Result	Tax Year		TOTAL
	2010	2011	
Negative	Inspectie, controle en toezicht		nt
Nihil			
Positive			
TOTAL			

Table 46: corrected counts for the random sample

% Audit Result	Tax Year		TOTAL
	2010	2011	
Negative	Inspectie, controle en toezicht		
Nihil			
Positive			
TOTAL			

Table 47: corrected percentages for the random sample

### 5.4 Audit Counts by Tax Year

# Audit Result	Tax Year						
	2006	2007	2008	2009	2010	2011	2012
Missing	Inspectie, controle en toezicht						
Negative							
Nihil							
Positive							
TOTAL							
TOTAL COMPLETED							

Table 48: Audit Counts by Tax Year

The current data for tax year 2012 only has approximately 200k completed audits. It is hoped that by the time we get to Lab Phase 2 we will have more data for 2012. We have no audit data for 2013 at this time.

# completed	Tax Year						
	2006	2007	2008	2009	2010	2011	2012
Winst	Inspectie, controle en toezicht						
Niet Winst							
Niet Winst btl							
TOTAL COMPLETED							

Table 49: Completed Audits by Tax Year and Type of Aangifte

## 6 Appendix B: Detail for the Baseline, Lab Phase 2



Counts come from the SAS dataset aang\_vord\_aans\_totaal\_new3.sas7bdat, stored in the AWS+ directory /AD005/shared/datasets/BI-A\_data/IH/IHriskmodel.

Counts as of 7th July 2014.

Tax year is defined by

```
belastingjaar = 2011.
```

Individuals aka Not Profit aka Niet Winst are defined by

```
ind_winst = 'Niet Winst'.
```

IH tax is defined by

```
middel = 'IH'.
```

Including only those aangifte that had a VA followed by a DA is defined by

```
type_afhandeling = 'VA, DA'.
```

Including only those aangifte that had a positive correction is defined by

```
bedrag_bel > 0.
```

Excluding those aangifte that were, or could be, corrected by MAC is defined by excluding where

```
UWB in ('H0225', 'H0288', 'H0290', 'H0500', 'H0501', 'H0516',  
'H0560', 'H0658').
```

Double check on excluding those aangifte that were corrected by MAC is defined by excluding where

```
wijze_afdoen = 'massaal automatisch corri'.
```

Excluding those aangifte that were part of the harde uitworp process is defined by excluding where

```
wijze_afdoen = 'harde-uitworpregeling'.
```





- 8.6 Appendix 6: Reference 7: Persoonsgegevens “Belastingdienst IH Risk Model: Requirements 2 (Pilot 1 and onwards),” /IH selectie/Project Managment/Quality/.

**Belastingdienst IH Risk Model  
Requirements 2  
(Pilot 1 and onwards)  
Version: 0.3**

**DRAFT**



**Table of Contents**

<b>1. OVERVIEW</b> .....	<b>129</b>
<b>2. PILOT 1</b> .....	<b>129</b>
2.1 BACKGROUND AND CONTEXT.....	129
<b>3. VALUE MODEL</b> .....	<b>129</b>
3.1 BACKGROUND AND CONTEXT.....	130
3.2 INITIAL VALUE MODEL .....	130
3.3 GENERATE INSIGHTS FROM THE VALUE MODEL.....	130
3.4 SECOND VALUE MODEL.....	130
<b>4. RE-BUILD SUBJECT MODELS</b> .....	<b>130</b>
4.1 BACKGROUND AND CONTEXT.....	131
4.2 METHOD.....	131
<b>5. BACKLOG 2013</b> .....	<b>131</b>
<b>6. PILOT 2</b> .....	<b>133</b>
<b>7. IMPLEMENTATION AKA ‘SHADOW RUN’</b> .....	<b>133</b>



**Revision History**

Date	Version	Description	Author
26/11/14	0.1	Creation	Persoonsgegevens
18/12/14	0.2	Value model details	
21/1/15	0.3	Sign off on value model thresholds	



# 1 Overview

This document tracks the evolving requirements for the 'IH Risk Model' analytics project. It starts from the Pilot 1 phase onwards and picks up where the previous requirements document, Belastingdienst\_IH\_Risk\_Business\_Requirements\_v0.7, stops.

# 2 Pilot 1

**Summary:** Pilot to test the whole aangifte and 10 subject models.

**Requirements Generated:** in the preceding requirements document Belastingdienst\_IH\_Risk\_Business\_Requirements\_v0.7

**Requirements Definition Agreed:** no date logged but the details in the preceding requirements document have been accepted by Persoonsgegevens

The 'whole aangifte' model hit rate needs to be greater than  $\frac{1}{2}$  at an absolute minimum. A hit rate of  $\frac{1}{3}$  however will fulfil the requirements.

**Requirements Completed and Signed Off:**

## 2.1 Background and Context

1800 cases were chosen for audit from the 2013 random sample.  
200 are used to test the work instructions.  
1400 were selected on the basis of risk.

Additionally 200 historical cases (i.e. previously audited) were chosen to check for defects in the model.

Pilot 1 takes place October to December 2014.

Analysis of the results will take place in January 2015.

# 3 Value Model

**Summary:** We will be building two value models, each with a single threshold.

**Requirements generated:**

The use of a cut-off threshold model: 25<sup>th</sup> November 2014 in a meeting attended by

- [Redacted]
- Persoonsgegevens
- Persoonsgegevens tns
- Persoonsgegevens

Refined to two threshold models, in a meeting 17<sup>th</sup> December 2014 attended by

- Persoonsgegevens
- Persoonsgegevens
- Persoonsgegevens evens
- Persoonsgegevens
- Persoonsgegevens

Two thresholds decided. A 'low value' threshold (of  $\frac{1}{2}$  the correction policy threshold), and a high value threshold (of  $\frac{1}{2}$  the median value).

**Requirements Definition Agreed:**

26<sup>th</sup> November 2014, by email from [Persoonsgegevens]

17<sup>th</sup> December 2014, verbally in the meeting by [Persoonsgegevens]

18<sup>th</sup> December 2014, agreed by email by [Persoonsgegevens]

**Requirements Completed and Signed Off:**

### 3.1 Background and Context

Generate descriptive statistics for value for IH Niet Winst cases: mean, median etc.

Consult widely as to at what value an audit is 'worth doing'.

A workshop for this will be taking place in the third week of December

### 3.2 Initial Value Model

Present the results of descriptive statistics and workshop to [Persoonsgegevens] to decide on a threshold for value.

Build a value model that generates the probability of the value of a positive correction being higher than the threshold

i.e.  $\text{Prob}(\text{value} \geq \text{€ threshold})$

This build is scheduled to take place in January.

### 3.3 Generate Insights from the Value Model

Generate a plot of  $\text{Prob}(\text{value} > \text{€X})$  against the actual value (for previously audited cases). Are the higher probabilities associated with higher values?

### 3.4 Second Value Model

If the higher probabilities are not associated with higher values, generate a second value model that has a higher threshold, such as value in the top 10% of values

i.e.  $\text{Prob}(\text{value} \geq \text{€ high value})$

We will attempt to build this at the same time as the low value model.

## 4 Re-Build Subject Models

**Summary:** We will investigate the current section to subject mapping, and if appropriate will re-build some or all of the subject models with a different section to subject mapping.

**Requirements Generated:** by the team of auditors that are collaborating with us in Pilot 1 November 2014.

**Requirements Definition Agreed:**

Stakeholder workshop on 4<sup>th</sup> December 2014. Attended by representatives from

- LTO
- The audit team
- Fiscal experts.

**Requirements Agreed Completed:****4.1 Background and Context**

The aangifte contains 57 sections, which we mapped to 10 subjects:

- Assets and Debts
- Expenses
- Foreign
- Healthcare
- Salary and Pension
- Other Income
- Own Home
- Tax Credits
- Travel
- Partner

A model was built for each subject. As Pilot 1 started it became clear that many of the auditors disagreed with the section to subject mapping we used.

**4.2 Method**

Each of the team of auditors that are collaborating with us in Pilot 1 have been given a document that allows them to record what they think the subjects should be and what sections should be in each subject.

Additionally we will hold a workshop with a broad range of stakeholders, at which we will come to a consensus as to what the new mapping should be. This workshop took place on the 4<sup>th</sup> December.

**5 Backlog 2013**

**Summary:** score the 2013 backlog

**Requirements Generated:**  Exact date not available

**Requirements Definition Agreed:**  Exact date not available

**Requirements Agreed Completed:**

**5.1 Background and Context**

There are approximately 1,500,000 aangiften from 2013 that have either not been processed by the IH selection module or that have been flagged by the IH selection module but have not yet been audited.

This team has been asked to score them for risk to help the audit team prioritise this workload.

The total 2013 backlog is approximately 1.5 million.

**5.2 Method**



We will score this backlog with the whole aangifte model, the re-built subject models and the two value models (high and low).

They will be scored at the end of January 2015.



## 6 Pilot 2

**Summary:** Pilot 2 will test the work instructions more fully. It will start at the beginning of February 2015.

**Requirements Generated:** Persoonsgegevens Exact date not available

**Requirements Definition Agreed:**

**Requirements Agreed Completed:**

### 6.1 Background and Context

The work instructions were tested on 200 aangifte as part of pilot 1. One of the outcomes of that pilot was that the auditors were unhappy with the make-up of the subject areas. In response to this a workshop was held to generate better subject areas and new subject models were built (see section 4).

Pilot 2 will test the new work instructions.

### 6.2 Method

## 7 Implementation aka 'Shadow Run'

**Summary:**

**Requirements Generated:**

**Requirements Definition Agreed:**

**Requirements Agreed Completed:**





8.7 Appendix 7: Reference 8: Persoonsgegevens "Belastingdienst IH Risk Model: Technical Definitions," /IH selectie/Project Management/Quality/.

**Belastingdienst IH Risk Model**  
**Technical Definitions**  
**Version: 0.7**

**DRAFT**

**Table of Contents**

<b>1. OVERVIEW .....</b>	<b>137</b>
<b>2. LAB PHASE 1 .....</b>	<b>138</b>
2.1 SELECTING THE AANGIFTEN .....	138
2.2 SETTING THE TARGET .....	138
<b>3. LAB PHASE 2 .....</b>	<b>139</b>
3.1 SETTING THE BASELINE .....	139
3.2 SELECTING THE AANGIFTEN .....	140
3.3 SETTING THE TARGET: WHOLE AANGIFTE MODEL .....	141
3.4 SETTING THE TARGET: SUBJECT LEVEL MODELS (EXCLUDING PARTNER) .....	142
3.4.1 <i>Subject Model Build</i> .....	142
3.4.2 <i>Subject Model Score</i> .....	143
3.5 SETTING THE TARGET: FISCAL PARTNER MODEL .....	143
3.6 SECTION TO SUBJECT MAPPING .....	143
<b>4. PILOT 1 AND BACKLOG 2012 .....</b>	<b>149</b>
<b>5. PILOT 2 .....</b>	<b>149</b>
5.1 WHOLE AANGIFTE MODEL .....	149
5.1.1 <i>Selecting the Aangiften: Whole Aangifte Model Build</i> .....	149
5.1.2 <i>Setting the Target: Whole Aangifte Model Build</i> .....	150
5.1.3 <i>Selecting the Aangiften: Whole Aangifte Score</i> .....	150
5.2 SUBJECT MODELS .....	150
5.2.1 <i>Selecting the Aangiften: Subject Models Build</i> .....	150
5.2.2 <i>Setting the Target: Subject Models Build</i> .....	151
5.2.3 <i>Selecting the Aangiften: Subject Model Score</i> .....	151
5.3 SECTION TO SUBJECT MAPPING .....	151
5.4 VALUE MODELS .....	156
5.4.1 <i>Selecting the Aangiften: Value Models Build</i> .....	156
5.4.2 <i>Setting the Targets : Value Models Build</i> .....	156
5.4.3 <i>Selecting the Aangiften: Value Models Score</i> .....	157
5.5 COMPLEXITY MEASURES .....	157
5.5.1 <i>Complexity Scores</i> .....	157
5.5.2 <i>Specialist Areas</i> .....	158
<b>APPENDIX A: STRUCTURE OF THE MODEL BUILD KEY TABLE .....</b>	<b>159</b>
<b>APPENDIX B: STRUCTURE OF THE MODEL SCORE KEY TABLE .....</b>	<b>160</b>
<b>APPENDIX C: STRUCTURE OF THE MODEL SCORE OUTPUT TABLE .....</b>	<b>160</b>
<b>APPENDIX D: VALUE MODEL STATISTICS .....</b>	<b>161</b>



Revision History

Date	Version	Description	Author
26/11/14	0.1	Creation	Persoonsgegevens
9/12/14	0.2	Added Pilot 2 section	
9/12/14	0.3	Added Value Model section, changed name of doc, started complexity section	
15/12/14	0.4	Subject definitions added	
18/12/14	0.5	Details of value models added	
19/12/14	0.6	More notes on the value model	
14/4/14	0.7	Baseline definition added	



## 1 Overview

This document

- tracks the evolving definition of the target for the 'IH Risk Model' analytics project.
- details the criteria for an aangifte to be included in a model build or scored dataset.
- gives the definition(s) for the complexity measures.

This document is deliberately specific and technical. The aim is to have a single repository of all the decisions made about the modelling.

Here we use the word 'target' in the technical analytics sense. That is, a variable whose value we are building a model to predict.

For this project the targets are binary:

- `target = 1`, means that there was a positive correction as the result of an audit
- `target = 0`, means that there was no correction, or there was negative correction.

Unless otherwise stated all variables come from a dataset called `aang_vord_aans_totaal`.

This dataset is built and owned by the BI&A team, and is a unified view of cases that have been worked.

Variable names, selection criteria, code snippets etc. are indicated by being in the `Courier` font.

`<subject>` means that the item should be repeated, once for each one of the subject models.



## 2 Lab Phase 1

### 2.1 Selecting the Aangiften

We used the following criteria to select aangiften for the model build:

- Tax year 2010 or 2011,
- Niet Winst (individuals),
- Random sample

i.e.

```
IF (  
    (belastingjaar = 2010 or 2011) AND  
    (aangiftesoort = 'Niet Winst') AND  
    (RISCIOCAT = '201000' or '201001') AND  
)  
THEN it's a valid aangifte to include
```

### 2.2 Setting the Target

We used the field POS\_NEG, where "P" indicates a positive correction, "N" a negative correction and "Z" no correction. Our target was hence defined by:

- POS\_NEG = "P" is a hit, that is target = 1
- POS\_NEG = "Z" or POS\_NEG="N" or POS\_NEG is missing, is a miss, that is target = 0.

We later found out this field was not exactly what we had thought.

A value of POS\_NEG = "P", that is a positive correction, includes fines and interest levied, typically for being late.

Since these fines are very common the frequency of them is liable to 'drown out' the tax correction itself. This means that in effect the initial model predicted the risk that an aangifte would be late. And since the model uses data in the aangifte itself we were predicting lateness of a form we had already received.



### 3 Lab Phase 2

We refined the target to make sure that we were only registering a hit when there was a positive correction as the result of an audit, and that fines and interest are not included.

This phase also introduced 10 subject level models, each of which needs an appropriate target.

#### 3.1 Setting the Baseline

The baseline is the hit-rate we compare our results to.

We used the following criteria:

- Tax years 2010 and 2011
  - As these are the most recent years for which most of the auditing is complete
- Niet Winst (individuals)
- IH tax
- Overige Uitworp
  - these are the 'risk' UWBs, and are manually audited
  - this excludes Harde Uitworp, the 'non-risk' UWBs, which are rules defined by political decisions to audit certain groups of taxpayers and so do not represent risk as used in this document.
- Exclude UWBs that will become MAC
  - To include these would be unfair to the risk model, as they are a simple comparison between a value on the aangifte and its matching contra data, and will be replaced by an automatic process in future years.
  - These UWBs are: H225,H288,H290,H500,H501,H516,H560,H658.
- Exclude the VIP list
  - To include these would be unfair to the UWB Risk Module as positive correction rate for this group is very low. This group does not represent a risk.
  - This is defined by the field ind\_vip.
- Exclude the random sample.
  - To include these would be unfair to the UWB Risk Module as positive correction rate for this group is low. This group does not represent a risk.
  - This is defined by the two files spk2010 and spk2011.
- We can calculate the change in tax as a result of an audit.
  - See the discussion on Voorlopige Aanslag (VA) and Definitieve Aanslag (DA) in the section on setting the target.
  - We refer to these as 'being trackable'.
- We have a result for the audit.

```

IF (
  (belastingjaar = 2010 OR belastingjaar = 2011) AND
  (ind_winst = 'Niet Winst') AND
  (MIDDEL = 'IH') AND
  (wijze_afdoen = 'overige-uitworpregeling') AND
  (type_afhandeling = 'VA, DA') AND
  (bedrag_bel_DA NE .) AND
  (ind_VIP = 0) AND
  (finr is not in SPK2010 or SPK2011) AND

```



```

(UWB NOT IN
 ('H0255', 'H0288', 'H0290', 'H0500', 'H0501', 'H0516', 'H0560', 'H0658')
)
}
THEN it's a valid aangifte to include in the baseline.

```

Counts for IH Niet Winst trackable aangiften, that were flagged by overige-uitworp:

	ALL IH Niet Winst trackable
negative	Inspectie, controle en toezicht
nil	
positive	
total	

Of those, the following are MAC, VIP and random sample, and so need to be excluded:

Inspectie, controle en toezicht
---------------------------------

This leaves us with:

Inspectie, controle en toezicht
---------------------------------

Further details can be found in the spreadsheet NietWinstBaselines\_March15v0.1

### 3.2 Selecting the Aangiften

We used the following criteria to select aangiften for this model build:



- Tax year 2011,
- Niet Winst (individuals),
- IH Tax<sup>33</sup>
- Manually audited, excluding MAC UWBs

We excluded cases that went through the MAC process as this is a very efficient automatic check and a risk model is not going to improve on it.

We also excluded manually audited aangiften that had just a single UWB where that UWB was one that is scheduled to become a MAC process in future years. The UWBs in this category are: H0255, H0288, H0290, H0500, H0501, H0516, H0560, H0658

```

IF (
  (belastingjaar = 2011) AND
  (ind_winst = 'Niet Winst') AND
  (MIDDEL = 'IH') AND
  (
    (wijze_afdoen = 'overige-uitworpregeling')
    OR
    (wijze_afdoen = 'harde-uitworpregeling')
  )
  AND
  (UWB NOT IN
   ('H0255', 'H0288', 'H0290', 'H0500', 'H0501', 'H0516', 'H0560', 'H0658')
  )
)
THEN it's a valid aangifte to include

```

### 3.3 Setting the Target: Whole Aangifte Model

Initially we then used `POS_NEG = "P"` to get the positive corrections, but as discussed above in section 1 we discovered that this incorporated fines and interest as well as tax corrections.

We can only calculate the extra tax due as the result of an audit if the taxpayer has a Voorlopige Aanslag (VA) followed by Definitieve Aanslag (DA). Since this is the expected process it is what happens most of the time, but we still need to exclude out-of-process cases to make sure that our target is not corrupted.

A new field, `type_afhandeling`, was created that classifies the path of an aangifte through the tax system. In-process aangiften are identified by the path

<sup>33</sup> This is just a 'belt and braces' check





```
type_afhandeling = 'VA, DA'
```

The new calculation of value is then in the new field `bedrag_bel_DA`. If

```
type_afhandeling = 'VA, DA'
```

then `bedrag_bel_DA` is the change in tax as result of the audit. So a positive correction is when

```
bedrag_bel_DA > 0.
```

This simplifies to:

```
IF (
  (it's a valid aangifte to include) AND
  (type_afhandeling = 'VA, DA') AND
  (bedrag_bel_DA > 0)
)
THEN whole_aangifte_model_target = 1
ELSE whole_aangifte_model_target = 0
```

### 3.4 Setting the Target: Subject Level Models (excluding partner)

For a subject on an aangifte to be included in a subject model build it must be looked at during the audit, which currently means that:

- There must be at least one UWB applied to that aangifte that looks at some of the data in that subject.
- Or the whole aangifte has been looked at. Which happens when the aangifte is part of the random sample or in the 'VIP' set of aangiften.

Initially, if the subject was looked at then the target for the subject is set to the whole aangifte model target, as discussed in section 3.3 above.

This was later refined. We extracted data from ABS that tells us:

- Is there any data for this subject<sup>34</sup> in this aangifte?
- Is any of the data in the subject<sup>35</sup> edited during the audit?

#### 3.4.1 Subject Model Build

- If there is no data for the subject then
  - `subject_target = 0`
- If the subject is not checked during the audit then
  - `subject_target = 0`
- If the subject has data and is checked during the audit but no changes are made, then
  - `subject_target = 0`
- If the subject has data, is checked during the audit and is edited during the audit then
  - `subject_target = whole_aangifte_model_target`

<sup>34</sup> We exclude section 1, the personal details, from this check as all subjects include this section, and we don't want to include subjects for which the personal details are the only data.

<sup>35</sup> As footnote 2



In practice this simplifies to

```
IF (
  (it's a valid aangifte to include) AND
  (the subject has data) AND
  (the subject is checked during the audit) AND
  (the subject is edited during the audit) AND
  (whole_aangifte_model_target = 1)
)
THEN subject_target = 1
ELSE subject_target = 0
```

### 3.4.2 Subject Model Score

In order for an aangifte to receive a score for a subject model then it must have data in that subject.

## 3.5 Setting the Target: Fiscal Partner Model

The subject area Fiscal Partner has to be handled slightly differently from the other subjects as it is not edited during the audit. So for this subject

```
IF (
  (it's a valid aangifte to include) AND
  (the subject has data) AND
  (the subject is checked during the audit) AND
  (whole_aangifte_model_target = 1)
)
THEN partner_target = 1
ELSE partner_target = 0
```

## 3.6 Section to Subject Mapping

This mapping was designed by the analytics team and validated with members of LTO. Every subject also contains section 1 – personal details.

Section Number	Description	
19	Resultaat uit ter beschikking gestelde vermogensbestanddelen	result of providing assets
20	Balansgegevens	balance sheet



28	Inkomen uit aanmerkelijk belang	substantial interest
29	Waarde bezittingen	assets
30	Schulden	debts
31	Voordeel uit sparen en beleggen	summary of 29-30
47	Afgezonderd particulier vermogen	separated private assets

Table 50: Assets and Debts

Section Number	Description	
25	(Terug)ontvangen bedragen eerder afgetrokken	(Back) amounts previously deducted received
32	Betaalde alimentatie en andere onderhoudsverplichtingen	alimony paid out
33	Kosten levensonderhoud kinderen < 30	expenses for a child under 30
34	Uitgaven voor specifieke zorgkosten	medical expenses
35	Uitg. weekendbez. ernstig gehand. kinderen	disabled kids
36	Studiekosten en andere scholingsuitgaven	educational costs
37	Kosten voor Rijksmonumentenpanden	listed building expenses
38	Kwijtgescholden durfkapitaal	lost venture capital
39	Giften	donations
40	Restant persoonsgebonden aftrek	previous year expenses

Table 51: Expenses

Section Number	Description	
15	Buitenlands inkomen tegenwoordige DB	salary from working abroad



16	Buitenlands inkomen vroegere DB	pension from abroad
48	Dividend, kansspelink, rente btl spaartegoed	dividend, foreign interest
52	Werken in loondienst in België	Belgium

Table 52: Foreign

Section Number	Description	
34	Uitgaven voor specifieke zorgkosten	medical expenses
53	Niet gehele jaar verz.plicht volksverz. en ZVW	Not all year verz.plicht volksverz. and health insurance Act

Table 53: Healthcare

Section Number	Description	
18	Resultaat uit overige werkzaamheden (arbeid)	income from other work
22	Ontv. Alimentatie/afkoopsom	alimony received
23	Periodieke uitkeringen	periodical income received
24	Ontvangen overige inkomsten	other income received
48	Dividend, kansspelink, rente btl spaartegoed	dividend, foreign interest
50	Te conserveren inkomen	Income to preserve
51	Voorkoming dubbele belasting	income taxed elsewhere

Table 54: Other income

Section Number	Description	
21	Eigen woning	owner occupied home

Table 55: Own home



Section Number	Description	
57	Overzicht drempelinkomen partner	selection of fields from the partner's aangifte

Table 56: Partner



Section Number	Description	
12	Loon of uitkering Ziektewet	salary from normal job
13	Pensioen, lijfrente of andere uitkering	income from pension
14	Vrijgesteld ink. Funct. Inter. Org	salary from an international organisation
15	Buitenlands inkomen tegenwoordige DB	salary from working abroad
16	Buitenlands inkomen vroegere DB	pension from abroad
26	Betaalde premies voor inkomensvoorzieningen	insurance contributions
27	Lijfrente afgekocht/neg. uitgaven inkomensvoorz	negative expenses
49	Revisierente	revisionary interest
52	Werken in loondienst in België	Belgium

Table 57: Salary and Pensions

Section Number	Description	
41	Uitbetaling Algemene Heffingskorting	low earner
42	Levensloopverlofkorting	life-discount
43	Heffingskortingen ouders met kinderen < 27 jaar	Tax credits parents with children under 27 years
44	Alleenstaande ouderenkorting	Single elderly discount
45	Jonggehandicaptenkorting	young disabled Discount
46	Heffingskortingen i.v.m. beleggingen	Tax i.v.m. investments

Table 58: Tax Credits



Section Number	Description	
17	Reisaf trek openbaar vervoer	Travel Deduction public transport

Table 59: Travel



## 4 Pilot 1 and Backlog 2012

For these projects we scored the aangiften using the models developed in Lab Phase 2.

## 5 Pilot 2

For this phase we introduced value models and re-built the subject models.

The models are built on the same 2011 data that was used for the Lab Phase 2 model build and will score 2013 data.

### 5.1 Whole Aangifte Model

We are not refreshing the whole aangifte model. So this is just a repeat of the pilot 1 model build.

#### 5.1.1 Selecting the Aangiften: Whole Aangifte Model Build

We used the following criteria to select aangiften for this model build:

- Tax year 2011,
- Niet Winst (individuals),
- IH Tax<sup>36</sup>
- Manually audited, excluding MAC UWBs

We excluded cases that went through the MAC process as this is a very efficient automatic check and a risk model is not going to improve on it.

We also excluded manually audited aangiften that had just a single UWB where that UWB was one that is scheduled to become a MAC process in future years. The UWBs in this category are: H0255, H0288, H0290, H0500, H0501, H0516, H0560, H0658

```
IF (
  (belastingjaar = 2011) AND
  (ind_winst = 'Niet Winst') AND
  (MIDDEL = 'IH') AND
  (
    (wijze_afdoen = 'overige-uitworpregeling')
    OR
    (wijze_afdoen = 'harde-uitworpregeling')
  )
  AND
  (UWB NOT IN
  ('H0255', 'H0288', 'H0290', 'H0500', 'H0501', 'H0516', 'H0560', 'H0658')
  )
)
```

---

<sup>36</sup> This is just a 'belt and braces' check





```
THEN whole_aangifte_model = 1  
ELSE whole_aangifte_model = 0
```

### 5.1.2 Setting the Target: Whole Aangifte Model Build

We are not refreshing the whole aangifte target, so this is a repeat of Pilot 1.

We can only calculate the extra tax due as the result of an audit if the taxpayer has a Voorlopige Aanslag (VA) followed by Definitieve Aanslag (DA). Since this is the expected process it is what happens most of the time, but we still need to exclude out-of-process cases to make sure that our target is not corrupted.

```
IF (  
    (whole_aangifte_model = 1) AND  
    (type_afhandeling = 'VA, DA') AND  
    (bedrag_bel_DA > 0)  
)  
THEN whole_aangifte_target = 1  
ELSE whole_aangifte_target = 0
```

### 5.1.3 Selecting the Aangiften: Whole Aangifte Score

We will score all Niet Winst aangiften in the 2013 backlog. Exact definition of which sofis is still being defined.

```
IF (aangifte in 2013 Niet Winst backlog)  
THEN whole_aangifte_score = 1  
ELSE whole_aangifte_score = 0
```

## 5.2 Subject Models

We are re-building the subjects as a result of the section-to-subject mapping workshop that ran in December 2014. The logic remains the same as pilot 1.

There is no longer a partner model.

### 5.2.1 Selecting the Aangiften: Subject Models Build

Include an aangifte is selected for the Whole Aangifte model and if it has data for the sections that make up the subject (excluding the personal details section).

```
IF (  
    (whole_aangifte_model = 1) AND  
    (there is non-zero data in at least one relevant section37)  
AND  
    (the subject is checked during the audit38)
```

---

<sup>37</sup> Excluding personal



```
)  
THEN <subject>_model = 1  
ELSE <subject>_model = 0
```

### 5.2.2 Setting the Target: Subject Models Build

We are re-building the subjects as a result of the section-to-subject mapping workshop.

```
IF (  
  (<subject>_model = 1) AND  
  (the subject is edited during the audit) AND  
  (whole_aangifte_target = 1)  
)  
THEN <subject>_target = 1  
ELSE <subject>_target = 0
```

### 5.2.3 Selecting the Aangiften: Subject Model Score

We will subject score aangiften that are selected for the whole aangifte score and where the aangifte has data for the subject.

```
IF (  
  (whole_aangifte_score = 1) AND  
  (there is non-zero data in at least one relevant section)  
)  
THEN <subject>_score = 1  
ELSE <subject>_score = 0
```

## 5.3 Section to Subject Mapping

This mapping was decided at a stakeholder workshop that was attended by representatives from LTO, the audit teams and fiscal experts.

The new sections were mapped using the 2013 aangifte. Some of section numbering has changed since the 2011 aangifte, which is what was used for the original model build. Since we are re-using the sections as building blocks we note the section number in 2011 and in 2013. All subjects also contain section 1 – personal details.

---

<sup>36</sup> That is, there is a UWB on the aangifte that points to somewhere in that subject



Section Number		Description	
2011	2013	Box 3 Bezittingen & Schulden	Box 3 Assets & Liabilities
29	29	Waarde bezittingen	value possessions
30	30	Schulden	debts
31	31	Voordeel uit sparen en beleggen	Benefit from savings and investments
47	49	Afgezonderd particulier vermogen	Separated private assets
48	50	Dividend, kansspelinkomen, rente betaling spaartegoed	Dividend, gambling income, interest payment savings

Table 60: Assets - new mapping

Section Number		Description	
2011	2013	Uitgaven en Budgetten	Spending and budgets
25	25	(Terug)ontvangen bedragen eerder afgetrokken	(Back) received amounts previously deducted
26	26	Betaalde premies voor inkomensvoorzieningen	Premiums paid on income
32	32	Betaalde alimentatie en andere onderhoudsverplichtingen	Alimony paid and other maintenance obligations
33	33	Kosten levensonderhoud kinderen < 21	Cost of living children <21
36	36	Studiekosten en andere scholingsuitgaven	Tuition and other educational expenses
37	37	Kosten voor Rijksmonumentenpanden	Costs for monument Properties
38	38	Kwijtgescholden durfkapitaal	Waived venture capital

Table 61: Expenses - new mapping

Section Number		Description	
----------------	--	-------------	--



2011	2013	Specifieke Zorgkosten & Giften	Specific care costs & gifts
25	25	(Terug)ontvangen bedragen eerder afgetrokken	(Back) received amounts previously deducted
34	34	Uitgaven voor specifieke zorgkosten	Expenditure on specific care costs
35	35	Uitgaven weekendbezoek ernstig gehandicapte kinderen	Spending a weekend visit severely disabled children
39	39	Giften	gifts
40	40	Restant persoonsgebonden aftrek	Remainder personal allowance

Table 62: Health- new mapping

Section Number		Description	
2011	2013	Buitenland	Foreign
15	15	Buitenlands inkomen tegenwoordige dienstbetrekking	Foreign income present employment
16	16	Buitenlands inkomen vroegere dienstbetrekking	Foreign income past employment
50	52	Te conserveren inkomen	Preserve income
51	53	Voorkoming dubbele belasting	Prevent double taxation
52	54	Werken in loondienst in België	Are employed in Belgium

Table 63: Foreign - new mapping

Section Number		Description	
2011	2013	Premieheffing & Zorgverzekeringswet	Contributions & Health Insurance



53	55	Niet gehele jaar verzekeringsplicht volksverzekering en ZvW	Not full year insurance national insurance and Insurance Act
54	56	Inkomensafhankelijke bijdrage ZvW	Income-related contribution ZvW

Table 64: Health - new mapping

Section Number		Description	
2011	2013	Overige Inkomsten	Other Income
22	22	Ontvangen alimentatie/afkoopsom	Receive alimony / indemnity
23	23	Periodieke uitkeringen	periodic benefits
24	24	Ontvangen overige inkomsten	Receive other income
27	27	Lijfrente afgekocht/negatieve uitgaven inkomensvoorziening	Annuity bought / negative expenditure Income Scheme
49	51	Revisierente	revision Interest

Table 65: Other Income - new mapping

Section Number		Description	
2011	2013	Eigen woning	Own home
21	21	Eigen woning	own home
37	37	Kosten voor Rijksmonumentenpanden	Costs for Monument Properties

Table 66: Expenses - new mapping

Section Number	Description
----------------	-------------



2011	2013	ROW & Aanmerkelijk Belang	ROW & Substantial interest
18	18	Resultaat uit overige werkzaamheden (arbeid)	Income from other activities (work)
19	19	Resultaat uit ter beschikking gestelde vermogensbestanddelen	Income from assets made available
20	20	Balans gegevens	balance sheet data
28	28	Inkomen uit aanmerkelijk belang	Income from a substantial interest

Table 67: ROW - new mapping

Section Number		Description	
2011	2013	Inkomsten die onder de loonheffing vallen, en reisafrek	Income under the income tax, and travel deduction
12	12	Loon of uitkering Ziektewet	Salary or benefits Sickness
13	13	Pensioen, lijfrente of andere uitkering	Pension, annuity or other benefit
14	14	Vrijgesteld inkomen functionele internationale organisatie	Exempt income functional international organization
15	15	Buitenlands inkomen tegenwoordige dienstbetrekking	Foreign income present employment
16	16	Buitenlands inkomen vroegere dienstbetrekking	Foreign income past employment
17	17	Reisafrek openbaar vervoer	Travel Deduction public transport
49	51	Revisierente	revision Interest

Table 68: Salary - new mapping

Section Number		Description	
2011	2013	Heffingskortingen	Tax credits
41	41	Uitbetaling Algemene Heffingskorting	Payment Terms Tax Credit



	42	Bijzondere verhoging heffingskortingen	Special tax increase
42	43	Levensloopverlofkorting	Life-Off
43	44	Heffingskortingen ouders met kinderen < 18	Tax credits parents with children <18
44	45	Alleenstaande ouderenkorting	Single elderly discount
45	46	Jonggehandicaptenkorting	young disabled Discount
46	47	Heffingskorting voor groene beleggingen	Tax credit for green investment
	48	Tijdelijke compensatie ZVW-bijdrage VUT/pre-pensioen	Temporary Compensation Insurance Act contribution early retirement / pre-retirement

Table 69: Tax credits - new mapping

## 5.4 Value Models

The value models will be applied to all aangiften selected by the whole aangifte model.

They will be built on the same 2011 data as was used for the subject model re-builds and was used for the 'whole aangifte' model.

### 5.4.1 Selecting the Aangiften: Value Models Build

Same rules as for selecting aangiften for the whole aangifte model: niet winst, tax year 2011, aangiften that have a VA followed by a DA, and manual audits excluding the UWBs that will become MAC. See section 5.1.1 for details.

```
IF whole_aangifte_model = 1
THEN value_model = 1
ELSE value_model = 0
```

### 5.4.2 Setting the Targets : Value Models Build

We are using two targets:

- 'low' or 'some' value: threshold is Inspectie, controle en toezicht the correction policy threshold
- 'high' value: threshold is Inspectie, controle en toezicht the median value.

Summary statistics for value are in Appendix D.

```
LET low_value_threshold = Inspectie, controle en toezicht
LET high_value_threshold = Inspectie, controle en toezicht
IF (value_model = 1)
```



```
THEN (  
    IF (bedrag_bel_DA > low_value_threshold)  
    THEN low_value_target = 1;  
    ELSE low_value_target = 0;  
  
    IF (bedrag_bel_DA > high_value_threshold)  
    THEN high_value_target = 1;  
    ELSE high_value_target = 0;  
)
```

#### 5.4.3 Selecting the Aangiften: Value Models Score

If the aangifte is scored by the whole aangifte model, then it is also scored by the value model.

```
IF whole_aangifte_score = 1  
THEN (  
    low_value_score = 1  
    high_value_score = 1  
)  
ELSE (  
    low_value_score = 0  
    high_value_score = 0  
)
```

The scores will be used to select between one of three strings that will be part of the information shown to the work divider:

- *'Little to No Value Predicted'*
  - Neither value model fires.
- *'Low to Medium Value Predicted'*
  - The low value model fires, but the high value model does not.
- *'Medium to High Value Predicted'*
  - The high value model fires.

Exactly what probability counts as 'firing' will be decided based on the model's performance metrics.

### 5.5 Complexity Measures

This is extra information added to the scored and ranked output to help the work dividers allocate the aangifte to the appropriate auditor.

#### 5.5.1 Complexity Scores





This is a set of 10 yes/no flags, one for each of the ten subjects. This tells the work divider and the auditor which parts of the aangifte have been filled in.

### 5.5.2 Specialist Areas

A list of words that highlight particular meta-codes which require a specialist auditor. This list is still being defined.



## 6 Appendix A: Structure of the Model Build Key Table

Columns in the key table extracted from aang\_vord\_aans\_totaal\_new3

- sofi
- belastingjaar
  - should be 2011
- ind\_winst
  - should be 'Niet Winst'
- MIDDEL
  - should be 'IH'
- wijze\_afdoen
- RISCIOCAT
- UWB
- type\_afhandeling
  - should be 'VA, DA'
- bedrag\_bel\_DA
  - should be > 0

Columns in the key table generated by the IH team

- 1) whole\_aangifte\_model
  - 0/1 flag, 1 means include in the whole aangifte model build
- 2) whole\_aangifte\_target
  - 0/1 flag, 1 means it's a hit for the whole aangifte model
- 3) value\_model
  - 0/1 flag, 1 means to include in the value model build
- 4) low\_value\_target
  - 0/1 flag, 1 if means it's a hit for the low value model
- 5) high\_value\_target
  - 0/1 flag, 1 if means it's a hit for the high value model
- 6) <subject>\_has\_data
  - 0/1 flag, 1 means there is non-zero data for this subject
- 7) <subject>\_checked
  - 0/1 flag, 1 means that there is a UWB pointing at that subject
- 8) <subject>\_edited
  - 0/1 flag, 1 means that there was a change made as a result of the audit
- 9) <subject>\_model
  - 0/1 flag, 1 means include in that subject model build
- 10) <subject>\_target
  - 0/1 flag, 1 means it's a hit for that subject

See Section 5 for details of how to create these.



## 7 Appendix B: Structure of the Model Score Key Table

- `sofi`
- `whole_aangifte_score`
  - 0/1 flag, 1 if the `sofi` should be scored by the whole aangifte model
- `low_value_score`
  - 0/1 flag, 1 if the `sofi` should be scored by the low value model
- `high_value_score`
  - 0/1 flag, 1 if the `sofi` should be scored by the high value model
- `<subject>_score`
  - 0/1 flag, 1 if the `sofi` should be scored by the `<subject>` model

## 8 Appendix C: Structure of the Model Score Output Table

- `sofi`
- `complexity_score`
  - still being defined
- `specialist_areas`
  - still being defined
- `whole_aangifte_score`
  - 0/1 flag, 1 if the `sofi` was scored by the whole aangifte model
- `whole_aangifte_model_prob`
  - value between 0 and 1, the result of scoring with the whole aangifte model
- `whole_aangifte_rank`
  - value between 0 and 99, the result of ranking the whole aangifte model scores
- `value_score`
  - 0/1 flag, 1 if the `sofi` was scored by the value model
- `low_value_prob`
  - value between 0 and 1, the result of scoring with the low value model
- `high_value_prob`
  - value between 0 and 1, the result of scoring with the high value model
- `value_string`
  - text that contains the value prediction
    - 'Little to No Value Predicted'
    - 'Some to Medium Value Predicted'
    - 'Medium to High Value Predicted'
- `<subject>_score`
  - 0/1 flag, 1 if the `sofi` was scored by the `<subject>` model
- `<subject>_prob`
  - value between 0 and 1, the result of scoring with the `<subject>` model
- `<subject>_rank`
  - value between 0 and 99, the result of ranking the `<subject>` model scores



### 9 Appendix D: Value Model Statistics

Details of which aangiften were included in the metrics:

- Tax year 2011.
- 'Overige Uitworp' manual audits only.
- Excluded 'Haarde Uitworp' as our model will not affect those cases.
- Positive corrections only.
- Where there is a DA followed by a VA as that's the only time we can calculate the change in tax as a result of the audit.

Notes	Centile	Value
<b>Lowest Percentile</b>	<b>1%</b>	Inspectie, controle en toezicht
	<b>5%</b>	
	<b>10%</b>	
<b>Correction Policy Threshold</b>	<b>15%</b>	
	<b>25%</b>	
<b>Median</b>	<b>50%</b>	
	<b>75%</b>	
	<b>90%</b>	
	<b>95%</b>	
<b>Highest Percentile</b>	<b>99%</b>	

Table 70: Value Statistics



8.8 Appendix 8: Reference 9: Persoonsgegevens “Belastingdienst IH Risk Model: Standardised Methods,” //IH selectie/Documentatie/Software process/.

**Belastingdienst IH Risk Model**  
**SAS Standard Methods**  
**Version: 2.0**

**DRAFT**



**Table of Contents**

**1. OVERVIEW ..... 165**

**2. METACODE MEASURE ..... 166**

    2.1 NUMERIC METACODES ..... 166

    2.2 CHARACTER METACODES ..... 168

    2.3 DATE METACODES ..... 169

    2.4 YESNO METACODES ..... 169

**3. STANDARD IMPLEMENTATION ..... 170**

    3.1 CONSISTENCY CHECK ..... 170

    3.2 COMPLEXITY CHECK ..... 170

    3.3 VA COMPARISON CHECK ..... 170

    3.4 PRI-YEAR COMPARISON CHECK ..... 171

    3.5 CONTRA COMPARISON CHECK ..... 172

**APPENDIX A: M\_STD\_CONSISTENCY.SAS ..... 173**

**APPENDIX B: M\_STD\_COMPLEXITY.SAS ..... 173**

**APPENDIX C: M\_NUM\_COMPARE.SAS ..... 173**

**APPENDIX D: M\_CNT\_COMPARE.SAS ..... 173**



### Revision History

Date	Version	Description	Author
23/05/2014	1.0	Creation	Persoonsgegevens
14/01/2015	2.0	Creation	Persoonsgegevens



## 1 Overview

In InkomstenHeffingen (IH) project, standard measures are defined and implemented. These standard measures are the inputs to the Analytical Base Tables (ABT) which are used in risk scoring.

This document tracks the definition of the standard measures and the implementation of them in the Statistical Analysis Software, i.e. SAS.





## 2 Metacode measure

Tax return form consists of metacodes. In general there are four categories of metacodes. For each category, standard measures are defined.

### 2.1 Numeric metacodes

An example metacode; salary income.

The following measures are defined:

Measure	Definition
COUNT	number of non-zero, non-missing volnummers (NB: grouped numeric fields only)
SUM	contents over volnummers
FLAG	if there is a pre-calculated total field then compare SUM to the pre-calculated total. Set flag to 1 if different OR 0 if totals are equal.

A tax return case, which is submitted by tax payer, goes through two phases before it is audited by a tax auditor. These two phases are called "Aangegeven" (A) and "Voorvastgesteld" (V). The COUNT, SUM and FLAG measures are created for each phase.

In addition, IH project uses tax return data from two consecutive years. The COUNT, SUM and FLAG measures are created for both years.



Measures below are defined when comparing numeric metacodes:

Measure	Definition
ABSOLUTE DIFFERENCE	<ul style="list-style-type: none"> <li>• 0 if V count = A count</li> <li>• Positive value if V count &gt; A count</li> <li>• Negative value if V count &lt; A count</li> </ul> <ul style="list-style-type: none"> <li>• 0 if V sum = A sum</li> <li>• Positive value if V sum &gt; A sum</li> <li>• Negative value if V sum &lt; A sum</li> </ul> <ul style="list-style-type: none"> <li>• 0 if current year count = previous year count</li> <li>• Positive value if current year count &gt; previous year count</li> <li>• Negative value if current year count &lt; previous year count</li> </ul> <ul style="list-style-type: none"> <li>• 0 if current year sum = previous year sum</li> <li>• Positive value if current year sum &gt; previous year sum</li> <li>• Negative value if current year sum &lt; previous year sum</li> </ul>
RATIO	<ul style="list-style-type: none"> <li>• 1 if V count = A count</li> <li>• &lt; 1 if V count &lt; A count</li> <li>• &gt; 1 if V count &gt; A count</li> </ul> <ul style="list-style-type: none"> <li>• 1 if V sum = A sum</li> <li>• &lt; 1 if V &lt; A sum</li> <li>• &gt; 1 if V sum &gt; A sum</li> </ul> <ul style="list-style-type: none"> <li>• 1 if current year count = previous year count</li> <li>• &lt; 1 if current year count &lt; previous year count</li> <li>• &gt; 1 if current year count &gt; previous year count</li> </ul> <ul style="list-style-type: none"> <li>• 1 if current year sum = previous year sum</li> <li>• &lt; 1 if current year sum &lt; previous year sum</li> <li>• &gt; 1 if current year sum &gt; previous year sum</li> </ul>
FLAG	<ul style="list-style-type: none"> <li>• Compare count V to count A Binary flag: 0 if counts are the same, 1 otherwise</li> <li>• Compare sum V to sum A Binary flag: 0 if counts are the same, 1 otherwise</li> <li>• Compare count current year to count previous year Binary flag: 0 if counts are the same, 1 otherwise</li> <li>• Compare sum current year to sum previous year Binary flag: 0 if counts are the same, 1 otherwise</li> </ul>



## 2.2 Character metacodes

An example metacode: employer.

NOTE: In our data all character metacodes are 0 if present as we don't have the actual character data.

The following measures are defined:

Measure	Definition
COUNT	number of non-zero, non-missing volnummers (NB: grouped numeric fields only)

The COUNT measure is created for both phase V and phase A and for two consecutive years.

Measures below are defined when comparing the COUNT measures of the character metacodes:

Measure	Definition
FLAG	<ul style="list-style-type: none"> <li>Compare count V to count A Binary flag: 0 if counts are the same, 1 otherwise</li> <li>Compare count two years Binary flag: 0 if counts are the same, 1 otherwise</li> </ul>
ABSOLUTE DIFFERENCE	<ul style="list-style-type: none"> <li>0 if V count = A count</li> <li>Positive value if V count &gt; A count</li> <li>Negative value if V count &lt; A count</li> </ul> <ul style="list-style-type: none"> <li>0 if current year count = previous year count</li> <li>Positive value if current year count &gt; previous year count</li> <li>Negative value if current year count &lt; previous year count</li> </ul>
RATIO	<ul style="list-style-type: none"> <li>1 if V count = A count</li> <li>&lt; 1 if V count &lt; A count</li> <li>&gt; 1 V count &gt; A count</li> </ul> <ul style="list-style-type: none"> <li>1 if current year count = previous year count</li> <li>&lt; 1 if current year count &lt; previous year count</li> <li>&gt; 1 current year count &gt; previous year count</li> </ul>



### 2.3 Date Metacodes

An example: date of birth.

The following measures are defined:

Measure	Definition
FLAG	Set flag to 1 if date filed is filled in, 0 otherwise.

The FLAG measure is created for two consecutive years.

Measures below are defined when comparing the FLAG measures of the character metacodes:

Measure	Definition
FLAG	<ul style="list-style-type: none"> <li>Compare flag two years</li> </ul> Binary flag: 0 if flags are the same, 1 otherwise

### 2.4 Yesno Metacodes

An example: whether have a fiscal partner.

The following measures are defined:

Measure	Definition
FLAG	Set flag to 1 if yesno filed is filled in, 0 otherwise.

The FLAG measure is created for two consecutive years.

Measures below are defined when comparing the FLAG measures of the character metacodes:

Measure	Definition
FLAG	<ul style="list-style-type: none"> <li>Compare flag two years</li> </ul> Binary flag: 0 if flags are the same, 1 otherwise



### 3 Standard implementation

The defined measures in chapter 2 are implemented as SAS Macros. SAS Macro is a piece of codes that deliver a functionality and can be called and executed repeatedly. SAS Macro alleviate the burden of rewriting SAS codes while ensures the consistency in data processing.

#### 3.1 Consistency check

In tax return form, some information can be grouped together using metacodes. For example, the following metacodes are grouped together because they are related to income from employment: employer name, withheld salary tax and salary. It is important to check whether the grouped metacodes have the same number of filled-in instances. Consistency check is implemented to perform the check.

Consistency check uses the pre-defined measures and check whether there is a consistency across the phase A data and phase V data, and across two consecutive data. Consistency checks are implemented to numeric, character metacodes.

m\_std\_consistency.sas (appendix a) is the SAS code file containing the actual implementation of the check. It starts with putting grouped metacodes in an array. Then counts of each metacode in a group are compared, e.g. number of employers versus number of salaries. In addition, calculated total is compared with the pre-calculated total, e.g. calculated total salaries versus total salary from employment.

Outcome of consistency check is binary. If phase A data is consistent with phase B data then flag is 0 else is 1.

#### 3.2 Complexity check

In order to assess the complexity across tax return cases, four measures are defined and implemented. These measures are based on numeric metacodes.

Complexity measures are:

Measure	Definition
Total metacodes	Number of used numeric metacodes
Total metacodes level volnummer	Number of used numeric metacodes at the level of volnummer
Total corrected metacodes level volnummer	Number of changed metacodes comparing phase V with phase A
Percentage of corrected metacodes level volnummer	Percentage of changed metacodes comparing phase V with phase A

m\_std\_complexity.sas (appendix b) is the SAS code file containing the actual implementation of the check. The outcomes of this check are the measures as described above.

#### 3.3 VA comparison check

Comparison between phase V data and phase A data generates insights about the quality of a tax return case. Phase A data are the information that are provided by tax payers. Phase V data are the information after the Belastingdienst fiscal check.



There are two types of VA comparison check: one is applied to numeric metacode, the other is to character metacode.

For each type of VA comparison check, it starts with merging phase A data with phase B data. Then comparison check is performed to the selected metacodes variables. `m_num_compare.sas` and `m_cnt_compare.sas` are the SAS code files containing the actual implementation of the check (appendix c and appendix d). The check takes always two streams data as input for comparison, namely phase A data and phase V data.

The outcomes of the VA comparison check are the ones described in chapter 2.1 and 2.2.

### 3.4 Pri-year comparison check

Als van een belastingplichtige de aangifte van het huidige jaar (sterk) verschilt van de aangifte van het vorige jaar, kan dit duiden op een fout in de aangifte. Als iemand bijvoorbeeld dit jaar een veel hoger salaris heeft aangegeven dan vorig jaar, is de kans groter dat de aangifte een fout bevat dan als het salaris beide jaren gelijk is. Daarom is het belangrijk om een vergelijking te maken tussen de aangifedata van het huidige jaar en de aangifedata van het vorige jaar.

Pri-year comparison check gebruikt de voorgedefinieerde measures en vergelijkt de aangifedata van het huidige jaar met de aangifedata van het vorige jaar. Pri-year comparison checks zijn geïmplementeerd voor numerieke en karakter metacodes. Daarbij wordt er onderscheid gemaakt tussen variabelen die gaan over de numerieke waarde van een metacode en variabelen die gaan over het aantal volgnummers (de count) van een metacode.

`m_cnt_compare.sas` is het SAS codebestand dat de implementatie van de check voor de variabelen die over de count van een metacode gaan, bevat. Bij de aanroep van de macro wordt een count-variabele van het huidige jaar en de corresponderende count-variabele van het vorige jaar meegegeven. De macro creëert nieuwe variabelen van het type flag, ratio en absolute verschil, op de wijze zoals voorgaand beschreven.

`m_num_compare.sas` is het SAS codebestand dat de implementatie van de check voor de variabelen die over de numerieke waarde van een metacode gaan, bevat. Bij de aanroep van de macro wordt een variabele van dit type van het huidige jaar en de corresponderende variabele van het vorige jaar meegegeven. De macro creëert nieuwe variabelen van het type flag, ratio en absolute verschil, op de wijze zoals voorgaand beschreven.



### 3.5 Contra comparison check

De Belastingdienst heeft niet alleen de beschikking over de data die wordt aangegeven door belastingplichtigen, maar ook over informatie vanuit verschillende externe bronnen, zoals bank- en loongegevens. Deze informatie wordt *contradata* genoemd. Het is belangrijk om een vergelijking te maken tussen deze *contradata* en de waarden op het aangifteformulier. Als bijvoorbeeld iemand een veel lager banksaldo heeft aangegeven dan het saldo dat volgt uit de bankgegevens van de belastingplichtige, dan is het waarschijnlijk dat het aangegeven saldo onjuist is.

In de *contra comparison check* worden de benodigde gegevens uit de *contradata* gehaald en middels het SAS codebestand `m_num_compare_contra.sas` vergeleken met de aangiftegegevens. Bij de aanroep van de macro wordt een voorvastgestelde waarde van het aangifteformulier en de corresponderende waarde uit de *contradata* meegegeven. De macro creëert nieuwe variabelen van het type *flag*, *ratio* en *absolute verschil*, op de wijze zoals voorgaand beschreven.

De *contra-gegevens* die in het IH-risicomodel gebruikt worden, komen uit verschillende bronnen, zoals RIS, FLG en RBG. Een exacte beschrijving van welke bronnen gebruikt zijn, is te vinden in het document `Belastingdienst_IH_Data_Source`.



4 Appendix A: m\_std\_consistency.sas



m\_std\_consistency.sas

5 Appendix B: m\_std\_complexity.sas



m\_std\_complexity.sas

6 Appendix C: m\_num\_compare.sas



m\_num\_compare.sas

7 Appendix D: m\_cnt\_compare.sas



m\_num\_compare.sas





**8.9 Appendix 9: Reference 10:** Persoonsgegevens **“Belastingdienst IH Risk Model: Data Sources,” //IH selectie/Documentie/Software process/.**

**Belastingdienst IH Risk Model  
Data Source  
Version: 1.2**



**Table of Contents**

<b>1. OVERVIEW .....</b>	<b>177</b>
<b>2. TYPES OF DATA .....</b>	<b>178</b>
2.1 DATA CATEGORY .....	178
2.2 DATA USAGE .....	179
<b>3. DATA DELIVERY .....</b>	<b>179</b>
3.1 DATA REQUESTER .....	179
3.2 DATA PROVIDER .....	179
3.3 DATA LOCATION.....	180
3.4 DATA UPDATE.....	180
<b>APPENDIX A: SAMPLE CONTRA DATA REQUEST LIST .....</b>	<b>181</b>
<b>APPENDIX B: CONTACTS DATA REQUESTERS .....</b>	<b>182</b>
<b>APPENDIX C: B/CA BICC CONTACTS TAX RETURN DATA.....</b>	<b>183</b>
<b>APPENDIX D: CONTRA DATA CONTACT .....</b>	<b>184</b>
<b>APPENDIX E: TAX RETURN DATA SPECIFICATION.....</b>	<b>185</b>



**Revision History**

Date	Version	Description	Author
22/01/2015	1.0	Creation	Persoonsgegevens



## 1 Overview

A data source is simply the source of the data. It can be a flat file, a SAS dataset, or a particular database on a DBMS.

The purpose of this document is to define the data sources and the guidelines that are (to be) applied during the analysis, design and build of the Inkomstenheffing (IH) risk model solution.



## 2 Types of Data

In IH risk model project, there are two categories of data being used for analysis. These are tax return data and contra data.

### 2.1 Data Category

#### Tax return data

Tax return data is the information about tax payers in the field of income, property, social benefits., etc. Tax return data is gathered by Belastingdienst digitally and in the form of paper.

IH project uses tax return data from the fiscal year 2011 and 2010. This is because that there are sufficient amount of data for the purpose of statistical modelling.

#### Contra data

Contra data are information about tax payer's insurance, bank deposit., etc. Contra data is delivered by third-party organization, for example, banks delivers savings information about tax payers. List below shows the contra data in IH project.

Contra name	Domain	Used (Y/N)
Eigen Woning Lijfrente Sparen	primary house, annuity	Y
Fiscal Partners	fiscal partner	Y
FLG	fiscal salary information	Y
OV-verklaringen	public transportation	Y
RBG	bank savings	Y
Verzekeringsproducten	insurance product	Y
Wajong	young disabled person benefit	Y
WEP	financial investment	Y
WOZ-base	primary house	Y
Convenant aangifte	metadata about tax return filing	N
Way of filling tax form (paper, system etc.)	metadata about tax return filing	N
VIA aangifte downloaden yes/no	metadata about tax return filing	N
Supplier of tax filling programme	metadata about tax return filing	N
Number of aangifte submitted until selection	metadata about tax return filing	N
Spontane aangifte yes/no	metadata about tax return filing	N



## 2.2 Data Usage

Tax return data and contra data are used to create Analytical Base Tables (ABT). In turn, ABTs are used for risk scoring. ABT contains variables which are derived from tax return data and contra data, as well as comparing tax return data with contra data.

Standard methods are applied to data to create the variables in ABT. Below are the listed usage of data:

- Aangegeven and Voorvastgesteld data comparison

Tax return data are stored in the phase “Aangegeven” and phase “Voorvastgesteld”. “Aangegeven” data tells what information are submitted by tax payers. “Voorvastgesteld” data tells what are the corrected information after the fiscal check. Fiscal check is performed by Belastingdienst. Comparing “Aangegeven” with “Voorvastgesteld” indicates the quality of a tax return case.

- Current-year and previous-year comparison

Comparing two years' tax return data in order to capture the changes between them.

- Tax return data and contra data comparison

Comparing tax payers' data with contra data in order to verify the accuracy of the data provided by tax payers.

## 3 Data Delivery

### 3.1 Data Requester

initiate data requests on behalf of IH project team to data provider.  are business analysts from the team Business Intelligence & Analytics in Belastingdienst.

In appendix B, more information about data requester is given.

### 3.2 Data Provider

#### **Tax return data**

Data is provided by the team Belastingdienst/Centrale Administratie Business Intelligence Competence Center (B/CA BICC ) Analytics & Datamining.

In appendix C, there is list of contacts of the team B/CA BICC.

#### **Contra data**

Contra data is provided by  from the team InformatieVoorziening.

In appendix D, contact information is given.



### 3.3 Data Location

#### **Tax return data**

Tax return data is located in Teradata database DL\_BEL\_BIA\_ONTW\_PMI\_IH and are stored in the following 4 tables:

- ABS\_WAARDE
- ABS\_HEFFINGSZAAK
- ABS\_INKOMEND\_BERICHT
- ABS\_BEHANDELDOSSIER

In Appendix E: Tax return data specification, specification information about the four tables such as column name, column type, column format., etc are attached.

#### **Contra data**

Contra data is located on the SAS server i.e. apandu05 and can be found in the following paths:

- /AD005/shared/datasets/PDV/FLG
  - FLG fiscal salary information
- /AD005/shared/datasets/PDV/RIS
  - Eigen Woning Lijfrente Sparen
  - OV-verklaringen
  - RBG
  - Verzekeringsproducten
  - Wajong
  - WEP
- /AD005/shared/datasets/PDV/RISWOZ
  - WOZ-base primary house

### 3.4 Data Update

When IH project solutions are used for a project, for example, scoring all backlog cases in the fiscal year of 2013, both tax return data and contra data are required to be refreshed. In general, two fiscal years' tax return data are required.

The data updating process is carried out between data requesters and data providers. Data requesters put data refreshing request formally in a mail and send it to data providers. Data providers answer the initial request with an estimation about lead time of data delivery. The mails from data providers are always assigned or associated with unique service codes so that data requesters can use them to communicate with data providers.

Data providers deliver data always to the data locations as described in section 3.3.



## 4 Appendix A: Sample Contra Data Request List



Overzicht uitvraag  
data v0.3.xlsx





## 5 Appendix B: Contacts Data Requesters Business Intelligence & Analytics

### Address:

Persoonsgegevens

Persoonsgegevens



## 6 Appendix C: B/CA BICC Contacts Tax Return Data

**Business Intelligence Competence Center  
Belastingdienst / Centrale Administratie**

**Address:**

Persoonsgegevens

**E-mail:**

Persoonsgegevens

**Team members:**

Persoonsgegevens



## 7 Appendix D: Contra Data Contact

Belastingdienst Centrale Administratie  
Belastingdienst Informatie Competence Center

Inspectie, controle en toezicht

Persoonsgegevens



## 8 Appendix E: Tax return data specification



teradata aangifte  
tables specifications.)



8.10 Appendix 10: Reference 11: Persoonsgegevens, "IH Risk Value Model: Descriptive Statistics," /IH selectie/Models/ValueModels/.



Programma Broedkamer  
*Innovatie Projectbureau*



# Metrics for the Value of an IH Niet Winst Case

Risicomodel IH

Versie: 0.2

3-12-2014, Utrecht



## Can The IH Team Replicate the Value in the Benefit Tracking?



### Which Cases?

- tax type = IH, niet winst
- path = VA followed by a DA
  - *otherwise we can't correctly calculate the change in tax due to an audit*
- positive corrections only
- audit type = other
  - *overige uitworp*
- remove the VIP cases
- using these selection criteria gets us close to the benefit tracking value

### Values

- number of cases  
Inspectie, controle en toezicht
- mean (average) value
  - Inspectie, controle en toezicht
- median value\*  
Inspectie, controle en toezicht
  - *this is much lower than the mean, telling us that the mean is not representative*
- maximum value
  - Inspectie, controle en toezicht
  - *this large value has a large effect on the mean, making it higher*

\* the median is value lying at the midpoint of a frequency distribution of observed values, such that there is an equal probability of falling above or below it



## Not All Years Are Equal

Tax year
2006
2007
2008
2009
2010
2011
2012
2013

Inspectie, controle en toezicht

- Early tax years have too small a sample size to be statistically meaningful and their values are biased (*see panel on the left*).
- 2012 and 2013 are not fully audited yet, so are also unrepresentative.
- This leaves **2010** and **2011** as sensible years to base a measure on.

### Cases Opened Prior to 1 Jan 2011

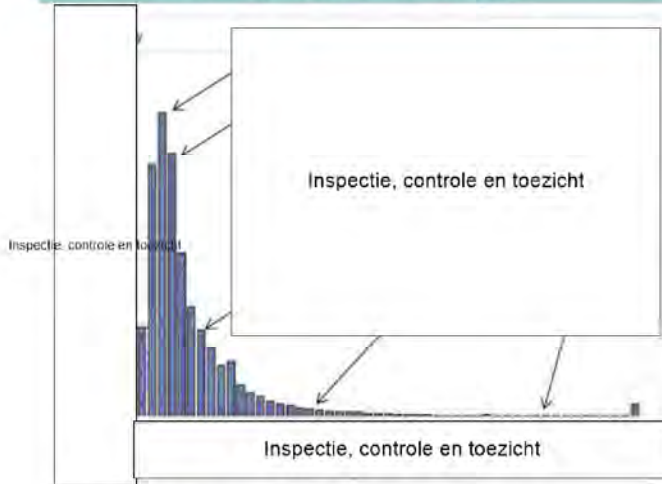
- The value of cases was derived from the collections domain, Incasso
- Incasso started recording cases from the 1st of January 2011
  - i.e. first full tax year is 2010
- All cases with status 'open' on or after 1 Jan 2011 were saved in Incasso.
- Relatively few cases from earlier years were still 'open' on 1 Jan 2011 as they were more likely to have been closed and paid for.
- Old cases that were 'open' when Incasso started tended to be cases with a high value, probably because the taxpayer might find these harder to pay off or might need more time.
- Smaller value cases are often easily paid, so are less likely to be open for several years.

These counts are just the overige uitworp





# Mean, Median, Mode and Outliers



### Which Cases?

- tax type = IH, niet winst
- path = VA followed by a DA
  - otherwise we can't correctly calculate the change in tax due to an audit
- tax years 2010 and 2011
- positive corrections only
- audit type = manual
  - **overige uitworp**
  - **harde uitworp**

### Values

- Number of cases:
- Mean value:
- Median value:
- Maximum value:

- The data has a very long tail in the high values
- There are some very, very high values, or outliers, that we can't expect to be repeated year after year

of the values are under   
 are under

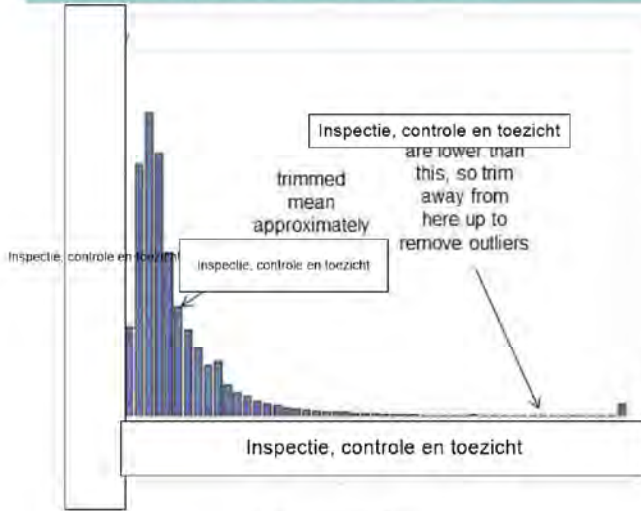
\* the mode is the value that occurs most frequently  
 † this is a larger number of cases than previously, as we're looking at all manually audited cases







# A More Representative Measure



### Problems with the Distribution

- The data has a very long tail in the high values
- There are some very, very high values, or outliers, that we can't expect to be repeated year after year
- 95% of the values are under
- 99% are under

### Solution:

#### A More Representative Measure

- Trim away the top
- That's cases with
- Then re-calculate the mean
- This is called a Trimmed Mean
- **Trimmed Mean value** 
  - or **approximate**

An average value of  is what we intend to present  
 as part of the process of deciding on the appropriate threshold for the IH niet winst value model.



**8.11 Appendix 11: Reference 15: Programma Broedkamer, "Programboard IH v0.99 21-8-2014," /IH Selectie/Communicatie/Program Board 21-8-2014**



Programma Broedkamer  
*Innovatie Projectbureau*



Programboard IH

Risicomodel IH

Versie: 1.0

21-8-2014, Utrecht

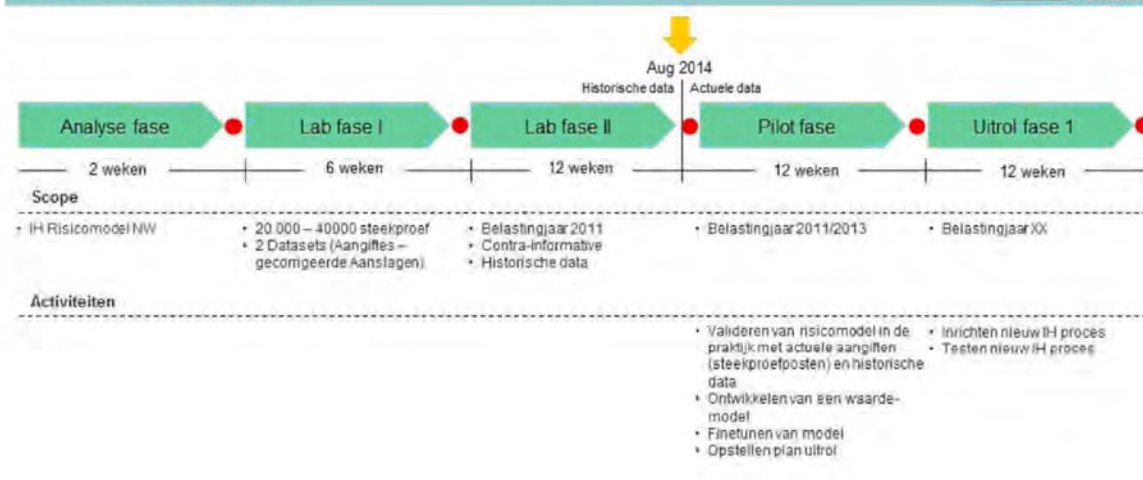


Agenda  

Status project Risicomodel IH
Resultaat Labfase 2 <input data-bbox="715 622 983 667" type="text" value="Persoonsgegevens"/> ns
Vervolg stappen
Stakeholders en betrokkenen
Gevraagde acties en besluiten



# Status en voortgang





## Heads up: Wat is het Risicomodul IH

Door gebruik te maken van statistische methoden kunnen we nieuwe inzichten/patronen vinden in de data. Hierdoor kunnen we gericht selecteren met een hogere hitrate op positieve correcties.

Heeft kinderen	Heeft partner	Heeft meer dan 1 baan	Is ondernemer	Vraagt korting alleenst ouder	Correctie
1	1	0	0	0	1
0	1	0	1	0	0
0	0	1	0	1	1
1	0	1	0	0	1
0	0	0	1	1	1
1	0	0	0	1	1
0	1	0	0	0	0
0	1	1	1	0	0
0	0	0	0	0	0
0	0	1	1	0	0
0	1	0	1	0	1
0	1	0	0	1	0
0	0	1	1	1	1
1	1	1	0	0	1

Specifiek

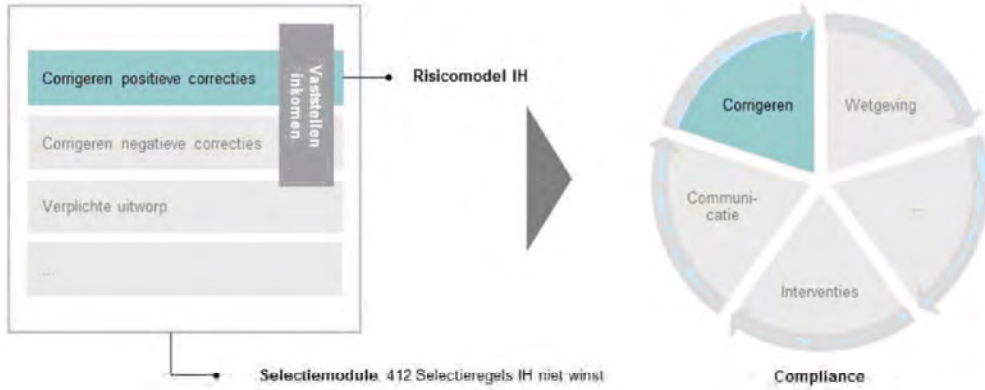
$$P(\text{hit}) = f_0^1 (\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n)$$



### Focus Risicomodel IH binnen de doelstellingen van IH Toezicht



Het is van belang om te onderstrepen op welke doelstelling binnen IH Toezicht het risicomodel zich op dit moment focust: het corrigeren van positieve correcties. Dit is slechts één van de instrumenten om fiscale regelnaleving te stimuleren.





Agenda

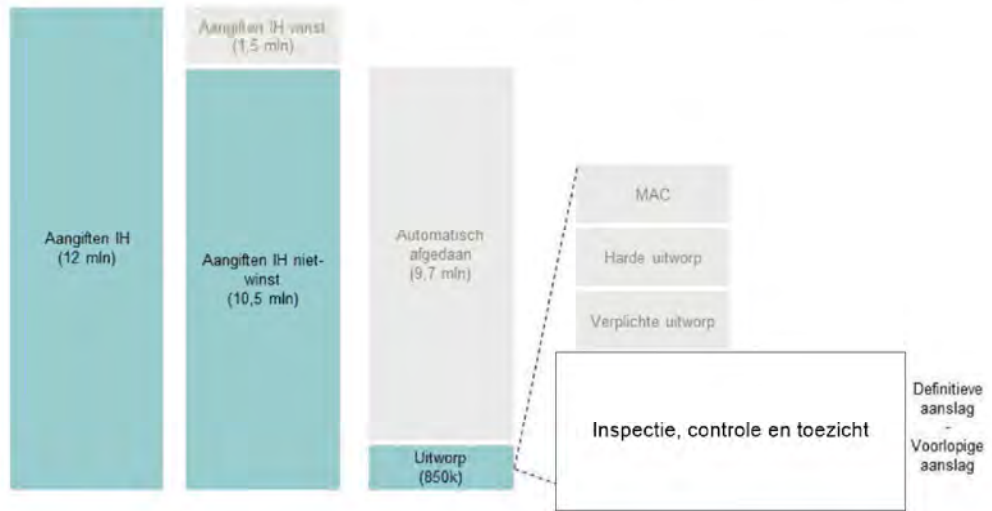
Status project Risicomodel IH	
Resultaat Labfase 2	Persoonsgegevens
Vervolg stappen	
Stakeholders en betrokkenen	
Gevraagde acties en besluiten	



### Inzoomen op de scope van het risicomodel



Om het resultaat van Labfase 2 te kunnen beoordelen is het belangrijk te begrijpen wat wel en niet in scope is van het risicomodel en wat de bijbehorende huidige hitrate is. In blauw is weergegeven wat in scope is.







## Ontwikkeling van risicomodellen voor aangiftesecties



In Labfase 2 is het Risicomodel IH verder ontwikkeld en is tevens gewerkt aan het toewijzen van risico scores op 10 aandachtsgebieden, die zijn opgesteld op basis van de aangiftesecties.

Gehele aangifte



Alle invulvelden van het aangifte formulier zijn verwerkt in het risicomodel voor de gehele aangifte.

Aandachtsgebieden



- Inkomen
- Reisaftrek
- Uitgaven & budgetten
- Zorgkosten
- Overige inkomsten
- Eigen woning
- Bezittingen & schulden
- Buitenland
- Fiscaal partnerschap
- Heffingskortingen

Alle invulvelden van het aangifte formulier zijn tevens verwerkt in één of meerdere modellen van de aandachtsgebieden.



### Risico score per aandachtsgebied

Iedere aangifte kan worden voorzien van een risico score voor de gehele aangifte en een risico score per aandachtsgebied. Op basis van deze scores wordt de behandelinstructie voor een aangifte bepaald.

BSN-nummer: XXXX	
Datum: XXXX	
Maximale score: 93	
<b>Gehele aangiften</b>	
Gehele aangiften	91
Aandachtsgebied 1: Inkomen	32
Aandachtsgebied 2: Reisaftrek	11
Aandachtsgebied 3: Uitgaven & budgetten	93
Aandachtsgebied 4: Zorgkosten	32
Aandachtsgebied 5: Overige inkomsten	50
Aandachtsgebied 6: Eigen woning	26
Aandachtsgebied 7: Bezittingen & schulden	72
Aandachtsgebied 8: Buitenland	88
Aandachtsgebied 9: Fiscaal partnerschap	13
Aandachtsgebied 10: Heffingskortingen	69

**Voorbeeld:**  
 Behandelinstructie: Bekijk  
 aandachtsgebied Uitgaven &  
 Budgetten en aandachtsgebied  
 Buitenland



### Prioriteren van de aangiften

De output van het IH risicomodel is een grote tabel, met daarin een risico score voor de gehele aangifte én een risico score per aandachtsgebied. De aangiften kunnen worden gesorteerd op risico score van hoog naar laag, waarna de behandeling op basis van deze sortering plaats kan vinden.

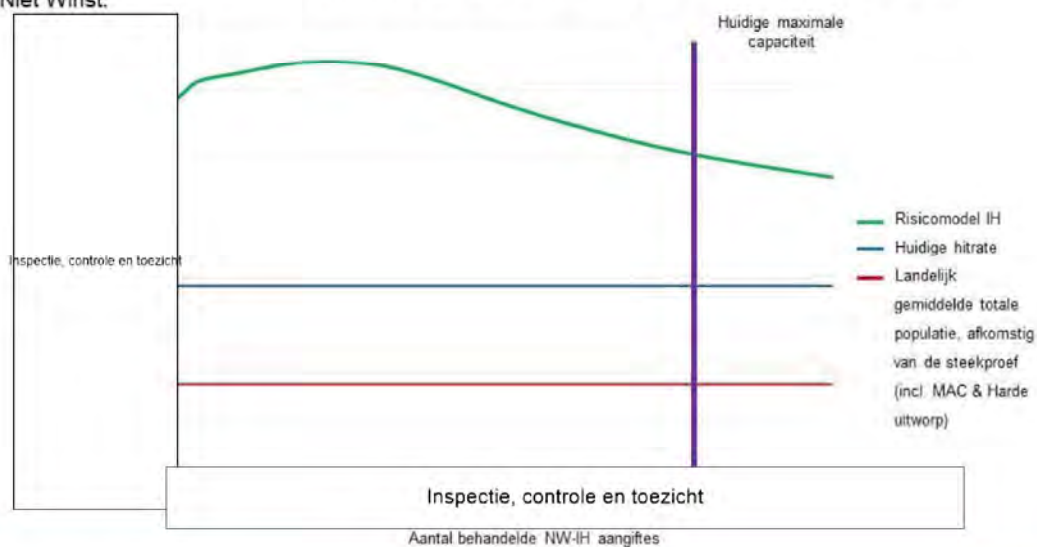
BSN	Maximale score	Gehele aangifte	Inkomen	Reisafrek	Uitgaven & Budgetten	Zorgkosten	Overige Inkomsten	Eigen woning	Bezittingen & schulden	Buitenland	Fiscaal partnerschap	Specifiek Heffingen korting
###	81	90	79	23	15	81	3	55	19	91	15	27
###	90	85	54	17	87	34	90	29	41	44	78	25
###	88	88	11	23	38	87	37	9	83	65	36	81
###	74	73	87	15	9	12	53	87	71	74	13	5
###	73	71	73	61	34	58	32	88	34	22	80	21
---	---	---	---	---	---	---	---	---	---	---	---	---
---	---	---	---	---	---	---	---	---	---	---	---	---
---	---	---	---	---	---	---	---	---	---	---	---	---
###	3	2	3	1	1	2	1	2	2	3	1	2



Resultaat hitrate IH op positieve correcties na Labfase 2



Het eindresultaat van Labfase 2 is de groene lijn. Deze lijn geeft aan wat de risico score voor de gehele aangifte op een positieve correctie is, gerelateerd aan het aantal te behandelen aangiftes IH Niet Winst.





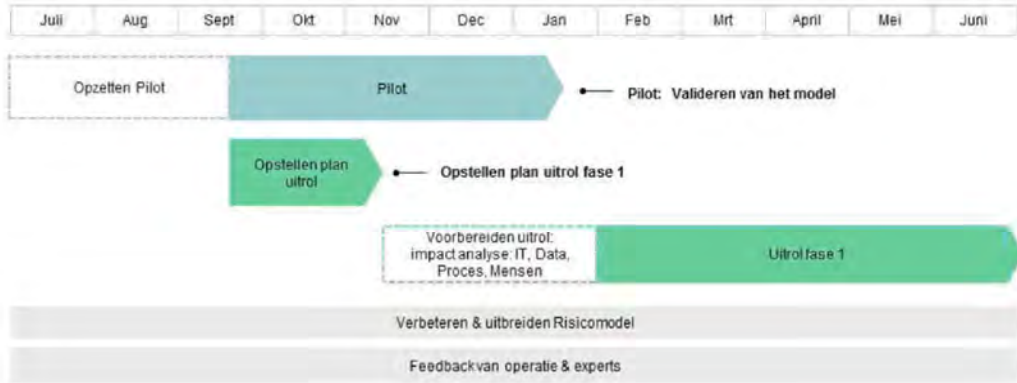
Agenda  

Status project Risicomodel IH	
Resultaat Labfase 2 <table border="1"><tr><td>Persoonsgegevens</td></tr></table>	Persoonsgegevens
Persoonsgegevens	
Vervolg stappen	
Stakeholders en betrokkenen	
Gevraagde acties en besluiten	



# Planning

Het risicomodel IH wordt gevalideerd in de praktijk, hiervoor starten we een Pilot met actuele aangiften en een aantal historische aangiften.







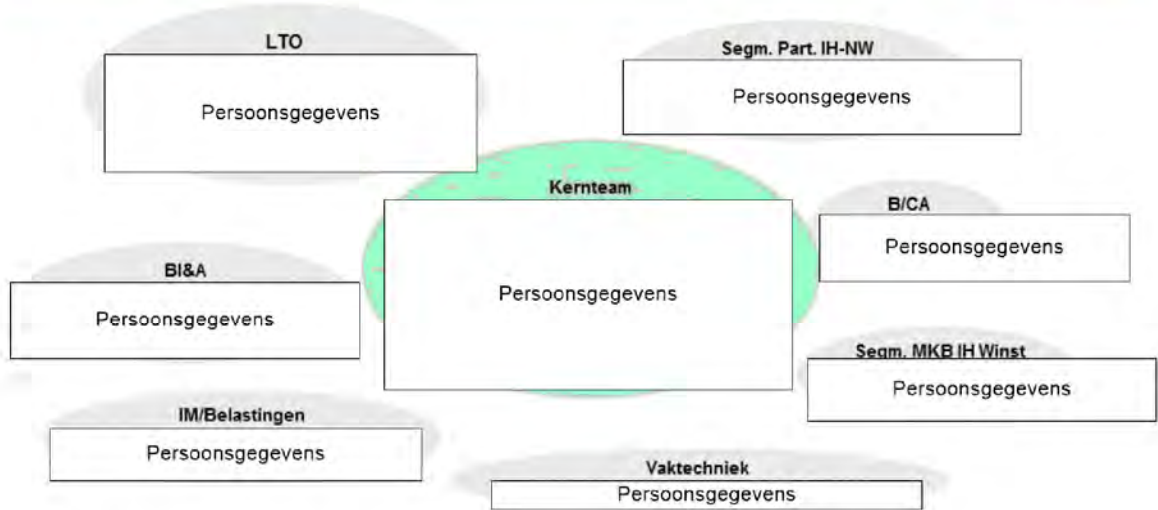
Agenda  

- Status project Risicomodel IH
- Resultaat Labfase 2  s
- Vervolg stappen
- Stakeholders en betrokkenen**
- Gevraagde acties en besluiten





Stakeholders





Agenda



Status project Risicomodel IH

Resultaat Labfase 2

Vervolg stappen

Stakeholders en betrokkenen

Gevraagde acties en besluiten



## Gevraagde acties en besluiten

### Gevraagd Besluit:

“Start met Pilot”

zijnde:

- Start met Pilot ter validatie van het risicomodel op 22 september met historische aangiften en actuele aangiften.

Voorwaarden waaraan wordt voldaan:

- De hitrate van het Risicomodel is significant hoger dan de huidige hitrate op positieve correcties.
- Er is een pilot team Particulieren IH Niet-Winst beschikbaar op 1 locatie.
- Er is een teamlead beschikbaar om samen de Pilot te coördineren.
- De steekproefaangiften zijn tijdig beschikbaar.
- De resultaten van de steekproef kunnen worden teruggevoerd in ABS.



## Gevraagde acties en besluiten

### Gevraagd Besluit:

“Start scoren huidige Uitworp voorraad”

---

#### zijnde;

- Het ontwikkelen en uitvoeren van een plan voor het voorzien van de huidige Uitworp voorraad van een risico score, om de prioritering te vergemakkelijken.

#### Benodigd;

- LTO expertise
- B/CA dataleverantie omtrent IH niet winst



## Gevraagde acties en besluiten

### Gevraagd Besluit:

"Start voorbereiding implementatie fase 1"

zijnde;

- Het opstellen van een implementatieplan voor uitrol.

Benodigd;

- LTO expertise
- B/CA dataleverantie omtrent IH niet winst
- Procesexpertise IH (IM)
- Vaktechniek expertise
- B/CA IT expertise



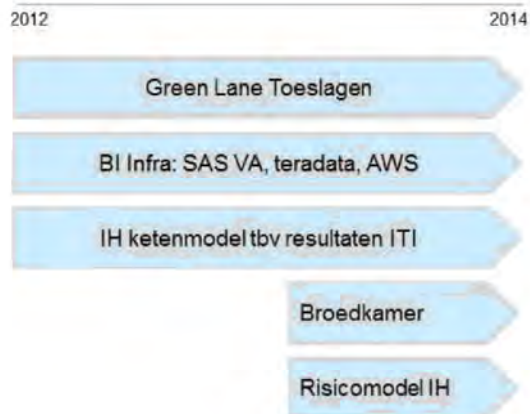
# Appendix



## Context &amp; Vraagstuk



## Context



## Waarom een Risicomodel IH?

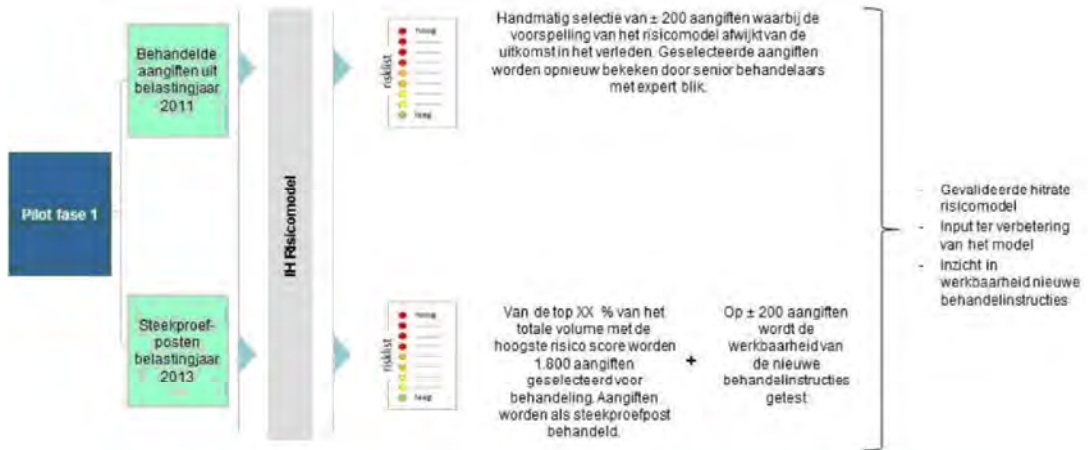
1. We willen gericht selecteren met een hogere hit-rate op positieve correcties.
2. We willen onze capaciteit efficiënter en effectiever inzetten; hiervoor is een pull v push method nodig.
3. We willen onze medewerkers waardevol en nuttig werk bieden.
4. We willen minder afhankelijkheid van de transactionele systemen (ABS)
5. We willen de kosten rondom selectie verlagen.
6. Omdat we het nu kunnen!



Opzet Pilot 1 Optioneel

**Doelstelling:**

- Valideren voorspelling risicomodel (hitrate %)
- Verbeteren risicomodel
- Testen werkbaarheid nieuwe behandelingsinstructies

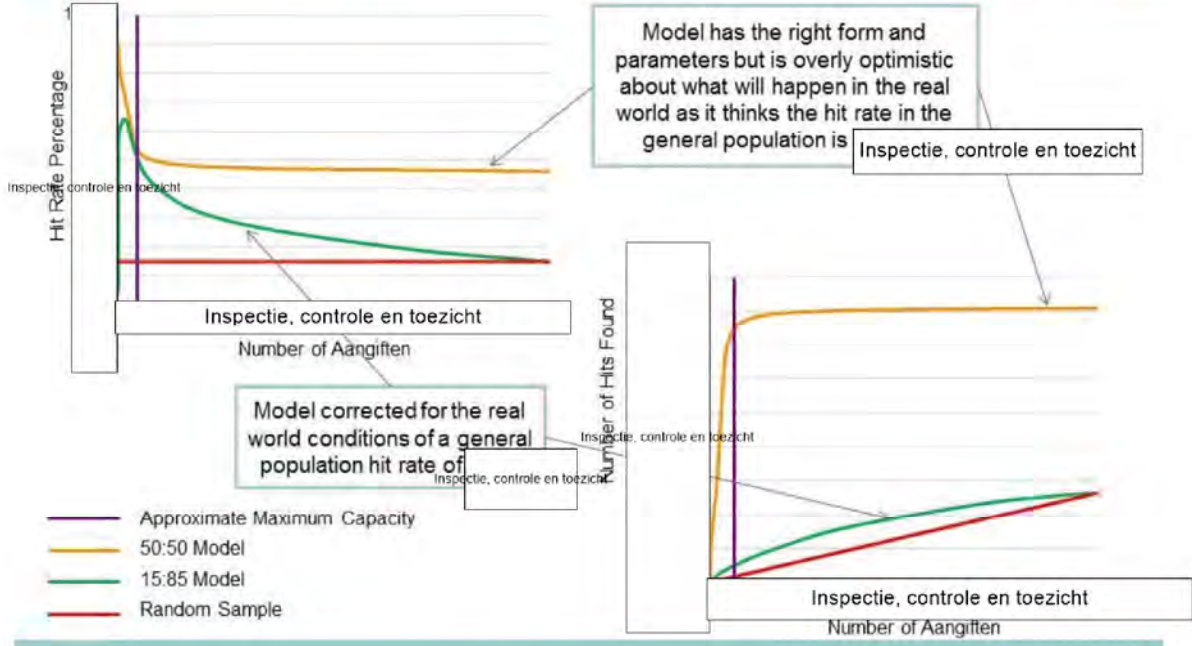


VERTROUWELIJK – NIET VERSPREIDEN 21-8-2014, Utrecht





### Risicomodel gehele aangifte



VERTROUWELIJK – NIET VERSPREIDEN 19/05/2015



8.12 Appendix 12: Reference 28: Persoonsgegevens “Belastingdienst IH Risk Model: Testing the New Work Instructions,” //IH Selectie/Pilot/Testen nieuwe behandelinstructies/140901 - Instructie presentatie testen nieuwe behandelinstructies.ppt.



Programma Broedkamer  
*Innovatie Projectbureau*



# Instructie testen nieuwe behandelinstructies

Risicomodel IH

Versie: 1.0

29-9-2014, Utrecht



## Agenda



Project Risicomodel IH

Opzet pilot

Planning

Instructies testen nieuwe behandelinstructies



# Wat is het Risicomodel IH

Door gebruik te maken van statistische methoden kunnen we nieuwe inzichten/patronen vinden in de data. Hierdoor kunnen we gericht aangiften selecteren met een hogere hitrate op positieve correcties.

Heeft kinderen	Heeft partner	Heeft meer dan 1 baan	Is ondernemer	Heeft een kind alleenstaande ouder	Erkenningscode
1	1	0	0	0	1
0	1	0	1	0	0
0	0	1	0	1	1
1	0	1	0	0	1
0	0	0	1	1	1
1	0	0	0	1	1
0	1	0	0	0	0
0	1	1	1	0	0
0	0	0	0	0	0
0	0	1	1	0	0
0	1	0	1	0	1
0	1	0	0	1	0
0	0	1	1	1	1
1	1	1	0	0	1

Specimen

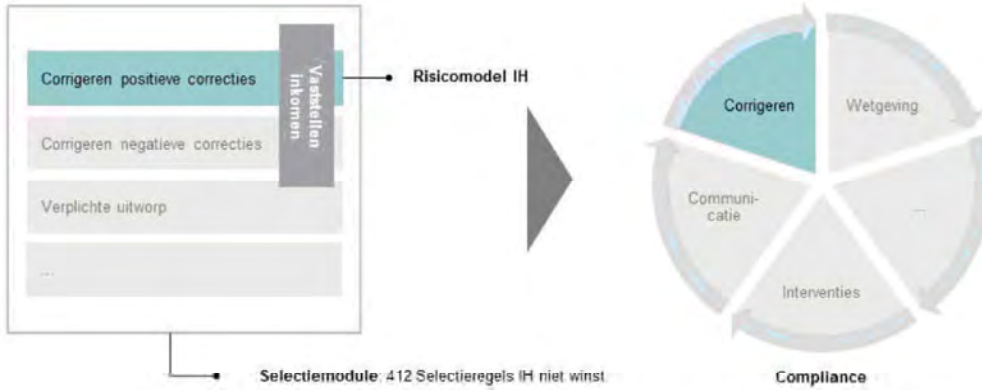
$$P(\text{hit}) = f \frac{1}{0} (\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n)$$



### Focus Risicomodel IH binnen de doelstellingen van IH Toezicht



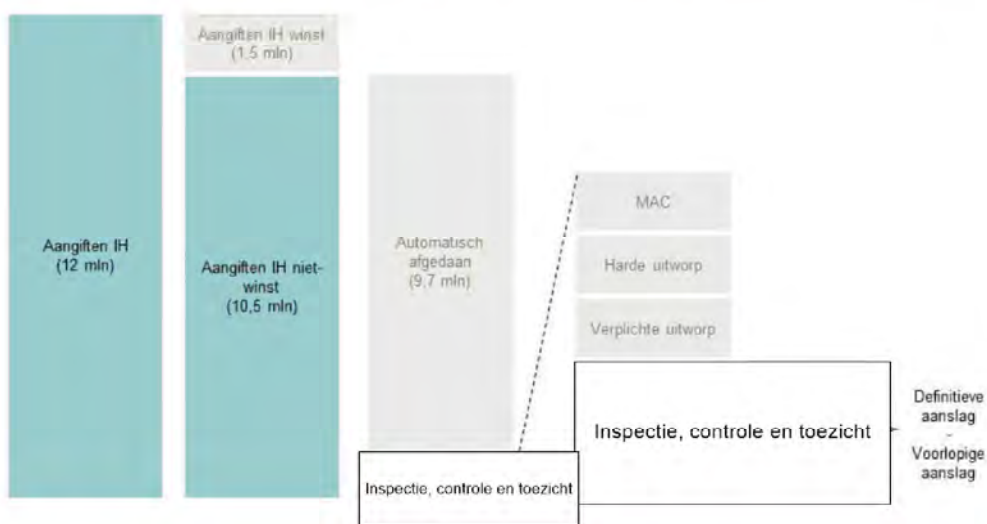
Het is van belang om te onderstrepen op welke doelstelling binnen IH Toezicht het risicomodel zich op dit moment focust: het corrigeren van positieve correcties. Dit is slechts één van de instrumenten om fiscale regel naleving te stimuleren.





# Inzoomen op de scope van het risicomodel

De scope van het huidige risicomodel is weergegeven in blauw. Het risicomodel richt zich op de wat op dit moment 'Overige uitwerp' wordt genoemd.

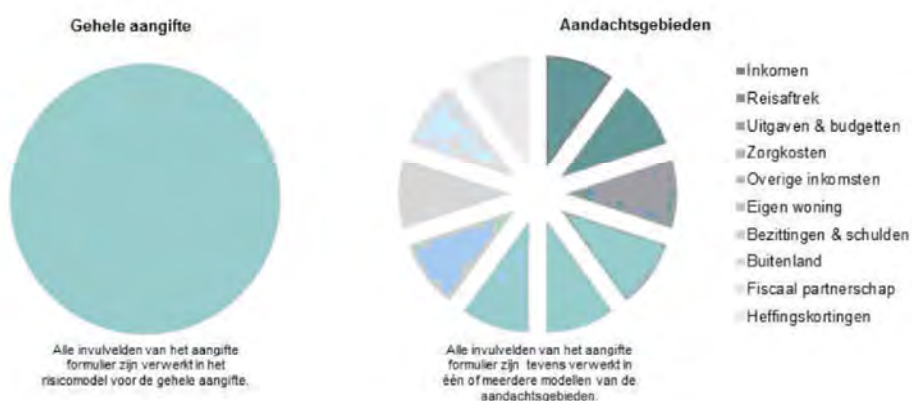




## Ontwikkeling van risicomodellen voor aangiftesecties



Op dit moment is er een Risicomodel ontwikkeld en is tevens gewerkt aan het toewijzen van risico scores op 10 aandachtsgebieden\*, die zijn opgesteld op basis van de aangiftesecties. Aandachtsgebieden zijn gegroepeerde rubrieken uit het aangifte formulier waarop een voorspelling kan worden gedaan.





### Prioriteren van de aangiften

De output van het IH risicomodel is een grote tabel, met daarin een risico score voor de gehele aangifte én een risico score per aandachtsgebied. De aangiften kunnen worden gesorteerd op risico score van hoog naar laag, waarna de behandeling op basis van deze sortering plaats kan vinden.

BSN	Maximale score	Gehele aangifte	Inkomen	Reisafrek.	Uitgaven & Budgetten	Zorgkosten	Overige inkomsten	Eigen woning	Bezittingen & schulden	Buitenland	Freuaal partnerschap	Hebba kortingen
***	87	50	76	23	15	81	3	55	19	91	15	27
***	91	89	54	17	87	34	90	29	41	44	78	25
***	56	80	11	23	39	87	37	9	83	65	38	91
***	74	73	87	15	9	12	53	87	71	74	13	5
***	73	71	73	81	34	58	32	88	34	22	80	21
---	---	---	---	---	---	---	---	---	---	---	---	---
---	---	---	---	---	---	---	---	---	---	---	---	---
---	---	---	---	---	---	---	---	---	---	---	---	---
***	3	2	3	1	1	2	1	2	2	3	1	2

Voorbeeld





# Risico score per aandachtsgebied

Iedere aangifte kan worden voorzien van een risico score voor de gehele aangifte en een risico score per aandachtsgebied. Op basis van deze scores wordt de behandelingsinstructie voor een aangifte bepaald.

BSN-nummer: XXXX	
Datum: XXXX	
Maximale score: 93	
<b>Gehele aangiften</b>	
Gehele aangiften	69
Aandachtsgebied 1 Inkomen	32
Aandachtsgebied 2 Reisaftrek	11
Aandachtsgebied 3 Uitgaven & budgetten	93
Aandachtsgebied 4 Zorgkosten	32
Aandachtsgebied 5 Overige inkomsten	50
Aandachtsgebied 6 Eigen woning	26
Aandachtsgebied 7 Bezittingen & schulden	72
Aandachtsgebied 8 Buitenland	88
Aandachtsgebied 9 Fiscaal partnerschap	13
Aandachtsgebied 10 Heffingskortingen	69

**Voorbeeld:**  
 Behandelingsinstructie: Bekijk  
 aandachtsgebied Uitgaven &  
 Budgetten en aandachtsgebied  
 Buitenland.



## Agenda



Project Risicomodel IH

Opzet pilot

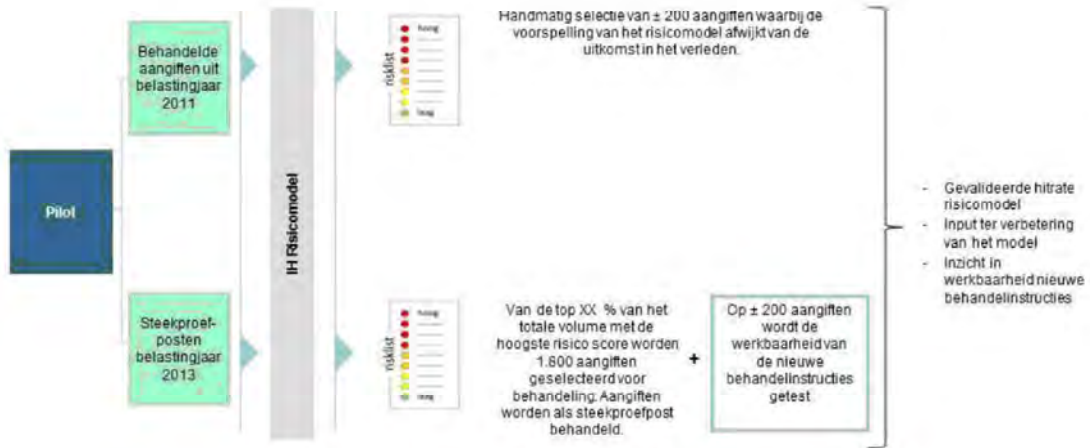
Planning

Instructies herbehandeling Historische aangiften



# Opzet Pilot

De pilot kent 3 doelstellingen: valideren van het risicomodel, het verbeteren van het risicomodel en het testen van de nieuwe behandelinstructies.



VERTROUWELIJK – NIET VERSPREIDEN 21-8-2014, Utrecht



## Testen van de nieuwe behandelinstructies



Voor 200 aangiften willen we de nieuwe behandelinstructies uitproberen om te achterhalen wat wel werkt en wat niet werkt.

### Doelstellingen testen nieuwe behandelinstructies:

- Testen van werkbaarheid behandelinstructies
- Inzicht verkrijgen in verschil in behandelijd met nieuwe behandelinstructies
- Inzicht verkrijgen in werkbaarheid groepering aandachtsgebieden
- Inzicht verkrijgen in de juistheid van de complexiteitsindicatie
- Analyseren van impact implementatie nieuwe behandelinstructies
- Inzicht verkrijgen in de mate waarin behandelinstructies worden gevolgd.

De behandelinstructies liggen nog niet vast, hier kan nog alles aan worden aangepast. Feedback is dus heel belangrijk!



## Testen van de nieuwe behandelinstructies



Er zijn 3 soorten nieuwe behandelinstructies. Tijdens de Pilot testen we 2 van de 3 behandelinstructies: nr. 2 en nr. 3.

### Behandelinstructies:





## Agenda



Project Risicomodel IH

Opzet pilot

Planning

Instructies herbehandeling Historische aangiften



# Planning

De planning is om na iedere 50-75 behandelde steekproefposten die zijn geselecteerd voor het testen van de nieuwe behandelinstructies een evaluatiemoment in te plannen, om de bevindingen en ervaringen te bespreken die zijn opgedaan.

Planning testen nieuwe behandelinstructies





## Agenda



Project Risicomodel IH

Opzet pilot

Planning

Instructies herbehandeling Historische aangiften





## Instructies



**Stap 1:** Bekijk het Excel-bestand met de lijst met BSN-nummers waarvoor de nieuwe behandelinstructies worden getest.

- **Ieder Excel-bestand bevat 2 tabbladen:**
  - **Tabblad 1 – Invulsheet:** Hierin staan de BSN-nummers vermeld, de behandelinstructies, de bijbehorende aandachtsgebieden en rubrieken en de moeilijkheidsgraad.
  - **Tabblad 2 – Risico scores:** Overzicht van de risico scores op de verschillende aandachtsgebieden en het gehele model.

Voorbeeld:





## Instructies



**Stap 2:** Bekijk de behandelinstructie voor het BSN-nummer dat je gaat behandelen op tabblad 1.

**VOORBEELD**

BSN	Behandelinstructie	Aandachtsgebied	Rubrieken	Moeilijkheidsgraad	Welke moeilijkheidsgraad zou jij deze aangifte geven?	Opmerking en
XXX	Begin bij de aangewezen aandachtsgebieden, besteed hier de meeste aandacht aan en bekijk vervolgens de gehele aangifte	Gehele aangifte	Alle rubrieken			
XXX	Bekijk de aangewezen aandachtsgebieden	Persoonlijke gegevens	1			
XXX	Begin bij de aangewezen aandachtsgebieden, besteed hier de meeste aandacht aan en bekijk vervolgens de gehele aangifte	Bezittingen en schulden	19, 20, 28, 29, 30, 31, 47			
XXX	Bekijk de aangewezen aandachtsgebieden	Uitgaven en budgetten	25, 32, 33, 34, 35, 36, 37, 38, 39, 40			
XXX	Bekijk de aangewezen aandachtsgebieden	Buitenland	15, 16, 48, 52			
XXX	Begin bij de aangewezen aandachtsgebieden, besteed hier de meeste aandacht aan en bekijk vervolgens de gehele aangifte	Zorgkosten	34, 53			

Behandelinstructie voor de aangifte.

Voor de namen van de verschillende rubrieken, zie het uitgereikte document met daarop de groepering van rubrieken naar aandachtsgebieden

De moeilijkheidsgraad is bepaald op basis van XXXX



## Instructies



### Stap 3: Bekijk vervolgens de aangifte in ABS en volg de nieuwe behandelinstructie

- **Er zijn twee soorten instructies:**
  1. Bekijk de aangewezen aandachtsgebieden
  2. Begin bij de aangewezen aandachtsgebieden, besteed hier de meeste aandacht aan en bekijk vervolgens de gehele aangifte
- Om te weten welke rubrieken onder het betreffende aandachtsgebied vallen, zie het uitgereikte document met daarop de groepering van rubrieken naar aandachtsgebieden.
- Mocht er eventueel extra informatie nodig zijn van de belastingplichtige voor het behandelen van de aangewezen aandachtsgebieden, wacht dan met het opvragen van deze informatie tot dat je de gehele aangifte hebt bekeken, zodat alles in één brief kan worden opgevraagd.
- Breng eventuele correcties aan in de aangewezen aandachtsgebieden in de behandelinstructie.



## Instructies



**Stap 4:** Indien je dit nog niet gedaan had, controleer de aangifte als een steekproefpost en breng eventuele correcties aan in ABS en SPAR applicatie.

- Volg de normale procedure voor het afhandelen van een steekproefpost.
- Indien behandelinstructie 3 van toepassing was, is de aangifte al volledig gecontroleerd en kan deze stap worden overgeslagen.



## Instructies



**Stap 5:** Bekijk na de behandeling de moeilijkheidsgraad in het Excel-bestand en geef aan welke moeilijkheidsgraad jij deze aangifte zou geven.

BSN	Behandelinstructie	Aandachtsgebied	Rubrieken	Moeilijkheidsgraad	Welke moeilijkheidsgraad zou jij deze aangifte geven?	Opmerking en
XXX	Begin bij de aangewezen aandachtsgebieden en bekijk vervolgens de gehele aangifte	Gehele aangifte	Alle rubrieken			
XXX	Bekijk de aangewezen aandachtsgebieden	Persoonlijke gegevens	1			
XXX	Begin bij de aangewezen aandachtsgebieden en bekijk vervolgens de gehele aangifte	Bezittingen en schulden	19, 20, 28, 29, 30, 31, 47			
XXX	Bekijk de aangewezen aandachtsgebieden	Uitgaven en budgetten	25, 32, 33, 34, 35, 36, 37, 38, 39, 40			
XXX	Bekijk de aangewezen aandachtsgebieden	Buitenland	15, 16, 48, 52			
XXX	Begin bij de aangewezen aandachtsgebieden en bekijk vervolgens de gehele aangifte	Zorgkosten	34, 53			
XXX	Begin bij de aangewezen aandachtsgebieden en bekijk vervolgens de gehele aangifte	Overige inkomsten	18, 22, 23, 24, 48, 50, 51			
XXX	Begin bij de aangewezen aandachtsgebieden en bekijk vervolgens de gehele aangifte	Eigen woning	21			

Dit is de verwachte moeilijkheidsgraad die we willen toetsen

Hoe moeilijk vind jij deze aangifte op de schaal XXX?



## Instructies



**Stap 6:** Vul na afloop van het behandelen van alle BSN-nummer in het Excel-bestand de eind-vragenlijst in.

- Het zijn in totaal 14 meerkeuze vragen.
- Deze vragen gaan onder andere over:
  - Behandeltijd met de nieuwe behandelinstructies
  - Werkbaarheid van de aandachtsgebieden
  - Inzicht in risicoscores
- Tevens is er ruimte voor opmerkingen. Elke opmerkingen is zeer waardevol!
- De resultaten worden besproken tijdens een van de evaluatiemomenten tussen de behandeling van de aangiften door.



**8.13 Appendix 13: Reference 44 IH Risk Model Team, “Programma Broedkamer, “Programboard IH 14-1-2015” .,” /IH Selectie/Communicatie/Program Board 14-1-2015**



Programma Broedkamer  
*Innovatie Projectbureau*



Program Board IH

Risicomodel IH

Versie: 1.0

14 januari 2015, Utrecht



## Agenda



Status project risicomodel IH

Resultaten Pilot 1

Vervolgstappen project risicomodel IH

Resultaten IMB – Scenario's integreren risicomodel in IH-keten

Gevraagde besluiten







## Agenda



Status project risicomodel IH

Resultaten Pilot 1

Vervolgstappen project risicomodel IH

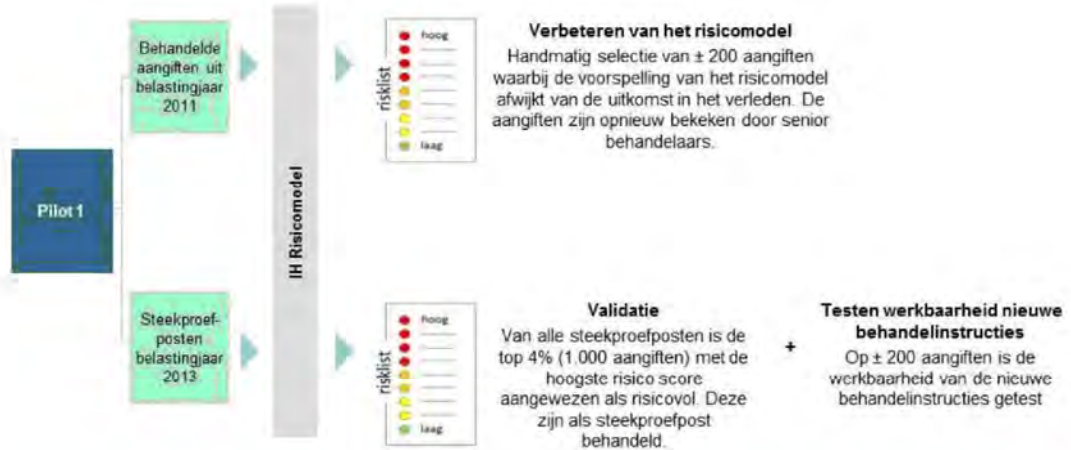
Resultaten IMB – Scenario's integreren risicomodel in IH-keten

Gevraagde besluiten



# Heads up Opzet Pilot 1

Pilot 1 bestond uit 3 onderdelen met als doel het valideren en verbeteren van het risicomodel en het opdoen van de eerste inzichten met betrekking tot de werkbaarheid van de nieuwe behandelinstructies.

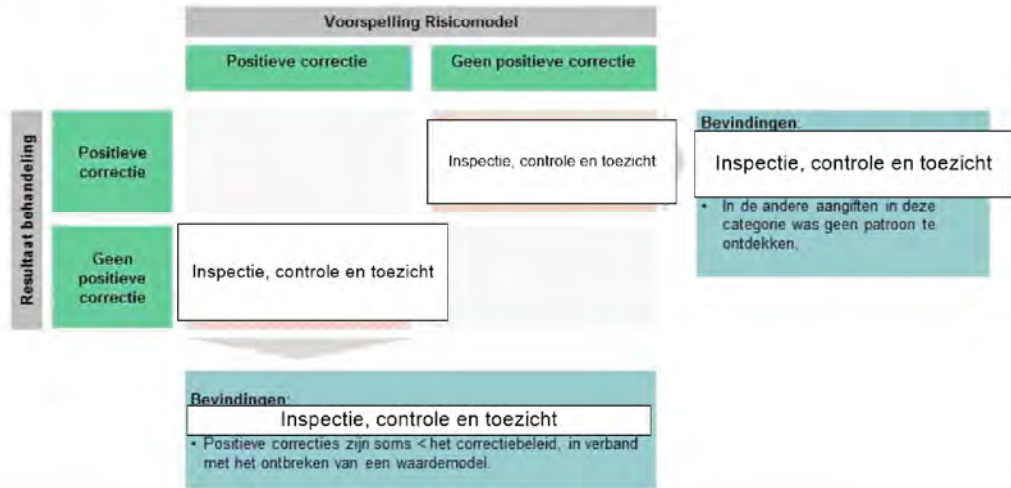




### Resultaten herbehandeling oude aangiften Belastingjaar 2011



Uit de herbehandeling is gebleken dat het risicomodel IH op dit moment geen belangrijke risicopatronen mist. Daarnaast heeft de input van de behandelaars geleid tot een aantal verbeteringen van het model.

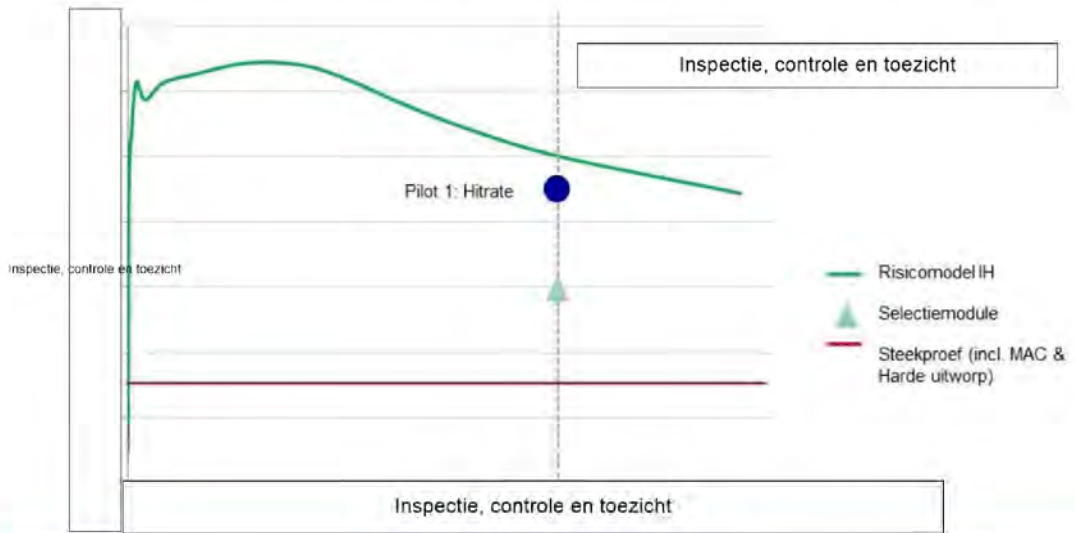




### Validatie Risicomodel IH: Resultaat Labfase 2



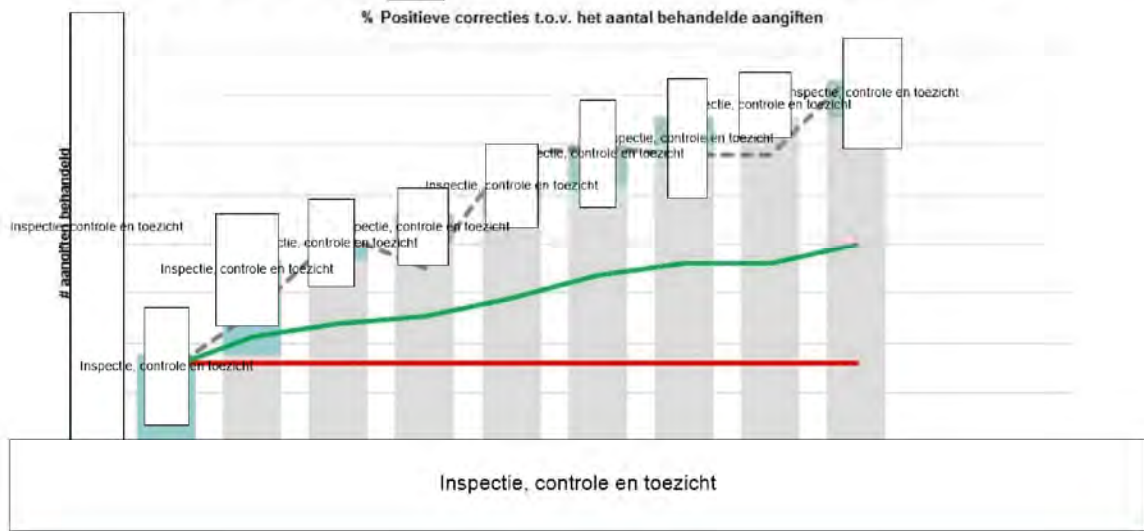
Labfase 2 is geëindigd met een verwachte risico score van  aangiften. Dit is een theoretische score, de ervaring leert dat risicomodellen in de praktijk iets lager scoren.





# Hitrate risicomodel IH

De hitrate van het Risicomodel IH komt naar verwachting uit binnen een range van **10% tot 20%**. Dit resultaat is bereikt met slechts **10%** van de contra-informatie t.o.v. de huidige selectiemodule.





## Resultaten Testen nieuwe behandelinstructies



De nieuwe behandelinstructies lijken aan te sluiten bij de huidige behandelwijze in de praktijk. Daarnaast heeft het testen van de nieuwe behandelinstructies geleid tot een nieuwe indeling van de aandachtsgebieden.

"Ik verwacht dat met het gebruik van deze aandachtsgebieden er meer gecorrigeerd gaat worden per aangifte, wat de compliantie ten goede komt"

*Behandelaar Pilot 1*

"Zonder overzicht van afwijkende contra-informatie, ben je continu aan het switchen tussen alle applicaties. De werkbaarheid van de nieuwe behandelinstructies van het risicomodel staat of valt met dit overzicht"

*Behandelaar Pilot 1*

"Ik heb de afgelopen jaren nog nooit zoveel correcties gevonden in steekproefposten"

*Behandelaar Pilot 1*



"Ik denk dat de meeste behandelaars liever integraal of op aandachtsgebied controleren dan enkel op uitwerpbewering"

*Behandelaar Pilot 1*



"Wanneer een behandelaar exact de huidige instructie van de uitwerpbewering opvolgt, zal het behandelen op aandachtsgebied iets meer tijd kosten. Aangezien niet iedereen exact deze instructie volgt, zal het minimaal zijn."

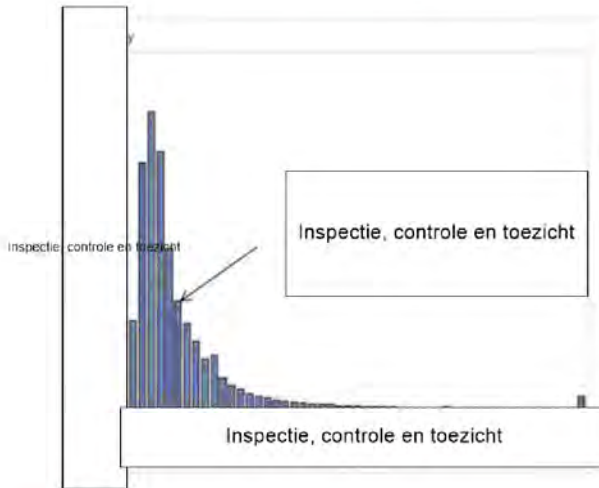
*Behandelaar Pilot 1*



Ontwikkeling van het waardemodel



Voor de ontwikkeling van het waardemodel is een verkenning van de data gedaan en als gevolg hiervan zijn twee drempelwaarden gekozen voor het waardemodel.



Cut-off percentage	Cut-off waarde	# aangiften
Geen cut-off		
Laagste waarde	1%	
	5%	
	10%	
	15%	
	25%	Inspectie, controle en toezicht
Mediaan	50%	
	75%	
	90%	
	95%	
Hoogste waarde	99%	





## Agenda



Status project risicomodel IH

Resultaten Pilot 1

Vervolgstappen project risicomodel IH

Resultaten IMB – Scenario's integreren risicomodel in IH-keten

Gevraagde besluiten



## Voorstel versnelling uitrol op enkelvoudig risico



Het voorstel is om op korte termijn al te starten met uitrol, naast Pilot 2, met aangiften met een enkelvoudig risico die fiscaal relatief makkelijk te behandelen zijn.

Risicomodel voor gehele aangifte



Alle invulvelden van het aangifte formulier zijn verwerkt in het risicomodel voor de gehele aangifte.

Submodellen voor onderdelen van de aangifte



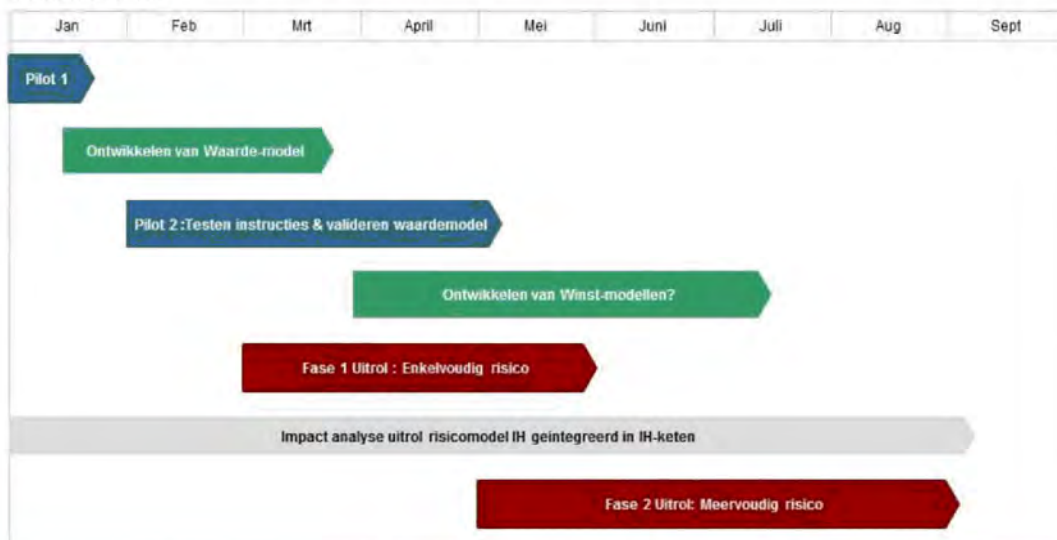
Alle invulvelden van het aangifte formulier zijn verwerkt in één of meerdere modellen van de aandachtsgebieden.



Tijdslijn implementatie & doorontwikkeling risicomodel IH



Onderstaande tijdslijn is de voorgestelde planning voor implementatie en doorontwikkeling van het risicomodel IH.





Agenda  

- Status project risicomodel IH
- Resultaten Pilot 1
- Vervolgstappen project risicomodel IH
- Resultaten IMB – Scenario's integreren risicomodel in IH-keten**
- Gevraagde besluiten



Voorgestelde werkwijze  
Belastingjaar 2014



Voor Belastingjaar 2014 is het voorstel om de uitworp van de selectiemodule te verrijken met de gegevens van het risicomodel en op basis daarvan keuzes te maken over wat te behandelen.

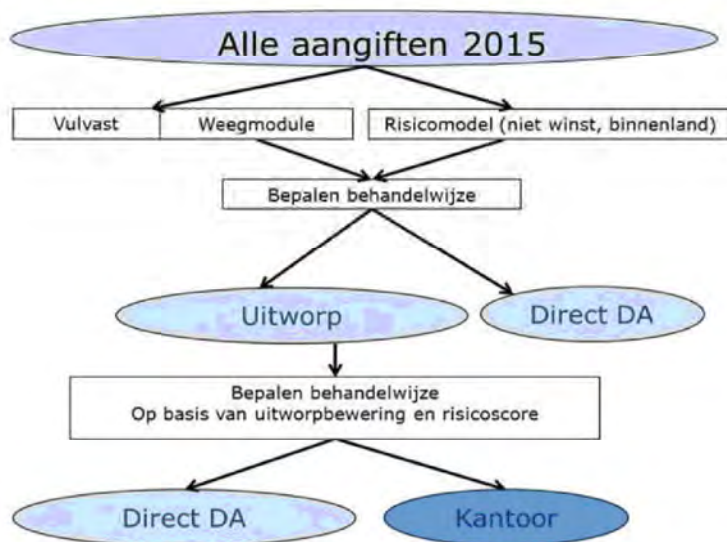




Voorgestelde werkwijze  
Belastingjaar 2015



Voor Belastingjaar 2015 is het voorstel om de selectiemodule en het risicomodel naast elkaar te laten draaien, waarbij deze elkaar aanvullen. Risico's worden voornamelijk uitgeleverd vanuit het risicomodel.





Agenda  

- Status project risicomodel IH
- Resultaten Pilot 1
- Vervolgstappen project risicomodel IH
- Resultaten IMB – Scenario's integreren risicomodel in IH-keten
- Gevraagde besluiten**



Gevraagd besluit # 1



Gevraagd Besluit:

"Akkoord voor start fase 1 uitrol: enkelvoudig risico op 1 maart 2015"

zijnde:

- Akkoord voor start uitrol op 1 of 2 locaties, bij voorkeur op de locatie Eindhoven, Den Haag en/of Arnhem
- Teamomvang van totaal 30 behandelaars met ervaring op 2 geselecteerde aandachtsgebieden
- Behandeling van 10.000 aangiften; er worden twee aandachtsgebieden met een hoge score geselecteerd waarbij van ieder aandachtsgebied 5.000 aangiften worden behandeld
- De uitrol van deze fase start op 1 maart 2015.





**8.14 Appendix 14: Reference 45 Persoonsgegevens "Belastingdienst IH Risk Model: Backlog 2013 - Definities en keuzes uitsluitingen LTO pilot2 uitrol1," //IH selectie/Backlog2013/.**

## **Belastingdienst IH Risk Model Backlog 2013 – Definities en keuzes uitsluitingen LTO pilot2 uitrol1**

### **Inhoud**

1. DEFINITIES EN KEUZES UITSLUITINGEN PER DEELPROJECT.....	255
2. LTO: PROCES SAMENSTELLING LIJST.....	257
3. UITROL 1: PROCES SAMENSTELLING LIJST.....	258



## 1 Definities en keuzes uitsluitingen per deelproject

Hieronder een overzicht van welke aangiften IH2013 er zijn uitgesloten of juist opgenomen in de diverse deelprojecten.\*

	<b>Definitie</b>	<b>LTO regulier 2013</b> <small>(Persoonsgegevens)</small>	<b>Pilot 2</b> <small>(Persoonsgegevens)</small>	<b>Uitrol 1</b> <small>(Persoonsgegevens)</small>
<b>Moeters (VIP / diverse anderen)</b>	Volgens de risicocategorieën die LTO heeft opgeleverd n.a.v. <a href="#">memo MT PDB</a> (ook te vinden in <a href="#">formele opdracht LTO</a> )	Uitsluiten van de selectie door het risicomodel	Niet uitsluiten, als er een hoog risico in zit pakken we hem mee	Niet uitsluiten, als er een hoog risico in zit pakken we hem mee
<b>Winst-posten</b>	<ul style="list-style-type: none"> <li>A) Volgens de risicocategorieën die LTO heeft opgeleverd OF</li> <li>B) Heeft winst rubriek (2 t/m 11) ingevuld in aangifte OF</li> <li>C) Heeft volgens ABS "winst aangifte" ingediend</li> </ul>	Uitsluiten van de selectie door het risicomodel, volgens definitie A	Uitgesloten volgens definitie B & C	Uitgesloten volgens definitie B & C
<b>C/M biljetten (buitenland)</b>	<ul style="list-style-type: none"> <li>A) Volgens de risicocategorieën die LTO heeft opgeleverd</li> <li>B) Volgens de uitworpbevestiging die hiervoor dient</li> <li>C) Heeft velden XXXX en YYYY van aangifteformulier ingevuld (deze velden hebben 'toevallig' dezelfde metacodes als velden van de winstrubrieken)</li> </ul>	Uitsluiten volgens definitie A (BETER IS VOLGENS DEFINITIE B ivm poetsen)	Uitgesloten volgens definitie C (maar zijn al uitgesloten, vanwege uitsluiten winst-posten)	Uitgesloten volgens definitie A en C (bij gebruik van uitsluitend definitie C blijken toch niet alle buitenland-aangiften eruit gehaald te worden)
<b>MAC</b>	A) Heeft UITSLUITEND nu geldende MAC uitworpbevestiging(en), eventueel in combinatie met UWB's die betrekking hebben op W-	Uitsluiten volgens definitie A	Momenteel niet uitgesloten	Uitsluiten volgens definitie B (BETER IS VOLGENS DEFINITIE A, maar helaas



	<p>en M-biljetten. De nu geldende MAC uitworpbeweringen zijn de volgende vijf:</p> <p>H0290 H0560 H0658 H0225 H0288 H0516</p> <p>B) Hetzelfde als A), behalve dat uitworpbewering H0516 niet als MAC wordt beschouwd (destijds dachten we dat het een potentieel MAC-UWB was, maar het bleek een 'echte' MAC-UWB te zijn. Deze vergissing is door LTO gemaakt).</p>			kwamen we hier te laat achter)
<b>Bekende fraude</b>	<p>A) Volgens de risicocategorieën die LTO heeft opgeleverd (deze maken ook deel uit van de risicocategorieën van Moeters)</p> <p>B) Volgens de uitworpbewering die hiervoor dient: H0912.</p>	Uitsluiten volgens definitie A (maar zijn al uitgesloten, vanwege uitsluiten Moeters)	Niet uitgesloten bij selectie, maar uiteindelijk worden deze posten (volgens definitie B) buiten de pilot om behandeld.	Uitgesloten volgens definitie B
<b>Posten 'niet behandelen' (VOW, steekproef, PGB en Box3)</b>	Volgens de risicocategorieën die LTO heeft opgeleverd.	Uitsluiten van de selectie door het risicomodel	Momenteel niet uitgesloten	Momenteel niet uitgesloten

\* Dit overzicht is nog niet volledig. De volgende acties staan nog open:

- 1 C/M biljetten buitenland scherp maken: allen
- 2 Winst afstemmen welke definitie we willen hanteren: allen
- 3 Afstemmen of we moeters in de pilot en uitrol mogen opnemen:

Persoonsgegevens



## 2 LTO: proces samenstelling lijst

- 1 LTO delivers a new (final) list to us of sofis that have been selected by the selection module.
- 2 We filter those sofis from our 3.1 million scored aangiftes.
- 3 We map per sofi the UWB's to the subjects and then exclude the cases where not any of the UWB's are pointing at a risky subject area.  
(RM\_advies\_geldig=1)
- 4 We also exclude 'Moeters' (for example fraud-cases), 'niet behandelen'-cases, Winst- and C- and M-forms, by the riskcategories that LTO says correspond with this.  
(ind\_bak12=0 AND  
ind\_winstLTO=0 AND  
ind\_CMbiljet=0 AND  
ind\_bak\_niet\_behandelen=0)
- 5 'Pure' MAC-cases, pilot2-cases, uitrol1-cases and the cases that haven't been scored (because they weren't in the 3.1 million list), are also excluded.  
(zuivere\_mac=0 AND  
niet\_gescoord=0 AND  
ind\_pilot=0 AND  
ind\_uitrol=0)
- 6 The final table will be saved and sent to LTO (containing only columns that are important for them).
- 7 Not only the final table will be saved, but also the complete table that's been created before excluding the cases mentioned above, consisting indicators that are used in excluding the cases (for example the indicator 'ind\_winstLTO' that indicates the winst-cases by LTO's definition). This table also contains indicators for the definitions used in Uitrol 1, so the table can also easily be used in selecting the cases for Uitrol 1.



### 3 Uitrol 1: proces samenstelling lijst

- 1) Use the table mentioned in step 7 of chapter 2.
- 2) Exclude the pilot2-cases, the cases that haven't been scored (because they weren't in the 3.1 million list), the 'pure' MAC-cases, the Winst-cases, the C- and M-forms and the fraud-cases, by the definitions described in chapter 1.

```
(ind_pilot=0 AND  
niet_gescoord=0 AND  
zuivere_mac=0 AND  
winst_btl_volgensBK=0 AND  
ind_CMbiljet=0 AND  
ind_fraude=0)
```

- 3) Exclude the cases that have specialisms in or have data in the subject 'foreign'.
- 4) Pick 11.000 cases (including 1.000 spare cases) by the method described in [this spreadsheet](#).
- 5) Adjust the indicator for uitrol (ind\_uitrol) in the table mentioned in step 7 of chapter 2, so that the indicator has value 1 if and only if a BSN is in the (new) uitrol-list\*.

\* The LTO-list (chapter 2) is created after the uitrol-list (but the table described in step 7 of chapter 2, is created before the uitrol-list).



**8.15 Appendix 15: Reference 46:** Persoonsgegevens **“Belastingdienst IH Risk Model: Uitwerking\_prioritering\_IH2013\_Concept\_v3\_20150215,” /IH selectie/Backlog2013/.**

**Aan: MT Toezicht IH-niet winst**

**Van: EHI**

**Nadere Uitwerking herprioritering uitworp IH2013**

**CONCEPT – 13-2-2015**

### **Inleiding**

In het MT Toezicht IH-niet winst is op 3 februari 2015 de memo van EHI “Prioritering voorraad IH2013 t.b.v. segment PDB” van 30 januari 2013 besproken. Aan EHI is gevraagd om samen met de Broedkamer de prioritering van de uitworp IH2013 verder uit te werken. In dit memo wordt hiertoe een voorstel gedaan.

### **Uitwerking**

Onder verwijzing naar de bovengenoemde memo richt de uitwerking zich alleen op de categorieën 3, 5 en 7. Over de andere categorieën is al besloten door het MT.

Bij de uitwerking dient rekening te worden gehouden met:

- Broedkamer pilot 2 (2.300 posten waarvan een deel al in behandeling is)
- Broedkamer uitrol 1 (10.000 posten welke door de Broedkamer geselecteerd zijn)
- De selectie van deze 12.300 posten staat dus al vast
- Voldoende aanbod voor alle behandel niveaus B, C en E
- Voldoende aandacht voor de populatie/ onderwerpen welke in Handhavingscommunicatie betrokken zijn.

Het is wenselijk rekening te houden met een zo breed mogelijk pakket van werkzaamheden.

De volgende werkwijze wordt voorgesteld:

De gehele voorraad IH2013 is door het risicomodel (RM) gescoord. Daarmee zijn de volgende criteria voor prioritering beschikbaar:

1. Risicoscore (een indicatie welke loopt van 1 tot en met 99)
2. Waarde (Indicatie laag/ gemiddeld/ hoog)
3. Indicatie of het RM hetzelfde aandachtsgebied aanwijst als de uitworp (nee =0 of ja = 1)
4. Indicatie behandelniveau (3, 4 of 5, resp. E, B of C)
5. De risicocategorie uit CTH (bepaald door de uitworp van de weegmodule).

De uitworp uit de categorieën 5 en 7 wordt met inachtneming van bovenstaande aandachtspunten gerangschikt. Hierbij worden alleen die posten meegenomen waarbij de onder 3 genoemde indicatie gelijk is aan 1.

Dit omdat daarmee de aangifte op basis van de uitworpbewering kan worden opgepakt. Dat is nl. de enige informatie die de aanslagregelaar ter beschikking heeft in ABS.

Vervolgens komen de posten met de hoogste score en hoogste waarde bovenaan te staan.

Over categorie 3 konden we geen consensus bereiken over de door EHI voorgestelde wijze. We stellen daarom voor deze categorie op te splitsen en dan als volgt:

Onderdelen “Correctie t-1” en “Revisierente” toe te voegen aan categorie 2, die daarmee geheel in behandeling wordt gegeven; Onderdeel “Vangnet div.” toe te voegen aan categorie 7, die daarmee meeloopt in de hierboven voorgestelde werkwijze.

Overigens hechten wij er waarde aan nog het volgende op te merken. De kennis van het risicomodel bij EHI is gebaseerd op presentaties en mondelinge toelichtingen van en door BI&A. Documentatie



van de specificaties van het risicomodel en documentatie van testresultaten van het risicomodel en waardenmodel zijn niet beschikbaar.

EHI heeft zich, door het ontbreken van de documentatie, een onvoldoende beeld kunnen vormen over de werking van het risicomodel en kan dan ook geen onderbouwde uitspraken doen over de inzet van het risicomodel; wel geeft de goede samenwerking en afstemming tussen BI&A en EHI het beeld dat het verantwoord is om het risicomodel in te zetten bij de prioritering van de uitworp IH-niet-winst 2013.

Het voorstel is gebaseerd op het maximale gebruik van het risicomodel.

**Gevraagd besluit:**

In te stemmen met voorgestelde uitwerking van de prioritering

# Model documentation

## DTA EWS model

---

30<sup>th</sup> May 2015

**Belastingdienst:**

Persoonsgegevens

**McKinsey and Company:**

Persoonsgegevens



# TABLE OF CONTENTS

<b>1</b>	<b>Model Name</b>	<b>4</b>
<b>2</b>	<b>Executive Summary</b>	<b>5</b>
2.1	Model purpose	5
2.2	Data description	5
2.3	Final model	5
2.4	Summary model results	5
2.5	Key limitations of the model	6
<b>3</b>	<b>Model Context</b>	<b>7</b>
3.1	Dutch tax authority strategic agenda (Investeringsagenda)	7
3.2	The “Broedkamer” Projects	7
<b>4</b>	<b>Model Setup</b>	<b>9</b>
4.1	Model Objectives	9
4.1.1	Need for EWS within DTA	9
4.1.2	Scoping of the EWS	9
4.2	Model Framework	10
4.2.1	Time frame	11
4.2.2	Population definition	12
4.2.2.1	Active companies definition	12
4.2.2.2	Healthy companies definition	12
4.3	Explored models	13
4.3.1	Regression model (Selected approach)	13
4.3.2	Regression model with segmentation	13
4.3.3	Random forest	13
4.3.4	Tree approach	13
<b>5</b>	<b>Inputs Data and Preparation</b>	<b>15</b>
5.1	Dependent Variable	15
5.1.1	Analysis	15
5.2	Independent Variables	16
5.2.1	Data exploration	16
5.2.1.1	Omzetbelasting (OB) Data (Value added tax data)	17
5.2.1.2	Vennootschapsbelasting (VPB) Data (Corporate tax data)	18
5.2.1.3	Loonheffingen (LH or LHN) Data (Income/wage tax data)	18
5.2.1.4	Balans Data (Balance sheet data)	18

5.2.1.5	Incasso Data (Collections data).....	18
5.2.1.6	Miscellaneous Data.....	19
5.2.1.7	IKB and ATK+ qualitative data.....	19
5.2.2	Expert analysis.....	19
5.2.3	Variable Creation.....	20
5.2.3.1	Omzetbelasting (OB) Data (Value added tax data).....	20
5.2.3.2	Vennootschapsbelasting (VPB) Data (Corporate tax data).....	20
5.2.3.3	Balans Data (Balance sheet data).....	21
5.2.3.4	Incasso Data (Collections data).....	21
5.2.3.5	Miscellaneous Data.....	22
5.2.4	Variable Transformation (Creating trend variables).....	23
<b>6</b>	<b>Model Development Process.....</b>	<b>25</b>
6.1	<i>Univariate Analysis.....</i>	<i>25</i>
6.1.1	Pre-transformation factor assessment.....	25
6.1.2	Fill rate assessment.....	25
6.2	<i>Correlation Assessment.....</i>	<i>26</i>
6.3	<i>Variable Treatment.....</i>	<i>27</i>
6.3.1	Bucketing.....	27
6.3.2	Weight of Evidence (WoE) calculation.....	27
6.4	<i>Multivariate Analysis.....</i>	<i>27</i>
6.4.1	Multifactor model development.....	27
6.5	<i>Model Results.....</i>	<i>28</i>
6.5.1	Model parameters.....	28
6.5.2	Model variables buckets.....	28
6.6	<i>Signal allocation.....</i>	<i>30</i>
<b>7</b>	<b>Model Validation.....</b>	<b>33</b>
7.1	<i>Diagnostic Tests for Potential Violations of Model Requirements/Assumptions.....</i>	<i>33</i>
7.1.1	Model predictive power.....	33
7.1.2	Plausibility.....	33
7.1.3	Statistical significance of model coefficients – model p-value.....	33
7.1.4	Consistent rank ordering.....	33
7.1.5	Factor contributions.....	34
7.1.6	Over fitting.....	34
7.1.7	Applicability.....	35
7.1.8	Segmented performance.....	35
7.2	<i>Business review of results.....</i>	<i>35</i>
<b>8</b>	<b>Appendices.....</b>	<b>37</b>
8.1	<i>Variable list.....</i>	<i>37</i>
8.2	<i>Univariate analysis.....</i>	<i>37</i>

8.3	<i>Correlation analysis</i> .....	37
8.4	<i>Scoring analysis</i> .....	37
8.5	<i>Model Results</i> .....	37
8.6	<i>Random forest results</i> .....	37
8.7	<i>Modular approach results</i> .....	37

## 1 Model Name

Model Name:	Early Warning System version 2
Model Owner:	<input type="text" value="Persoonsgegevens"/>
Model Developer:	Belastingdienst: <input type="text" value="Persoonsgegevens"/> McKinsey and Company: <input type="text" value="Persoonsgegevens"/>
Date Documentation:	May 31, 2016

## 2 Executive Summary

### 2.1 Model purpose

The purpose of this model is to predict companies heading for payment problems. The reason for building an early warning system comes from a need for DTA to take action at an early stage in the process. In the current state DTA works in a reactive manner and action is taken when triggered by an event. This means that in cases of bankruptcy there is significant loss that could be decreased if action would be taken at an earlier point in time. Thus the purpose of this model is to predict companies heading for payment problems before the problems are visible. The dependent variables which the model predicts are:

- Insolvency (companies going bankrupt)
- “Oninbaar” (booked losses on filings)
- Company getting a warrant

### 2.2 Data description

#### Scope

The model has been built on the selection of the population of all companies. The companies that have been selected to run cover the majority of the losses and are define as all the large corporates (GO) and the large companies in SME (MKB) segment, which are defined as those with assets or turnover greater than 1M EUR.

The time frame for which the model is built is on all the available data at DTA on terradata (as of 2011). An observation period is selected of 2 years and a performance period of 12 months. The model is trained on four observation point one on each quarter of 2013. The first two quarters in 2014 are used as out of time validation sets.

#### Data sources

Before variable creation there has been an extensive data exploration in which both internal tax data and external data sources have been explored. The sources that have been selected for the variable creation are VAT (“OB”), corporate tax (“OB”), balance sheet info (“Balans”), collections data (“Incasso”), Bank data and a few pieces of miscellaneous company info such as years in business or industry.

### 2.3 Final model

The final model selected for the EWS is a logistic regression model. It is built up of the following 6 variables and related estimate value and the weight it contributes to the score.

Variable	Estimate	Weights
Maximum delay (days) in last 2 years across all tax types	-0.6346	46%
Quarter since last warrant in last 2 year in any of the 3 tax types (OB,VPB, LH)	-0.427	21%
# of quarters of delay OB filing in last 2 years	-0.5758	18%
Years in business	-1.5994	6%
Ratio of retained earnings over assets for last year	-0.6299	4%
# of quarters for which there was a delay in payment across all tax types in last 2 years	-0.2166	4%

Table 2.1: Final model variables

### 2.4 Model validation results

Table below shows the GINI value on development, validation and out of time dataset. The GINI is consistent across all three datasets that proves that model is quite stable and robust

Dataset	GINI
Development	55%
Validation	54%
Out of time	54%

Table 2.2: Model validation results

Below charts shows the comparison between average predicted probability and average bad rate across different buckets. The trend is quite similar that proves the prediction error is low and model is stable.

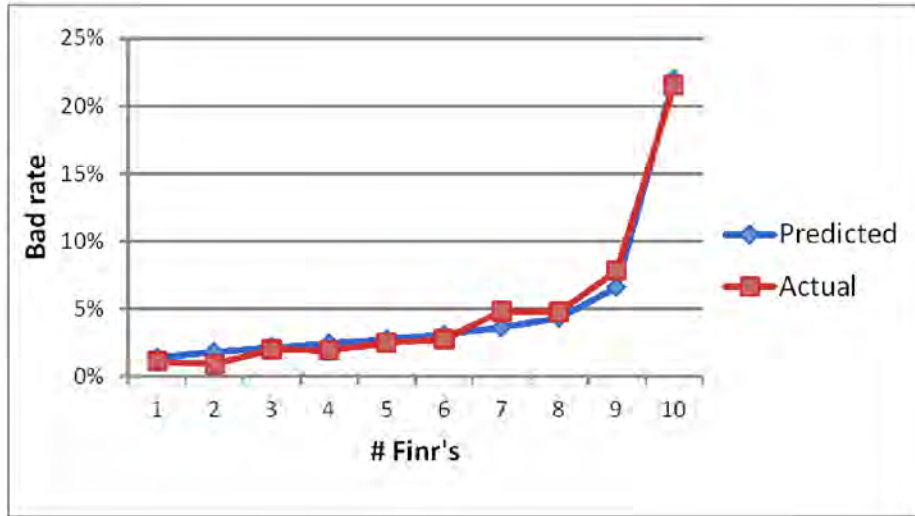


Figure 2.1: Actual and predicted badrates by population bucket

## 3 Model Context

### 3.1 Dutch tax authority strategic agenda (Investeringsagenda)

The Ministry of Finance has outlined a strategic investment agenda<sup>1</sup> in which the course of action for the tax authority is set. This agenda is built upon three main tracks: reducing complexity, making processes robust and objectifying performance. The fourth track is an underlying one to the others and is to ensure that there is traction within the organization. These are portrayed in figure 3.1.

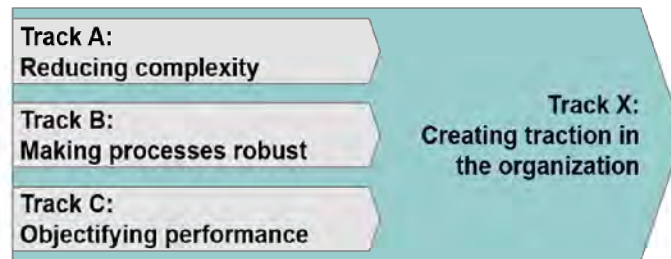


Figure 3.1: Investeringsagenda

The efforts in Track B are reflected in the creation of the “Broedkamer”<sup>2</sup> and also the relevance of creating an early warning system. Track B outlines that the IT processes need to become robust so that they can be simplified and improved, by making incremental improvement rather than revising everything at once, as drastic change is too risky at the moment for the organization. In an additional letter<sup>3</sup> from the State Secretary of Finance to the House of Representatives, he indicates that the newly created “Broedkamer” has an important role in improving the current processes through data-driven approaches.

### 3.2 The “Broedkamer” Projects

**In order to increase efficiency of the processes and to reduce the tax gap (i.e., the difference between taxes that need to be paid to the government and which ones are actually paid).** The Belastingdienst has started to develop more data-driven approaches. To achieve these goals, the “Broedkamer” was established, a project where data scientists and business analysts work together to develop innovative, data-driven ideas and models. **An important part of this project is to collect, structure and standardize the numerous datasets available from the different departments (‘Datafundamenten’) and make it accessible for analytical purposes.** Based on this data, models are developed to optimize processes and to increase the amount of paid taxes. **One of these projects monitors debtors and collections processes in order to optimize collections of claims (‘Dynamisch monitoren’), while another project aims at predicting and identifying companies that are committing VAT fraud based on network analysis (‘OB Carroussel fraude’).** Other projects develop risk models that predict and prioritize the potentially incorrect tax returns and which should increase the amount of positive corrections of tax returns (‘IH Risicomodel’) or reduce the amount of incorrect VAT tax returns (‘OB Negatief’). EWS builds on the

<sup>1</sup> Letter to the House of Representatives of the Netherlands : 19 May 2014: Brede agenda Belastingdienst

<sup>2</sup> Broedkamer in Dutch can be translated to incubator

<sup>3</sup> Letter to the House of Representatives of the Netherlands : 20 May 2015: Investeringsagenda Belastingdienst



ideas developed in these projects and aims to predict which companies are likely to head into financial difficulties.

## 4 Model Setup

### 4.1 Model Objective

#### 4.1.1 Need for EWS within DTA

The Belastingdienst works in a way that is of a reactive (or correctional) nature, meaning that action is undertaken in reaction to an event to correct the negative result. In tax collections for companies this means the process is often only started once a debt occurs, or becomes of a significant size. **Due to this reactive nature the losses in cases where companies are unable to pay are often higher than they would have been if action was taken earlier.** This is why there is a need for an Early Warning System; a system which can give a warning if a company is heading towards payment troubles at an early stage in the process, so that losses can be minimized.

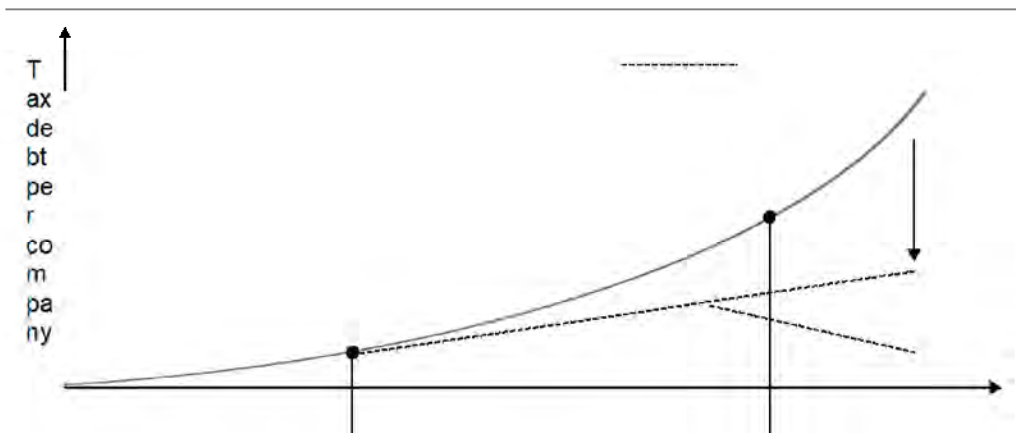


Figure 4.1: Conceptual effect of EWS on tax debt

#### 4.1.2 Scoping of the EWS

Segmentation currently used at DTA distinguishes large corporates (“Grote Ondernemingen”, GO) and small and medium size enterprises (“Midden- en Kleinbedrijf”, MKB). In first version of the EWS model all GO<sup>4</sup> and the larger companies of the MKB segment (“MKB-midden”) are prioritized.

**Selecting only GOs and MKB-midden allowed narrowing the scope and create a manageable sample of companies to work with in initial phase of Early Warning System.**

Inspectie, controle en toezicht	Inspectie, controle en toezicht
Inspectie, controle en toezicht	Focusing on this smaller sample

will allow creating processes faster and focusing on significant improvements and loss reduction.

MKB companies were selected based on their turnover and assets – average of turnover OR assets had to be greater than 1 million EUR at entity level in last three available years. In initial analysis we also considered equity and number of employees but due to data quality and completeness these two dimensions were not considered in final selection. Trend analysis is based on “Balans” data set, data

<sup>4</sup> Within GO a sub-selection is made to exclude APA/ATR

stored in P\_PMI.BALANSINFO table. This contains yearly information about companies about their balance sheet, cash flows and other financial aspects. Based on that we have calculated average turnover and average assets in last 3 available years, 2012-2014, 2011-2013 or 2010-2012, depending what last 3 years were available for each company (FINR).

## 4.2 Model Framework

The framework of the model is depicted in Figure 3 and at the core consists of four elements; (1) input in the form of a selection of companies that are in scope, (2) input in the form of data and variables that characterize the model, (3) the model consisting of a logistic regression to select the most predictive variables for “bad”<sup>5</sup> and (4) the output of the model giving each company a score based on the probability that the company will become “bad”. The creation of this model is divided into 9 different phases which form an iterative process and are depicted in Figure 4.

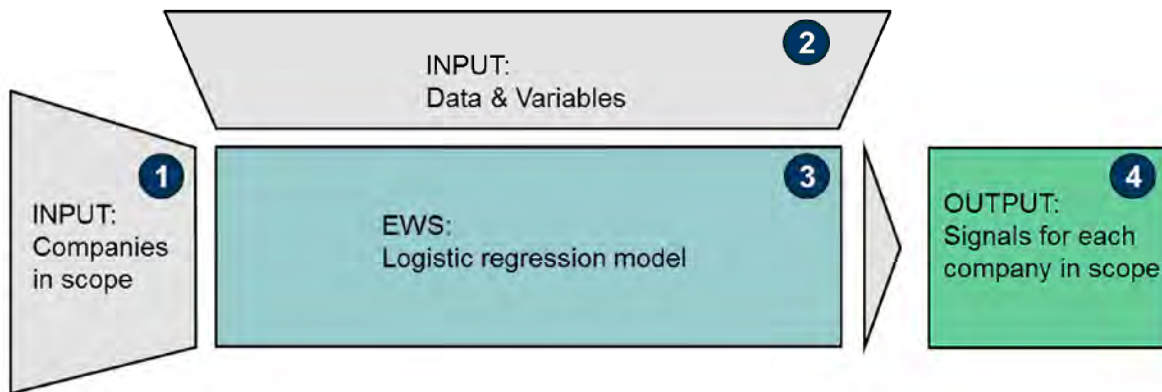


Figure 4.2: Conceptual illustration of an early warning system

<sup>5</sup> “Bad” for this model refers to companies which are “insolvent”, have “oninbaar” or are more than 365 days past due.

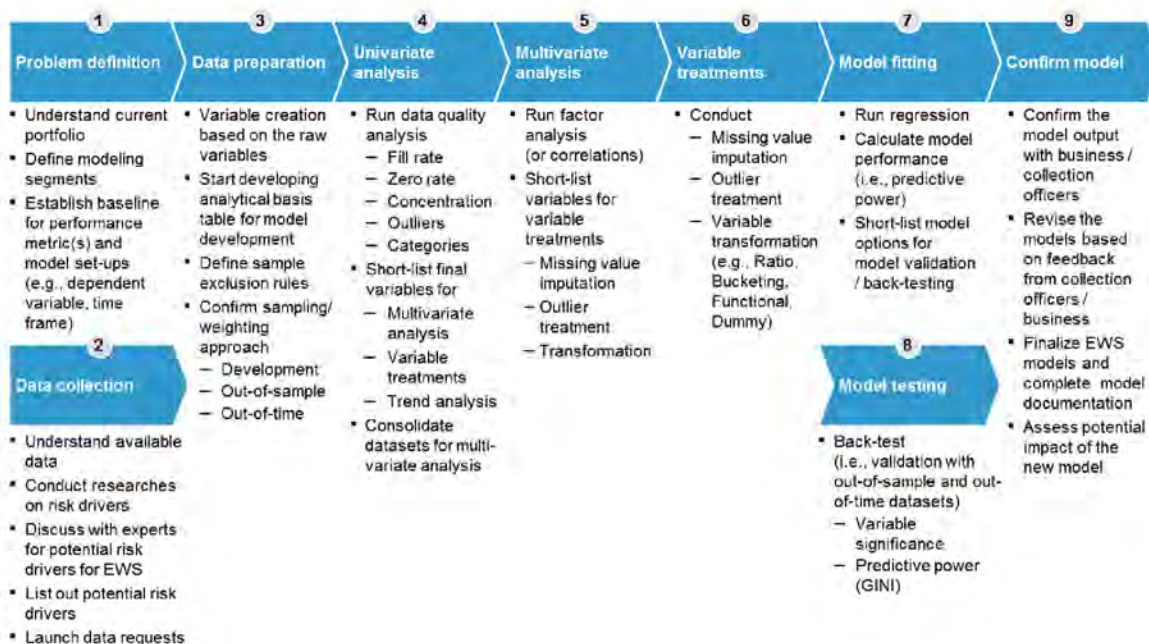


Figure 4.3: Overview of phases in the model development, it should be noted that the development is an iterative process.

#### 4.2.1 Time frame

The data available of adequate quality to build the model, is available as of 2011 and up to the current date<sup>6</sup>. When training the model an observation point must be selected at a certain point in the past. From this observation point two time periods are defined:

1. **A performance period** towards the future to observe whether or not a certain company has gone “bad”. The performance period is define as 12 months; the reason for this because factors are not forward looking for long period.
2. **An observation period** is a period in past, to observe the independent variables of all the companies, so that a distinction of the characteristics can be observed of companies that have and have not gone bad. This period is defined as two years. The reason a two year period is chosen is because this can be seen as a period which can be seen as an adequate period to be able to predict a 12 months further, as selecting a longer period would mean that the number of observation points that could be chosen would be limited.

To train and test the model multiple observation points are chosen. These are spread out by quarters because  of all companies file their OB taxes on quarterly basis, and selecting a higher frequency will thus not significantly add value. Given the time period for which data is available and the above selected performance and observation periods it has been possible to select 6 observation points.

<sup>6</sup> Current data at the start of the project is in October 2014

When creating a model it is of importance that the model is not built and tested on the same data set and not on the same time period. For this reason the data is split into several sections; the last two observation points are kept separate for testing, this is called the out-of-time sample. The first four observation points are the so-called in-time sample, this data set is randomly split 70%-30%, where the 30% is called the out-of-sample and kept aside to test the model and the 70% percent is the so-called in-sample selection upon which the model is trained.

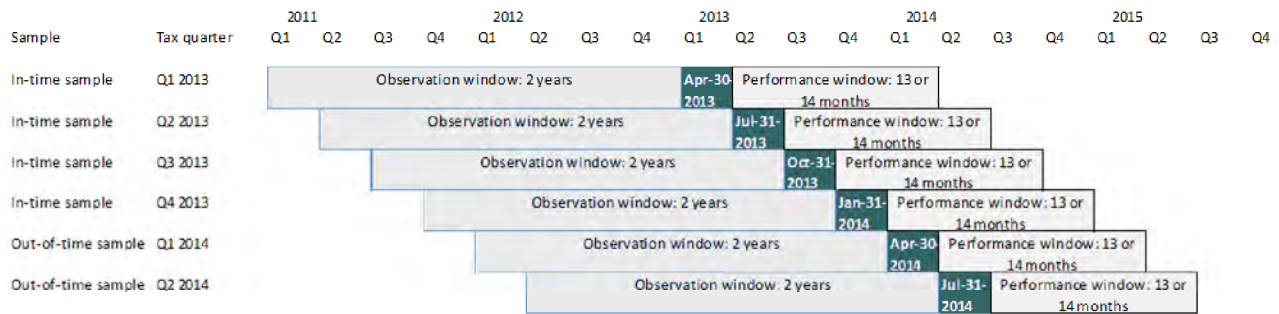


Figure 4.4: Time-line showing quarterly observation points, observation and performance windows, and in- and out-of-time samples.

#### 4.2.2 Population definition

The total population as defined in Section 4.1.2 needs to be filtered upon three criteria, before being able to select the company for training the model; firstly the company should be “active”, secondly the company should be “healthy” and third it should not be OB fiscal unit. In the model development, natural persons were excluded, because they were out of the original scope of our model. This means one-man-businesses were not part of the model development. In the model scoring code, natural persons are included.

##### 4.2.2.1 Active companies definition

To define the companies which are active at a certain point in time, and more specifically at the 6 defined observation points an active company flag has been created. This is based on two dates from D\_P\_CLC.G\_CLC\_ECON\_ACT table, start date of a company: BVR\_EA2\_ING\_D and date a company has ceased to exist: BVR\_EA2\_VER\_D. If start date was lower or equal to the observation point and end date was equal or greater than observation point then company was flagged as active for that point in time.

##### 4.2.2.2 Healthy companies definition

In addition to “active” companies a criteria has been created to make sure that companies at observation point have zero days past due. This is done so that the model can be ensured to be forward looking and predict the “bad” flag that has not yet started to happen. For each of six observation points the companies have been selected that did not have any claim with oninbaar status, were not in insolvent state and did not have a ‘dwangbevel’ (=warrant) or ‘aanmaning’ (=first reminder ob or second reminder vpb) at the observation point. This provided a healthy population of companies for which we measured probability of significant delay and inability to pay tax.

Based on that we created 6 tables (named: O\_A\_EWS.data\_all\_ph2\_<quarter\_year>) which included only healthy companies at each observation point, which is at a month after the end of the quarter. (eg.

Observation point 30 apr 2013 is <quarter\_year> = q1\_2013). This allowed to create relevant independent and dependent variables for model input.

#### 4.2.2.3 *Fiscal Units*

Individual companies under an entity could file OB tax together under a fiscal unit. We dropped the finr of the head of the OB fiscal unit from the dataset as they don't have any physical presence. For both the OB fiscal units and VPB fiscal units (where the vpb fiscal unit finr is also a company), we used a division key to divide the tax amount of the fiscal unit between the underlying finrs. In addition, for the variables as used in model 2, the underlying finrs are given the same behavior as the head fiscal unit OB or VPB (regarding these tax types). For further detailed logic on finr-fiscal unit split of tax filing amounts and other characteristics, please refer to section 5.2.5.

### 4.3 Explored models

#### 4.1.1 Regression model (Selected approach)

Logistic regression methodology was used for the model development as the outcome was binary (0/1). The model was developed on the full population having independent variables from all available data sources (e.g. OB, VPB, incasso etc.). Model validation was done on in-sample data (30% of the total sample) and out of time data. For definition of in-sample and out of sample data, please refer to Section 4.2.1 . This approach of model development was preferred to capture the interaction between variables from different data sources.

## 5 Inputs Data and Preparation

### 5.1 Dependent Variable

#### 5.1.1 Analysis

We have defined 3 types of dependent variable:

1. Insolvency based on Chamber of Commerce data (“KvK”) – this data provides bankruptcy as registered legal statuses. These are defined as any records in `sdv_frau.b_kv_k_nnp` table with field `kvk_brt_rechtstoestand` equal to “FAILLISSEMENT”.
2. “Oninbaar” based on “incasso” data; this provides information about companies where it is asses that the debt will not be collected. It is defined as any record in `D_P_INC.n_inc_vordering_verrijkt` table with field `onincode` not empty. This criterion has some overlap with the first, but there is still a significant part of “oninbaar”, which is not insolvent by “KvK”.
3. Warrant issued in next 12 months

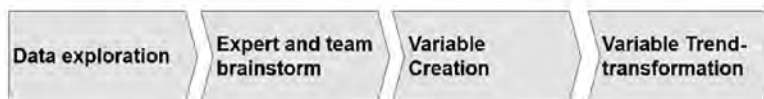
The condition for the dependent variable is an OR condition where the dependent variable requires at least one of the three criteria to be met. The dependent variable has been validated with discussion with process experts; and yields  of the population to be flagged bad after 12 months.

The decision to add a third dependent variable based on warrant is made because:

1. Using only the criteria for insolvency and “oninbaar” results in a dependent variable that covers a too small selection of the total population . Thus the option is explored to add warrant as a criterion, as  of “oninbaar” cases has also received a warrant.

### 5.2 Independent Variables

The decision for which variable to create has been based on the following process; each step is explained in detail in this chapter



#### 5.2.1 Data exploration

Prior to describing the databases it is important to understand the data preparation process in place at the tax authority. There is a vast amount of data collected at DTA and this is done in very robust but often very old systems. To make sure that the data becomes readily available the “datafundament” projects have been set up by the “Broedkamer” to clean and make the data ready for use for numerous projects. In the process of data preparation there are three different layers of data:

1. Brown: Raw unprocessed data, this is the roughest level of data and before it can be properly used in a model it still needs processing and cleaning.

2. Blue: Cleaned data, blue level data can be used in models or for research purposes.
3. Purple: Information level data, purple level data is also called information level data and consists of data from the blue level which combines variables to give more insightful information.

Before the creation and selection of the variables the available data is explored. This section describes per data source the finding of the type of data in the database. Also the usefulness for the model is assessed and whether the data is used moving forward in the variable creation. The data sources explored are VAT (“OB”), corporate tax (“OB”), wage tax (“LH”), balance sheet info (“Balans”), bank data, collections data (“Incasso”) and miscellaneous data (Misc)

The remainder of this section focuses on describing some of the data sources which have been explored. Each subsection describes the type of data available, the coverage and frequency. An overview of the data sources is given in the following figure, depicting the layer each data source is in (January 2016):

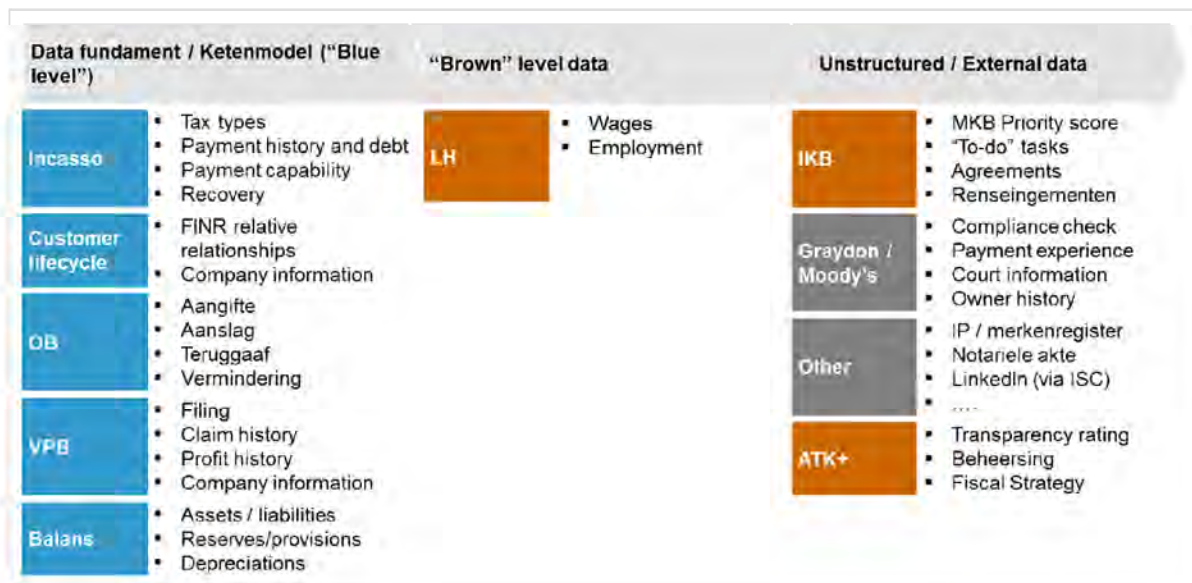


Figure 5.2: Representation of the data sources and the quality of data preparation

### 5.2.1.1 Omzetbelasting (OB) Data (Value added tax data)

The OB database is currently purple level data and contains information related to OB tax. It contains information on companies OB filing amount, VAT paid to suppliers, OB paid amount and turnover for each tax year. The data is present either at monthly, quarterly or yearly level depending on when the company is required to file their OB tax. It also contains information on companies that have exceeded the OB negative norm in the past and, also for which manual assessment of OB filing was done for each tax year. The quality of data is high and is expected to give a strong prediction on a company heading for default.



### 5.2.1.2 *Vennootschapsbelasting (VPB) Data (Corporate tax data)*

The data from the so-called Corporate Tax datapillar (VPB-datafundament) contains data from the claims on the corporate tax. By this is meant the information concerning claims that are sent to the taxpayer -potentially after a file from the taxpayer. There are several kinds of claims: preliminary claims/voorlopige aanslagen (during or shortly after the taxyear) and final claims/ definitieve aanslagen (three years after the taxyear). There are also extra claims /navorderingen, but due to the usually long time after the final claim they are not considered in this model. The processing and the possible correction of the file coming from the taxpayer are also contained in this datafundament. The complete files are in the ABS-system and the actual cashflow is in the COA/DACAS (the collections-system).

### 5.2.1.3 *Loonheffingen (LH or LHN) Data (Income/wage tax data)*

The Loonheffingen data is currently brown level data, this brings a limitation that the data cannot be used directly for modelling purposes. Given the scope and deadline of this project; at this point in time the LH data source cannot be used.

Once the data in LH does become available in blue or purple level it is a potentially valuable data source for an EWS. The data which is included in the LH data is information concerning companies' employment i.e. the number of employees, the hours worked, salaries, tenure, etc.

### 5.2.1.4 *Balans Data (Balance sheet data)*

The balance data originates from the fact that every year companies are required to submit their balance sheet to the tax authority. The system in which this is registered is ABS and the fields which are (or can be) filled in are be found in ABS field file.

The balance sheet info is submitted yearly and is *done several months after* the end of the tax year. Rendering the frequency and recency of the data to be quite low. The fields which are reported on do give strong insights into companies' performance.

### 5.2.1.5 *Incasso Data (Collections data)*

The Incasso database contains data from the collections process and consists of a wide range of information all related to a company's payment behavior. The type of information which can be retrieved from this data base consists of statuses in the process such as reminders, warrants and fines which are sent, contains payment date and the deadlines for payments as well as the amount of time postponed or appealed. Using this data it will be possible to create insights into a company's payment behavior.

The quality of the incasso data is high; the data is updated frequently and is expected to give a strong prediction on a company heading for default.

	Inspectie, controle en toezicht
Inspectie, controle en toezicht	
Inspectie, controle en toezicht	

<sup>7</sup> Exception of the preliminary claims of corporate tax

#### 5.2.1.6 *Miscellaneous Data*

The variables in the EWS that are described as miscellaneous are variables related to employees and variables related to static<sup>8</sup> characteristics of a company, where the static characteristics are those that do not change regularly, such as the region or the industry. The both data types are updated monthly and the employee related information has a two month delay. The fill rates for each of the static characteristics are high, but for employee information is low.

#### 5.2.1.7 *Bank data*

Bank data contains information on the saving account, mortgages, equity value and loans at firm level. However the information is missing for majority of the companies present in DTA system.

#### 5.2.1.8 *IKB and ATK+ qualitative data*

IKB and ATK+ are two systems in which, among others, qualitative data is stored. ATK+ is the system used by GO and has several promising pieces of information where a KC'er (client coordinator) judges a company upon several aspects with a score of 1-5

- Transparency
  - Open communication
  - Openness about disruptions
  - Does the organization understand DTA expectations
- Internal monitoring
  - Monitoring of risks (currently not in use)
  - Availability for questions/checks
  - Expertise/competence
- Fiscal Strategy
  - Positive outlook
  - Shared interest

The limitation of the ATK+ data is there are no date stamps for each separate field only for the entire scoring (so it is not known when a specific field was added or changed). Furthermore there is no historic information available; only the latest updated score is stored. This means that it is not possible to train a model on historic data, meaning that at this point in time it is not possible to incorporate ATK+ data into the EWS.

The other system used for qualitative data is IKB and is used by the MKB (SME segment). This system contains largely text fields or separate documents with information which is difficult to implement into the model, but it also contains several variables. One of which is the “aandachtscategorie”, a score given to each company bi-annually on two axes. The first axis is based on the total taxable amount and given a score of 1-3. The second axis is a score which is risk category which consist of numerous elements and result in a letter A-E where B is for beginning/new company.

### 5.2.2 **Expert analysis**

After an initial exploration of the available data at the Tax Authority the team has consulted multiple sources of expert analysis. With the goal to discover if there could be additional data source still to be explored and whether these could be feasible for use in the early warning system. The experts which

---

<sup>8</sup> Static, in this context refers to not changing often

have been consulted consist of numerous interviews and brainstorming within DTA, with experts from the broedkamer and experts from the field providing information which, given their experience, could be a signal for a company heading towards default. Furthermore the team has consulted several persons within McKinsey who are experts in the field of this type of predictive modeling. Gathering all this information provides the team with a comprehensive list of possible variables depicted below. The figure also indicates which were and were not further explored.

### 5.2.3 Variable Creation

This section describes the manner in which the variables have been created. Each subsection describes the creation of the variables in a single data source, and the final subsection describes the transformations which have been applied to the variables to create a variables.

#### 5.2.3.1 Omzetbelasting (OB) Data (Value added tax data)

For OB, variables were created using information on OB filing amount, companies turnover, OB paid amount and OB amount paid to suppliers. The information on filing, turnover and OB paid to supplier is present in the dataset "l\_ob\_aangiften ". As filing for OB tax can be done monthly, quarterly or yearly, and given approximately 70% of the OB filings are done quarterly, all the amounts were first rolled up at a quarterly level (e.g. Jan, Feb and Mar amounts were rolled up to get quarter 1 amount). The amount and flags variables were first calculated for each quarter and then rolled up to get the yearly variables. Similar approach was used for the creation of variable related to paid amount that is present in the dataset "n\_ob\_aangiften\_inc".

The dataset "n\_ob\_nnosmob\_keten" was used for the creation of flags variable related to OB negative norm and manual assessment. First, flag were created for each quarter in historical period and then rolled up to get the yearly variables. The complete list of variables is present in Appendix XXX.

#### 5.2.3.2 Vennootschapsbelasting (VPB) Data (Corporate tax data)

The variable creation considering the claims of the corporate tax takes place in several steps. First we create from an input file of companies – represented by their fiscal number (finr)- a table consisting of those companies that actually are obliged to pay corporate tax (in the given timeframe), meaning are present in the corporate tax system (VBN).

Next we create (once) a temporary file of all preliminary claims and reductions on preliminary claims in order to refer later to the second last signal.

Subsequently we collect information from the (reductions on the) preliminary claim and the final claim per tax year. A tax year is the year on which the tax is paid; this is not necessarily a calendar year although it usually is. This variable collection happens in three steps:

- records that have a reduction on a preliminary claim,
- records that do not have a reduction on a preliminary claim but do have a preliminary claim
- records that do not have a preliminary claim (and therefore no reduction either)

Then we determine what the second last signal is from the file described in the above paragraph and add the information to the file in two separate procedures. This is all in a defined macro which takes as input parameters the date of observation (to filter out information after the observation point), the tax

year and a character variable that distinguishes between the table names of the temporary files for the different tax years.

Then the information of all the desired tax years is combined into a single file. For the rare cases that there are multiple tax years in a calendar year we only take the last one.

After this, we transpose the table such that the information for different tax years is ordered as variables instead of records per tax year. Now the table is unique on finr level.

In the next two procedures, calculations are being done on the selected variables. The name of the output file is `ews_btb_vpb_OP` excluding a postfix to distinguish between the different observation points and the output library is `o_a_ews`. The macrovariables that need to be put in every time are:

- The source file containing the finrs of the companies (taken to be in library `o_a_ews`)
- The date of the observation point
- The 4 taxyears included descending from the year of the observation point
- The character postfix given to the table to distinguish between different observation points

### 5.2.3.3 *Balans Data (Balance sheet data)*

De data in de tabellen van het ABS zijn gestructureerd middels een finr dat aangifte doet namens de hele fiscale eenheid VPB. Er is dan een record per fiscale eenheid met geaggregeerde waarden en per individueel finr dat onderdeel uitmaakt van de fiscale eenheid (FE).

De creatie van balansvariabelen gebeurt in meerdere stappen. Eerst worden de dochterbedrijven geselecteerd (volgnummer van de onderneming binnen de FE is niet gelijk aan 1 en verbind de tabel via het finr van de dochterbedrijven), vervolgens de individuele gegevens van het moederbedrijf (volgnummer van de onderneming binnen de FE is niet gelijk aan 1, er is geen finr van een dochter en verbind de tabel via het finr van de moeder). Ten slotte zijn er nog gegevens van bedrijven die VPB-aangifte doen, maar geen onderdeel zijn van een FE. Deze finrs zijn dus niet aanwezig in de voorgaande tabellen per boekjaar en worden in de laatste stap geselecteerd.

De records worden per observatiepunt geselecteerd t.o.v. de aangifte datum uit de tabel `d_p_vpb.i_vpb_btv` en er bestaat een record per elementnummer uit de balans. Deze wordt vervolgens getransponeerd en er worden variabelen afgeleid (ratio's etc.)

### 5.2.3.4 *Incasso Data (Collections data)*

To create the variables from the Incasso data several tables from the `D_P_INC` database in the terradata library are used. Each of the tables provides a piece of information required for a certain type of variable. All of the tables are built up on "vorderingsnummer" a unique identifier for tax filing; some of the tables also include finr, the unique identifier used in this model for companies. The tables that do not contain finr need to be matched to the ones that do so that variables can be linked to finr.

Table	Description	Finr incl.	Variable examples
<code>i_inc_basisvordering</code>	At level of the vorderingsnummer contains general information about that vordering	Yes	Fine flags Fine amounts Interest amount
<code>n_inc_vordering_verrijkt</code>	Aggregated at a vordering contains information about about the vordering on many	Yes	The vordering amount, the amount still due, the start, due and end dates.

	fields		
<b>i_inc_vooraad_financieel</b>	At regular intervals in time contains information about the vordering (2014 onward)	No	Amount that has been payed, and is still due on a vordering
<b>ic_vooraad_financieel_hist</b>	Same as above but for 2011-2013, and is joined with more current data.	No	Same as above (note for use that fields in table are different)
<b>i_inc_levenscyclus_v_vw</b>	Aggregated at a vordering contains information about the statuses the vordering has been in	Yes	First reminder Second reminder Postponement Warrant
<b>g_inc_vord_stat</b>	Lower level prepped data used to recover the status of a request on a vordering	No	Statuses such as a declined postponement

The appendix contain a full list of the Incasso variables created, but the variables can be divided into two different types one-off events which can be attributed to one quarter (e.g., a fine, a reminder or an amount still due at one point in time) and event which has a duration which extends beyond the length of one quarter (in incasso this is the maximum delay in a certain period of 1 or 2 years for different tax types).

To create the variables attributable to a quarter several steps are taken

1. First each of the variables (events) is attributed to a certain quarter in which it has occurred.
2. A. Second the variables (events) are aggregated at a quarter level for each tax type and company (finr)  
B. Events are aggregated at a quarter level and company (finr) across tax types
3. The table is transposed creating a new column for each quarter and variable (event) for each tax type (as well as across all tax types)
4. The trends are calculated by comparing quarters (see section 5.2.4)

The variables not attributable to a quarter are the variables related to the maximum delay. These are calculated for a certain observation period in which the following steps are taken:

1. The period upon which the maximum delay is observed is defined
2. Each vordering attributed a delay in days based on certain conditions:
  - a. If the vordering has been payed before observation date and payment date is later than due date then time between payment and due date is calculated
  - b. If the vordering has been payed after observation date and payment date is later than due date then time between observation and due date is calculated
  - c. If the vordering has not been payed and the due date is before the observation date then time between observation and due date is calculated
3. The maximum delay is then selected for each finr for that time period for each tax type as well as overall.

### 5.2.3.5 *Miscellaneous Data*

#### 5.2.3.5.1 *Employee variables*

Data included in employee metrics comes from internal database P\_PMI.INKOMEN\_CYCLUS. It includes employees and their salary, for EWS model data is used covering the 2011-2014.

A count is made of the employees who were employed by company at the reporting date and their type of income (TYPE\_INKOMEN) was recorded as “loon”, referring to wage. These counts are done on a quarterly basis. Average salary is calculated as aggregations of monthly net income per employee (NETTO\_MAANDINKOMEN\_YYYY) are aggregated at a finr for each quarter. From there all trend variables are created.

#### 5.2.3.5.2 *Static characteristics variables*

The static characteristic variables created are a company’s region, industry and years in business are created for each observation point.

Years in business is created by subtracting the establishment year (bvr\_pso\_oprichting\_d from D\_P\_CLC.g\_clc\_nnp) from the observation point year.

A company’s region is created from the table D\_P\_CLC.i\_clc\_lc\_vestadr\_ind. Firstly entries registered outside of the Netherlands are given region “OTHER\_COUNTRY” by checking the field bvr\_adr\_land not equal to 1 or missing. Next the region is selected for the correct point in time by selecting the latest date for the registration, by selecting the entries with the maximum bvr\_pad\_verval\_d. Then lastly the province can be extracted for each finr from the field geo\_pcz\_provincie.

For the industry the same procedure is followed as for the region, but the table and fields are different. The table containing the information is D\_P\_CLC.g\_clc\_econ\_ent\_relaties and the fields for industry and date are kvk\_brn\_cur\_branche\_3 and bvr\_ece\_verval\_d respectively.

### 5.2.3.6 *Bank Data*

To create variable related to the bank account information of companies, we used d\_P\_Bank.i\_bank\_nnp\_sparen for saving accounts, d\_p\_bank.i\_bank\_nnp\_beleggen for variables related to equity. We also explored dataset d\_p\_bank.i\_bank\_nnp\_LENEN for mortgages and loan account information, however because of low fill rate didn’t tested the variables.

## 5.2.4 **Variable Transformation (Creating trend variables)**

Each of the variables is created for a specific quarter in a specific year. This in many cases a trend or an aggregation over several quarters can be more predictive than a single quarter absolute value. Variables are transformed based on the type of variable that is created. The types of variables created can be amounts, categorical, flags or sums of flags. The following matrix gives an overview of the transformations:

---

<sup>9</sup> Included in flags are also one-off amounts, such as fine amount. The reason these are not included in amounts because in amount growth variables are created and this only makes sense if they are continually filled.

	Amount variables	Flag variables <sup>9</sup>	Categorical variables
<i>Transformation 1</i>	(Last quarter)	(Last quarter)	(No transformation)
<i>Transformation 2</i>	(last quarter) / (quarter before)	(Sum of last half year)	
<i>Transformation 3</i>	(Last half year) / (half year before)	(Sum of last year)	
<i>Transformation 4</i>	(Last year) / (year before)	(Sum of last 2 years)	
<i>Transformation 5</i>	(Last quarter)/(same quarter a year earlier)	(Quarters since last occurrence)	
<i>Transformation 6</i>	(Last half year)/(same half year a year earlier)		

Table 5.1: Variable transformations

### 5.2.5 Fiscal unit and FINR logic

Individual companies under an entity could file OB or VPB tax together under a fiscal unit. At the same time, they can file one tax type together and other separately. Because of this filing arrangement, it is difficult to separate out the tax amount corresponding to individual companies within a fiscal unit. To counter this problem, we decided to use separate approach for variable creation and value creation which is outlined below.

#### 5.2.5.1 Variable creation

Based on the discussion with experts, for variables related to delay payment, warrant or first reminder, we decided to assign the same value of fiscal unit to individual companies (e.g. If the fiscal unit has 1 delay payment for OB in a period, the individual companies will also be considered to have 1 delay payment for OB in the same period)

#### 5.2.5.2 Value creation

Uit de Verlies-en Winstrekening kunnen de winst en omzet worden verkregen. Per finr wordt het finr van de FE OB gevonden als deze van toepassing is. Als deze er niet is dan is het finr verantwoordelijk voor de hele omzet (100%). Als er geen omzet-gegevens op de balans te vinden zijn én het finr maakt geen onderdeel uit van een FE dan laten we ook het percentage veld leeg. Als het wel onderdeel is van een FE en er zijn geen omzetgegevens bekend in de hele FE, dan vullen we nul (0%) in. Anders bepalen we de ratio van de omzet t.o.v. de hele FE.

Hetzelfde geldt voor FE VPB. Deze is anders gestructureerd in de bron. Eerst worden het finr, finr van de moeder, jaar en winst van de verlies & winst rekening, winst van de aangifte en de kasstroom van de VPB bepaald van dochterbedrijven en vervolgens van alle moederbedrijven. In een volgende procedure worden de aangiftegegevens van individuele bedrijven toegevoegd.

In de laatste procedure worden de loonaangiftegegevens uit het boekjaar toegevoegd.

<sup>9</sup> Included in flags are also one-off amounts, such as fine amount. The reason these are not included in amounts because in amount growth variables are created and this only makes sense if they are continually filled.





## 6 Model Development Process

### 6.1 Univariate Analysis

The univariate analysis aims to make a selection of the variables which are most predictive towards a company going bad. The univariate analysis consists of two parts; the first is the pre-transformation factor assessment which assesses the predictive power of a single variables and the second is the fill rate assessment. Within each dataset explored both assessments are done and thresholds for both are selected to choose the best variables within each dataset.

#### 6.1.1 Pre-transformation factor assessment

Before starting the variable treatment, an assessment is performed on the raw predictive power of each variable. This requires calculation of the un-transformed gini (measure of predictive power) of each factor and then filtering out the variable which do not have sufficient predictive power. The primary metric used for predictive power assessment is gini. It is an indication of the position of the predictive power as compared to completely random prediction (0% prediction) to a perfect prediction (100% prediction). The following figure illustrates the concept of gini.

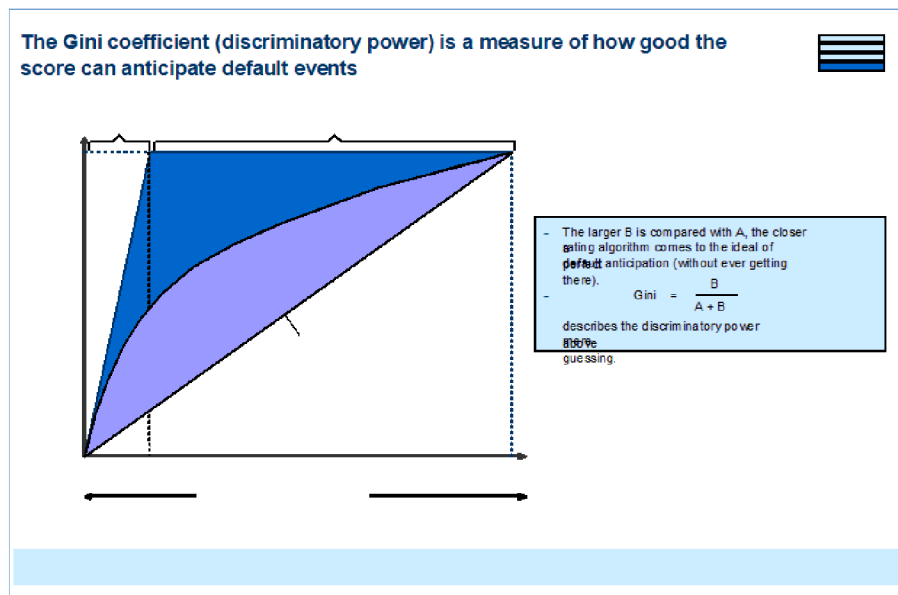


Figure 6.1: Gini explanation

Table in Section 6.1.2 shows the thresholds which have been set for each dataset for “gini with missing” and “gini non missing”, which are the measures of gini including missing values and gini excluding missing values respectively. A full list of univariate gini for all variables is provided in appendix 8.2.

### 6.1.2 Fill rate assessment

Within each dataset a selection has been made of the variables which show to have an adequate fill rate. Similar to the selection of predictive power described in the previous section, it is also necessary to select variables which have a minimum fill rate. The reason for this is because variables where fill rate is too low will not be able to provide valuable predictive information for the full range of the model. Table 6.1 shows the thresholds which have been set for each dataset for fill rate, as to select an acceptable set of variables with acceptable fill rate. A full list of fill rates for all variables is provided in appendix 8.2.

Dataset	Gini with missing threshold	Gini non missing threshold	Fill rate
OB	5%	5%	35%
VPB	3%	3%	20%
Balans	5%	5%	35%
Incasso	5%	5%	50%
Incasso (subset of vars)	10%	10%	25%
Miscellaneous	5%	5%	40%

Table 6.1: Fill rate assessment

## 6.2 Correlation Assessment

Given that all the variables would be transformed to predict the same final outcome, a certain variable correlation is expected. However this correlation should be constrained to lower values, as it has been noted that statistical procedures used to develop multivariate models underperform in cases of large factor correlation (e.g., over-fitting, making the two correlated factors statistically insignificant). The correlation between variables was measured using Pearson's correlation coefficient. The variables were considered as correlated if the Pearson's correlation coefficient value is greater than 0.6. However because of the high missing rate, the correlation was high among variables across different data sources, as a result this cut-off was adjusted for each data sources to have at-least 5 variables from each data sources for next step. Below table shows the correlation coefficient cut-off for each data sources. Among the correlated variables, variable having the highest Univariate GINI were selected for the next step. A full list of correlation analysis for all variables across different data sources is provided in Appendix 8.3.

<i>Data sources</i>	<i>Pearson's coefficient cut-off</i>
<i>Incasso</i>	0.6
<i>OB</i>	0.6
<i>Balans</i>	0.6
<i>Miscellaneous</i>	0.6
<i>VPB</i>	0.6

Table 6.2: Correlation thresholds

## 6.3 Variable Treatment

### 6.3.1 Bucketing

Post correlation analysis, variables were grouped into different buckets based on its value and each bucket was then checked on the following constraints (i) each bucket should have sufficient records (ii) the bad rate trend across buckets should be monotonous (iii) the bad rate trend should be in line with the business sense. The variables satisfying the above constraints were selected for the next step. The automatic created buckets can be manually adjusted to fit these constraints.

After bucketing the variables, buckets are checked again for consistency in trend (with post\_binning script). Especially because missing values are treated as a separate bucket, or merged with another group if this makes sense. For example, missing values in delay in payment should be merged with the group with 0 days delay in payment. Post-binning allows to check the trend with these changes included.

### 6.3.2 Weight of Evidence (WoE) calculation

After post bucketing, WoE is calculated for each bucket of the variable. The WoE value is a widely used measure of the “strength” of a grouping for separating good and bad companies. The WoE is calculated for each bucket using the below formulae:

The WoE recoding of variables is well suited for subsequent modeling using Logistic Regression. Specifically, logistic regression will fit a linear regression equation of WoE-coded continuous variables to predict the logit-transformed binary Bads over Goods dependent variable. The Logit transformation is simply the log of the odds, i.e.,  $\ln(p(\text{Bads})/p(\text{Goods}))$ . Therefore, by using WoE-coded predictors in logistic regression, the predictors are all prepared and coded to the same (WoE) scale, and the parameters in the logistic regression equation can be directly compared.

## 6.4 Multivariate Analysis

### 6.4.1 Multifactor model development

Post WoE calculation, the variables were then passed through an iterative procedure wherein a stepwise approach is utilized to develop the model. At every step, the least predicted variables **not** satisfying the following constraints were removed :

- (i) variable coefficients having negative sign as WoE values are being used
- (ii) variables being significant at 1% significance level
- (iii) VIF (Variance inflation factor) value for all variables is less than 2 (this removes any Multicollinearity between variables)
- (iv) the weightage of individual variable in the model is less than 50% (this removes any overdependence of the model on that particular variable)

- (v) variable significant in validation and out of time data. The model iteration passing all tests is the final model developed. The model construct is designed to maximize the likelihood function, which also guarantees the maximum predictive power. The following figure illustrates the optimization used to develop the model.

$$\text{Maximize } \prod \text{predicted probability}^{\text{actual outcome}} * (1 - \text{predicted probability})^{(1 - \text{actual outcome})}$$

$$\text{predicted probability } y = \frac{1}{1 + e^{d_1 + d_2 * \text{predicted score}}}$$

$$\text{predicted score} = w_0 + \sum_{i=1}^{\text{chosen factors}} w_i * \text{transformed factor } i$$

## 6.5 Model Results

### 6.5.1 Model parameters

Below table shows the final model parameters. It includes the final model variables, the model coefficients corresponding to each variable, significance value (p-value), weights of each variable in the model and variance inflation factor (VIF) value.

Variables	Description	Estimate	P-value	Weights	VIF
Intercept	Constant	-2.9243	<.0001		
total_maxdelay_fufn_lst2yr_WOE	Maximum delay (days) in last 2 years across all tax types	-0.6346	<.0001	46%	1.7
total_warrant_qlast_WOE	Quarter since last warrant in last 2 year in any of the 3 tax types (OB,VPB, LH)	-0.427	<.0001	21%	1.4
ob_flg_dl_flg_fufn_lst2yr_WOE	# of quarters of delay OB filing in last 2 years	-0.5758	<.0001	18%	1.1
years_in_business_WOE	Years in business	-1.5994	<.0001	6%	1.0
zscore_t2_lstyr_WOE	Ratio of retained earnings over assets for last year	-0.6299	<.0001	4%	1.0
total_delaypayment_lst2yr_WOE	# of quarter for which there was a delay in payment across all tax types in last 2 years	-0.2166	<.0001	4%	1.6

Table 6.3: Model specifications

### 6.5.2 Model variables buckets

The following tables show the description of buckets created for each variable, and the corresponding bad rate and % population. The buckets were created keeping in mind the constraints mentioned in chapter 6.3.1

- Quarter since last warrant in last 2 year in any of the 3 tax types (OB,VPB, LH)

Variable	Bins	Good	Bad	Bad rate	% population
Quarter since last Warrant in last 2 year in any of the 3 tax types (OB,VPB, LH)	warrant in last 8 quarters	52,635	21,237	29%	5%
	No warrant in last 8 quarters	1,339,192	52,413	4%	95%

- # of quarters for which there was a delay in payment across all tax types in last 2 years

Variable	Bins	Good	Bad	Bad rate	% population
# of quarters for which there was a delay in payment across all tax types in last 2 years	1(0-1)	1,185,581	42,422	3%	84%
	2(2-2)	124,266	10,723	8%	9%
	3(3-3)	47,413	6,719	12%	4%
	4(4-19)	34,567	13,786	29%	3%

- Maximum delay (days) in last 2 years across all tax types

Variable	Bins	Good	Bad	Bad rate	% population
Maximum delay (days) in last 2 yrs across all tax types	delay between 1 to 7 days	1,213,580	39,366	3%	85%
	Delay between 8 - 30 days	98,855	9,913	9%	7%
	Delay between 30 - 61 days	43,672	9,280	18%	4%
	Delay more than 61 days	35,720	15,091	30%	3%

- Years in Business

Variable	Bins	Good	Bad	Bad rate	%
----------	------	------	-----	----------	---

					population
years in business	Missing	448	23	4.9%	0%
	0-2 years	162,515	10,480	6.1%	12%
	3-5 years	196,600	11,150	5.4%	14%
	6-10 years	319,237	17,419	5.2%	23%
	11-15 years	233,581	12,449	5.1%	17%
	16-20 years	177,360	8,517	4.6%	13%
	>20 years	302,086	13,612	4.3%	22%

- # of quarters for which there was delay in OB filing in last 2 years

Variable	Bins	Good	Bad	Bad rate	% population
# of quarters in which there was delay in filing OB in last 2 years	Missing	343,172	8,115	2%	24%
	0	907,802	47,992	5%	65%
	1	108,154	10,387	9%	8%
	more than 1	32,699	7,156	18%	3%

- Ratio of Retained earnings over assets for last year

Variable	Bins	Good	Bad	bad_rate	% population
Ratio of retained earnings over asset of last year	Missing	533,506	25,999	5%	38%
	2(<-0.05)	213,629	16,898	7%	16%
	3(-0.05-0.3)	248,174	14,285	5%	18%
	4(0.3-0.7)	224,650	10,403	4%	16%
	>0.7	170,748	6,023	3%	12%
	0 value for previous year assets	1,120	42	4%	0%

## 6.6 Signal allocation

Giving a signal on a certain company is of central importance to the success of an Early Warning System, this section describes the methodology used to determine which companies to give which signals. Risk is commonly defined as “the probability of something happening multiplied by the effect by the resulting effect when it does”. The methodology used is one which stays near to this principle in signaling the companies with the largest risk. On one axis the probability of going bad is given as the score generated by the model and on the other axis the collectable tax is plotted as the measure of severity when a company actually goes bad. 6.2 shows conceptually where red yellow and green signals are located along the two axis of probability and value.

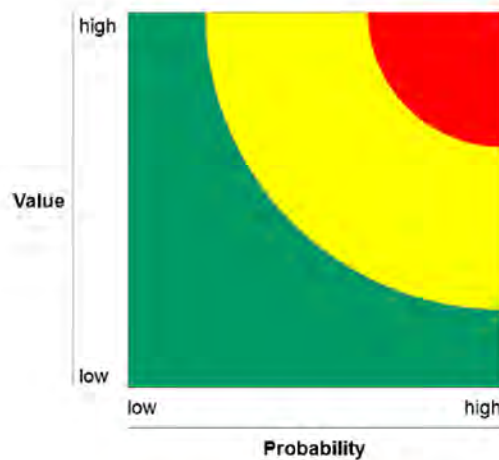


Figure 6.2: Conceptual signaling

In the EWS the probability is generated as the score from the model and value is defined as the net sum of the taxable amounts of “OB”, “VPB” and “LH”, the three main taxes for companies. To be able to properly analyze where the threshold should for each of the signals should be it is necessary to identify 1. the number of signals in each bucket, 2. the portion of the value in each bucket and 3. the number of companies that actually go bad.

To be able to do this in a user friendly way in Excel the score and value have been split into buckets, the scores are split into buckets according a step size of 0.04, and the value has been split into 10 percentile groups; this is done for both GO and MKB. Each of the buckets is then analyzed on how many finr there are in the bucket, how many actual bad and how much value is in the bucket. After the bucketing of the scores has been done several option are presented to the process experts and in an iterative process an option is chosen.

This step is followed by translating the red, yellow and green signals back to the finrs, this is done by examining the bucket in which a finr has been placed and assigning the according signal<sup>10</sup>.

The final step in the signaling is rolling up the signal from finr to entity; for this several options are examined and in collaboration with the process experts the highest score option is chosen:

<sup>10</sup> In the scoring code, 1 refers to red, 2 to yellow and 3 to green.

- **Value based:** summing the amount of taxable value which is in each signal of an entity (amount in red, yellow and green) and giving the entity the signal with the highest value.
- **Count based:** count the number of signals in red, yellow and green and giving the entity the signal with the most finrs.
- **Importance based:** give the signal that the head of the entity or highest value finr has received.
- **Highest score:** give the entity the highest (most risky) signal that is detected in all the underlying finrs.



## 7 Model Validation

This chapter includes the various analysis conducted to check the model stability and robustness

### 7.1 Diagnostic Tests for Potential Violations of Model Requirements/Assumptions

#### 7.1.1 Model predictive power

Table below shows the GINI value on development, validation and out of time dataset. The GINI is consistent across all three datasets that proves that model is quite stable and robust

Dataset	GINI
Development	55%
Validation	54%
Out of time	54%

Table 7.1: GINI across samples

#### 7.1.2 Plausibility

The final coefficients for each variable are consistent in signs

#### 7.1.3 Statistical significance of model coefficients – model p-value

As mentioned in chapter 6.5.1, all variable coefficients are statistically significant at 99% confidence interval

#### 7.1.4 Consistent rank ordering

Below figure shows the comparison between average predicted probability and average bad rate across different buckets. The trend is quite similar that proves the prediction error is low and model is stable.

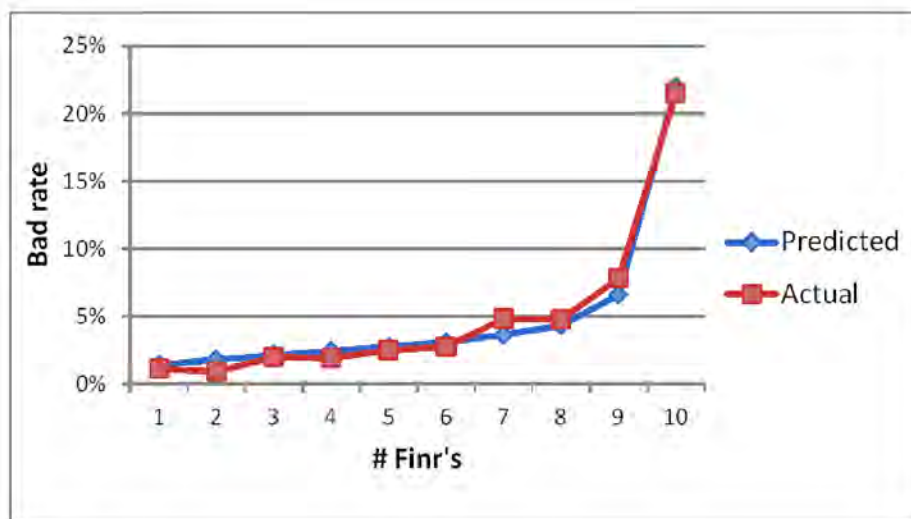


Figure 7.1: Actual and predicted bad rates by population bucket

### 7.1.5 Factor contributions

Below table shows that factor maximum delay is the most important factor in the model

Model variable	Weight in the model
Maximum delay (days) in last 2 years across all tax types	46%
Quarter since last warrant in last 2 year in any of the 3 tax types (OB,VPB, LH)	21%
# of quarters of delay OB filing in last 2 years	18%
Years in business	6%
Ratio of retained earnings over assets for last year	4%
# of quarter for which there was a delay in payment across all tax types in last 2 years	4%

Table 7.2: Factors contribution

### 7.1.6 Over fitting

Below table contains the model coefficient on development and validation dataset. The model coefficients are quite similar across development and validation sample that proves that there is no over fitting present.

Model variables	Development	Validation
Maximum delay (days) in last 2 years across all tax types	-0.6346	-0.6339
Quarter since last warrant in last 2 year in any of the 3 tax types (OB,VPB, LH)	-0.427	-0.4168
# of quarters of delay OB filing in last 2 years	-0.5758	-0.3624
Years in business	-1.5994	-1.5111
Ratio of retained earnings over assets for last year	-0.6299	-0.586
# of quarter for which there was a delay in payment across all tax types in last 2 years	-0.2166	-0.2259

Table 7.3: Coefficients comparison

### 7.1.7 Applicability

The model has a GINI of 54% in out of time data that proves that the model have good performance in recent dataset

### 7.1.8 Segment performance

The model performance is consistent across different segment that shows that model is quite stable. Below exhibit shows the GINI for GO /MKB segments and across different industries

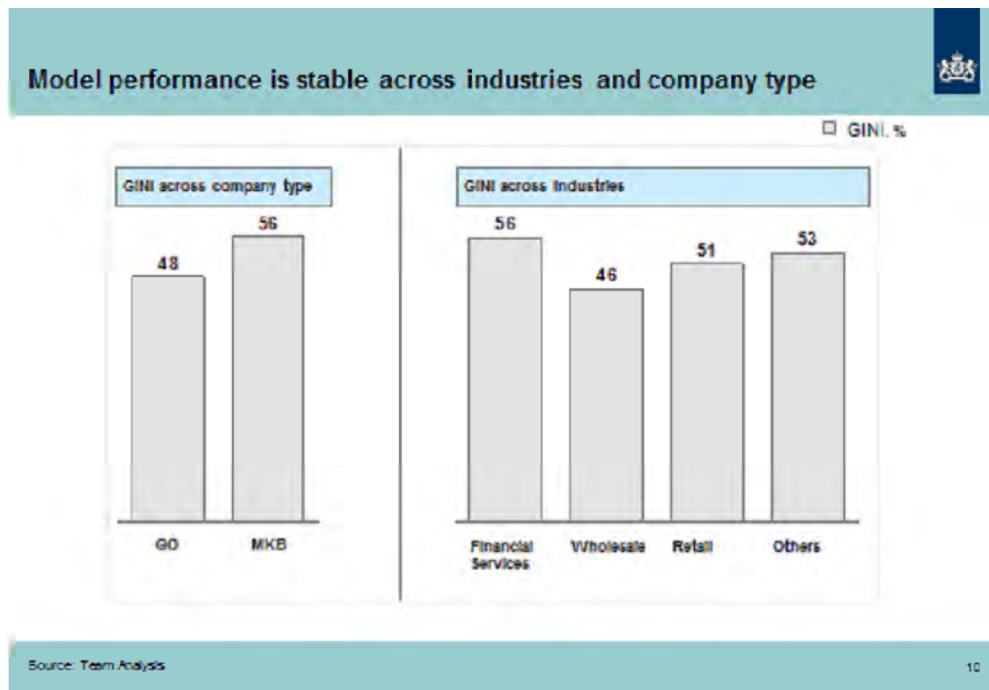


Figure 7.2: Segment performance

## 8 Appendices

### Table of Exhibits

Figure 2.1: Actual and predicted badrates by population bucket.....	8
Figure 3.1: Investeringsagenda.....	9
Figure 4.1: Conceptual effect of EWS on tax debt.....	11
Figure 4.2: Conceptual illustration of an early warning system.....	12
Figure 4.3: Overview of phases in the model development, it should be noted that the development is an iterative process.....	13
Figure 4.4: Time-line showing quarterly observation points, observation and performance windows, and in- and out-of-time samples.....	14
Figure 5.2: Representation of the data sources and the quality of data preparation.....	17
Figure 6.1: Gini explanation.....	26
Figure 6.2: Conceptual signaling.....	32
Figure 7.1: Actual and predicted bad rates by population bucket.....	33
Figure 7.2: Segment performance.....	35
Table 2.1: Final model variables.....	7
Table 2.2: Model validation results.....	7
Table 5.1: Variable transformations.....	24
Table 6.1: Fill rate assessment.....	27
Table 6.2: Correlation thresholds.....	27
Table 6.3: Model specifications.....	29
Table 7.1: GINI across samples.....	33
Table 7.2: Factors contribution.....	34
Table 7.3: Coefficients comparison.....	34

### 8.1 Variable list



Variable list .xlsx

### 8.2 Univariate analysis



20162204\_Univariate\_Gini\_overview4.xls

### 8.3 Correlation analysis



Correlation  
Analysis.xlsx

### 8.4 Scoring analysis



20160524  
ScoringAnalysisV3.xls



20160524  
Counts\_of\_score\_bu



20160524 Bucketing  
CodeV3.xlsx



20161905\_Determin  
ation of value treshol

### 8.5 Model Results



20160517\_Logistisch  
e\_regressie\_v4\_for\_

### 8.6 Scoring code



scoring\_code\_phase  
2\_vf V2 30052016.sa

## Draaiboek behandeling starters in 1<sup>e</sup> fase van ondernemerschap

### Inleiding

De negatieve aangifte van een onderneming wordt op basis van een risicoweging wel of niet geselecteerd voor een kantoortoetsing. Het risicomodel is nu nog niet in staat om de eerste fase van de onderneming goed te wegen. Het model heeft minimaal twee systeemaangiften in de voorliggende periode nodig om de huidige aangifte te beoordelen. Is er nog geen historie van minimaal 2 systeemaangiften aanwezig dan wordt de aangifte niet gewogen. De aangifte krijgt wel in de NnoBox het kenmerk "STARTER" mee.

De aangiften met de signalering STARTER moeten op het kantoor beoordeeld worden op risico. Wij ontvangen verhoudingsgewijs veel negatieve aangiften met het kenmerk "STARTER". Dat is op zich niet vreemd wat ondernemers moeten in het begin flink investeren. Toch moeten wij wel alert zijn omdat juist bij de nieuwe inschrijvingen wij regelmatig stuiten op fake ondernemingen.

De risico's bij de behandeling van negatieve aangiften zijn divers en de startersproblematiek is er maar een van. Wij hebben niet voldoende capaciteit om alle aangiften met het kenmerk "STARTER" te toetsen. Wij moeten op dit kenmerk dus door selecteren. Landelijk is voor de beoordeling van de 1<sup>e</sup> fase van ondernemers een percentage van inspectie, controle en toezicht vastgesteld. Dat betekent dat inspectie, controle en toezicht aangiften met het kenmerk STARTER in aanmerking kan komen voor toetsing.

### Nadere beoordeling ondernemer

De aangifte geeft feitelijk niet al te veel informatie over de onderneming. Via het telefonisch contact kan een "eerste indruk" worden verkregen en kunnen indien nodig facturen opgevraagd worden. Geadviseerd wordt om voorafgaand aan het bellen eerst de volgende punten door te lopen.

- - 
  - 
  - 
  - 
  - 
  - 
  -
- Vertrouwelijke bedrijfsinformatie

Bij het doornemen van bovenstaande punten zal een deel van de aangiften met het kenmerk "Starter" in de groep "Laag risico" terechtkomen. Een deel zal zoveel twijfel oproepen dat deze direct in toetsing worden genomen.

Probeer voor de overige aangiften via telefonisch contact een beter beeld te krijgen van de investeringen. Aandachtspunten tijdens het telefonisch contact:

- Zijn de investeringen ook in samenspraak met de activiteiten van de onderneming.
- Is er geen of weinig omzet vraag dan ook naar de huidige activiteiten. Wanneer is er omzet te verwachten?
- Is men bezig met een verbouwing vraag dan ook naar de tijdslijn voor oplevering.

Leg de verkregen informatie in het algemene scherm van de NnoBox vast. Vraag eventueel de facturen op als je tijdens het gesprek nog de nodige twijfel hebt.

**Selecteren blijft altijd een afweging en is soms best lastig. Het uitgangspunt dat in de week maar inspectie, controle en toezicht met het kenmerk "Starter" getoetst kan worden maakt het niet makkelijker. Durf te beslissen en weet dat er een vangnet bij verdere teruggaven aanwezig is.**





**Broedkamer OB**  
**Pilot OB negatief**  
**Kantoren Almere, Amsterdam en Zwolle**

### Versiebeheer

Versie	Datum	Auteur	Omschrijving
0.1	23 maart 2015	Persoonsgegevens	conceptversie
0.2	27 maart 2015		Aanpassing conceptversie
1.0	3 april 2015		Definitieve versie



### Inhoudsopgave :

- 1 INLEIDING**
- 2 Doel van de pilot**
- 3 Werkwijze algemeen**
- 4 Werkzaamheden pilot OB negatief**
  - 4.1 Selectie van normoverschrijdingen
  - 4.2 Voortraject Broedkamer OB
  - 4.3 Bijwerken NnoBox
  - 4.4 Controle op niet geselecteerde normoverschrijdingen
- 5 Behandeling op de pilotkantoren**
  - 5.1 Opvragen van bescheiden
    - 5.1.1 Telefonisch opvragen van bescheiden
    - 5.1.2 Opvragen van de bescheiden in briefvorm
    - 5.1.3 Teruggaaf van de aangifte voorkomen
  - 5.2 Bewaken tijdige ontvangst van bescheiden
  - 5.3 Geen reactie van de ondernemer
  - 5.4 Nadere beoordeling van de ondernemer
  - 5.5 Ontvangst reactie
  - 5.6 Beoordeling van de bescheiden
  - 5.7 Derdenonderzoek
  - 5.8 Uitgeven voor veldtoetsing
  - 5.9 Afwerken beslissing conform/correctie
  - 5.10 Aandachtspunt 2<sup>e</sup> ASD
  - 5.11 Tegentekenen / goedkeuren
  - 5.12 Teruggaven met een normoverschrijding van Inspectie, controle en toezicht

# Werkinstructie

---

## 1 INLEIDING

De resultaten van de landelijke behandeling van normoverschrijdingen OB blijft al langere tijd achter bij de verwachtingen. De resultaten zowel in de kantoortoetsing als in de veldtoetsing zijn lager dan bijvoorbeeld bij de behandeling van steekproefposten. Dat is vreemd want de posten worden geselecteerd voor behandeling toezicht.

Tijdens een bijeenkomst met de behandelaars is vastgesteld dat een deel van de normoverschrijdingen op basis van het aangiftepatroon geen risico's bevat. Daar staat tegenover dat risico's niet zichtbaar worden omdat de gestelde norm niet wordt overschreven.

De Broedkamer heeft de opdracht gekregen om op basis van een breder aanbod aan informatie en verbeterde technologie een selectiesysteem te bouwen. Het ontwikkelde systeem "OB negatief" is bij de Belastingdienst van Denemarken al in gebruik. Op basis van eerder verkregen informatie wordt een inschatting gemaakt van de nog te ontvangen aangifte. Na de vergelijking worden de grootste risico's doorgegeven. De handmatige selectie wordt hiermee grotendeels uitgeschakeld.

Begin april wordt een pilot gestart met de aanlevering van signalen vanuit het risicomodel OB negatief. De pilot heeft een tijdslijn van 4 maanden en de kantoren Zwolle (incl. Almere) en Amsterdam nemen hier aan deel.

In deze instructie is niet alleen de aanlevering van de signalen vanuit het risicomodel beschreven maar ook de werkafspraken met betrekking tot de behandeling van de aangiften.

## 2 Doel van de pilot

Inspectie, controle en toezicht Tijdens de pilot zal vastgesteld moeten worden of het model goed is "ingeregeld". Ongetwijfeld zal de instelling van parameters nog aangepast moeten worden.

Daarbij is de opgedane ervaring van kantoortoetsers van groot belang. Naast de harde informatie over wel of niet correctie is vooral de softe informatie van belang. Ook in de situatie dat er geen correctie is kan het signaal heel terecht zijn afgegeven en/of bruikbaar zijn geweest bijvoorbeeld voor een derden onderzoek.

Naast het verder verhogen van de kwaliteit van de risico's heeft ook het logistieke proces de volle aandacht. Taken en rollen worden op beide pilotkantoren nu nog anders belegd. De capaciteit wordt afgestemd op de ontvangst van aangiften. Dat vraagt in de piekperiode bij de ontvangst van kwartaalaangiften een bredere bezetting aan kantoortoetsers dan bij de maandaangevers. Behandelaars krijgen een breder werkaanbod dan alleen NNO werkzaamheden en dat vraagt de nodige flexibiliteit.

In de pilot gaan wij het logistieke traject, waaronder het werk op het juiste niveau beleggen, goed inregelen. Met behulp van een procesplaat wordt gewerkt aan een eenduidige werkmethode die op termijn ook landelijk uitgerold kan worden.

### 3 Werkwijze algemeen

In deze pilot richten wij ons op de normoverschrijdingen. Het is de bedoeling dat uiteindelijk alle aangiften zowel positief als negatief via een risicomodel worden beoordeeld maar dat traject krijgt gefaseerd vorm.

De normoverschrijdingen worden digitaal ingelezen en voorzien van een score. De hoogste risico's worden gesignaleerd. De uitworp zal echter niet alleen bestaan uit hoge risico's. Om het risicomodel goed te trainen worden er ook een aantal lagere risico's in behandeling gegeven. De totale uitworp aan de kantoren zal om en nabij de <sup>inspectie, controle en toezicht</sup> van het aantal normoverschrijdingen zijn. Dat is conform de landelijke productiecijfers NNO 2013/2014.

De geselecteerde aangiften worden via de NnoBox opgeleverd voor behandeling. De overige normoverschrijdingen worden dicht geboekt in de box. De geselecteerde aangiften hebben een harde uitworp. Dat betekent dat iedere aangifte wordt getoetst en dat hiervoor stukken opgevraagd moeten worden.

Men kan zowel telefonisch als schriftelijk stukken opvragen. Het schriftelijk opvragen wordt zoveel mogelijk centraal belegd. Mocht de reactie uitblijven dan wordt ook het herhaald verzoek centraal verzonden.

Na ontvangst van de stukken wordt de aangifte voor de beoordeling uitgegeven aan de kantoortoetsers.

## 4 Werkzaamheden pilot OB negatief

### 4.1 Selectie van normoverschrijdingen

Dagelijks worden de negatieve aangiften met een normoverschrijding in het Inspob-verslag opgenomen. <sup>inspectie, controle en toezicht</sup> verslag is ook in digitale vorm beschikbaar en het bestand met de <sup>inspectie, controle en toezicht</sup> wordt 's nachts ingelezen in de NnoBox. Voor de pilotkantoren 146(A'dam), 147(A'dam), 114(Almere), 115(Apeldoorn) en 121(Zwolle) wordt het automatisch dichtmelden voor aangiften met een risico indicatie van minder dan 10 punten met betrekking tot het segment MKB <sup>inspectie, controle en toezicht</sup> Alle MKB posten worden tot nader bericht "geblokkeerd" voor verdere behandeling. Alle velden worden hiervoor dicht geboekt. Ook de aangiften voor het segment G-O worden nog niet vrijgegeven voor selectie en behandeling.

De aangiftegegevens van de normoverschrijdingen worden dagelijks aangeleverd aan het risicomodel. Het risicomodel heeft op basis van de historie <sup>inspectie, controle en toezicht</sup> een inschatting van het aangiftebeeld. Hiervoor dient de ondernemer minimaal <sup>inspectie, controle en toezicht</sup> ingeleverd te hebben. Bij starters in het 1<sup>e</sup> jaar zal het model nog geen inschatting kunnen maken. N.B. De handmatige aangifte die in het basissysteem wordt verwerkt wordt niet als systeemaangifte herkent.

Het systeem levert tijdens de pilot zowel hoge als minder hoge risico's op. <sup>inspectie, controle en toezicht</sup> <sup>inspectie, controle en toezicht</sup> Het is dus geen test voor de behandelaars maar het goed valideren (=trainen) van het risicomodel. Het leeuwendeel van de uit te leveren risico's zal een inschatting "hoog risico" krijgen. De indicatie is niet zichtbaar voor de behandelaars. In de NnoBox wordt geen toekenning vanuit het risicomodel in punten of omschrijving zichtbaar.

## Werkinstructie

---

In de loop van de dag worden de gegevens van de geselecteerde aangiften doorgezeten naar de NnoBox. Het tijdstip is nu nog niet precies in te schatten. In de box worden de geselecteerde aangiften weer vrijgegeven voor behandeling. De overige aangiften blijven in het systeem dicht geboekt staan.

### 4.2 Voortraject Broedkamer OB

Dagelijks worden de aangiftegegevens van de inspectie, controle en toezicht in het segment MKB opgehaald en verwerkt in het risicomodel. De aangiften van het segment G-O worden niet meegenomen in de pilot. Het model selecteert op basis van risico's een aantal posten uit voor de verplichte toetsing. Er wordt geen risico indicatie mee gegeven.

Wel is het zo dat een aantal posten een kenmerk mee krijgen om de verdere behandeling te ondersteunen. De volgende kenmerken worden meegeleverd:

- 1 Gewijzigd
- 2 Starter
- 3 **Onderhand** inspectie, controle en toezicht
- Inspectie, controle en toezicht
- 5 toetsing correctie
- 6 toetsing conform
- Inspectie, controle en toezicht

Het is mogelijk dat er meerdere kenmerken van toepassing kunnen zijn en in die situaties wordt bovenstaande volgorde in de vermelding aangehouden. Er komt maar een kenmerk in de vermelding te staan.

De rubriek kenmerken kan zowel gevuld zijn bij de posten die geselecteerd zijn voor verplichte toetsing als bij posten die daar buiten blijven.

De kenmerken worden in de NnoBox geplaatst in een aparte kolom. Door op het veld KCer te drukken wordt een kolom "kenmerken" zichtbaar. Er zijn maximaal 7 posities beschikbaar. Daarom is de rubricering van het kenmerk aangepast.

Een korte toelichting op de kenmerken:

Gewijz.	Er is al eerder een aangifte over hetzelfde tijdvak in de NnoBox geplaatst. De aangifte moet beoordeeld worden op de vraag of de verbeterde aangifte wel of niet van invloed is op de keuze van behandeling.
Starter	Het risicomodel maakt onder meer gebruik van de eerder ingebrachte <u>aangiften om het risico in te schatten</u> . Als over voorliggende perioden nog geen <span style="border: 1px solid black; padding: 2px;">inspectie, controle en toezicht</span> aanwezig zijn dan wordt de nieuwe aangifte niet meegenomen in het risicomodel. De aangiften worden dan voorzien van het kenmerk starter.
OHW	Deze aangiften moeten geselecteerd worden op risico. Een eerdere toetsing staat in de NnoBox nog open. De aangifte moet beoordeeld worden op de vraag of de aangifte opgeleverd moet worden <u>aan de behandelaar en zo ja, of de uitbetaling ongeschort moet worden</u>

Inspectie, controle en toezicht

Vorcorr	Als in de pilot al eerder het OB nummer als risico is aangeleverd en de toetsing is met correctie afgedaan dan wordt het kenmerk meegegeven. Deze aangiften moeten beoordeeld worden op de vorm van verdere behandeling c.q. soort toetsing.
---------	--

## Werkinstructie

---

Vorconf Als in de pilot al eerder het OB nummer als risico is aangeleverd en de toetsing is afgedaan zonder correctie dan wordt dit kenmerk meegegeven. Deze aangiften moeten beoordeeld worden op wel of niet in toetsing nemen.

Inspectie, controle en toezicht.

G-O Alle aangiften worden dicht geboekt en pas na het inlezen van de selectie worden posten voor behandeling vrijgegeven. De posten van het segment Grote Ondernemingen worden niet meegenomen in deze pilot en worden daarom allemaal vrijgegeven met dit kenmerk.

### 4.3 Bijwerken NnoBox

Het risicomodel selecteert de posten op risico indicaties. De grootste risico's worden gekenmerkt als verplichte uitworp. De programmeur van de NnoBox krijgt dagelijks vanuit het risicomodel de terugmelding in score van alle inspectie, controle en toezicht In het overzicht is aangegeven of de aangifte wel of niet moet worden behandeld en/of de post op basis van kenmerken eerst beoordeeld moet worden.

Het ontvangen bestand wordt via een voorgeprogrammeerd traject in de NnoBox ingelezen. Het kenmerk wordt opgenomen in een aparte kolom. Daarna worden de normoverschrijdingen in de box verder bijgewerkt. De aangiften met het behandeladvies Nee worden in het systeem automatisch als conform dicht geboekt. Zowel de rubriek selecteur, behandelaar als de goedkeurder wordt afgevuld met "systeem". De aangiften met het behandeladvies Ja worden vrijgegeven voor verdere beoordeling.

Inspectie, controle en toezicht

### 4.4 Controle op niet geselecteerde normoverschrijdingen

Het risicomodel heeft op basis van de huidige instelling een inschatting gemaakt. In de pilot moet worden vastgesteld of de inschatting van risico's voldoende verantwoord is. Dat betekent dat niet alleen gekeken moet worden naar de geselecteerde aangiften maar ook naar de aangiften die een lager risico hebben.

De medewerkers van de kantoren Zwolle (Almere) en Amsterdam krijgen de harde uitworp opgeleverd. De focus moet bij de collega's vooral liggen bij de beoordeling van deze signalen en moet zo min mogelijk vertroebeld worden door de selectie van de overige normoverschrijdingen.

Vandaar dat de reeds dicht geboekte signalen bij een beperkte groep selecteurs (centaal team) worden neergelegd. Deze selecteurs vormen een virtueel werkverband waarbij wekelijks de posten worden toebedeeld voor beoordeling en de resultaten in de daarop volgende week worden besproken met de programmeurs van het risicomodel. Het kenmerk ondersteunt de beoordelaar in het maken van de afweging. Op het feedback formulier wordt aangegeven of er wel of niet een kenmerk is meegegeven.

## Werkinstructie

---

Alleen bij een groot risico op teruggaaf wordt de aangifte alsnog in behandeling genomen. De selecteur geeft hiervoor advies via een feedbackformulier. Dit formulier heeft als doel om de beoordeling van de selecteur nadien goed te kunnen vergelijken met de uitworpcriteria van het risicomodel. Als het voorstel wordt gevolgd dan wordt het signaal in de NnoBox weer ter behandeling gegeven.

Deze posten worden op een fictieve behandelaar (user-id wordt nog bekend gemaakt) gezet zodat men op de locaties dagelijks kan controleren of aangiften alsnog vrij zijn gegeven. Het in te vullen feedbackformulier wordt als bijlage opgenomen in de NnoBox.

## 5 Behandeling op de pilotkantoren

### 5.1 Opvragen van bescheiden

Het risicomodel heeft op basis van de historie dagelijks een aantal aangiften geselecteerd. De aangiften met een kenmerk worden door een ervaren kantoortoetsers beoordeeld op de wijze van behandeling. Een beperkt deel van de aangiften wordt dus nog steeds geselecteerd voor verdere behandeling.

Op basis van deze selectie worden signalen of (1) direct afgedaan als conform, of (2) op naam gezet van de huidige behandelaar of (3) opgeleverd voor overige klantbehandeling. De aangiften met een kenmerk hebben na de selectie voor kantoortoetsing ook de verplichte toetsing. De aangiften worden samen met alle aangiften zonder kenmerk getoetst. Daarvoor moeten er bescheiden worden opgevraagd. Dat kan zowel telefonisch als schriftelijk.

Bij het telefonisch contact is de kans aanwezig dat de ondernemer c.q. adviseur inhoudelijk op de aangifte in gaat. Het telefonisch opvragen moet dan ook gebeuren door de medewerkers die ook de kantoortoetsing uitvoeren.

Inspectie, controle en toezicht	
Inspectie, controle en toezicht	De signalen die niet telefonisch zijn benaderd worden administratief behandeld.

#### 5.1.1 Telefonisch opvragen van bescheiden

Het telefonisch opvragen heeft meerdere voordelen. Het is snel in het maken van contact, de rechtstreekse benadering bij een verzoek om teruggave wordt op prijs gesteld en het voorkomt een administratief traject van aangifte verwijderen en uitstel inbrengen.

Bij het telefonisch contact zal men niet in moeten gaan op verklaringen van de ondernemer of adviseur. De toetsing moet tenslotte plaatsvinden op basis van de nog toe sturen bescheiden. Dat kan aangeleverd worden via de post en als de ondernemer dat wil via de mail. Toezending kan het beste plaats vinden via een in Lotus hiervoor ingerichte postbus.

In de NNO box wordt het telefonisch contact (wie en wanneer) vastgelegd en in IKB wordt in een verslag van heffing het toetsen van de aangifte en de telefonische afspraken opgenomen.

## Werkinstructie

---

Als in een oogopslag het al duidelijk is dat de aangifte direct in veldtoetsing moet worden gegeven dan wordt de post via het formulier "uitgifte voor veldtoetsing" voorbereid voor uitgifte en het traject opvragen van stukken overgeslagen. Omdat wordt afgezien van selectie zal dit maar in een beperkt aantal gevallen van toepassing kunnen zijn. In deze situaties wordt eerst uitstel ingebracht om uitbetaling te voorkomen.

### Adviezen bij het bellen:

- Gebruik indien aanwezig het telefoonnummer wat in de NnoBox staat.
- Noteren wie je aan de telefoon krijgt, zeker als je een adviseur belt. Leg zowel de naam van het adviesbureau/accountantskantoor als de naam van de contactpersoon vast.
- Blijkt het telefoonnummer in de NnoBox niet te kloppen, dan kan je deze corrigeren d.m.v. het juiste nummer. Handig voor toekomstige telefonische contacten.
- De ondernemer /adviseur zal waarschijnlijk direct een toelichting willen geven. Maak in het begin van het gesprek duidelijk dat je de betreffende bescheiden toegezonden wilt hebben.
- Facturen laten mailen naar de postbus ..... (actie voor het kantoor) en aangeven dat de afzender het BTW nummer als kenmerk in de mail zet. Verder het verzoek om de gegevens bij voorkeur in één PDF file laten scannen en niet voor iedere factuur een aparte PDF laten maken. Dring aan op spoedige inzending en wijs op de consequenties als de stukken niet tijdig binnen zijn (blokkering/vertraging).
- Indien de ondernemer niet wil of kan mailen, dan is toezending per post of fax of bezorging vanzelfsprekend mogelijk. Wijs wel op de deadline en laat bel.pl weten dat tijdige toezending van invloed is op de termijn van blokkering van de teruggaaf.
- Noteer in de NnoBox en IKB goed wat je hebt afgesproken, wanneer je hebt gebeld enz. Een collega moet de post naadloos kunnen overnemen als je er zelf niet bent. Dus hoe meer je vastlegt, hoe beter en duidelijker het voor die ander is. Kijk wel uit met eventuele bedenkingen tegen de bel.pl ook schriftelijk vast te leggen. Dit kan problemen geven als de ondernemer deze informatie te zien krijgt, bijvoorbeeld bij een boekenonderzoek!
- Binnen de NnoBox kan een aparte gebruikersnaam aangemaakt worden. Wijzig de naam van de behandelaar in deze gebruiker om de post in beeld te kunnen houden als zijnde een lopende post waarvoor gebeld is en stukken zijn opgevraagd.

De tijdige ontvangst van stukken moet wel goed bewaakt worden. In principe wordt bij telefonisch contact geen uitstel ingebracht. Bij de tweede teruggaafselectie wordt de aangifte uitbetaald. Wekelijks moeten de teruggaafverzoeken van de vorige week doorgelopen worden en indien nodig moet er alsnog uitstel (zie par. 4.4) ingebracht worden.

### 5.1.2 Opvragen van de bescheiden in briefvorm

In de NnoBox is voor het kantoor een standaardbrief (zie bijlage 1) opgenomen om de bescheiden op te vragen. In de brief zijn alle aangifterubrieken met betrekking tot de omzet en omzetbelasting belasting apart opgenomen. De gegevens in de aangifte zijn bepalend om wel of niet de gegevens op te vragen.

## Werkinstructie

Inspectie, controle en toezicht	
aangifte	Inspectie, controle en toezicht
Rubriek 1a leveringen/diensten belast met hoog tarief	
Rubriek 1b leveringen/diensten belast met laag tarief.	
Rubriek 1c leveringen/diensten belast met overige tarieven behalve 0%	
Rubriek 1e leveringen/diensten belast met 0% of niet bij u belast	
Rubriek 2a leveringen/diensten waarbij de heffing van OB naar u is verlegd	
Rubriek 3a leveringen/diensten naar landen buiten de EU (uitvoer)	
Rubriek 3b leveringen naar of diensten in landen binnen de EU	
Rubriek 3c installatie / afstandsverkopen binnen de EU	
Rubriek 4a leveringen/diensten uit landen buiten de EU	
Rubriek 4b leveringen/diensten uit landen binnen de EU	
Rubriek 5b Voorbelasting	

In de NnoBox is de standaardbrief opgenomen. Men klikt de rubriek brief aan, selecteert het adres en de juiste persoon waarna d.m.v. dubbelklikken de brief wordt gegenereerd.

Inspectie, controle en toezicht

De brief verschijnt op het beeldscherm en kan hierna aangepast worden. Op de achterzijde van de brief staan de op te vragen rubrieken vermeld. De rubrieken die niet voldoen aan de criteria worden in het overzicht verwijderd.

De brief wordt hierna opgeslagen en tweemaal dubbelzijdig uitgeprint. Eenmaal voor verzending naar de ondernemer/adviseur en eenmaal voor deponering in het blanco omslagvel. Het omslagvel wordt gebruikt om de beslissing op het verzoek om teruggaaf vast te leggen en de mutatie aan de collega's van SMP door te geven. Op het omslagvel wordt het OB nummer, tijdvak en de naam van de ondernemer vermeld.

De opgeslagen brief wordt ook digitaal als bijlage in de NnoBox vastgelegd. Daarnaast wordt het toetsen van de aangifte direct in IKB onder verslag van heffing opgenomen en de brief wordt als bijlage toegevoegd.

### 5.1.3 Teruggaaf van de aangifte voorkomen

De negatieve aangifte wordt in principe bij de eerstvolgende teruggaafselectie meegenomen voor de beschikking teruggaaf.

Inspectie, controle en toezicht

Inspectie, controle en toezicht

Dat zal bij het opvragen van bescheiden via een brief bijna altijd een tekort tijdsbestek zijn. Vandaar dat uitstel van uitbetaling moet worden ingebracht. Uitstel inbrengen kan alleen als er geen aangiftegegevens vermeld staan. De aangifte moet dus eerst verwijderd worden. Als het een eerder aangiftetijdvak betreft waar de tweede aanslagselectie al van is verwerkt dan hoeft de niet ingebracht te worden. De uitbetaling is dan al via het systeem geblokkeerd.

#### Stap 1

Maak een printscreen van de ingebrachte aangifte.

#### Stap 2

Het verwijderen van de aangifte gebeurt door middel van de waarbij alleen het aangiftenummer op het mutatieformulier wordt vermeld.



## Werkinstructie

---

### Stap 3

Hierna kan de inspectie, controle en toezicht met uitstel voor onbepaalde tijd inspectie, controle en toezicht ingebracht worden. Op het formulier komt de inspectie, controle en toezicht het aangiftenummer en de inspectie, controle en toezicht

### Stap 4

De aangiftegegevens moeten hierna wel weer opnieuw in het basissysteem OB ingebracht worden. Dat kan via de inspectie, controle en toezicht met de volgende gegevens:

Aangiftenummer		Datum binnenkomst
Controlegetal		Dagnummer
CAP	0	Volgnummer
	Omzetbedrag	Bel.bedrag
1a	Belast met %	
1b	Belast met .,%	
1c	Belast met overige	
1d	Privegebruik	
1e	0%/Niet belast	
2a	Heffing verlegd naar aangever	
3a	Export buiten de EG	
3b	Leveringen binnen de EG	
3c	Instal./Televerk. EG	
4a	Import buiten de EG	
4b	Totaal verwervingen uit de EG	
	T O T A A L Omzetbelasting	
	Voorbelasting	
	KO-regeling	
	Schatting vorige aangifte(n)	
	Schatting deze aangifte	
	T O T A A L te betalen/terug te ontvangen	

Het opnieuw opnemen van de aangiftegegevens leidt de volgende dag tot een normoverschrijding. Deze overschrijding wordt door het risicomodel herkend en niet meegenomen in de selectie van die dag. Deze posten worden automatisch in de NnoBox dicht geboekt.

Geadviseerd wordt om eerst het traject van aangifte verwijderen, uitstel en opnieuw inbrengen voor alle posten te doorlopen en daarna de brieven samen te stellen en het verslag van heffing op te voeren in IKB.

## 5.2 Bewaking tijdige ontvangst van bescheiden

Na het verlenen van uitstel en verzending van de brieven wordt het omslagvel met hierin de verzonden brief en de print van de aangifte centraal in een kast per week op dagnummer (eventueel verder op BSN) bewaard. De ondernemer / adviseur heeft 3 weken de tijd om de stukken toe te sturen.

Na vier weken wordt een herinneringsbrief (zie bijlage 2) verzonden. Het rappel staat in de NnoBox opgenomen en kan gelijk aan de eerste brief via de box vervaardigd worden. Verzending van het rappel vindt eens in de week plaats. Voorafgaand wordt via BVR de juistheid van het toezendadres gecontroleerd. In het rappel moet ook aangegeven worden voor welke rubrieken de gegevens opgevraagd worden en dat is vanzelfsprekend gelijk aan de eerder verzonden brief.

## Werkinstructie

---

De rappelbrief wordt tweemaal geprint (voor verzending en het omslagvel) en daarnaast digitaal opgeslagen in de NnoBox en in het verslag van heffing in IKB

### 5.3 Geen reactie van de ondernemer

Als na vier weken nog geen reactie wordt ontvangen dan gaan wij het verzoek om teruggaaf afwijzen. In de herinneringsbrief is aangegeven dat op basis van de ingediende aangifte de kans aanwezig is dat er naheffingsaanslag wordt opgelegd met een verzuimboete van 10%.

De aangifterubrieken met een terug te vragen bedrag worden op € 0 gesteld omdat men niet heeft aangetoond recht te hebben op deze bedragen. Als in de aangifte ook positieve bedragen vermeld staan dan is de ondernemer feitelijk een bedrag verschuldigd. De basis voor deze berekening is de ingediende aangifte. Op het printscreen worden alle negatieve bedragen op € 0 gezet. Hierna wordt vastgesteld wat het saldo van de aangifte is geworden.

Bij het saldo € 0 wordt het verzoek om teruggaaf afgewezen met de standaardoverweging. Bij een positief saldo wordt het bedrag opgelegd door middel van een naheffingsaanslag. De aanslag is voorzien van een niet standaardoverweging en een verzuimboete van 10%.

tekstblok van de overweging:

*Op grond van artikel 47 en 49 van de Algemene wet inzake rijksbelastingen is een ieder gehouden desgevraagd aan de inspecteur gegevens en inlichtingen te verstrekken. Op informatieverzoeken van ..... en ..... is tot op heden geen reactie ontvangen. Het verzoek om teruggaaf op grond van artikel 17 Wet OB 1968 is niet gemotiveerd en wordt daarom voor de terug te vragen bedragen gecorrigeerd.*

Zowel de NnoBox als het verslag van heffing wordt verder bijgewerkt waarna het signaal zelf is afgedaan.

### 5.4 Nadere beoordeling van de ondernemer

Het is op zich vreemd dat de ondernemer / adviseur een verzoek om teruggaaf in stuurt en dan niet meer reageert. Het verzoek zelf is administratief afgehandeld maar een nadere beoordeling van de ondernemer is wel op zijn plaats. Vandaar dat het omslagvel met de brief, herinnering en de eventuele aanslag opgeleverd wordt aan de werkverdeler.

De werkverdeler beoordeelt op basis van het eerdere aangiftgedrag of Inspectie, controle en toezicht of eerdere aangiftetijdvakken alsnog in kantoortoetsing moeten worden genomen of dat er aanleiding is om een boekenonderzoeken te starten. Voor de aandachtsgebieden "midden" wordt de nadere beoordeling ook voor andere middelen uitgevoerd.

Bij nadere toetsing wordt opnieuw een verslag van heffing opgevoerd om de contacten met de ondernemer vast te leggen. Als er direct contact is wordt nagevraagd wat de reden is geweest van het niet reageren.

### 5.5 Ontvangst reactie

De meeste reacties worden door middel van de retourenveloppe direct op naam/afdeling opgeleverd. Maak indien nodig nog afspraken met de collega's van de CFD om onnodige registratie in GBV te voorkomen.

De ontvangen bescheiden worden toegevoegd aan het omslagvel. De ontvangst van de reactie wordt aan de hand van de eerder verzonden brief gecontroleerd op volledigheid van alle stukken. Mocht dat niet zo zijn dan wordt telefonisch / schriftelijk (evt. met tussenkomst van de werkverdeler) contact opgenomen met de ondernemer / adviseur. De werkverdeler beslist in deze situaties of de nog niet complete stukken wel of niet opgeleverd kunnen worden aan de kantoortoetsers.

Als de aangifte met de bescheiden afgegeven kan worden dan wordt in het verslag van heffing de ontvangst aangetekend en hierna wordt het omslag met alle bescheiden opgeleverd aan de werkverdeler.

De werkverdeler geeft op basis van de zwaarte van de te beoordelen posten in de NnoBox de aangiften uit aan kantoortoetsers. In principe is het beoordelen van de aangiften belegd op E niveau. Bij aspecten zoals vrijgestelde prestaties en onroerende goederen is opschaling naar F niveau soms wenselijk. Als de toetsing alleen betrekking heeft op de voorbelasting dan kan dit ook bij ervaren medewerker op C niveau belegd worden.

### 5.6 Beoordeling van de bescheiden

Inspectie, controle en toezicht

## Werkinstructie

---

### 5.7 Derdenonderzoek

Bij twijfel worden de ontvangen facturen bij de leverende onderneming beoordeeld op aangeven van omzet. Het vermelde bedrag moet over dezelfde periode dan minimaal op de aangifte vermeld staan. Als het bedrag niet is aangegeven dan kan de leverende ondernemer in onderzoek  worden gegeven.

Als de behandelaar twijfels heeft over de echtheid van facturen kan er ook een derdenonderzoek ingesteld worden bij de ondernemer die de facturen zou hebben uitgereikt.

Bij een derdenonderzoek komen ook ondernemers in beeld die niet tot het ambtsgebied van het betreffende kantoor behoren. Heeft het onderzoek betrekking op een onderneming van het kantoor dan kan intern afstemming plaatsvinden. Heeft het onderzoek betrekking op een onderneming buiten het ambtsgebied van het kantoor dan wordt voor de betreffende onderneming een notitie opgenomen onder de rubriek renseignements. De betreffende facturen worden gescand en als bijlage aan de notitie toegevoegd. Daarna wordt contact gezocht met het competente kantoor via het bestand contactpersonen (zie NnoBox handleidingen).

### 5.8 Uitgeven voor veldtoetsing

Als bij de kantoortoetsing wordt vastgesteld dat de toetsing beter ter plaatse kan worden uitgevoerd en/of ook voorliggende tijdvakken in onderzoek moeten worden betrokken dan wordt de toetsing van de aangifte voor veldtoetsing aangeboden. Bij uitgifte voor veldtoetsing wordt aan de hand van de vastleggingen in de administratie van de ondernemingen getoetst of de aangifte correct is ingevuld.

De reden van voordracht wordt vastgelegd op het formulier "INN veldtoets" (zie bijlage 3). Het formulier heeft bij onderdelen een keuzemenu wat is afgestemd op de toetsing van negatieve aangiften omzetbelasting. Bij het onderdeel aanleiding moet duidelijk omschreven worden waarom de aangifte over gaat van kantoortoetsing naar veldtoetsing. Het formulier wordt als bijlage geplaatst zowel in het verslag van heffing als in de NnoBox.

De omslag met bijbehorende stukken en het geprinte formulier worden opgeleverd aan de werkverdelers.

De coördinator veldtoetsing beoordeelt het signaal en voor de aandachtsgebieden "Midden" wordt vastgesteld of er op basis van andere middelen/processen aanleiding is om de opdracht verder uit breiden. De coördinator zorgt voor de toedeling aan de collega veldtoetsing en de oplevering van de stukken via het Toezicht ondersteuningsteam (TOT).

Het TOT voert de controleopdracht op  en zorgt voor zowel de digitale oplevering als de oplevering van het controledossier met stukken.

**NB. Het is een toetsing van een verzoek om teruggaaf waar termijnen aan verbonden zijn. Belangrijk om hier goede werkafspraken over te maken. Dat is ook van toepassing op de afhandeling omdat de afwerking van het  dossier soms wat langer op zich laat wachten en dat geeft een onnodige vertraging.**

De beslissing van de veldtoetsers wordt inhoudelijk tegen gelezen voordat het verzoek om teruggaaf wordt afgehandeld.

## Werkinstructie

---

### 5.9 Afwerken beslissing conform / correctie

Als de kantoortoetsers de bescheiden heeft beoordeeld en de eventuele vraagpunten zijn afgehandeld dan wordt in IKB het verslag van heffing gevuld. Bij afwijking op de aangifte wordt de correctie in een brief toegelicht en voorzien van een standpunt. De brief wordt zowel in IKB als in de NnoBox als bijlage opgenomen. De beslissing wordt duidelijk in IKB en de NnoBox vastgelegd. Het omslagvel wordt (evt. digitaal) ingevuld en voorzien van paraaf en de naam.

In de NnoBox wordt ook het resultaat van de correctie opgenomen en hierna wordt de aangifte afgedaan (1<sup>e</sup> vink).

De kantoortoetsers vult daarnaast het digitale feedback formulier in. Met deze informatie wordt een beter beeld verkregen van het risico. Met deze (softe) informatie kan het risicomodel beter ingeregeld worden. Het feedbackformulier wordt als bijlage opgenomen in de NnoBox. Daarnaast wordt het formulier opgeslagen

Inspectie, controle en toezicht

Hierna gaat het ingevulde omslagvel met de bescheiden naar de goedkeurder.

### 5.10 Aandachtspunt 2<sup>e</sup> ASD

In het systeem worden ook aangiften ingelezen over verstreken tijdvakken. Het systeem verwerkt de aangiften niet automatisch en teruggaafverzoeken over deze perioden moeten dus "handmatig" verwerkt worden.

De verwerking kan pas plaatsvinden als de teruggaaf akkoord is bevonden. De posten die centraal worden beoordeeld komen  open te staan. Bij de aantekening staat vermeld wat de reden van openstelling is. Het is dus mogelijk dat men geen risico heeft geconstateerd maar wel heeft vastgesteld dat de aangifte nog handmatig verwerkt moet worden.

Voor de overige signalen is in principe uitstel voor verleend. Daar zal altijd een handmatige verwerking moeten plaatsvinden. Wel graag alert zijn bij het telefonisch opvragen van stukken.

### 5.11 Tegentekenen / goedkeuren

De goedkeurder krijgt het omslagvel ter goedkeuring opgeleverd. Hij beoordeelt of de vastlegging van de behandeling volgens afspraken is verlopen en of dit ook voor derden goed is te begrijpen. De goedkeurder controleert of de beslissing juist op het omslagvel is vastgelegd.

Steekproefsgewijs worden bescheiden doorgenomen om ook inhoudelijk de kwaliteit te toetsen. Het is echter niet zo dat de goedkeurder alle signalen ook inhoudelijk zou moeten toetsen. Bij het inwerken van nieuwe collega's kan de inhoudelijke beoordeling ook verlegd worden naar degene die de collega inwerkt.

Na akkoordbevinding wordt het signaal door de werkverdelers in de NnoBox goedgekeurd (2<sup>e</sup> vink). Het omslagvel wordt met de bescheiden opgeleverd voor de mutatie in het basissysteem OB. Alle bescheiden blijven in het omslagvel en worden als onderdeel van het intoetswerk gearhiveerd.

## Werkinstructie

---

### 5.12 Teruggaven met een

Inspectie, controle en toezicht

Inspectie, controle en toezicht

Met Auditdienst van het Rijk (ADR) zijn afspraken gemaakt over normoverschrijdingen.

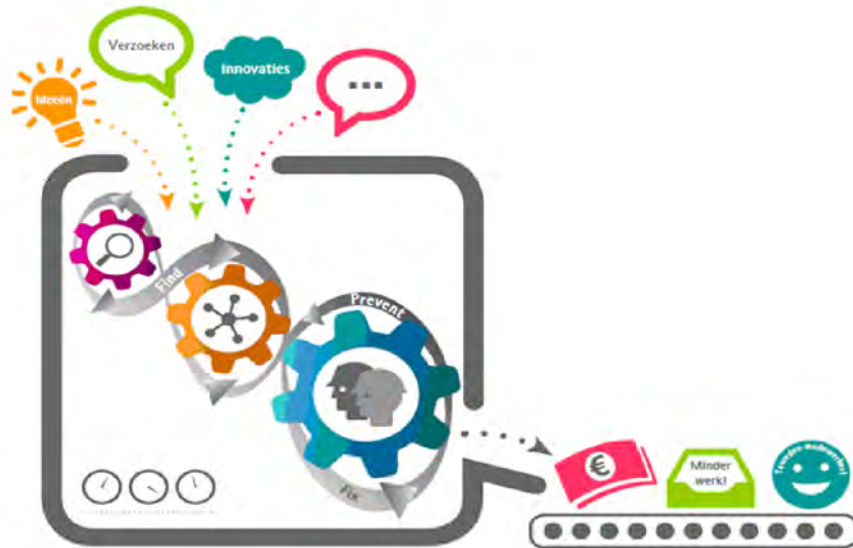
Inspectie, controle en toezicht

Inspectie, controle en toezicht

De lijnverantwoordelijke controleert of de inhoudelijke toetsing door de juiste niveau van medewerker en tegentekenaar heeft plaatsgevonden en of de vastlegging conform de instructie heeft plaatsgevonden.



Belastingdienst



## Werkinstructie OB negatief

*Datum: 12-08-2015*

*Versie: 0.5*

*Status: concept*



1. Inleiding
2. Doel & scope van de pilot
3. Werkwijze algemeen
4. Selectie aangiften door risico model
5. Controle op niet geselecteerde normoverschrijdingen door centrale team
6. Behandeling op de pilotkantoren



# 1. Inleiding



De Broedkamer heeft de opdracht gekregen om op basis van een breder aanbod aan informatie en verbeterde technologie een selectiesysteem te bouwen. Het ontwikkelde systeem “OB negatief” is bij de Belastingdienst van Denemarken al in gebruik.

Inspectie, controle en toezicht

Inspectie, controle en toezicht

Om vast te stellen of het model goed is ‘ingeregeld’, zal er een pilot gedaan worden op de pilotkantoren 146(A’dam), 147(A’dam), 114(Almere), 115(Apeldoorn) en 121(Zwolle). Ongetwijfeld zal de instelling van parameters nog aangepast moeten worden. Daarbij is naast de harde informatie over wel of niet correctie,

Inspectie, controle en toezicht

Inspectie, controle en toezicht

Ook in de situatie dat er geen correctie is kan het signaal heel terecht zijn afgegeven en/of bruikbaar zijn geweest bijvoorbeeld voor een derden onderzoek.

## 2. Doel & scope van de pilot

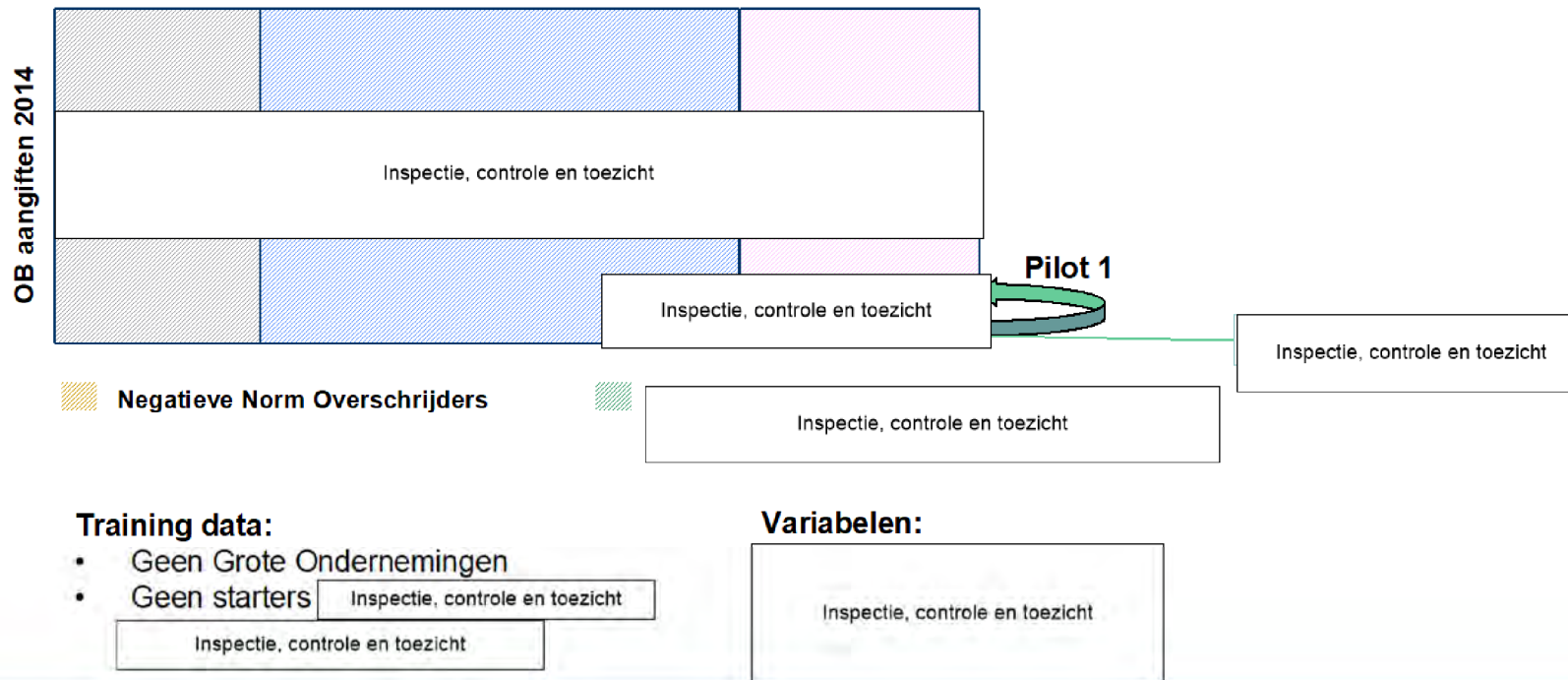


### Doel

Het valideren van een voorspellend risicomodel, vertrouwd en begrepen door de business, dat het OB negatief (door-)selectie proces verbetert, gebaseerd op goede en gevalideerde OB- en klantdata.

### Scope

Pilot 1 is gericht op de normoverschrijdingen



# 3. Werkwijze algemeen



## 1. Aangiften inlezen en blokkeren

Inspectie, controle en toezicht alle MKB posten in de NNO box op "Hold". Het risico model heeft Inspectie, controle en toezicht gedraaid vandaar Inspectie, controle en toezicht  
Inspectie, controle en toezicht posten. De aangiften van het segment G-O worden niet meegenomen in de pilot.

## 2. Risicomodel draaien

De normoverschrijdingen worden digitaal ingelezen en voorzien van een score. De hoogste risico's worden gesignaleerd. De uitworps zal echter niet alleen bestaan uit hoge risico's. Om het risicomodel goed te trainen worden er ook een aantal lagere risico's in behandeling gegeven. De totale uitworp aan de kantoren Inspectie, controle en toezicht van het aantal normoverschrijdingen zijn. Dat is conform de landelijke productiecijfers NNO 2013/2014.

## 3. Aangiften vrijgeven

In de loop van de dag worden de gegevens van Inspectie, controle en toezicht doorgezet naar de NnoBox. In de box worden de geselecteerde aangiften weer vrijgegeven voor behandeling, Inspectie, controle en toezicht  
Inspectie, controle en toezicht

## 4. Stukken opvragen

De geselecteerde aangiften hebben een harde uitworp. Dat betekent dat iedere aangifte wordt getoetst en dat hiervoor stukken opgevraagd moeten worden. Men kan zowel telefonisch als schriftelijk stukken opvragen. Het schriftelijk opvragen wordt zoveel mogelijk centraal belegd. Mocht de reactie uitblijven dan wordt ook het herhaald verzoek centraal verzonden.

## 5. Beoordeling

Na ontvangst van de stukken wordt de aangifte voor de beoordeling uitgegeven aan de kantoortoetsers.

# 4. Selectie aangiften



De aangiftegegevens van de normoverschrijdingen worden dagelijks aangeleverd aan het risicomodel. Het risicomodel heeft op basis van de historie een inschatting van het aangiftebeeld. Hiervoor dient de ondernemer

Inspectie, controle en toezicht

N.B. De handmatige aangifte die in het basissysteem wordt verwerkt wordt niet als systeem aangifte herkend.

Het systeem levert tijdens de pilot

Inspectie, controle en toezicht

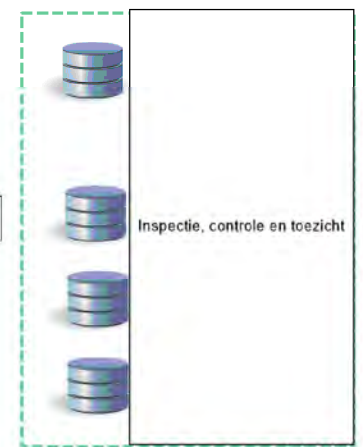
- Inspectie, controle en toezicht

- In de NnoBox wordt geen toekenning vanuit het risicomodel zichtbaar.

Inspectie, controle en toezicht

## Risicomodel

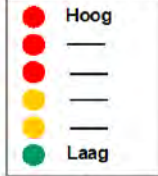
4 gekoppelde data bronnen



Risicoscore model



Risico lijst



• Hoogste kans op correctie

NB: Voor de pilotkantoren wordt het automatisch dichtmelden voor aangiften met een risico indicatie van

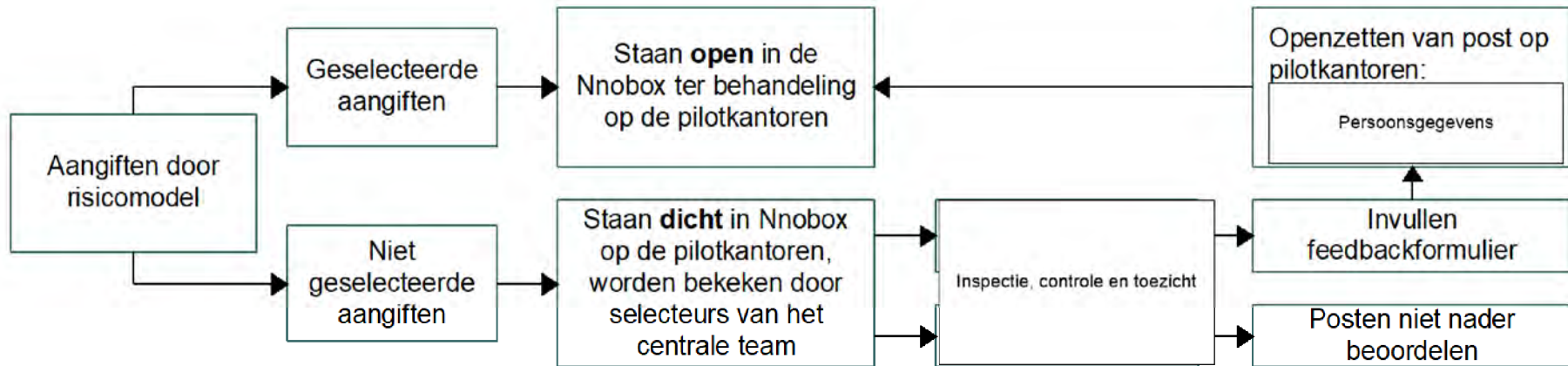
Inspectie, controle en toezicht

# 5. Controle op niet-geselecteerde aangiften



Het model heeft op basis van de huidige instelling  In de pilot moet worden vastgesteld of de inschatting van risico's  Dat betekent dat niet alleen gekeken moet worden naar de geselecteerde aangiften, die als harde uitworp voor de medewerkers van de pilot kantoren zichtbaar zijn

De aangiften  worden door een groep selecteurs (centaal team) beoordeeld.  aangifte alsnog in de Nnobox open gezet voor de pilot locaties   De selecteur geeft hiervoor advies via een feedbackformulier als bijlage in de Nnobox, om de beoordeling nadien goed te kunnen vergelijken met de uitworpcriteria van het risicomodel.





## 6. Behandeling op de pilot kantoren

6.1. Geselecteerde aangiften verdelen obv kenmerken

6.1.1. Proces aangiften verdelen

6.2. Opvragen bescheiden en teruggaaf voorkomen

6.3. Bewaking termijnen

6.4. Geen reactie

6.5. Nadere beoordeling

6.6. Ontvangst bescheiden

6.7. Beoordeling

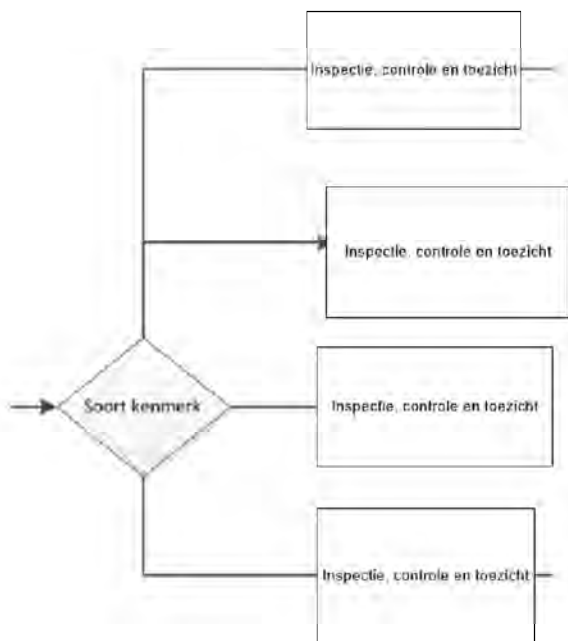
6.8. Verwerken beslissing

6.9. Veldtoetsing

# 6.1. Geselecteerde aangiften verdelen o.b.v. kenmerken

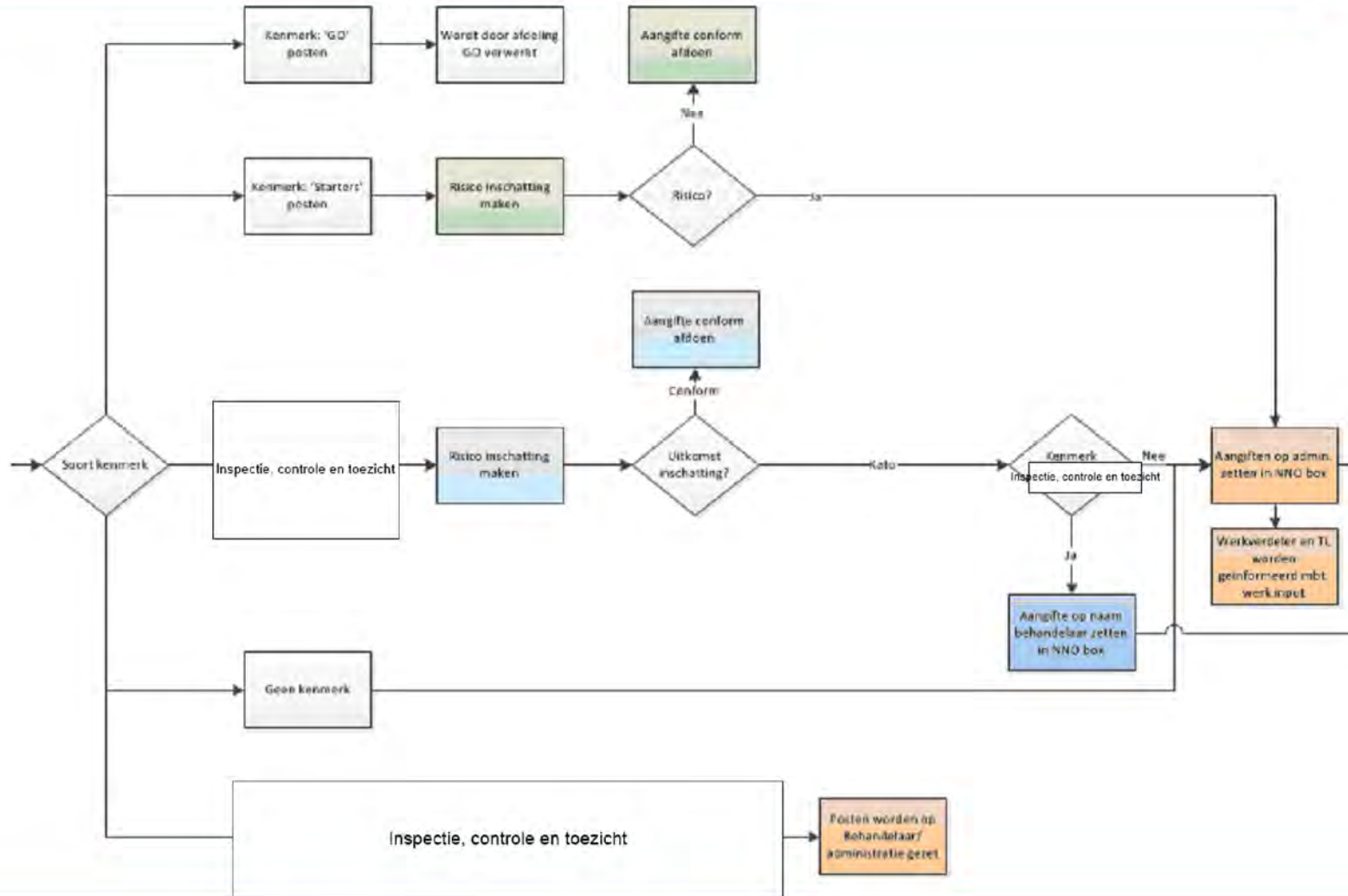


Het risicomodel heeft op basis van de historie dagelijks een aantal aangiften geselecteerd. Er komen aangiften zonder en met een kenmerk binnen. De aangiften met een kenmerk worden door een ervaren kantoortoetser beoordeeld op de wijze van behandeling. De aangiften zonder kenmerk worden direct doorgezet naar administratie. Administratie blokkeert de aangifte en vraagt bescheiden op middels een brief.



Kenmerk	Omschrijving
	Inspectie, controle en toezicht

# 6.1.1. Proces aangiftes verdelen



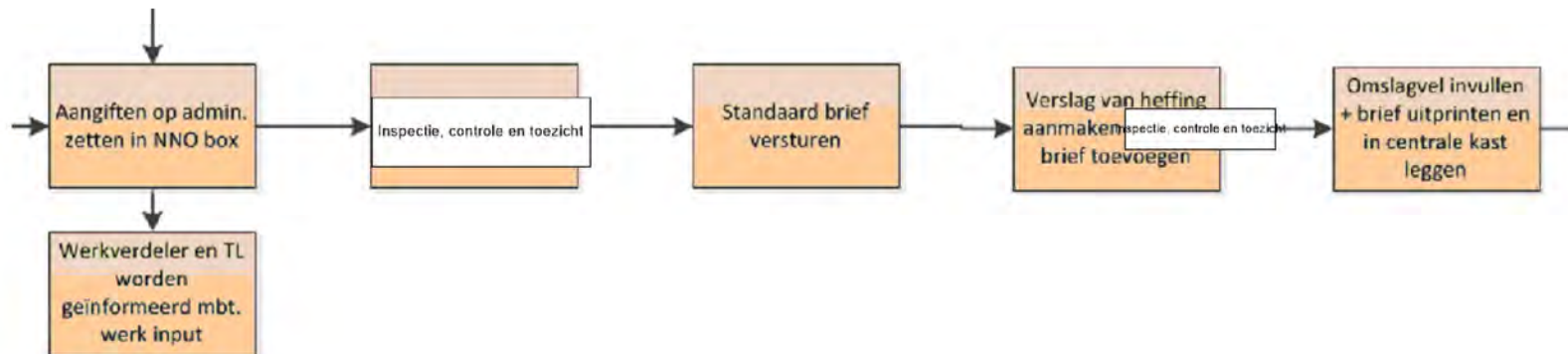


## 6.2. Opvragen bescheiden en teruggaaf voorkomen



Opvragen van bescheiden gebeurt door administratie:

- Aangifte wordt op ID van administratie gezet.
- Administratie blokkeert de aangifte doormiddel van
- Vervolgens wordt een standaardbrief opgesteld en richting klant verstuurd. Belangrijk is dat de adresgegevens van te voren worden gecontroleerd. Deze brief wordt vervolgens opgeslagen in verslag van heffing  en in de NNO box.
- De gegevens van de aangifte worden op een omslagvel genoteerd en in de centrale kast gelegd in afwachting van stukken.



## 6.3. Bewaking termijnen

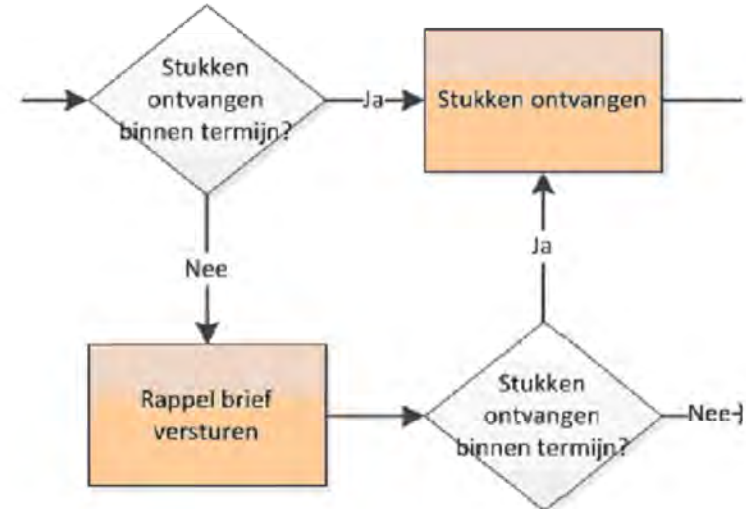


Wekelijks wordt gecontroleerd door administratie of bescheiden van opgevraagde posten binnen de wettelijke termijnen van 4 weken binnen zijn:

Bescheiden **niet** binnen: in de 5<sup>de</sup> week rappelbrief versturen

- Check of het toezendadres juist is
- Vraag de juiste rubrieken aan (zie bijlage)
- Print de brief 2x uit (voor verzending en het omslagvel)
- Sla de brief op in de Nnobox
- Plak in verslag van heffing (VvH)

Inspectie, controle en toezicht



- Bescheiden binnen: zie volgende slide

## 6.4. Geen reactie



Als na vier weken nog geen reactie wordt ontvangen (op rappelbrief) dan gaan wij het verzoek om teruggaaf afwijzen:

- De aangifterubrieken met een terug te vragen bedrag worden op € 0. Bij het saldo € 0 wordt het verzoek om teruggaaf afgewezen met de standaardoverweging.
- Bij een positief saldo wordt het bedrag opgelegd door middel van een naheffingsaanslag. De aanslag is voorzien van een niet standaardoverweging en een verzuimboete van Inspectie, controle en toezicht

tekstblok van de overweging:

*Op grond van artikel 47 en 49 van de Algemene wet inzake rijksbelastingen is een ieder gehouden desgevraagd aan de inspecteur gegevens en inlichtingen te verstrekken. Op informatieverzoeken van ..... en ..... is tot op heden geen reactie ontvangen. Het verzoek om teruggaaf op grond van artikel 17 Wet OB 1968 is niet gemotiveerd en wordt daarom voor de terug te vragen bedragen gecorrigeerd.*

Zowel de NnoBox als het verslag van heffing wordt verder bijgewerkt waarna het signaal zelf is afgedaan.

