

# VU Research Portal

## Happy@Work

Geraedts, A.S.

2014

### **document version**

Publisher's PDF, also known as Version of record

[Link to publication in VU Research Portal](#)

### **citation for published version (APA)**

Geraedts, A. S. (2014). *Happy@Work: E-mental health in occupational health care*. [PhD-Thesis - Research and graduation internal, Vrije Universiteit Amsterdam].

### **General rights**

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

### **Take down policy**

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

### **E-mail address:**

[vuresearchportal.ub@vu.nl](mailto:vuresearchportal.ub@vu.nl)

# Chapter 7

---

***The longitudinal prediction of costs due to health care uptake and productivity losses in a cohort of employees with and without depression or anxiety***

---

Anna S. Geraedts, Marjolein Fokkema, Annet M. Kleiboer, Filip Smit, Noortje M. Wiezer, Maria Cristina Majo, Willem van Mechelen, Pim Cuijpers, Brenda W. J. H. Penninx

*Journal of Occupational and Environmental Medicine 2014, 56(8):794-801*

## **Abstract**

**Objective:** To examine how various predictors and subgroups of respondents contribute to the prediction of health care and productivity costs in a cohort of employees.

**Method:** We selected 1548 employed people from a cohort study with and without depressive and anxiety symptoms or disorders. Prediction rules, using the RuleFit program, were applied to identify predictors and subgroups of respondents, and to predict estimations of subsequent 1-year health care and productivity costs.

**Results:** Symptom severity and diagnosis of depression and anxiety were the most important predictors of health care costs. Depressive symptom severity was the most important predictor for productivity costs. Several demographic, social, and work predictors did not predict economic costs.

**Conclusions:** Our data suggests that from a business perspective it can be beneficial to offer interventions aimed at prevention of depression and anxiety.

**Keywords:** Depression; Anxiety; Economic costs; Prediction rules

## Introduction

Depressive and anxiety disorders are highly prevalent in both the general [1-3] and the working [4-5] population. For example, in the Netherlands, the lifetime prevalence of depressive disorders is 19% and similar rates are reported for anxiety disorders (19.3%) [6]. Both disorders have a major impact on social relations, family life and activities at work. The adverse impact of depressive and anxiety disorders on daily life functioning has shown to result in substantial economic costs [7-10]. Smit et al. [9] calculated and compared the total costs of common mental disorders, such as mood disorders, anxiety disorders, and alcohol disorders, in a Dutch population-based cohort study. Mood disorders showed to have the highest total costs of all common mental disorders. The per-person excess costs for a mood disorder were €5009 annually, which sums to €311 million per one million persons aged 18-65 years. Anxiety disorders ranked second, with €3587 per person, totalling to €404 million per one million people aged 18-65 years. Comparable costs have been reported in other studies and the higher total costs of anxiety disorders can be explained by the high prevalence of this disorder [7, 8, 11-20].

Productivity losses stemming from absenteeism and decreased productivity while at work (presenteeism) [9, 21-25] account for the largest share of the total costs of depressive and anxiety disorders, with estimates ranging between 70-85% for depression [9, 26, 27]. One study [15] on the prevalence and costs of disorders of the brain within the European Union (EU) showed that the estimated annual costs of work absenteeism and lost productivity while at work (henceforth called 'productivity costs') in the EU due to mood disorders were €72 billion for depression and €28 billion for anxiety disorders. The costs for anxiety disorders were lower, because anxiety disorders have less impact on activities at work. These productivity costs are directly paid for by employers.

Since the economic burden of depression and anxiety disorders is substantial, it is important to shed light on the main drivers that contribute to these costs. This would provide a more precise understanding of predictors of the costs and when significant cost drivers can be identified, it might be possible to ameliorate these factors before costs are generated. It is thereby important to distinguish between health care costs and productivity costs, since both cost groups might reveal different predictors. However, little is known about predictors of health care and productivity costs and only few studies have been published on this topic. Chisholm et al. [28] found a positive association between health care costs and depression severity on cross-sectional data of a multi-national cohort-study (i.e. Longitudinal Investigation of Depression Outcomes). These costs were amplified by comorbid somatic illnesses. They also found an association between depression severity and productivity costs, but this relationship was not observed in all countries that took part in the cohort study. Knapp et al. [29] have studied predictors of costs of community care

for former psychiatric hospital in-patients in the UK. These costs were calculated over the first twelve months after discharge and the influence of several predictor variables, such as gender and age, that were measured at discharge were studied in order to predict the costs. They found that costs of community care could partly be predicted by characteristics such as age, marital status and psychiatric history. Although these studies have established the importance of some predictors of health care and productivity costs, knowledge in this field is still limited and only available for individual predictors.

The current study contributes to the further understanding of cost drivers from an occupational medicine perspective. In line with the previous studies in this field, we will examine the contribution of several individual predictors of health care and productivity costs. In addition, this study will identify subgroups of respondents, which has to our knowledge not been attempted before, and will utilize a broad range of (psychiatric, demographic, social, and work) predictors. To identify important predictors of health care and productivity costs, we will use prediction rule ensembles. Prediction rules describe the characteristics of subgroups in a sample that have a markedly higher or lower value for the outcome variable (economic costs, in this study). For example, one could find that the subgroup of men above the age of 45 have estimated €350 higher health care costs than others. The description of this subgroup of respondents provides us with a more specific prediction of costs than would be provided by age and gender, as separate predictors. In addition, combining multiple prediction rules into an ensemble improves predictive accuracy [30]. Prediction rule ensembles provide an estimate of the relative importance of individual predictor variables for prediction of the outcome variable, as well.

The analyses in this study were based on a Dutch cohort study [31] and were studied using a longitudinal design; the predictors were measured at one assessment and the outcome variable (the economic costs) was measured a year later and covered the costs between the first and second assessment. This longitudinal design allowed for prediction of future economic costs. In order to make an accurate estimation of the productivity costs and to be able to explain the productivity costs, we selected respondents from the cohort study with a paid job for at least eight hours per week.

We have looked at several predictor variables; diagnosis and symptom severity of depression and anxiety, demographic variables, psychosocial working conditions, and social conditions. We have selected these variables as they have been shown to influence economic costs and/or the severity or development of depressive and anxiety disorders in various cohort studies [6, 28, 29, 32-36]. The current study is considered exploratory as not much is known about predictors of health care and productivity costs. Based on previous studies [28, 29], we postulate that the severity of symptoms of depression and anxiety, older age, male gender, and being single are related to higher health care and productivity costs.

## Methods

### Research population

The Netherlands Study of Depression and Anxiety (NESDA) is a multisite naturalistic cohort study (N=2,981, age 18–65 years) examining the long-term course and consequences of depressive and anxiety disorders. The study was designed to include persons with depressive and anxiety disorders in different healthcare settings and stages of developmental history. Therefore, in addition to respondents with current disorders (n=1,550), respondents with remitted disorders (n=656), respondents at risk (due to family history, n=261), and healthy controls (n=514) were included. Participants with and without depressive and anxiety disorders were recruited from the general population (n=564), general practices (n=1,610), and in outpatient mental health organizations (n=807). Participants who were recruited from the general population were selected from two previous cohort studies in the Netherlands and these participants were asked to take part in the NESDA cohort study. Participants who were recruited from general practices were randomly selected from 65 general practices and included through a screening procedure. Participants who were recruited from outpatient mental health organizations were approached by the researchers after their intake at the mental health organization if their primary diagnosis was depressive or anxiety disorder. Across recruitment settings, uniform exclusion criteria were used; those with a primary diagnosis of obsessive compulsive disorder, bipolar disorder, psychotic disorder, or severe addiction disorder and those who were not fluent in Dutch were excluded from the cohort study. The NESDA study protocol was approved by the Ethical Review Board of the participating institutes, and all participants provided written informed consent. The objectives, rationale and methods of NESDA have been described elsewhere [31].

In this study we use the data of the baseline ( $t_0$ ) and one-year follow-up ( $t_1$ ) assessments from the NESDA cohort. Predictor variables were measured at baseline and health care and productivity costs were measured at one-year follow-up. The baseline assessment included a face-to-face interview and written questionnaires to assess demographic and contextual characteristics, psychological and health outcomes, and clinical diagnosis. At  $t_1$  a written questionnaire was used which included changes in contextual circumstances, such as a change of address or marital status, and the severity of psychological and health outcomes. All assessments were conducted by trained research assistants and monitored.

For the current study, only respondents with a paid job for  $\geq 8$  hours per week at baseline (n= 1,873) were selected. Of those 1,548 (82.7%) participated in the one-year follow-up assessment and were included in the analysis. The NESDA cohort is a representative sample of persons with a history or current depressive and anxiety disorder, but is not a representative sample of the Dutch (working) population.

### **Predictor variables**

#### *Depressive symptoms*

Severity of depressive symptoms was assessed with the Inventory of Depressive Symptomatology (IDS) self-report questionnaire [37, 38]. This questionnaire consists of 30 items scored between 0 and 3. Of the items 11-14 only two items are answered by the respondent leading to a summed score between 0 and 84 [37] with higher scores indicating more symptoms. This questionnaire has good psychometric properties [37, 38].

#### *Anxiety symptoms*

Anxiety symptoms were measured with the use of the 21-item Beck Anxiety Inventory [39]. Items are scored on a 4-point Likert scale and the total score ranges between 0 and 63, with higher scores indicating more symptoms. The psychometric properties of the questionnaire are good [39].

#### *Psychiatric diagnoses*

Diagnoses of depressive and anxiety disorders were assessed at the baseline interview, by means of the Composite International Diagnostic Interview lifetime interview, version 2.1 [40]. The Composite International Diagnostic Interview is a structured interview to assess psychiatric diagnoses, as defined in the *Diagnostic and Statistical Manual of the American Psychiatric Association, Fourth Edition (DSM-IV)* [41]. The Composite International Diagnostic Interview lifetime version assesses the presence of a current disorder, but also the history of (remitted) depressive and anxiety disorders. For the present study we used current (past month), past 6 months, past year, and lifetime diagnosis of the following depressive disorders: major depressive disorder, minor depressive disorder and dysthymia, and of the following anxiety disorders: social phobia, panic disorder (with or without agoraphobia), agoraphobia, and generalized anxiety disorder. The numbers of depressive disorders in the past six months and during lifetime were also calculated by summing the number of episodes of major depression, minor depression, and dysthymia. The same calculation was made for anxiety disorders.

#### *Psychosocial working conditions*

Psychosocial working conditions were measured with the Job Content Questionnaire (JCQ) [42]. This questionnaire consists of four scales; job demands (work fast, enough time to do work, work hard, excessive work, conflicting demands; 5 items), decision latitude (e.g. job requires creativity, learn new things, have freedom to plan tasks; 13 items), job support (e.g. friendly coworkers, supportive management; 8 items), and job insecurity (steady work, job security, future layoff; 3 items). All items were rated with positive answers scored as one and negative answers scored as zero. The items of each scale were summed and divided by the number of items, which resulted in a scale score between zero and one.

*Social conditions*

Social conditions were assessed in two different ways. First, we asked for the frequency of the respondent's social activities. These were scored dichotomously; at least one social activity per month versus no frequent social activities per month. Second, we used the Close Person Inventory (CPI) [43] to determine the amount of social support received by the respondent's partner and up to two friends (confidants) with whom the respondent had contact on a regular basis. A total of four questions about emotional support provided by the respondent's partner were scored on a 5-point Likert scale ranging from "never" (1) to "very often" (5). Total scores ranged between 5 and 20 and the score 0 was used if no partner was present. The same questions were used to ask for emotional support received from the respondent's confidants. The scores on the questions were then summed and divided by two to generate a total emotional support score from confidants. Total scores ranged between 5 and 20 and the score 0 was used if no confidants were present.

*Demographic variables*

The demographic characteristics we assessed were: gender, age, educational level, social economic status, nationality (Dutch or non-Dutch), marital status (presence of a partner or not), educational level of partner, employment status of partner, and household income. The educational level was determined by the highest education finished with a diploma and the social economic status in accordance with the Dutch Central Bureau of Statistics (CBS) [44]. Net household income was categorized in a net household income of <€2000, €2000-3600, and >€3600 per month.

**Outcome measure: Economic costs**

Costs stemming from health care uptake and productivity losses were assessed at one-year follow-up ( $t_1$ ) with the Trimbos and iMTA Questionnaire on Costs Associated with Psychiatric Illness (TiC-P) [45] and the Dutch Manual for Cost Research [46]. This questionnaire contains two parts; Part One consists of questions related to healthcare utilization by respondents and Part Two consists of questions related to absenteeism and lesser productivity while at work (presenteeism). Part One was used to determine the total health care costs and Part Two was used to compute the costs stemming from productivity losses (productivity costs). Costs for hospital stays and medication use were not calculated, because hospitalization is not expected for common mental disorders and the pharmacy costs form an ignorable fraction of the total costs that are dominated by productivity losses and health care uptake, other than pharmacy use in European countries such as the Netherlands [18, 19]. The recall period of the questions was six months and the annualized costs were calculated by multiplying the estimated costs by two. All costs are expressed in Euro (€) for the fiscal year 2011.



### Analysis

To identify important predictors we used models consisting of prediction rules with use of the RuleFit program [47, 48]. A more detailed description of the RuleFit program can be found in a Supplemental Digital Content (see document Supplemental Digital Content 1, a detailed description of RuleFit) and in this section we will briefly describe the program.

Prediction rules provide two major advantages, compared to linear regression models, traditionally used in econometric prediction: applicability and flexibility. The applicability of the results of multiple linear regression analyses for practitioners working in applied settings has been questioned by several authors [49-51]. Instead, they propose the use of non-linear heuristics, which are more suited to the requirements and limits of real-world decision-making processes. At the same time, rule-based methods are flexible, in that they do not rely on distributional assumptions about the data, and have the ability to detect non-linear effects of predictor variables, without the need for explicit specification of those effects, prior to analysis [52].

RuleFit is a so-called ensemble method: a method that combines the predictions of a large number of simple models [30, 53]. By combining predictions of a large number of simple models, also known as base learners [54], ensemble methods perform better than any of their constituent members [55, 56].

The base learners in a RuleFit model are prediction rules: statements of the form *if [condition], then [prediction]*. The condition specifies a set of values of predictor variables, and the prediction specifies the expected increase or decrease in the criterion variable, when an observation satisfies the specified condition. For example, the prediction rule: *if [age  $\geq 35$  & IDS  $\geq 45$ ], then [+250 health care costs]* would indicate that the expected health care costs in the next year made by a person at least 35 years of age, with a IDS score of at least 45, are €250 higher, compared to persons who do not meet those conditions (i.e., persons who are aged below 35, and/or have IDS scores below 45).

Rules in a RuleFit model are derived from CART trees [57] and collected in an initial ensemble. Furthermore, each of the predictor variables is added to the initial ensemble, to allow for estimation of linear functions. The final ensemble is formed by applying a regularized regression of the response variable on all prediction rules and predictor variables. By default, the RuleFit program uses the lasso penalty [58] for determining the coefficients for the prediction rules and variables. Application of the lasso penalty forces most coefficients to zero, resulting in ensembles that are both stable and interpretable [57].

Besides coefficients, the output of the RuleFit program provides a measure of importance for every prediction function and input variable in the ensemble. The importance of an input variable is calculated as a weighted sum of the importances of the prediction functions in the ensemble in which the variable appears. Importances are rescaled, to have a maximum of 100 (indicating the most important predictor variable) and a minimum of 0 (indicating

a variable which does not contribute to the prediction of the ensemble). Inspection of predictor variables and their importances provides an overview of the relative importance of each predictor variable to the predictive model.

Calculations of the predictor variables and total health care and productivity costs (outcome measure) were performed in Stata [59] and then imported in the R program for statistical computing [60]. Due to missing values in the TiC-P, health care costs could be calculated for 1,339 respondents, and productivity costs could be calculated for 1,377 participants. All analyses were then performed in R. RuleFit [47] is implemented in R, as well. The default settings for RuleFit were used, with exception of the approximate maximum number of prediction functions in the final ensemble, which was set to ten. The RuleFit program includes all cases with non-missing values on the outcome variable for estimation of the model.

## Results

### Baseline characteristics

Baseline demographic characteristics, baseline scores, and diagnoses status of the 1,548 employed respondents are shown in Table 1.

### Health care costs

The first RuleFit model was built for the prediction of health care costs. These health care costs were made in the year following baseline assessment. The RuleFit ensemble for prediction of health care costs consisted of 16 prediction rules and linear functions, based on 11 predictor variables. This RuleFit ensemble explained about 14% of health care cost variance ( $R^2 = 0.140$ ).

The five most important predictors of health care costs are presented in Table 2. The number of depressive disorder diagnoses in the past six months was the most important predictor of health care costs in the next year. Depressive symptom severity, anxiety symptom severity, and a current diagnosis of major depression were also important predictors.

**Table 1.** Baseline information; demographics, severity scores, diagnosis (n=1548)

<b>Gender</b>		
Female	N (%)	1018 (65.8)
Male	N (%)	530 (34.2)
<b>Age</b>	Mean (SD)	41.9 (11.6)
<b>Educational level<sup>1</sup></b>		
Low	N (%)	259 (16.7)
Middle	N (%)	605 (39.1)
High	N (%)	684 (44.2)
<b>Nationality</b>		
Dutch	N (%)	1515 (97.9)
Other	N (%)	33 (2.1)
<b>Marital status</b>		
Partner	N (%)	1141 (73.7)
No partner	N (%)	407 (26.3)
<b>Educational level partner<sup>1*</sup></b>		
Low	N (%)	210 (18.6)
Middle	N (%)	424 (37.5)
High	N (%)	497 (43.9)
<b>Employment partner</b>		
Job	N (%)	908 (79.6)
No job	N (%)	233 (20.4)
<b>Nett household income**</b>		
< €2000	N (%)	569 (37.0)
€2000 - €3600	N (%)	784 (50.9)
> €3600	N (%)	186 (12.1)
<b>Socioeconomic status</b>		
Low	N (%)	67 (4.3)
Middle	N (%)	797 (51.5)
High	N (%)	684 (44.2)
<b>Depression</b>		
Severity (IDS)	Mean (SD)	19.2 (13.2)
Depressive disorder diagnosis in past 6 months	N (%)	540 (34.9)
One depressive disorder	N (%)	436 (80.7)
Two depressive disorders	N (%)	104 (19.3)
Lifetime depressive disorder diagnosis	N (%)	975 (63.0)
Dysthymia only	N (%)	18 (1.9)
MDD only	N (%)	690 (70.8)
MDD and Dysthymia	N (%)	267 (27.4)
<b>Anxiety</b>		
Severity (BAI)	Mean (SD)	10.6 (9.7)
Anxiety disorder diagnosis in past 6 months	N (%)	594 (38.4)
One anxiety disorder	N (%)	353 (59.4)
Two anxiety disorders	N (%)	191 (32.2)
Three anxiety disorders	N (%)	50 (8.4)
Lifetime anxiety disorder diagnosis	N (%)	859 (55.5)

<b>Psychosocial working conditions (JCQ)</b>		
Job demands	Mean (SD)	0.5 (0.3)
Decision latitude	Mean (SD)	0.7 (0.3)
Job support	Mean (SD)	0.7 (0.3)
Job insecurity	Mean (SD)	0.6 (0.2)
<b>Social conditions</b>		
Frequent social activities	N (%)	679 (44.2)
Social support partner (CPI)	Mean (SD)	15.5 (2.6)
Social support confidants (CPI)	Mean (SD)	13.4 (3.7)

Note: BAI= Beck Anxiety Inventory; CPI= Close Person Inventory; IDS= Inventory of Depressive Symptomatology; JCQ= Job Content Questionnaire; MDD= Major Depression Diagnosis. <sup>1</sup> low = lower vocational education or less, middle = general intermediate education or high school, high = higher vocational education or university.

\* n=10 missing \*\* n=9 missing

**Table 2.** Importances from the five most important variables appearing in the model for the prediction of health care and productivity costs.

<b>Importance</b>	<b>Variable</b>
<b>Health care costs</b>	
100.00	Number of depressive disorder diagnoses in past six months
92.04	Depressive symptom severity
89.75	Anxiety symptom severity
80.18	Current major depression diagnosis
60.32	Number of major depression diagnoses in lifetime
<b>Productivity costs</b>	
100.00	Depressive symptoms severity
39.32	Number of depressive disorder diagnoses in past six months
35.44	Gender
35.34	Age
16.34	Anxiety disorder diagnosis in lifetime

The ten most important prediction functions for health care costs are presented in Table 3, giving a more precise and detailed impression of the ensemble. As shown in Table 3, the expected health care costs of a person with an anxiety symptom severity score of 15 or less and without a current major depression diagnosis were €312 lower in the following year, than a person who did not meet these conditions. This prediction rule indicates that the expected health care costs of people without, or with only mild symptoms of depression and anxiety are lower compared to people with depression and anxiety symptoms. The second prediction rule in Table 3 further supports this; the expected health care costs of a person with a depressive symptom severity score of 35 or less and an anxiety symptom severity score of 27 or less were €343 lower in the year that followed compared to a person who did not meet these conditions, i.e. persons with a depressive symptom severity score of at

least 36 and/or an anxiety symptom severity score of at least 28. The third most important prediction function is a linear function and indicates that a unit increase in the number of depressive disorder diagnoses in the past six months resulted in an estimated €183 increase in health care costs in the next year.

**Table 3.** Ten most important prediction functions for health care costs.

Type	Description	Coefficient <sup>1</sup>	Importance
Rule	Anxiety symptom severity ≤ 15 No current major depression diagnosis	-312	100.00
Rule	Depressive symptom severity ≤ 35 Anxiety symptom severity ≤ 27	-343	80.08
Linear	Number of depressive disorder diagnoses in past six months	183	74.42
Rule	No current major depression diagnosis No current social phobia diagnosis	-195	60.87
Rule	Depressive symptom severity ≥ 25	176	53.23
Rule	Net household income < €2000 One or more depressive disorder diagnoses in past six months	211	49.14
Rule	No major depression diagnosis in past six months ≤ 2 anxiety disorder diagnoses in past six months	-118	36.97
Rule	Middle or high educational level of partner ≤ 6 episodes of major depression diagnoses in lifetime No panic disorder (without agoraphobia) in past six months	-85	28.45
Rule	≤ 6 episodes of major depression diagnoses in lifetime No panic disorder (without agoraphobia) in past six months	-107	28.42
Linear	Depressive symptom severity	2	18.47

<sup>1</sup>Rounded and expressed in Euro (€)

### Productivity costs

The second RuleFit model was built for prediction of productivity costs as they occurred in the year after baseline assessment. These are the costs stemming from work absenteeism and productivity losses. The RuleFit model consisted of 11 prediction rules and linear functions, based on 11 predictor variables. The RuleFit model explained about 11% of productivity cost variance ( $R^2 = 0.114$ ).

The five most important predictors of productivity costs are presented in Table 2. Depressive symptom severity was the most important predictor of productivity costs. The number of depressive disorder diagnoses in the past six months, gender, age, and a lifetime anxiety disorder diagnosis were important predictors as well. However, these predictors were of less importance than depressive symptom severity as indicated by the importance rates in Table 2.

The ten most important prediction functions of productivity costs are presented in Table 4 giving a more precise and detailed impression of the ensemble. The most important

prediction function is a linear function and indicates that with every point increase on the depressive symptom severity scale (IDS) the estimated increase in productivity costs in the year that followed was €91. The second most important prediction function in Table 4 is a prediction rule, indicating that the expected productivity costs of a person aged 30 or older with one or more diagnoses of a depressive disorder in the past six months were €2387 higher in the next year compared to a person who did not meet these conditions, i.e. people younger than 30 years of age and/or without a depressive disorder diagnosis in the past six months. The third prediction function in Table 4 is a prediction rule of lesser importance (importance= 41.95), indicating that the expected productivity costs of a man with a lifetime anxiety disorder diagnosis were €1054 higher in the year that followed than females and/or people without a lifetime anxiety disorder diagnosis.

**Table 4.** Ten most important prediction functions for productivity costs.

Type	Description	Coefficient <sup>1</sup>	Importance
Linear	Depressive symptoms severity	91	100.00
Rule	Age ≥30 ≥ 1 depressive disorder diagnosis in past six months	2387	90.72
Rule	Gender: male Anxiety disorder diagnosis in lifetime	1279	41.95
Rule	Gender: male Major depression diagnosis in past six months	1054	27.97
Rule	Depressive symptoms severity ≤ 23 Anxiety severity ≤15	-454	19.08
Rule	Middle or high educational level Depressive symptom severity ≤ 41 ≤1 depressive disorder diagnosis in past six months	-417	15.33
Rule	Depressive symptoms severity ≤ 38 No depressive disorder diagnosis in lifetime	-396	14.04
Rule	Gender: female Depressive symptoms severity ≤ 46	-274	11.48
Rule	Gender: female No depressive disorder diagnosis in lifetime	-223	9.59
Rule	Job control Depressive symptoms severity ≤ 25 No diagnosis of generalized anxiety disorder in lifetime	-66	2.84

<sup>1</sup>Rounded and expressed in Euro (€)

## Discussion

This study explored the role of several psychiatric, demographic, and psychosocial variables in the prediction of health care and productivity costs, in a cohort of employed respondents with and without depressive and anxiety disorders. Using ensembles of prediction rules both predictors and an estimation of the costs could be calculated.

The results of this study showed that for the prediction of health care costs both symptom severity and depressive and anxiety disorders were most important. Conversely, the absence of an anxiety or depressive disorder and mild symptom severity of anxiety and depression predicted lower health care costs. These findings are in line with the previous studies of Knapp et al. [29] and Chisholm et al. [28]. Other variables were of lesser importance for the prediction of health care costs. However, the prediction rules identified several subgroups of respondents, with markedly higher or lower health care costs. For example, the sixth prediction rule indicated that the estimated health care costs of a person with a net household income lower than €2000 and one or more depressive disorder diagnosis in the past six months were €211 higher compared to a person who did not meet these conditions. This finding indicates that the predictor net household income only affected health care costs in interaction with the predictor depressive disorder diagnosis.

The results of this study further showed that for productivity costs, depressive symptom severity was the most important predictor. The importance of depressive symptom severity in the prediction of productivity costs is in line with the findings of Chisholm et al. [28]. Age and gender also predicted productivity costs, but were less important than depressive symptom severity. Educational level and psychosocial working conditions were of lesser importance for the prediction of productivity costs, but both affected the productivity costs in interaction with other predictors, according to the sixth and tenth prediction rule.

The hypothesis that health care and productivity costs would be related to depressive and anxiety symptom severity, older age, male gender, and being single was partly supported by the results. Depressive and anxiety symptom severity were both important in the prediction of health care costs. For the prediction of productivity costs only depressive symptoms were important, as were older age and male gender. Marital status, as in being single, appeared not to be important in the prediction of costs, which is not in line with the findings of Knapp et al. [29]. The psychosocial working conditions were only of some importance in the prediction of productivity costs. There were also predictors that appeared not to be important; social conditions, social economic status, marital status and nationality. The finding of irrelevancy of marital status is not in line with the findings of Knapp et al. [29], but the other predictors have not been studied before and could therefore not be compared to other studies. These findings indicate that some variables, such as social conditions, are not related to healthcare or productivity costs, and may be irrelevant for the reduction of costs. However, the sample of this cohort study is a sample with an over representation of respondents with (a history of) depressive and anxiety disorders. It could therefore be possible that the presence of depressive and anxiety disorders are dominant and suppress the importance of other predictors. Furthermore, we were able to explain only a small percentage of the observed variation in health care and productivity costs; 14% of health care costs and 11% of productivity costs. This indicates that the larger share of the variance

could not be explained by the prediction model and that several other variables, not included in our dataset, can be important predictors of health care and productivity costs as well. Other potential predictors could be different psychiatric disorders, chronic medical diseases and other social circumstances, such as loneliness, as studies of Knapp et al. [29] and Chisholm et al. [28] have shown before.

### **Strengths and limitations**

This study has an important limitation that needs to be taken into account when interpreting the results. The NESDA sample is not a representative sample of the general Dutch population. The purpose of this cohort study is to examine the long-term course and consequences of depressive and anxiety disorders and focus on (statistical) relationships, rather than obtaining point estimates such as means and proportions, which will be biased. In this study, people with depressive and anxiety disorders were oversampled and therefore the prevalence rates were much higher compared to other cohort studies [2, 6, 8]. Furthermore, we selected those respondents from the total cohort who had a paid job for more than 8 hours per week. This group represented about two thirds of the entire sample but this study is not a cohort study with a specific focus on the working population, such as the Maastricht Cohort Study [61]. The findings of this study have therefore limited generalizability. As a consequence, we recommend against making very precise interpretations of the costs, but believe that the prediction functions have validity.

This study also has two important strengths. First, the results of this study are based on longitudinal data; predictor variables were studied at one assessment and the outcome variable was studied a year later. This results in more accurate predictions of the contributions of the predictors. Second, in this study we used the RuleFit program for the identification of important predictors, and the estimation of the health care and productivity costs. The RuleFit program has the unique possibility to identify non-linear effects (i.e., interaction effects) of predictors on health care and productivity costs. The advantage of using a rule-based method is that it does not rely on distributional assumptions. This is a major advantage in the use of cost data since they are frequently not normally distributed. In contrast to linear regression, prediction rules provide direct identification of subgroups that are likely to have high or low health care and productivity costs. To our knowledge we were the first to use a rule based method for the prediction of economic costs. Furthermore, despite the fact that we did not use a representative sample in this study, but a sample with an overrepresentation of respondents with a history or current presence of a depressive or anxiety disorder, the results should not be neglected as the NESDA sample helps to “zoom in” on the importance of mental disorders on productivity and health care costs. The results of this study can contribute to the identification of predictors of health care and productivity costs within a population of employees with a history or presence of mental disorders.



### **Implications and future research**

The results of this study showed that both symptom severity and a diagnosis of depression and anxiety were the most important predictors of health care and productivity costs. The average productivity costs were twice as high as the health care costs, implying that the employer pays for the larger share of the costs of depression and anxiety. In the Netherlands, employers are obliged to contribute to the employee's health insurance and therefore they pay an even larger part of the total costs. It could therefore be of interest to the employer to keep both health care and productivity costs as low as possible. The results of this study suggest that if the employer would invest in the treatment and prevention of depression and anxiety this might result in cost reductions for the employer, since the prediction rules indicated that the absence of depression and anxiety resulted in estimated lower health care and productivity costs. The importance of facilitating treatment for depression and anxiety by the employer in order to reduce costs has been suggested by other researchers as well, but more studies on cost-effective treatments for depression and anxiety are needed [62].

This study provides a first impression of the added value of using a rule based method for the prediction of economic costs. The rule based method identifies both predictor variables and subgroups of respondents, which gives a more precise view on predictors of health care and productivity costs. For example, the second prediction rule for the prediction of productivity costs indicated that the expected productivity costs of a person aged 30 or older with one or more depressive disorders in the past six months would be €2387 higher compared to a person who did not meet these conditions. None of the previous studies in this field [28, 29] have used a data-analytic method that incorporates interaction effects, and could therefore only show the effect of individual predictor variables. However, since this is an explorative study, with limited generalizability, it is not possible to give a precise and accurate answer to the question which groups will have high health care and productivity costs in the future and consequently should be monitored by, for example, the employer.

This study should be replicated within a representative cohort of the working population. For the prediction of economic costs the RuleFit program [47], or a similar program which uses a rule based method, should be used to give a more precise answer to the question which groups will make high health care and productivity costs in the future. This replication study should also include several other predictors, such as other psychiatric disorders, chronic medical diseases and other social conditions, since this study could only explain a small percentage of the variance in the health care and productivity costs.

### **Conclusion**

This study showed that symptom severity and diagnosis of depression and anxiety were the most important predictors of health care costs. Depressive symptom severity was the most important predictor for productivity costs. Several other predictors were identified but these were of lesser importance in the prediction of health care and productivity costs. This study also revealed several subgroups of respondents with markedly higher and lower estimated costs. The results of this study suggest that it could be beneficial for employers to facilitate treatment and prevention of depression and anxiety to reduce costs. A replication of this study based on a representative cohort study should be conducted in order to provide us with more precise answers to the question what groups of people are most likely to generate high costs stemming from health care uptake and productivity losses.

## References

1. Kessler RC, McGonagle KA, Zhao S, et al. Lifetime and 12-month prevalence of DSM-III-R psychiatric disorders in the United States: results from the National Comorbidity Survey. *Arch Gen Psychiatry*. 1994;51:8-19.
2. Alonso J, Angermeyer MC, Bernert S, et al. Prevalence of mental disorders in Europe: results from the European Study of the Epidemiology of Mental Disorders (ESEMeD) project. *Acta Psychiatr Scand*. 2004;109:21-27.
3. Waraich P, Goldner EM, Somers JM, Hsu L. Prevalence and incidence studies of mood disorders: a systematic review of the literature. *Can J Psychiatry*. 2004;49:124-138.
4. Wang JL, Adair CE, Patten SB. Mental health and related disability among workers: A population-based study. *Am J Ind Med*. 2006;49:514-522.
5. Organisation for Economic Co-operation and Development (OECD). *Sick on the job? Myths and realities about mental health and work*. Paris: OECD Publishing; 2012.
6. Bijl RV, Ravelli A, van Zessen G. Prevalence of psychiatric disorders in the general population: results of the Netherlands Mental Health Survey and Incidence Study (NEMESIS). *Soc Psychiatry Psychiatr Epidemiol*. 1998;33:587-595.
7. Berto P, D'Ilario D, Ruffo P, Di Virgilio RF. Depression: cost-of illness studies in the international literature: a review. *J Ment Health Policy Econ*. 2000;3:3-10.
8. de Graaf R, Tuithof M, van Dorsselaer S, ten Have M. *Absenteeism due to psychological or somatic disorders in workers. Results of the 'Netherlands Mental Health Survey and Incidence Study-2' (NEMESIS-2) - Verzuim door psychische en somatische aandoeningen bij werkenden. Resultaten van de 'Netherlands Mental Health Survey and Incidence Study-2' (NEMESIS-2)*. Utrecht: Trimbos-Instituut; 2011. [in Dutch]
9. Smit F, Cuijpers P, Oostenbrink J, Batelaan N, de Graaf R, Beekman A. Costs of nine common mental disorders: implications for curative and preventive psychiatry. *J Ment Health Policy Econ*. 2006;9:193-200.
10. Verow P, Hargreaves C. Healthy workplace indicators: costing reasons for sickness absence within the UK National Health Service. *Occup Med (Lond)*. 2000;50:251-257.
11. Andrews G, Issakidis C, Sanderson K, Corry J, Lapsley H. Utilising survey data to inform public policy: comparison of the cost-effectiveness of treatment of ten mental disorders. *BJP*. 2004;184:526-533.
12. Greenberg PE, Birnbaum HG. The economic burden of depression in the US: societal and patient perspectives. *Expert Opin Pharmacother*. 2005;6:369-376.
13. Greenberg PE, Kessler RC, Birnbaum HG, et al. The economic burden of depression in the United States: how did it change between 1990 and 2000? *J Clin Psychiatry*. 2003;64:1465-1475.
14. Greenberg PE, Sisitsky T, Kessler RC, et al. The economic burden of anxiety disorders in the 1990s. *J Clin Psychiatry*. 1999;60:427-435.
15. Gustavsson A, Svensson M, Jacobi F, et al. Costs of disorders of the brain in Europe 2010. *Eur Neuropsychopharmacol*. 2011;21:718-779.
16. Hawthorne G, Cheok F, Goldney R, Fisher L. The excess cost of depression in South Australia: a population-based study. *Aust N Z J Psychiatry*. 2003;37:362-373.
17. Hunkeler EM, Spector WD, Fireman B, Rice DP, Weisner C. Psychiatric symptoms, impaired function, and medical care costs in an HMO setting. *Gen Hosp Psychiatry*. 2003;25:178-184.
18. Löthgren M. Economic evidence in affective disorders: a review. *Eur J Health Econom*. 2004;Suppl 1:S20-S24.

19. Marciniak M, Lage MJ, Landbloom RP, Dunayevich E, Bowman L. Medical and productivity costs of anxiety disorders: case control study. *Depress Anxiety*. 2004;19:112-120.
20. Stewart WF, Ricci JA, Chee E, et al. Cost of lost productive work time among US workers with depression. *JAMA*. 2003;289:3135-3144.
21. Henderson M, Glozier N, Elliott KH. Long term sickness absence is caused by common conditions and needs managing. *BMJ*. 2005;330:802-803.
22. Lerner D, Henke RD. What does research tell us about depression, job performance, and work productivity? *J Occup Environ Med*. 2008;50:401-410.
23. Lim D, Sanderson K, Andrews G. Lost productivity among full-time workers with mental disorders. *J Mental Health Policy Econ*. 2000;3:139-146.
24. Kessler RC, Frank RG. The impact of psychiatric disorders on work loss days. *Psychol Med*. 1997;27:861-873.
25. Stansfeld S, Feeney A, Head J, Canner R, North F, Marmot MG. Sickness absence for psychiatric illness: the Whitehall II Study. *Soc Sci Med*. 1995;40:189-197.
26. Cuijpers P, Smit F, Oostenbrink J, de Graaf R, ten Have M, Beekman A. Economic costs of minor depression: a population-based study. *Acta Psychiatr Scand*. 2007;115:229-236.
27. Thomas CM, Morris S. Cost of depression among adults in England in 2000. *Br J Psychiatry*. 2003;183:514-519.
28. Chisholm D, Diehr P, Knapp M, Patrick D, Treglia M, Simon G. Depression status, medical comorbidity and resource cost. *Br J Psychiatry*. 2003;183:121-131.
29. Knapp M, Beecham J, Fenyo A, Hallam A. Community mental health care for former hospital in-patients. Predicting costs from needs and diagnoses. *Br J Psychiatry Suppl*. 1995;166:10-18.
30. Berk R. An introduction to ensemble methods for data analysis. *Sociol Methods Res*. 2006;34:263-295.
31. Penninx BWJH, Beekman ATF, Smit JH, et al. The Netherlands Study of Depression and Anxiety (NESDA): rationale, objectives and methods. *Int J Methods Psychiatr Res*. 2008;173(3):121-140.
32. Plaisier I, de Bruijn JGM, de Graaf R, ten Have M, Beekman ATF, Penninx BWJH. The contribution of working conditions and social support to the onset of depressive and anxiety disorders among male and female employees. *Soc Sci Med*. 2007;64:401-410.
33. Plaisier I, de Bruijn JGM, Smit JH, et al. Work and family roles and the association with depressive and anxiety disorders: differences between men and women. *J Affect Disord*. 2008;105:63-72.
34. Plaisier I, Beekman ATF, de Graaf R, Smit JH, van Dyck R, Penninx BWJH. Work functioning in persons with depressive and anxiety disorders: The role of specific psychopathological characteristics. *J Affect Disord*. 2010;125:198-206.
35. Plaisier I, de Graaf R, de Bruijn JGM, et al. Depressive and anxiety disorders on-the-job: the importance of job characteristics for good work functioning in persons with depressive and anxiety disorders. *Psychiatry Res*. 2012;200:382-388.
36. Teo AR, Choi H, Valenstein M. Social relationships and depression: Ten-year follow-up from a nationally representative study. *PLoS One*. 2013;4:e6.
37. Rush AJ, Gullion CM, Basco MR, Jarrett RB, Trivedi MH. The Inventory of Depressive Symptomatology (IDS): psychometric properties. *Psychol Med*. 1996;26:477-486.
38. Trivedi MH, Rush AJ, Ibrahim HM, et al. The Inventory of Depressive Symptomatology Clinician Rating (IDS-C) and Self-Report (IDS-SR), and the Quick Inventory of Depressive Symptomatology, Clinician Rating (QIDS-C) and Self-Report (QIDS-SR) in public sector patients with mood disorders: a psychometric evaluation. *Psychol Med*. 2004;34:73-82.

39. Beck AT, Epstein N, Brown G, Steer RA. An inventory for measuring clinical anxiety: psychometric properties. *J Consult Clin Psychol.* 1988;56:893-897.
40. World Health Organization. *Composite International Diagnostic Interview (CIDI)*. Geneva: WHO; 1990.
41. American Psychiatric Association. *Diagnostic and Statistical Manual of Mental Disorders - Fourth Edition (DSM-IV)*. Washington: American Psychiatric Association; 2001.
42. Karasek R, Brisson C, Kawakami N, Houtman ILD, Bongers PM, Amick B. The Job Content Questionnaire (JCQ): An instrument for internationally comparative assessments of psychosocial job characteristics. *J Occup Health Psychol.* 1998;3:322-355.
43. Stansfeld S, Marmot M. Deriving a survey measure of social support: the reliability and validity of the Close Person Questionnaire. *Soc Sci Med.* 1992;35:1027-1035.
44. Kardal M, Lobber BJH. *Healthy life expectancy in social economic status – De gezonde levensverwachting naar social economische status*. Den Haag: Centraal Bureau voor de Statistiek; 2008. [in Dutch]
45. Hakkaart-van Roijen L, van Straten A, Donker M, Tiemens B. *Manual Trimbos/iMTA questionnaire for costs associated with psychiatric illness (TIC-P)*. Rotterdam: Institute for Medical Technology Assessment; 2002.
46. Hakkaart-van Roijen L, Tan SS, BOuwmans CAM. *Dutch Manual for Cost Research. Methods and standard prizes for economic evaluations in health care – Handleiding voor kostenonderzoek. Methoden en standard kostprijzen voor economische evaluaties in de gezondheidszorg*. Diemen: College voor Zorgverzekeringen; 2011. [in Dutch]
47. Friedman J (2012). RuleFit with R [version 3]. [<http://www-stat.stanford.edu/~jhf/R-RuleFit.html>]
48. Friedman J, Popescu B. Predictive learning via rule ensembles. *Ann Appl Stat.* 2008;2:916-954.
49. Gigerenzer G, Goldstein D. Reasoning the fast and frugal way: models of bounded rationality. *Psychol Rev.* 1996;103:650-669.
50. Katsikopoulos K, Pachur T, Machery E, Wallin A. From meehl to fast and frugal heuristics (and back) new insights into how to bridge the clinical - actuarial divide. *Theory Psychol.* 2008;18:443-464.
51. Martignon L, Vitouch O, Takezawa M, Forster M. Naive and yet enlightened: From natural frequencies to fast and frugal decision trees. In: Hardman D, Macchi L, ed. *Thinking: Psychological perspective on reasoning, judgment, and decision making*. West Sussex: John Wiley & Sons; 2003:189-211.
52. Breiman L. Statistical modeling: The two cultures (with comments and a rejoinder by the author). *Stat Sci.* 2001;16:199-231.
53. Dietterich T. Ensemble methods in machine learning. In: Kittler J, Roli F, ed. *Multiple classifier systems 2000*. Heidelberg: Springer-Verlag Berlin; 2000:1-15.
54. Schapire RE. The strength of weak learnability. *Mach Learn.* 1990;5:197-227.
55. Opitz D, Maclin R. Popular Ensemble Methods: An Empirical Study. *J Artif Intell Res.* 1999;11:169-198.
56. Zhou ZH, Wu J, Tang W. Ensembling neural networks: many could be better than all. *Artif Intell.* 2002;137:239-263.
57. Breiman L, Friedman J, Olshen, R, Stone C. *Classification and Regression Trees*. New York: Wadsworth; 1984.
58. Tibshirani R. Regression shrinkage and selection via the Lasso. *J R Stat Soc Series B Stat Methodol.* 1996;58:267-288.
59. StataCorp. *Stata Statistical Software: Release 8*. College Station: StataCorp LP; 2003.

60. R Development Core Team. *R: A Language and Environment for Statistical Computing*. Vienna: R Foundation for Statistical Computing; 2010.
61. Kant IJ, Bültmann U, Schröder KAP, Beurskens AJHM, van Amelsvoort JGPM, Swaen GMH. An epidemiological approach to study fatigue in the working population: the Maastricht Cohort Study. *Occup Environ Med*. 2003;60(suppl 1):i32-i39.
62. Burnett- Zeigler I, Ilgen MA, Bohnert K, Miller E, Islam K, Zivin K. The impact of psychiatric disorders on employment: Results from a National Survey (NESARC). *Community Ment Health J*. 2013;49:303-310.

## Supplemental Digital File 1

### Rule ensembles

To identify important predictors we derived a model consisting of prediction rules by application of the RuleFit algorithm [1, 2]. RuleFit is a so-called ensemble method: a method that combines the predictions of a large number of simple models [3, 4]. Those simple models are referred to as base learners or weak learners, as their predictions often perform only marginally better than random guessing [5]. By combining predictions of a large number of base learners, ensemble methods perform better than any of their constituent members [6, 7]. The base learners in a RuleFit model are prediction rules: statements of the form *if [condition], then [prediction]*. The condition specifies a set of values of predictor variables, and the prediction specifies the expected increase or decrease in the criterion variable, when an observation satisfies the specified condition. Rules in a RuleFit model are derived from classification and regression trees (CARTs) [8].

### Classification and regression trees

The CART algorithm creates decision trees, which can be used for the prediction of binary outcome variables (classification trees) or continuous outcome variables (regression trees). The predictor variables in CART analysis can be categorical, ordinal or continuous. The CART algorithm recursively splits observations into increasingly smaller subgroups, whose members are increasingly similar with respect to the outcome variable. The algorithm starts with the root node, containing all observations in the dataset. Next, the observations are split using one predictor variable at a time: In every node, the variable and splitting point are selected, which split the observations into subgroups, for which the distributions of the outcome variable are most different. The CART algorithm produces binary splits only: observations are split into two subgroups at every split. The result is a decision tree consisting of branches and nodes. Every node of a decision tree can be turned into a decision rule, where the path from the root node to the node specifies the conditions of the rule.

### RuleFit algorithm

#### *Deriving rules*

In the RuleFit algorithm, a large number of CART trees are grown on randomly drawn from subsets of the data. By default, the fraction of observations drawn for growing each tree is a function of the sample size, with the fraction decreasing for increasing sample sizes. In the current study, the default setting was used, resulting in a fraction of randomly drawn observations of 0.28. The average number of terminal nodes in the CART trees can be controlled as well. In the current study, the default value of four terminal nodes on average was used. When a tree is grown, a rule is derived from every node of the tree. These rules

are functions of the input variables, taking a value of 1 when the rule applies, and a value of 0 when the rule does not apply. The trees are grown on random subsets of the data, until a prespecified number of rules has been generated and collected in the initial ensemble. In the current study, the prespecified number of rules was set to the default value of 2000. Also, by default, each of the predictor variables is added to the initial ensemble, to allow for estimation of linear functions, because linear functions are most difficult to approximate by prediction rules.

*Estimation of coefficients*

The final ensemble is formed by application of regularized regression of the response variable on all prediction rules and predictor variables. Whereas with ordinary least squares (OLS) regression, the coefficients of prediction functions are estimated by minimizing the residual sum of squares, with penalized regression, an additional penalty is placed on the coefficient. The RuleFit algorithm uses the lasso penalty [9] by default, and coefficients are estimated with the following function [2, equation 26]:

$$(\{\hat{a}_k\}_0^K, \{\hat{b}_j\}_1^n) = \arg \min_{\{a_k\}_0^K, \{b_j\}_1^n} \sum_{i=1}^N L(y_i, a_0 + \sum_{k=1}^K a_k r_k(x_i) + \sum_{j=1}^n b_j l_j(x_i)) + \lambda (\sum_{k=1}^K |a_k| + \sum_{j=1}^n |b_j|)$$

Where  $K$  is the number of prediction rules,  $n$  is the number of linear predictors, and  $N$  is the number of observations in the sample. Further,  $a_k$  is the coefficient of the  $k$ th prediction rule,  $b_j$  is the coefficient of the  $j$ th linear predictor variable,  $L$  is a loss function,  $\lambda$  is the penalty parameter,  $r_k$  is the  $k$ th prediction rule, and  $l_j(x_i)$  is a so-called Winsorized version of the  $j$ th linear predictor variable.

The value for the penalty parameter  $\lambda$  is determined by default by using a separate subset of the data, which is not used for estimation of the coefficients. This fraction of the dataset used for determining the penalty parameter is .2, by default. For smaller datasets ( $N < 1486$ ), k-fold cross validation is used to determine the value for  $\lambda$ , with smaller values for  $N$  resulting in a larger number of folds. Higher values for  $\lambda$  result in higher penalties for the coefficients, and therefore sparser models. The loss function  $L$  defaults to squared error loss, which was used in the current study, as well. Winsorized versions of predictor variables,  $l_j(x_i)$ , are used for robustness against outliers. Outliers are defined to be values smaller than the  $\beta$ , and larger than the  $(1-\beta)$  quantiles of the predictor variables. In the Winsorized version of variables, values smaller than the  $\beta$  quantile are replaced by the  $\beta$  quantile, and values larger than the  $(1-\beta)$  quantile are replaced the  $(1-\beta)$  quantile. In the current study,  $\beta$  is set to the default value of .025.

Due to application of the lasso penalty in the estimation function provided above, coefficient estimates are shrunken towards zero. This introduces a slight bias in the estimates, but greatly reduces their variance [9, 10]. The lasso penalty produces models which are more



sparse (including less prediction functions than models estimated with OLS), and both interpretable and stable, as a result [8, 9].

#### **Interpretation of the RuleFit ensemble**

The RuleFit program outputs all prediction functions (i.e., rules and linear functions) included in the ensemble, with their respective coefficients. The coefficients can be interpreted as regular regression coefficients: they represent the increase in the predicted value for a unit increase in the prediction function. Besides coefficients, the output of the RuleFit program provides a measure of importance for every prediction function and input variable in the ensemble. Importances for prediction functions are given by the absolute value of the coefficient, multiplied by their support. The support for a prediction rule is given by  $\sqrt{s_k \cdot (1-s_k)}$ , in which  $s_k$  is the proportion of observations in the full dataset to which the rule applies. The support for a linear function is the standard deviation of the corresponding predictor variable in the full dataset. The importance of an input variable is calculated as a weighted sum of the importances of the prediction functions in the ensemble in which the variable appears. Importances are rescaled, to have a maximum of 100 (indicating the most important predictor variable) and a minimum of 0 (indicating a variable which does not contribute to the prediction of the ensemble). Inspection of importances provides an overview of the relative importance of each predictor variables to the predictive model.

## References

1. Friedman J (2012). RuleFit with R [version 3]. [<http://www-stat.stanford.edu/~jhf/R-RuleFit.html>]
2. Friedman J, Popescu B. Predictive learning via rule ensembles. *Ann Appl Stat.* 2008;2:916-954.
3. Berk R. An introduction to ensemble methods for data analysis. *Sociol Methods Res.* 2006;34:263-295.
4. Dietterich T. Ensemble methods in machine learning. In: Kittler J, Roli F, ed. *Multiple classifier systems 2000*. Heidelberg: Springer-Verlag Berlin; 2000:1-15.
5. Schapire RE. The strength of weak learnability. *Mach Learn.* 1990;5:197-227.
6. Opitz D, Maclin R. Popular Ensemble Methods: An Empirical Study. *J Artif Intell Res.* 1999;11:169-198.
7. Zhou ZH, Wu J, Tang W. Ensembling neural networks: many could be better than all. *Artif Intell.* 2002;137:239-263.
8. Breiman L, Friedman J, Olshen, R, Stone C. *Classification and Regression Trees*. New York: Wadsworth; 1984.
9. Tibshirani R. Regression shrinkage and selection via the Lasso. *J R Stat Soc Series B Stat Methodol.* 1996;58:267-288.
10. Zou H, Hastie T, Tibshirani R. (2007). On the “degrees of freedom” of the lasso. *The Annals of Statistics.* 2007;35:2173-2192.