

# Sick dataset analysis

*Bogdan Jastrzębski*

*5 kwietnia 2020 r.*

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Initial Data Mining</b>	<b>2</b>
2.1	Balance of the “Class” Variable . . . . .	2
2.2	Distributions of Categorical Variables . . . . .	3
2.3	Missing values reduction . . . . .	5
2.4	Solution . . . . .	5
2.5	Skewness reduction . . . . .	6
<b>3</b>	<b>Prediction models</b>	<b>7</b>
3.1	Naive Bayes . . . . .	7
3.2	Logistic Regression . . . . .	8
3.3	Tree . . . . .	9
3.4	The KNN . . . . .	10
3.5	C-Tree . . . . .	11
<b>4</b>	<b>Conclusions and the last benchmark</b>	<b>12</b>

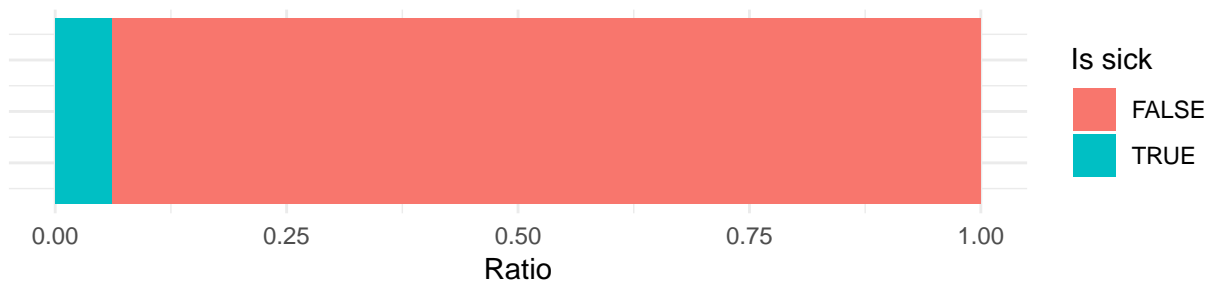
# 1 Introduction

In the following paper, I present an analysis of the “Sick” dataset, along with the strategy for predicting “Class” in an interpretable manner and it’s results.

## 2 Initial Data Mining

In this section I will address all the major issues with the dataset and describe the way to face them.

### 2.1 Balance of the “Class” Variable



The “Class” variable is not balanced, which indicates, that we need to measure model performance in more sophisticated way than calculating accuracy of a given model. I will use the following two measures:

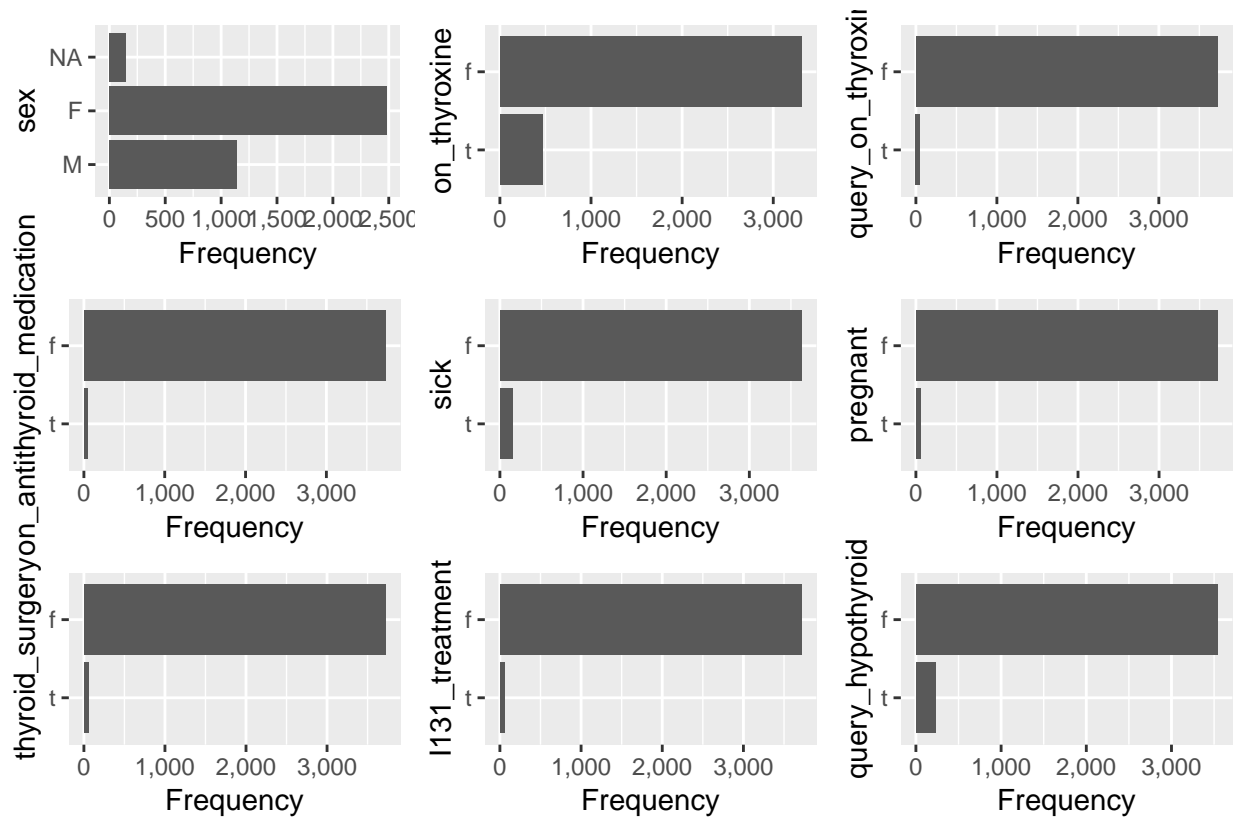
- auc
- auprc

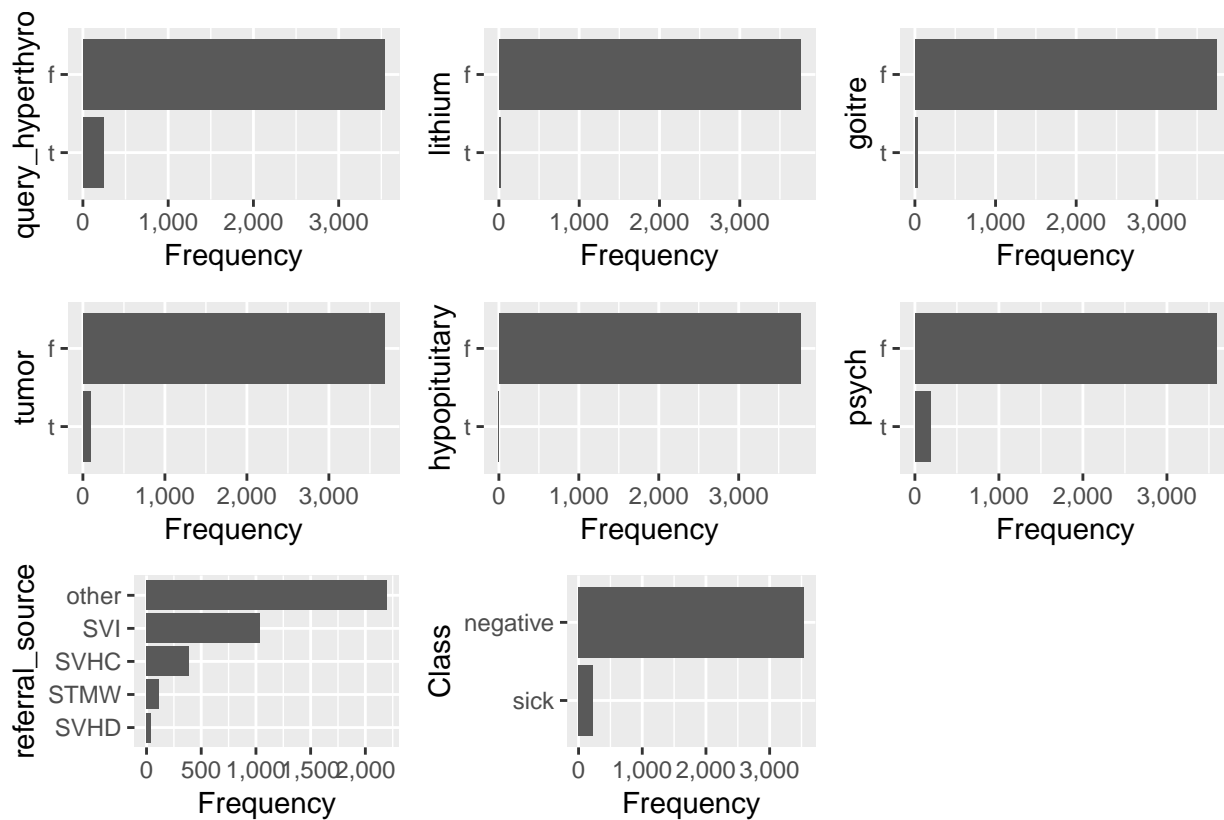
## 2.2 Distributions of Categorical Variables

In this section, I will deliberate on some of the variables in our dataset.

### 2.2.1 Exploration

Distributions of random variables are as follow:





## Page 2

There is a number of variables, that are nearly constant, namely:

- query\_on\_thyroxine
- on\_antithyroid\_medication
- pregnant
- thyroid\_surgery
- I131\_treatment
- lithium
- goitre
- hypopituitary

Do these variables have an impact on “Class”?

$\chi^2$  test results:

Names	p.values
query_on_thyroxine	0.7681159
on_antithyroid_medication	0.1099450
pregnant	0.0779610
thyroid_surgery	0.0744628
I131_treatment	0.1854073
lithium	1.0000000
goitre	1.0000000
hypopituitary	0.0694653

As we can see, some of those variables may in fact be connected with “Class”. I will exclude the “hypopituitary” variable due to a very little information it provides (distribution of “hypopituitary” is 1:3771). Even if this variable is important, we have no statistical certainty to say so, given that only one observation is positive.

### 2.2.2 Solution

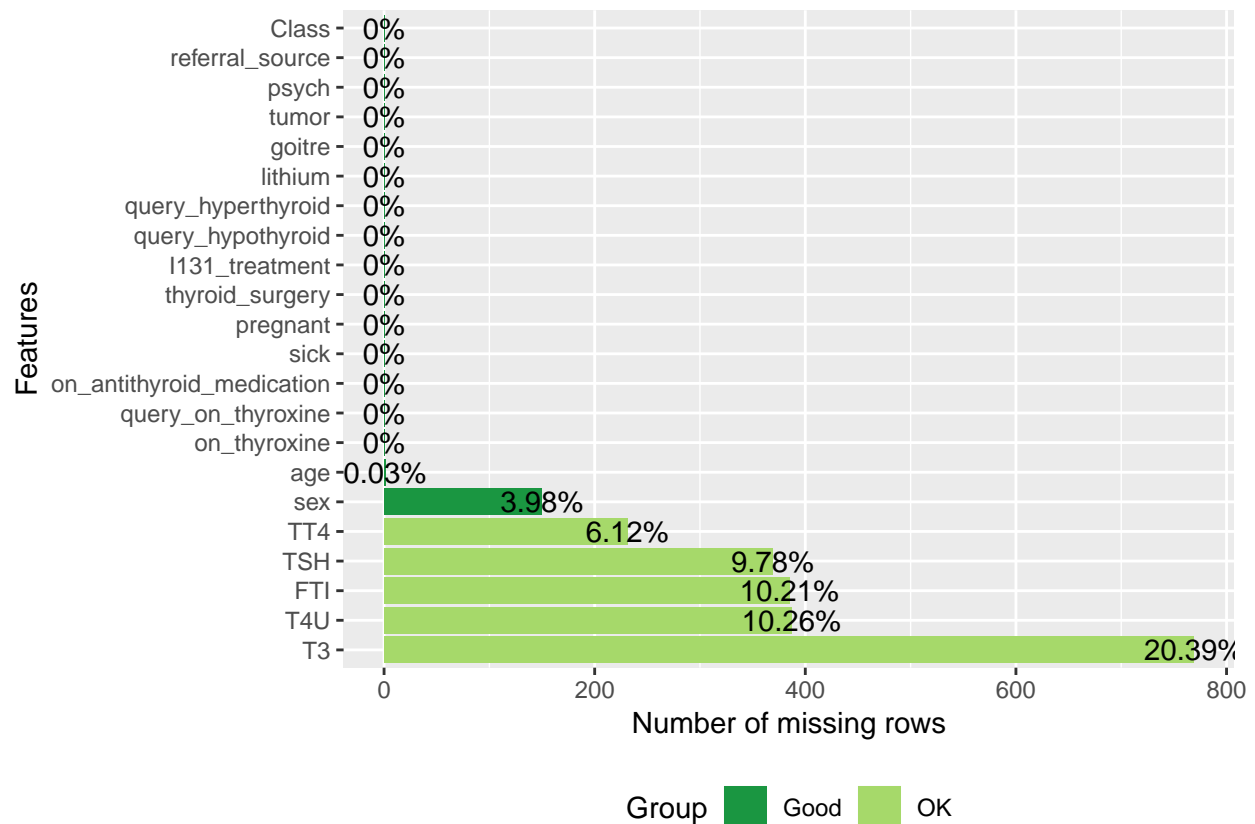
```
sick_tidy <- sick_tidy %>%  
  dplyr::select(-hypopituitary)
```

## 2.3 Missing values reduction

In this section, I will discuss the missing values in the dataset.

### 2.3.1 Exploration

There are a lot of missing values in the dataset.



However, there is no problem with removing observations with missing values, for the dataset size is very large. Given that we focus on interpretable models, which are generally not complex, there’s no need to impute missing values. Such techniques may cause bias in parameter estimation, i.e. lead to the assignment of an inaccurate level of importance to some features.

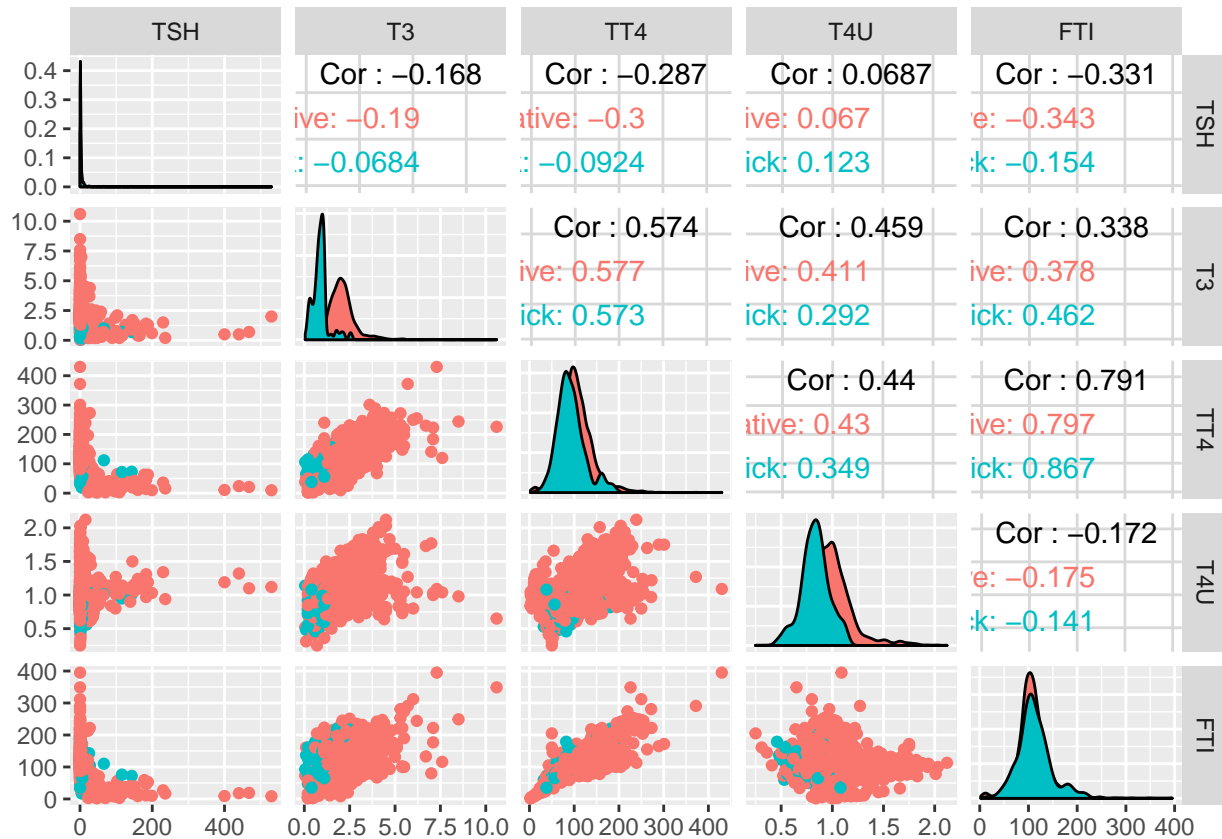
## 2.4 Solution

```
sick_tidy <- sick_tidy %>% na.exclude()
```

## 2.5 Skewness reduction

Some of the numeric variables are skewed. I will try to fix that transforming those variables.

### 2.5.1 Exploration



Skewness of these variables:

column	skewness	skewness_log	skewness_sqrt
TSH	13.231931	-0.4938557	5.4889020
T3	1.768918	-1.7744449	0.0254847
TT4	1.186118	-2.8400853	-0.2653533
T4U	1.234600	0.0157377	0.6624722
FTI	1.102262	-3.5446807	-0.5991264

As we can see, we can fix the skewness pretty easily, by taking logarithm or square root of these variables.

### 2.5.2 Solution

```
sick_t <- sick_tidy %>%
  mutate(log_TSH = log(TSH),
         sqrt_T3 = sqrt(T3),
         sqrt_TT4 = sqrt(TT4),
         log_T4U = log(T4U),
         sqrt_FTI = sqrt(FTI)) %>%
  dplyr::select(-TSH, -T3, -TT4, -T4U, -FTI)
after <- ggpairs(sick_t[, 18:22], aes(colour=sick_t$Class))
```

After the transformation:



### 3 Prediction models

In this section I will compare five different interpretable models:

- naive biases
- logistic regression
- basic tree
- knn

and the winning model.

#### 3.1 Naive Bayes

Let's perform CV resample with naive Bayes.

iter	auc	auprc
1	0.9205259	0.5531770
2	0.9458474	0.5759507
3	0.8823529	0.5464028
4	0.9328093	0.5010373
5	0.9461538	0.6309614

As we can see, naive bayes model performs quite well on auc, but auprc shows a room for improvement.

### 3.2 Logistic Regression

Results of logistic regression:

iter	auc	auprc
1	0.9674427	0.6588224
2	0.9204713	0.6459353
3	0.9294430	0.7720567
4	0.9233440	0.6791901
5	0.8759783	0.3883609

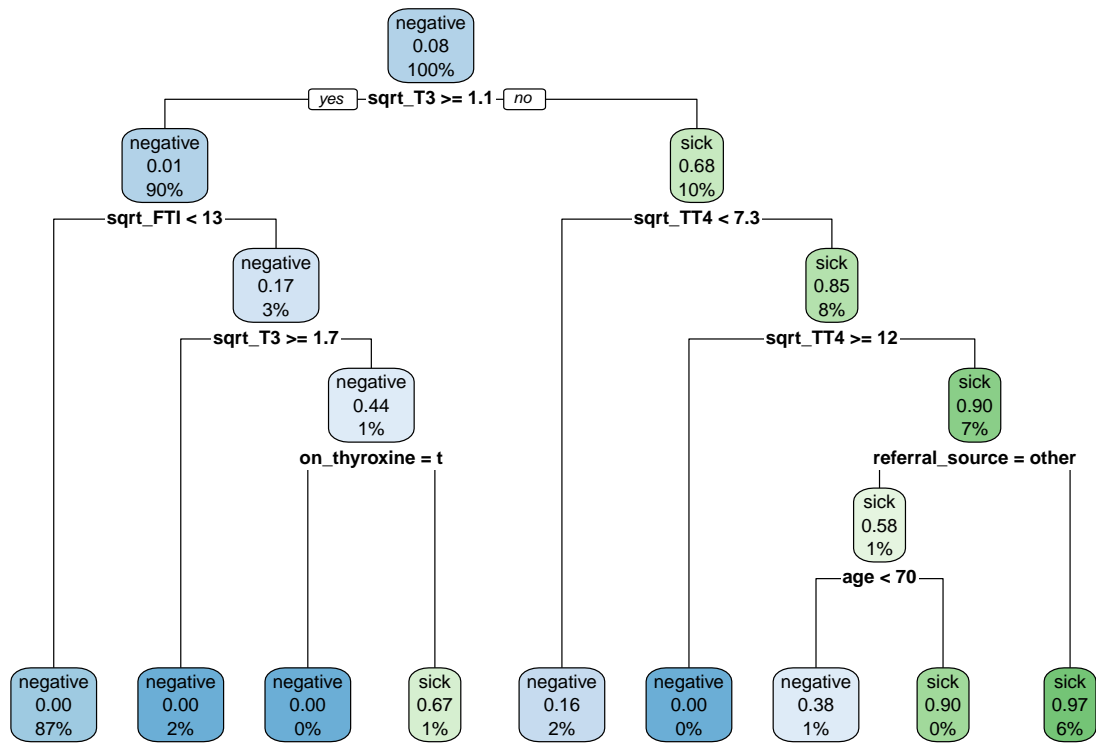
	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	4.7290675	1.1088622	4.2647928	0.0000200
age	0.0030426	0.0062927	0.4835039	0.6287380
sexM	0.2147685	0.2506418	0.8568744	0.3915143
on_thyroxinet	-1.0271864	0.5456642	-1.8824515	0.0597747
query_on_thyroxinet	-0.3059204	1.0116070	-0.3024103	0.7623393
on_antithyroid_medicationt	-12.8808241	1691.6382951	-0.0076144	0.9939246
sickt	0.9931992	0.4137108	2.4007092	0.0163633
pregnantt	-9.4753250	1243.1154542	-0.0076222	0.9939184
thyroid_surgeryt	-15.7367500	1673.1140581	-0.0094057	0.9924955
I131_treatmentt	1.1833063	1.2224896	0.9679480	0.3330703
query_hypothyroidt	1.2000409	0.4151496	2.8906227	0.0038448
query_hyperthyroidt	0.6085589	0.5827506	1.0442872	0.2963526
lithiumt	1.6547161	1.5527483	1.0656692	0.2865732
goitret	0.1079875	1.6831403	0.0641583	0.9488442
tumort	0.9673015	1.0293900	0.9396842	0.3473796
psycht	-0.3337404	0.7755397	-0.4303331	0.6669533
referral_sourceother	-1.0872804	0.5702644	-1.9066251	0.0565692
referral_sourceSVI	0.2897272	0.5085669	0.5696935	0.5688856
referral_sourceSTMW	-12.9323781	1029.2876974	-0.0125644	0.9899753
referral_sourceSVHD	-0.0143378	1.0104054	-0.0141902	0.9886783
log_TSH	-0.2351134	0.0846293	-2.7781554	0.0054668
sqrt_T3	-9.8717508	0.7497800	-13.1661965	0.0000000
sqrt_TT4	-1.1310958	0.7421060	-1.5241701	0.1274662
log_T4U	6.3909899	3.5394193	1.8056606	0.0709714
sqrt_FTI	1.5410437	0.7186901	2.1442395	0.0320137

Logistic regression is already much better in both measures. What if this problem is nonlinear?



### 3.3 Tree

The most basic nonlinear model is regression tree.



Tree results:

iter	auc	auprc
1	0.9980308	0.9279514
2	0.9626802	0.8441369
3	0.9168363	0.8632732
4	0.9285190	0.7507804
5	0.9558445	0.6912679

This basic tree achieves astonishing 20% improvement over naive bayes classifier in auprc.

### 3.4 The KNN

K-nearest neighbors is one of the oldest classifiers. Providing enough data, might make it very robust. Just like other “shallow”, nonlinear models, like svm with gaussian kernel for instance, it has a scalability problem. However, unlike SVM, it’s highly interpretable, despite that it doesn’t generalise knowledge.

KNN Performance:

iter	auc	auprc
1	0.9738698	0.8589582
2	0.9077860	0.6858856
3	0.9066434	0.7272863
4	0.9425751	0.8108829
5	0.9021739	0.7512860

Performance is good. What if knn was trained only on variables significant in logistic regression and numeric data?

iter	auc	auprc
1	0.9929554	0.9113244
2	0.9666243	0.6985563
3	0.9742813	0.8376024
4	0.9942935	0.8633175
5	0.9918079	0.8905217

Results are a bit better. The KNN model achieves about 84% in auprc benchmark, beating all other models. Note that dataset size is large enough, for it work good - there are about 2600 observations.

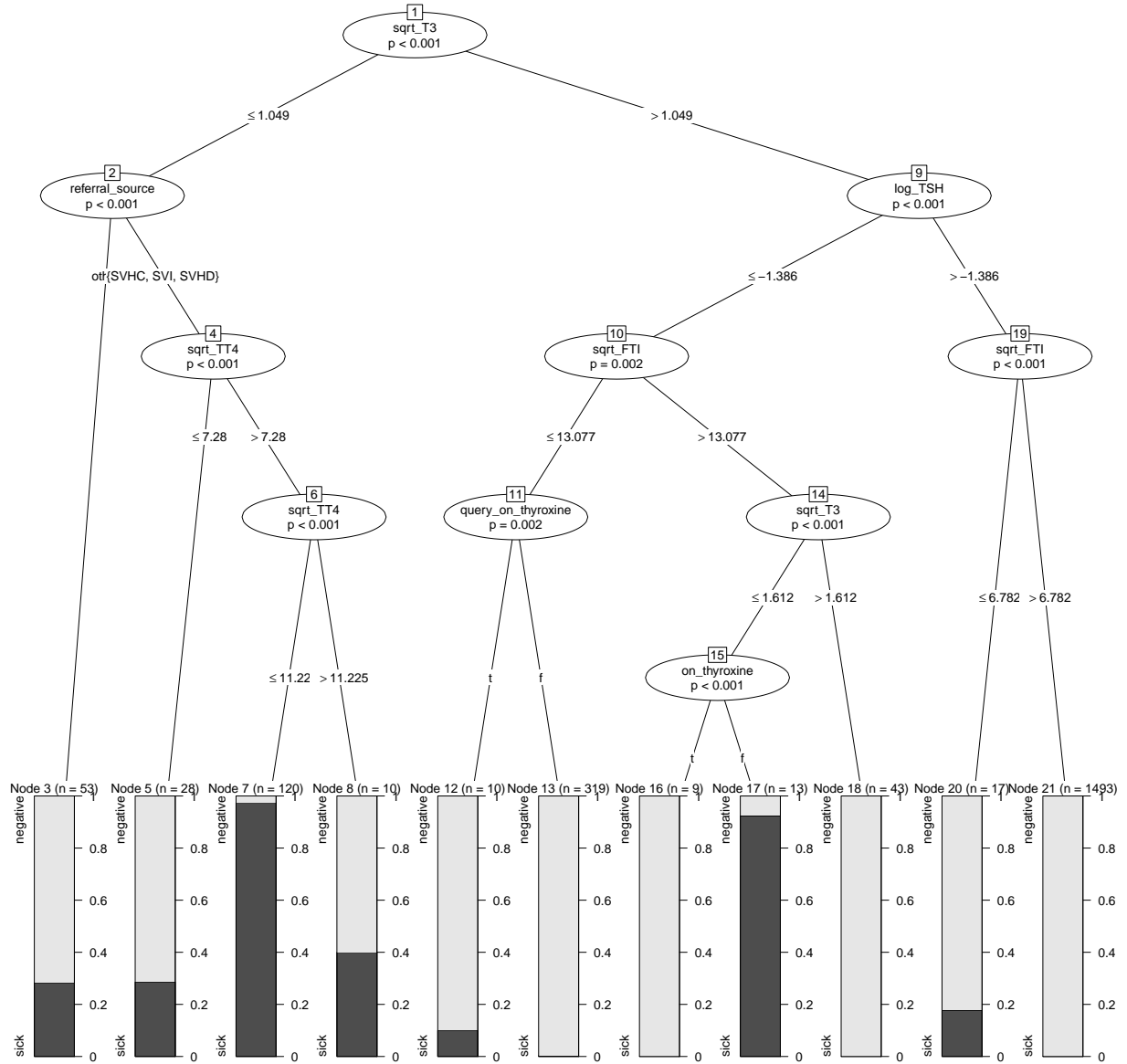
Can this result be improved? KNN relies on data, and we deleted about a third of the dataset. Let’s interpolate missing values and check, if it improves performance. It’s worth noticing, that test set should only include observations without missing values, since they are more trustworthy.

KNN Performance after imputation:

iter	auc.auc	auprc
1	0.9346216	0.7576713
2	0.8990817	0.7089261
3	0.9305470	0.7542170
4	0.9351010	0.7701688
5	0.9338911	0.7387038

As we can see, results are worst. The idea of interpolating data turned out to be not effective.

### 3.5 C-Tree



C-Tree results:

iter	auc	auprc
1	0.9856435	0.7028711
2	0.9541640	0.7133251
3	0.9116920	0.7847331
4	0.9719577	0.8848178
5	0.9941077	0.8069158

Results are a lot better. The tree is not overly complicated. It mostly uses the T3 variable and TSH.

## 4 Conclusions and the last benchmark

The C-tree model turned out to be the best. Perhaps the KNN could be improved by learning metric, but this goes beyond the scope of this paper.

Here's comparison of different models on a test dataset:

classif.naiveBayes	0.517524401027005
classif.binomial	0.716829507076891
classif.rpart	0.930269804414643
classif.kknn	0.75201453305437
classif.ctree	0.879151506258267
classif.kknn subset	0.83029979007248