

Machine Learning – Final Report

Volker Strobel (4524187)

August 9, 2016

Abstract

This report presents the techniques and results of a classification problem involving missing data as well as a mixture of categorical and continuous data. The problem of missing data points is addressed by Multiple Imputation Chained Equations (MICE). The classification is conducted using a support vector machine (SVM). The method was evaluated locally using cross-validation and remotely on a hold-out test-set using the Kaggle platform. The best Kaggle submission resulted in an accuracy of 85.494 %, which corresponds to the 3rd position out of 23 participants in the competition.

1 The Challenge

The Kaggle competition “Final Assignment IN4320”¹ asks participants to predict whether a person earns more than EUR 40k a year from $D = 13$ predictor variables. The provided dataset was created by a census bureau. The task can be divided into the following three challenges:

- Imputation of missing data points
- Handling a mixture of categorical and continuous independent variables
- Classification of the output variable

In this report, the terms *wealthy* or “label 1” describe a sample with income EUR > 40k and *non-wealthy* or “label 0” are used otherwise.

In Section 2, I analyze and visualize the structure of the data to set the stage for the classification pipeline. In Section 3, the used methods for imputation, classification, and transductive learning are detailed. Section 4 presents and discusses the results obtained during cross-validation and on Kaggle.

2 Analysis & Ideas

To motivate the used classifier, design, and technique choices, I start with an analysis of the given dataset and introduce possible concepts and ideas. Additionally, important patterns dependence structures are visualized.

In total, 13 features are given, of which 8 are categorical and 5 continuous (Table 1). The training set consists of 10500 samples and the test set has 38342 samples, resulting in a ratio of approx. 1 : 3.7. The test set is thus considerably larger than the training set. Since transductive learning is permitted, the information in the test set—information about the distribution of the variables—could help to increase the classification performance.

10500 samples constitute the training set of which 2500 (23 %) were labeled as wealthy and 8000 (77 %) as non-wealthy. If we assume that the training and test set were randomly sampled from the entire dataset, therefore, the predictions on the test set should reflect this ratio. Additionally always predicting non-wealthy should give a 0-1 loss of 77 % under this assumption, which can be used as a baseline for

¹<https://inclass.kaggle.com/c/final-assignment-in43202>

Categorical	Continuous
work class	age
education	number of years of education
marital status	income from investment sources
occupation	losses from investment sources
relationship	working hours per week
race	
sex	
native country	

Table 1: Overview of the used features

the classifier performance. The assumption could be tested by probing the test set with a “0-only” submission. This ratio might be useful for setting the class-weights of a classifier or for determining the threshold θ in majority voting: the final prediction after voting is “label 1” if at least θ preliminary predictions are “label 1”. The parameter θ can be modified to increase or decrease the amount of “label 1” predictions and to make the final predictions more or less conservative (see Section 3.4).

In total, approx. 20% of the data values are missing, with approx. 23% for the variables workclass and occupation, 20% for country of origin, and 19% for the remaining variables. There is no apparent difference in the percentage and pattern of missing data between the training and the test set. Only 6% of samples in the training set and in the test set are complete cases (i.e., have no missing values), making handling missing data a crucial step. I assume that the higher values for the variables workclass, occupation, and country stem from missing values in the original dataset and the missing values for the remaining variables were just introduced.

Given the description of the dataset and the total number of samples (48842), the dataset is potentially the UCI Adult Data set [5]. This would underline the pattern of missing data (higher values for workclass, occupation, and country of origin). A possible method would be to use the labeled samples of the UCI dataset and match them with the given test set, which should give an accuracy of 1.0. However, this method is not as straight-forward as it may seem. The UCI dataset is split in a different manner into training and test set (amount of training samples: 32561, test samples: 16281). Additionally, the categorical variables in the UCI dataset are encoded as strings, while in the given dataset, they are encoded as integers. Therefore, one would have to find a mapping from strings to integers. The large amount of missing data might impede this endeavor. Since using the UCI Adult dataset might defeat the goal of this assignment, I did not take any steps in this direction. A comparison on <http://www.cs.toronto.edu/~delve/data/adult/adultDetail.html> shows that the best performing classifier is a Forward Sequential Selection (FSS) naive Bayes model with an accuracy of 85.95%. The classifiers were trained after removing unknown values (7% of values had missing values; training set size: 30162, test set size: 15060). This accuracy can be used as an indicator of a good performance on the given dataset. However, there are two differences to the given dataset: (i) the original UCI dataset had a lower amount of missing data, and (ii) the UCI dataset had a higher number of training samples.

3 Methods

3.1 Overview of the Classification Pipeline

In Figure 1, the complete classification pipeline is visualized. In the following subsections, each step will be explained in more detail. The process starts with combining the training and the test set to yield a large dataset with missing values. Using MICE, the missing values are imputed and five training and test sets are obtained. Five classifiers are trained using the different training sets; for each training set all test sets are used for obtaining the “preliminary hypotheses”. Using majority voting with a modifiable threshold θ , the single hypotheses are pooled to yield the final prediction vector.

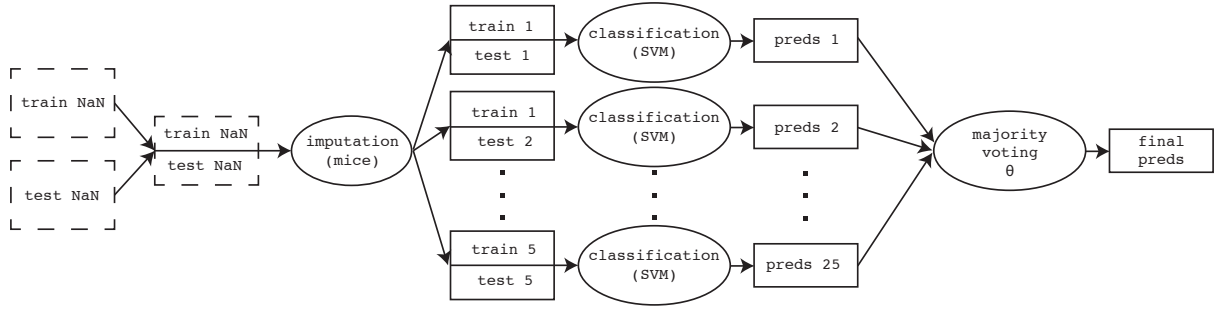


Figure 1: The figure illustrates the classification pipeline. Matrices are displayed as boxes and methods as ellipses; dashed lines indicate missing values.

3.2 Dummy Variables

Seven of the twelve measured variables are qualitative (workclass, education, marital status, occupation, relationship, sex, native country)—they are measured only at the nominal level. Since measurements on the nominal level do not allow for a particular ordering, a proxy method has to be used to make them suitable for statistical models. To this end, I defined dummy variables that take the value 0 or 1 to indicate if a certain category is present. Defining n dummy variables for the n different possible values of a categorical feature allows for capturing the full information in the original unmodified dataset and for using qualitative data in a straight-forward manner. Transforming the dataset greatly increased its size from $D = 13$ to $\tilde{D} = 104$ predictor variables.

3.3 Missing Data Values

Missing data values are a common problem in statistics. The failure of sensors, or the concealment of data impede machine learning accuracy. Thus, methods have been put forth for the imputation of missing data points, such as case deletion or single mean imputation [9]. However, such simple methods often discard useful information and are likely to give a rather poor result if a large part of the data is missing. The general goal is to find unbiased estimates of the missing values. The bias depends on the type of “missingness”:

- Missing Completely At Random (MCAR)
In the case of MCAR missingness, there is no dependence structure between the missing data and any other values. Listwise deletion leads to unbiased analyzes.
- Missing At Random (MAR)
In the case of MAR, the reason for missingness depends on the value of other non-missing predictor variables. Unbiased estimates can be obtained using the multiple imputation method [4], which creates multiple versions of filled-in datasets.
- Missing Not At Random (MNAR)
In the case of MNAR, the reason for missingness depends on non-available information. Such a case might for example occur, if people with high income refuse to reveal their income. Imputation of MNAR requires explicit modeling of the reason for missingness to yield unbiased estimates.

For getting unbiased estimates with listwise deletion, the data must be MCAR. Little’s MCAR test [6] is a statistical test with the null hypothesis that the data are MCAR. I used it at the significance level $\alpha = 0.05$ to analyze the interaction structure of the variables. The statistic was significant for the combined dataset consisting of the training and the test set ($\chi^2 = 28018.84$, $p < 0.001$, $df = 15639$). Therefore, there is evidence that the missing values in the dataset are not MCAR; the null hypothesis is rejected. It is not possible to distinguish between MAR and MNAR data using observed information only [10]. Consequently, I assumed that the data are MAR, since imputation of MNAR data would need further information about the reason for missingness, which was not given. Therefore, data points

were imputed using the multiple imputation method by a technique called Multivariate Imputation by Chained Equations (MICE) [2] using the R package `mice` (Version 2.25). Due to the multiple imputations, the algorithm incorporates the statistical uncertainty in the imputed values. MICE is a flexible approach that can also handle continuous and categorical variables. The general idea is to statistically model the conditional probability of each missing variable given the remaining variables. In detail, the MICE algorithm works as follows [2, 1]:

1. In the beginning, a simple imputation is carried out. To this end, all missing data points are replaced by random sampling with replacement from the observed datapoints.
2. One variable is selected at random, for example *occupation*. An intermediate regression model is built, using the remaining variables as predictors and *occupation* as target value. The chosen model is dependent on the target value. I used a logistic regression categorical data and linear regression without regularization for numerical data. The missing values in the variable *occupation* are replaced by the predictions of the model.
3. The previous step is executed for all variables with missing data. For each variable, the model is trained using both the already imputed values and the existing values. Once all variables have been predicted, one cycle is complete.
4. Several cycles are performed to stabilize the imputation results. I used $c = 5$ cycles.
5. The entire procedure is executed multiple times to yield several imputed datasets. Due to the random factors, the imputed values will be different, while the non-missing data entries will be the same in all datasets. I used $m = 5$ imputations.

The choices $c = 5$ and $m = 5$ were made due to practical limitations, considering time limitations and avoiding overflow errors. A value of m between 5 and 10 is a common value for the used method [7] and a value of $c = 5$ is considered “adequate” [11]. However, recent research shows that a higher value of m might be beneficial [12], which could be tested, if a more powerful machine would be at hand.

Importantly, the MICE algorithm was used on the total dataset, consisting of the training and the test set. This should increase the quality of the imputed values since more training examples can be used for building the intermediate models (Step 2 in the algorithm outline).

In the following, the two statistical models—linear regression and logistic regression—used for the imputations are shortly described:

Linear regression: For the linear regression (continuous variables), a design matrix \mathbf{X} is built from the predictor variables x_{ij} (j th feature of the i th sample):

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} & \dots & x_{1d} \\ 1 & x_{21} & x_{22} & x_{23} & \dots & x_{2d} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & x_{n3} & \dots & x_{nd} \end{bmatrix} \quad (1)$$

The first column in the matrix enables fitting the intercept in the regression. The solution for the weights $\hat{\beta}$ —which minimizes the squared-error loss function—is obtained using the normal equation $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$, where \mathbf{y} is the vector consisting of the target values. Predictions for one variable y_i are made with $\hat{y}_i = X_i\hat{\beta}$, where X_i stands for the features of the i th sample.

Logistic regression: Logistic regression (binary variables) uses a logit link function to link the linear predictor $X_i\beta$ to the probability of “label 1”:

$$\Pr(y_i = 1|X_i) = \text{logit}^{-1}(X_i\beta) = \frac{1}{1 + e^{-X_i\beta}} \quad (2)$$

The training of the logistic regression can be done via stochastic gradient descent (SGD). For this, the following update function is used—with a sufficiently small learning rate α and $p = \text{logit}^{-1}(X_i\beta_t)$:

$$\beta_{t+1} := \beta_t + \alpha(y - p)X_i \quad (3)$$

The difference $y - p$ is the error made at time step t by predicting y from X_i using the weight vector β_t . For predicting the target value (0 or 1) of the sample X_i , the hypothesis function $h(X_i)$ is used, which transforms the class probabilities into target values:

$$h(X_i) = \begin{cases} 1 & \text{if } \Pr(y_i = 1|X_i) > 0.5 \\ 0 & \text{otherwise} \end{cases} \quad (4)$$

3.4 Pooling

I trained models and made predictions on all possible combinations of imputed training and test sets. This results in 25 “preliminary hypotheses” \tilde{h}_{ij} for the i th sample: $\tilde{\mathbf{h}}_i = (\tilde{h}_{i,j})_{j \in \{1,2,\dots,25\}} = [\tilde{h}_{i,1}, \tilde{h}_{i,2}, \dots, \tilde{h}_{i,25}]$. These predictions have to be aggregated to obtain one final submission. For the aggregation, I used the following majority voting function on the preliminary hypotheses:

$$\mathbf{h}_i = \text{vote}_\theta(\tilde{\mathbf{h}}_i) = \begin{cases} 1 & \text{if } \sum_{j=1}^{25} \tilde{h}_{ij} \geq \theta \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

The function predicts the class *wealthy*, if at least $\theta \in \{1, 2, \dots, 26\}$ preliminary hypotheses are 1. Introducing θ in the function was motivated by the local cross-validation (next section): the default behavior of classifier was conservative and only 54.0 % of actual “wealthy” samples were correctly classified. The modifiable threshold θ serves as modulator for the “conservatism” of the classifier: the higher θ is, the less likely it is that the final prediction *wealthy* will be made. The threshold $\theta = 13$ would represent standard majority function with a 50 % majority. Figure 2 shows the trade-off between specificity and sensitivity. It can be seen that the accuracy is not greatly affected by the conservatism θ of the approach: a higher sensitivity is directly traded-off a smaller specificity. The maximum accuracy was achieved at $\theta = 9$.

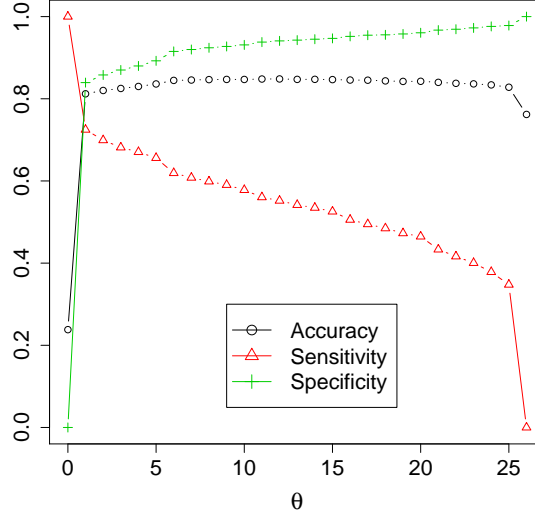


Figure 2: Accuracy, sensitivity, and specificity in dependence of the threshold θ . The accuracy stays roughly the same for $6 \leq \theta \leq 25$. In contrast, the sensitivity falls and the specificity rises for increasing values of θ .

3.5 Cross-Validation

To evaluate submissions and determine the ranks of participants, the Kaggle challenge used the 0-1 loss. While the score on the public leaderboard is the best validation of the methods employed, only one

submission could be submitted per participant and day. In order to evaluate the used methods more frequently, I implemented a local 5-fold cross-validation.

TODO: The cross-validation was performed per dataset of the imputation result, that is, no cross-dataset cross-validation was performed. Therefore, in the case of five imputations, 25 cross-validation results were obtained.

The cross-validation achieved similar results to the public Kaggle evaluation (Table 2).

Method	Local Score	Global Score
AdaBoost with 1 imputation	0.83886	0.84205
Random Forest with 5 imputations	0.82726	0.85056
SVM with 5 imputations	0.84413	0.85494
AdaBoost with 5 imputations		

Table 2: Comparing local and public scores

Each fold in the cross-validation contained 2100 test samples and 8400 training samples. Since the classes *wealthy* and *non-wealthy* are not equally distributed, a confusion matrix can give a more detailed picture of the classifier performance. The following confusion matrix was obtained:

		Predicted class		Total
		wealthy	non-wealthy	
Actual class	wealthy	$TP = 1351$	$FN = 1149$	$TP + FN = 2500$
	non-wealthy	$FP = 546$	$TN = 7454$	$FP + TN = 8000$
Total		$TP + FP = 1897$	$FN + TN = 8603$	$N = 10500$

From the confusion matrix, sensitivity and specificity can be calculated:

$$\text{Sensitivity} = \frac{TP}{TP + FN} = 54.0\% \quad (6)$$

$$\text{Specificity} = \frac{TN}{TN + FP} = 93.2\% \quad (7)$$

$$(8)$$

These statistics show that almost all non-wealthy samples are correctly classified, while only slightly more than half of the wealthy samples are correctly classified. The confusion matrix implies that the classifier is “conservative”: it avoids classifying samples as wealthy. The high overall accuracy is mainly caused by correctly classifying non-wealthy samples. I tried to make the classifier less conservative by predicting probabilities of the target values and modifying the threshold θ , in the ranges 0.3 – 0.5. This improved the sensitivity, but, in turn, led to a reduced specificity, thus not improving the classification accuracy in the local cross-validation.

3.6 The Classifier

As can be seen in Section 3.5, different classifiers were tested for the given problem. The best performing classifier was a support vector machine (SVM). For running the algorithm `Python 2.7.11` was used with the package `scikit-learn` in Version 0.17.0.

The choice of an SVM was motivated by the following points:

1. SVMs work “out-of-the-box” on a large number of different problems [8].
2. SVMs work in high-dimensional spaces ($\tilde{D} = 104$)
3. SVMs aim to minimize the generalization loss, instead of the empirical loss, by finding a decision boundary that maximizes the margin. Since the given training set contains only 10,500 samples, while the test set contains 38,342 samples, generalization is especially important.

4. SVMs allow for *transductive learning* using transductive SVMs. Therefore, the information in the unlabeled test data can be incorporated in the training process, possibly improving the accuracy.
5. SVMs showed the best performance in the local cross-validation.

If D is the number of features, Support vector machines build a $(D-1)$ -dimensional hyperplane that tries to maximize the margin between the two classes. The support vectors are the samples on the margin for the positive class and the negative class. The maximum-margin hyperplane is the hyperplane “right in the middle” of the support vectors.

The learning a linear classifier

$$f(\mathbf{x}) = w^T \mathbf{x} + b \quad (9)$$

The classification function is $h(\mathbf{x}) = \text{sign}(w^T \mathbf{x} + b)$, with \mathbf{x} being the feature vector, w the coefficient vector and b the intercept term. The goal is to find w, b such that the resulting decision boundary maximizes the margin.

A soft-margin is used with the hinge-loss function:

$$\max(0, 1 - y_i(\vec{w} \cdot \vec{x}_i + b)) \quad (10)$$

The primal form of the The solution to find w and b , starts with the *dual representation*, which can be derived from the primal form [8]:

$$\arg \max_{\alpha} \sum_j \alpha_j - \frac{1}{2} \sum_{j,k} \alpha_j \alpha_k y_j y_k (\mathbf{x}_j \cdot \mathbf{x}_k) \quad (11)$$

subject to the constraints: $\alpha_j \geq 0$ and $\sum_j \alpha_j y_j = 0$. The variable $y_i \in \{-1, 1\}$ represents the class of x_i and α is the weights vector that assigns a weight α_j to each data point. The *support vectors* are the samples for which the weights α_j are non-zero. This optimization problem can be solved using quadratic programming, for which several libraries exist. I used the `libsvm` library [3].

3.6.1 The kernel

To extend SVMs to non-linear classification, the *kernel trick* can be used, with which data can be embedded in a higher-dimensional space. The kernel $K(\mathbf{x}_j, \mathbf{x}_k)$ is then applied to the dot products of feature vectors in Equation 11. Therefore, the equation becomes:

$$\arg \max_{\alpha} \sum_j \alpha_j - \frac{1}{2} \sum_{j,k} \alpha_j \alpha_k y_j y_k K(\mathbf{x}_j, \mathbf{x}_k) \quad (12)$$

Practically, this means that the problem of finding a linear separator is shifted to a higher dimensional feature space. Mapping back this linear separator to the original feature space results in a non-linear decision boundary.

For the given problem, a Gaussian radial basis function kernel was used:

$$K(\mathbf{x}_j, \mathbf{x}_k) = \exp(-\gamma \|\mathbf{x}_j - \mathbf{x}_k\|^2) \quad (13)$$

The parameter γ was set using $1/\text{number of features}$, that is: $\gamma = 1/104$.

3.7 Transductive Learning

Since the test data is known at training time, this information can be used for using semi-supervised methods than can better capture the underlying distribution. In a first step, “naive” self-learning was tried, using an iterative approach: In the first iteration, the unlabeled data was labeled by predicting their labels using the training dataset. In the next iteration.

The approach did not significantly improve the classification performance.

The second approach used an implementation based on `scikit-learn` and the Contrastive Pessimistic Likelihood Estimation (CPLE) (<https://github.com/tmadl/semisup-learn>).

4 Results & Discussion

In total, I made six submissions during the competition. The best submission resulted in an accuracy of 85.494 %, which corresponds to the 3rd position out of 25 participants at the time of finishing this report.

In summary, the main challenge lay in handling the missing data values. Afterward, standard classification techniques could be used. The use of different classifiers had no substantial influence on the achieved accuracy: using a support vector machine, random forest regression, and AdaBoost classifier resulted in very similar results, as can be seen in Table 2.

The use of semi-supervised methods did not have a substantial effect on the classification performance; this is in line with findings other studies, showing that semi-supervised methods do not necessarily lead to improved performances and can even lead to a lower accuracy [13]. It might be that semi-supervised methods would have had a greater effect, if the number of training samples would have been substantially smaller.

Since the ground truth labels of the test set were not revealed, error statistics could only be computed based on the 0-1 loss on the public leaderboard and on the local cross-validation. The restriction of one submission per day increased the difficulty of the validation. My public score of 85.494 %—which is only 0.00364 below the best performing score in the competition (85.854 %) and 0.00456 below the best performing method of the Delve project ²—is an indicator that little improvement will be possible beyond my approach. Compared to the UCI dataset, the given dataset had a smaller amount of labeled data and more missing values, which made the Kaggle classification task more challenging. Using standard state-of-the-art classifiers, the maximum possible performance of the given dataset will be possibly capped clearly below 100 % due to mistakes during data acquisition, such as deliberate misinformation or transcription errors. Moreover, the used variables might only have limited explanatory power: additional variables, such as political view, morale, or number of children may be needed to increase the amount of explainable variation.

Due to the concealment of categorical variables, I did not consider manual feature engineering such as grouping certain levels of the categorical (dummy) variables.

Many steps in the classification pipeline were time-consuming, taking several hours to complete. In the future, more processing power could help to increase the classification performance. For example, using a larger amount of imputed datasets could capture the uncertainty in the imputed values to a greater degree and another method for transductive learning could help to increase the classification performance.

My code for this competition can be found at:
<https://github.com/Pold87/ml-final-ass>

Word Count: 3164

References

- [1] Melissa J. Azur et al. “Multiple imputation by chained equations: What is it and how does it work?” In: *International journal of methods in psychiatric research* 20.1 (2011), pp. 40–49.
- [2] Stef Buuren and Karin Groothuis-Oudshoorn. “mice: Multivariate imputation by chained equations in R”. In: *Journal of statistical software* 45.3 (2011).

²<http://www.cs.toronto.edu/~delve/data/adult/adultDetail.html>

- [3] Chih-Chung Chang and Chih-Jen Lin. “libsvm: A library for support vector machines”. In: *ACM Transactions on Intelligent Systems and Technology* 2 (3 2011). Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>, 27:1–27:27.
- [4] Moniek C.M. de Goeij et al. “Multiple imputation: dealing with missing data”. In: *Nephrology Dialysis Transplantation* 28.10 (2013), pp. 2415–2420.
- [5] M. Lichman. *UCI Machine Learning Repository*. 2013. URL: <http://archive.ics.uci.edu/ml>.
- [6] Roderick JA Little. “A test of missing completely at random for multivariate data with missing values”. In: *Journal of the American Statistical Association* 83.404 (1988), pp. 1198–1202.
- [7] Patrick Royston, Ian R White, et al. “Multiple imputation by chained equations (MICE): implementation in Stata”. In: *Journal of Statistical Software* 45.4 (2011), pp. 1–20.
- [8] Stuart Jonathan Russell et al. *Artificial intelligence: a modern approach*. Vol. 2. Prentice hall Upper Saddle River, 2003.
- [9] Joseph L Schafer and John W Graham. “Missing data: our view of the state of the art.” In: *Psychological methods* 7.2 (2002), p. 147.
- [10] Jonathan AC Sterne et al. “Multiple imputation for missing data in epidemiological and clinical research: potential and pitfalls”. In: *Bmj* 338 (2009), b2393.
- [11] Stef Van Buuren, Hendriek C Boshuizen, Dick L Knook, et al. “Multiple imputation of missing blood pressure covariates in survival analysis”. In: *Statistics in medicine* 18.6 (1999), pp. 681–694.
- [12] Ian R. White, Patrick Royston, and Angela M. Wood. “Multiple imputation using chained equations: issues and guidance for practice”. In: *Statistics in medicine* 30.4 (2011), pp. 377–399.
- [13] Xiaojin Zhu. *Semi-supervised learning literature survey*. 2005.