

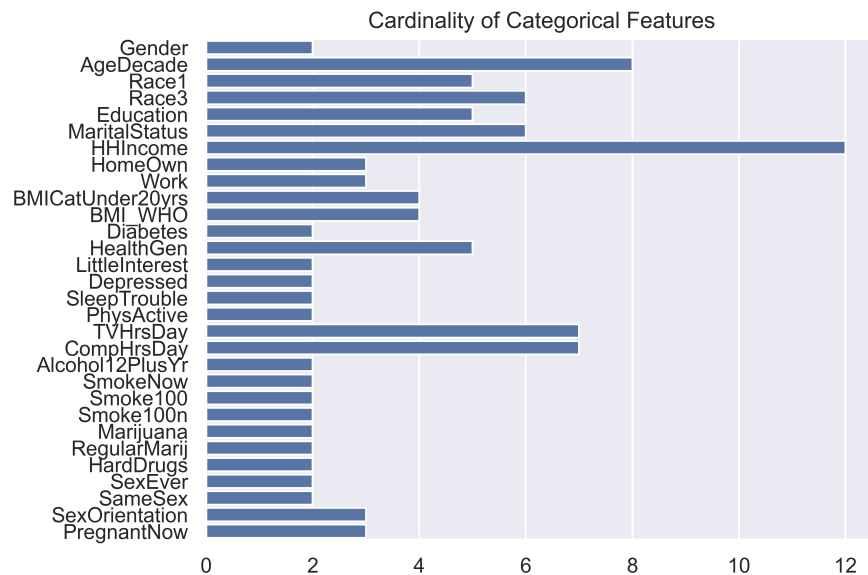
Machine Learning: Project 1

Johannes Misensky, Dejan Prvulovic & David Deket

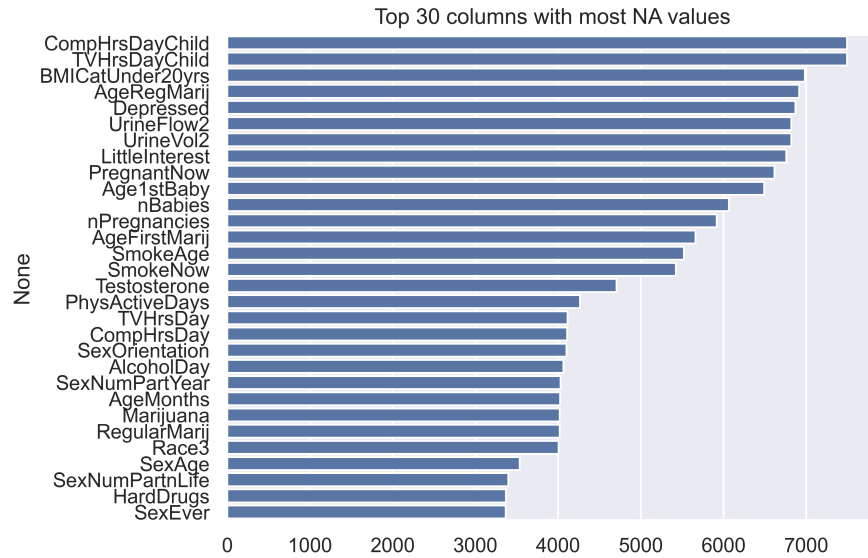
2023-12-29

Exploratory Data Analysis

Our first step was to explore the data, e.g. shape of the data, column types, missing values, etc. The dataset contains 8000 rows and 71 columns. 41 columns are numeric, 30 are categorical. Most categorical columns are binary, but some have more than two categories. The highest cardinality is 12, which means that one-hot encoding is feasible and won't lead to too many new columns.

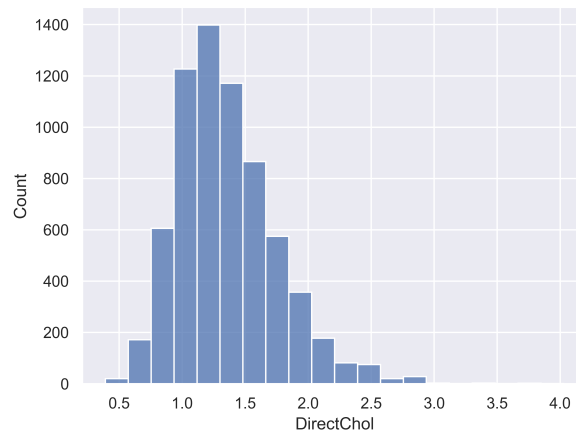


Next, we have a look at the missing values. 26 columns have at least 4000 missing values, which is 50% of the data. Usually, we would impute missing values with the median or mean, however with so many missing values, this would lead to a lot of bias. Instead, we drop all columns with more than 1500 missing values. This leaves us with 28 columns.



Regression

The target variable `DirectChol` for the regression has the following distribution:



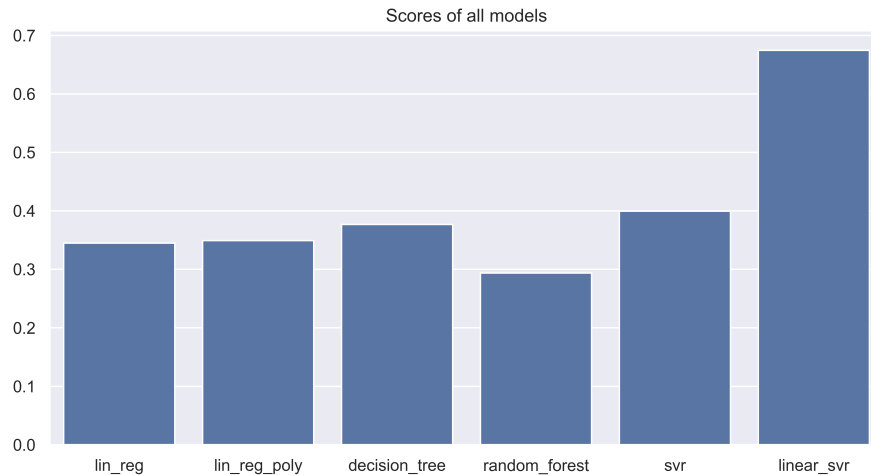
Our pre-processing pipeline for the regression task consists of the following steps:

1. Drop columns with more than 1500 missing values
2. For numerical columns, impute missing values with the median and use standard scaling
3. For categorical columns, use one-hot encoding with `handle_unknown='ignore'` to avoid errors when encountering unknown categories in the data

We use the following models for the regression task:

- Linear Regression (no and polynomial features)
- Decision Trees
- Random Forests
- Support Vector Regressions (with and without kernel trick)

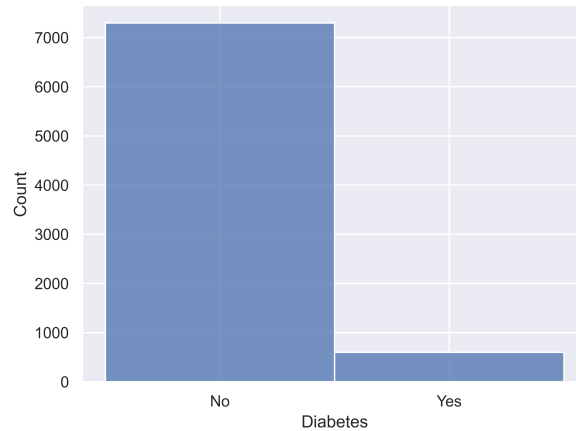
Models that have adjustable hyperparameters are trained using 5-fold cross-validation and scored on root mean squared error (RMSE).



The best model is the Random Forest with an RMSE of about 0.30. Interestingly, the linear regression models perform second best, indicating that the relationship between the features and the target variable is approximately linear and that data is more important than model choice. The worst model is the linear support vector regression, though the SVR with the kernel trick still performs worse than the linear regression models.

Classification

The target variable **Diabetes** for the classification has the following distribution:



Clearly, a class imbalance making the classification task more difficult. Therefore, we upsample the minority class by randomly duplicating rows until both classes have the same number of rows. There are more advanced techniques to deal with class imbalance (e.g. SMOTE), but we don't use them here in order to save another dependency.

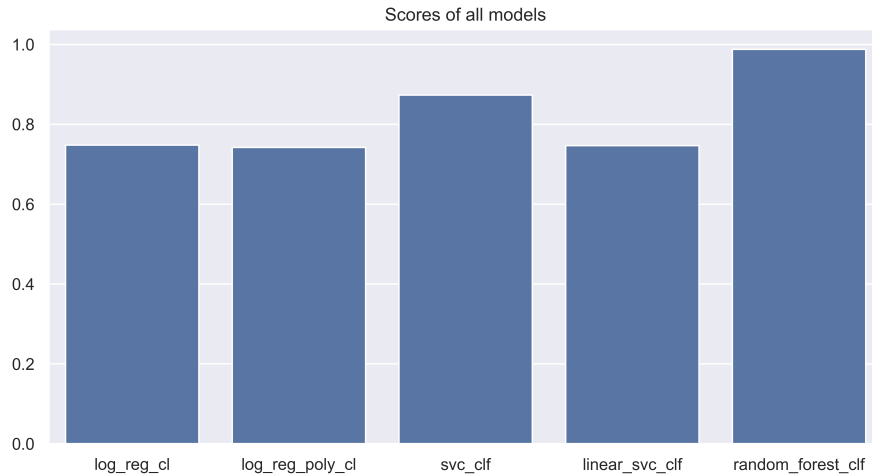
Our pre-processing pipeline for the classification task consists of the following steps (slightly different from the regression task to try out different techniques):

1. Drop columns with more than 1500 missing values
2. For numerical columns, impute missing values with the median, apply min-max scaling and PCA
3. For categorical columns, use one-hot encoding with `handle_unknown='infrequent_if_exist'` and set `sparse_output=True` to avoid memory errors

We use the following models for the classification task:

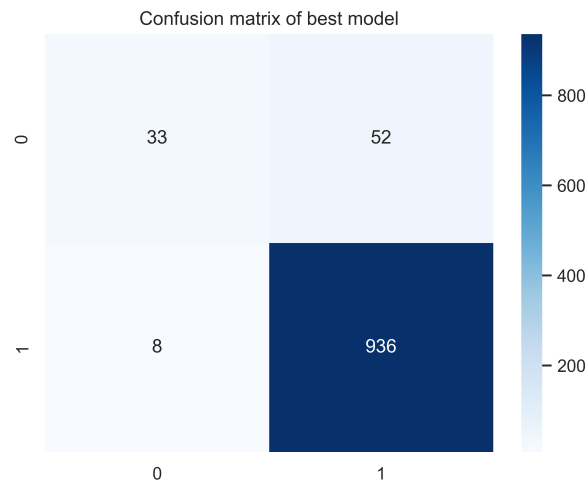
- Logistic Regression (no and polynomial features)
- Support Vector Classifiers (with and without kernel trick)
- Random Forest

Models that have adjustable hyperparameters are trained using 5-fold cross-validation and scored on the weighted F1 score.



The best model is again the Random Forest with an F1 score of close to 1. This time, the second best model is the support vector classifier, whereas the linear support vector classifier performs worst (kernel trick definitely helps). The logistic regressions with regularization perform similar to the SVCs, but the polynomial features don't help.

Lastly, we have a look at the confusion matrix of the Random Forest on the test set:



No diabetes = 1 is predicted correctly most of the time, which is easy to achieve given the class imbalance. Diabetes is predicted correctly 80% of the time (33 out of 33+8 true cases), however the model suffers from a low recall of 39%. This means we get many false positives (33 cases are truly diabetes from the predicted 33+52=85 cases), which is not ideal for a medical diagnosis.