

For data mining and machine learning tasks, data preprocessing is the most important and time-consuming process before the deployment of algorithms. It is a classical topic that how to select relevant features to fit to the algorithm. The goal for the project is to compare between the modern **feature selection techniques** in data preprocessing.

Current Progress

Currently, our task is using the diabetes data we found on the UCI repository to predict the probability of readmission by a patient to the hospital within 30 days. As far as progress towards this task, we started by cleaning the data. Initially the data had 50 variables (features). To reduce the number of variables that we will be using, we started by removing features that will have no effect on our desired outcome, including the “encounter id” and “patient number” that essentially act as indices. As we looked at each variable, we also noticed some missing data. The variables that had missing data were categorical, so we imputed the missing values with the mode of the non-missing values. This can introduce some bias into the data, but it is still better than removing the observations with missing information. Our dataset is quite large (101,766 observations), so doing a complete-case analysis was considered. However, we thought that the best model was to preserve all of the observations, while reducing the number of features as much as possible.

Despite what we had said in our data collection report, it turns out that the cleaning of the data was more involved than we initially thought it would be. Along with the removal of unnecessary variables and imputing missing data, most of the variables needed a mapping of some sort to assign each level of the categorical information to a number. Some mappings were already provided in the folder where the dataset was found. An example of a common mapping that we chose dealt with several drug variables. The data said “Up” if the dosage was increased during the patient’s visit, “Down” if it was decreased, “Steady” if the dosage was not changed, and “No” if that particular drug was not prescribed. We ended up using the mapping “No” = 0, “Down” = 1, “Steady” = 2, and “Up” = 3. This or something similar needed to be done for most of the variables. In addition, we normalized the variables so that they were all at the same scale. Some variables like “readmitted” were binary, meaning there was only 0 or 1. However, for a variable like the No-Down-Steady-Up example, the numbers 0, 1, 2 and 3 were converted to 0,

0.333, 0.667, and 1. This turned out to be a tedious process to do this for what ended up being 40 variables that we will use in our initial models (39 used to predict the “readmitted” outcome).

Once the data was cleaned, we knew that we wanted to find some model to form a baseline for the accuracy of predicting the probability of readmission by a patient. To achieve an initial accuracy percentage, we first split the data into a training set (90000 observations chosen randomly) and a testing set (the remaining observations). Because the response variable in question is binary (they were readmitted within 30 days, or they were not), we chose to initially perform a logistic regression to predict the probability of getting a “yes”. After doing the logistic regression using all of the variables in the normalized and cleaned data, the training set was able to correctly predict 88.89% of the results in the testing set.

Future Plan

Logistic regression is not necessarily our main technique that we will use. We want to experiment with different feature selection techniques to start seeing which features would be significant. We will try to improve the prediction accuracy as we discover all the possible feature selection techniques.

- **Embedded Methods**

Some techniques like Regularized Logistic Regression, Regularized Linear Regression (Lasso, Ridge, ElasticNet), and Decision Tree with max-depth have their own embedded approaches to implement feature selection while fitting the model to the data. We are not going to implement them but treat them as references.

- **Filter Methods**

Filter Methods refer to approaches which use statistical metrics to evaluate each feature/feature combinations. Common measures include the mutual information, Pearson correlation coefficient, inter/intra class distance or the scores of significance tests. These approaches not only can reflect the significance of a feature but easy to compute as well. Those approaches are supported by most of the data processing and statistical libraries and are widely used by data scientists, we are not going to implement them but treat them as references.

- **Wrapper Methods**

Wrapper Methods are less computationally efficient since they recursively run the data mining/machine learning algorithms with different feature sets, especially for much more complicated models like Kernel SVM, and Artificial Neural Networks. However, they produce a feature set which is tuned to a specific type of predictive model.

There are three most popular methods – Forward Feature Selection, Backward Feature Elimination and Recursive Feature elimination. Forward Feature Selection starts with an empty set and iteratively add the feature which best improves the performance. Backward Feature Elimination is the reverse version of Forward Feature Selection which starts with a set with all the features. Recursive Feature elimination is a greedy algorithm which aims to find the best performing feature subset. We are going to implement at least one the them.

- **Feature Clustering**

Clustering is always being used to group data points into different clusters. However, we can also treat each feature as a data point in a feature space. We want to group closed feature subspaces together to perform feature selection.

For assignment-based clustering (k-means, k-medians, k-medoids), users have to provide the number of clusters, which is not applicable in feature selection.

For density-based clustering (DBSCAN), outliers will be ignored by the algorithm. For feature selection, an outlier should be kept because it represents a feature that cannot be represented by other features.

For hierarchical clustering, it follows the strategy that points are grouped so that within each group there are no points that are very far from each other. Moreover, it will not let the user to specify the number of clusters at first. These properties fit the feature selection objective very well. The famous machine learning toolkit *scikit-learn* has also added a method `sklearn.cluster.FeatureAgglomeration`, which use the idea of hierarchical clustering. Thus, we are going to implement hierarchical clustering for feature selection.

With the help with those techniques, hopefully we can improve the prediction accuracy with full feature space.