

Feature Selection

Tushar B. Kute,
<http://tusharkute.com>



Feature Selection

- The training time and performance of a machine learning algorithm depends heavily on the features in the dataset. Ideally, we should only retain those features in the dataset that actually help our machine learning model learn something.
- Unnecessary and redundant features not only slow down the training time of an algorithm, but they also affect the performance of the algorithm. The process of selecting the most suitable features for training the machine learning model is called "feature selection".

Feature Selection

All Features



Feature Selection



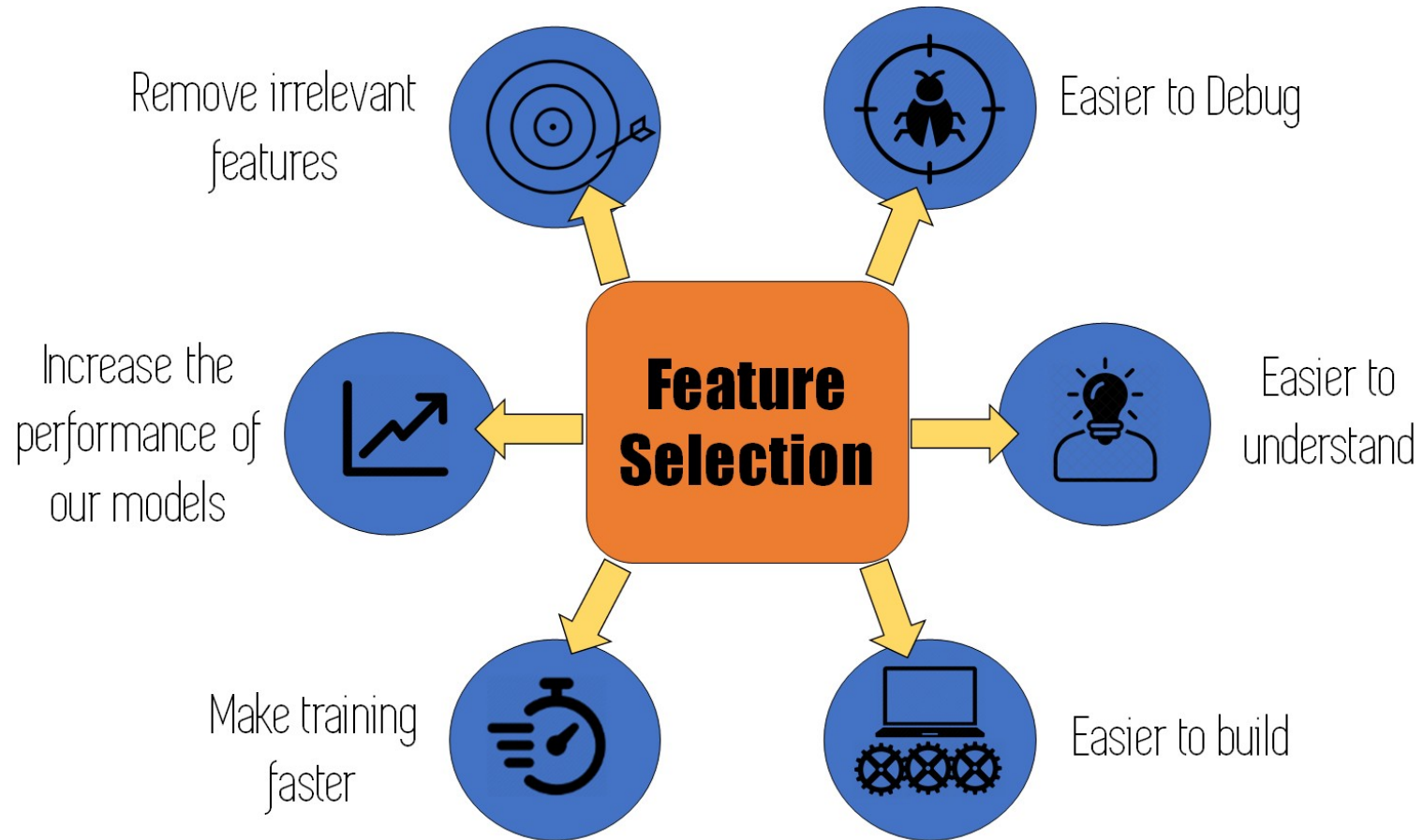
Final Features



Why feature selection ?

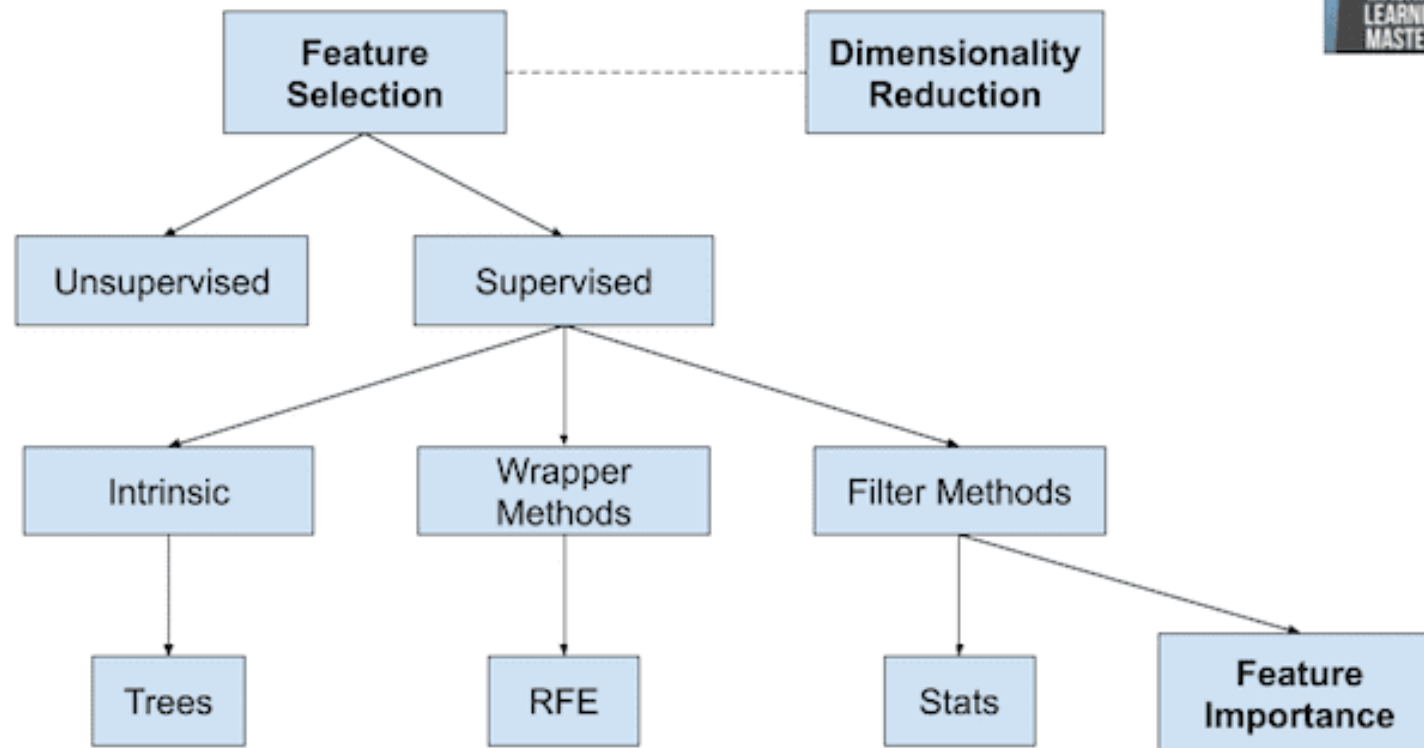
- There are several advantages of performing feature selection before training machine learning models, some of which have been enlisted below:
 - Models with less number of features have higher explainability
 - It is easier to implement machine learning models with reduced features
 - Fewer features lead to enhanced generalization which in turn reduces overfitting
 - Feature selection removes data redundancy
 - Training time of models with fewer features is significantly lower
 - Models with fewer features are less prone to errors

Why feature selection ?



Types

Overview of Feature Selection Techniques



Copyright © MachineLearningMastery.com

Filter Methods



- Filter method relies on the general uniqueness of the data to be evaluated and pick feature subset, not including any mining algorithm.
- Filter method uses the exact assessment criterion which includes distance, information, dependency, and consistency.
- The filter method uses the principal criteria of ranking technique and uses the rank ordering method for variable selection.

Filter Methods

- Filters methods belong to the category of feature selection methods that select features independently of the machine learning algorithm model. This is one of the biggest advantages of filter methods.
- Features selected using filter methods can be used as an input to any machine learning models.
- Another advantage of filter methods is that they are very fast. Filter methods are generally the first step in any feature selection pipeline.
- They are broadly categorized into two categories:
 - Univariate Filter Methods
 - Multivariate Filter Methods.

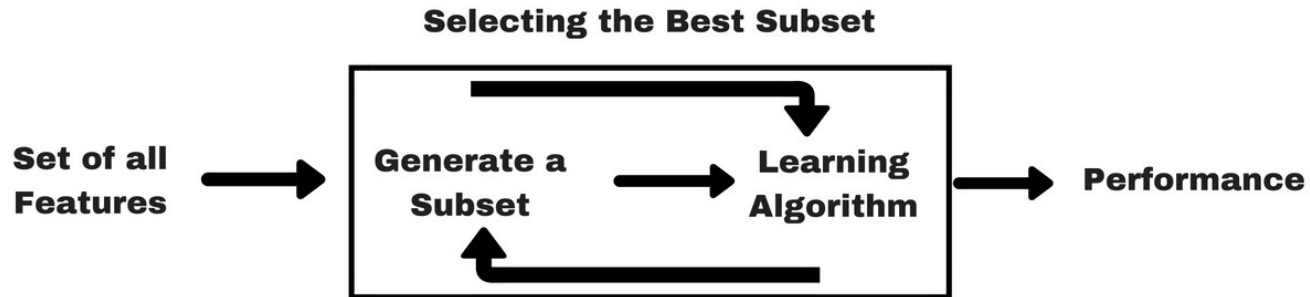
Univariate Filter Methods

- The univariate filter methods are the type of methods where individual features are ranked according to specific criteria.
- The top N features are then selected. Different types of ranking criteria are used for univariate filter methods, for example fisher score, mutual information, and variance of the feature.
- One of the major disadvantage of univariate filter methods is that they may select redundant features because the relationship between individual features is not taken into account while making decisions.
- Univariate filter methods are ideal for removing constant and quasi-constant features from the data.

Multivariate Filter Methods

- Multivariate filter methods are capable of removing redundant features from the data since they take the mutual relationship between the features into account.
- Multivariate filter methods can be used to remove duplicate and correlated features from the data.

Wrapper Methods



- A wrapper method needs one machine learning algorithm and uses its performance as evaluation criteria.
- This method searches for a feature which is best-suited for the machine learning algorithm and aims to improve the mining performance.
- To evaluate the features, the predictive accuracy used for classification tasks and goodness of cluster is evaluated using clustering.

Common Filter Methods

- Remove constant features
- Remove duplicate features and
- Remove correlated features
- Feature Importance: Information gain
- Feature Importance: chi-square test

Remove Constant Features

- Constant features are the type of features that contain only one value for all the outputs in the dataset.
- Constant features provide no information that can help in classification of the record at hand. Therefore, it is advisable to remove all the constant features from the dataset.
- Let's see how we can remove constant features from a dataset. The dataset that we are going to use for this example is the Santander Customer Satisfaction dataset, that can be downloaded from Kaggle.

Variance Threshold

- VarianceThreshold(threshold=0.0)
 - Feature selector that removes all low-variance features.
 - threshold
 - Features with a training-set variance lower than this threshold will be removed.
 - The default is to keep all features with non-zero variance, i.e. remove the features that have the same value in all samples.

Remove duplicate features

- Duplicate features are the features that have similar values.
- Duplicate features do not add any value to algorithm training, rather they add overhead and unnecessary delay to the training time.
- Therefore, it is always recommended to remove the duplicate features from the dataset before training.
- Removing duplicate columns can be computationally costly since we have to take the transpose of the data matrix before we can remove duplicate features.

Remove correlated features

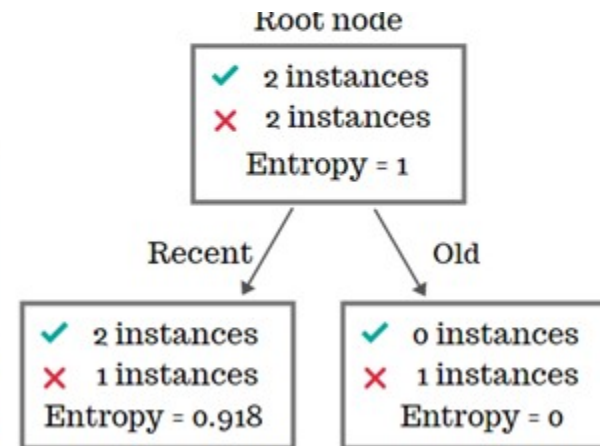
- A dataset can also contain correlated features. Two or more than two features are correlated if they are close to each other in the linear space.
- Take the example of the feature set for a fruit basket, the weight of the fruit basket is normally correlated with the price. The more the weight, the higher the price.
- Correlation between the output observations and the input features is very important and such features should be retained.
- However, if two or more than two features are mutually correlated, they convey redundant information to the model and hence only one of the correlated features should be retained to reduce the number of features.

Information Gain

- Information gain calculates the reduction in entropy from the transformation of a dataset.
- It can be used for feature selection by evaluating the Information gain of each variable in the context of the target variable.

Information gain for Age

Age	Mileage	Road Tested	Buy
Recent	Low	Yes	Buy ✓
Recent	High	Yes	Buy ✓
Old	Low	No	Don't buy ✗
Recent	High	No	Don't buy ✗



Chi-square test

- The Chi-square test is used for categorical features in a dataset. We calculate Chi-square between each feature and the target and select the desired number of features with the best Chi-square scores.
- In order to correctly apply the chi-squared in order to test the relation between various features in the dataset and the target variable, the following conditions have to be met:
 - the variables have to be categorical, sampled independently and values should have an expected frequency greater than 5.

Chi-square test

$$\chi^2 = \frac{(\text{Observed frequency} - \text{Expected frequency})^2}{\text{Expected frequency}}$$

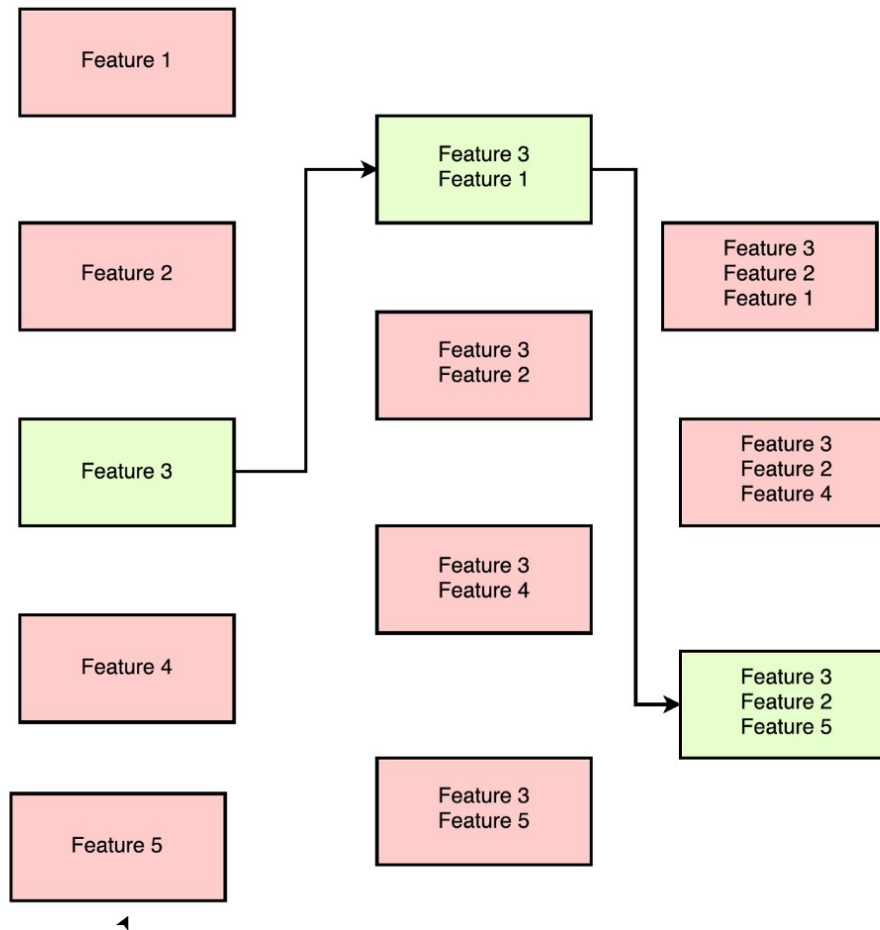
Recursive Feature Elimination

- The Recursive Feature Elimination (or RFE) works by recursively removing attributes and building a model on those attributes that remain.
- It uses the model accuracy to identify which attributes (and combination of attributes) contribute the most to predicting the target attribute.

Forward Feature Selection

- This is an iterative method wherein we start with the best performing variable against the target.
- Next, we select another variable that gives the best performance in combination with the first selected variable.
- This process continues until the preset criterion is achieved.

Forward Feature Selection



Backward Feature Elimination

- This method works exactly opposite to the Forward Feature Selection method.
- Here, we start with all the features available and build a model. Next, we remove the variable from the model which gives the best evaluation measure value.
- This process is continued until the preset criterion is achieved.

Comparing

Filter methods	Wrapper methods	Embedded methods
Generic set of methods which do not incorporate a specific machine learning algorithm .	Evaluates on a specific machine learning algorithm to find optimal features.	Embeds (fix) features during model building process . Feature selection is done by observing each iteration of model training phase.
Much faster compared to Wrapper methods in terms of time complexity	High computation time for a dataset with many features	Sits between Filter methods and Wrapper methods in terms of time complexity
Less prone to over-fitting	High chances of over-fitting because it involves training of machine learning models with different combination of features	Generally used to reduce over-fitting by penalizing the coefficients of a model being too large.
Examples – Correlation, Chi-Square test, ANOVA, Information gain etc.	Examples - Forward Selection, Backward elimination, Stepwise selection etc.	Examples - LASSO, Elastic Net, Ridge Regression etc.

Useful resources

- <https://stackabuse.com>
- <https://datacamp.com>
- <https://scikit-learn.org>
- www.towardsdatascience.com
- www.medium.com

Thank you

This presentation is created using LibreOffice Impress 5.1.6.2, can be used freely as per GNU General Public License



@mitu_skillologies



/mITuSkillologies



@mitu_group



/company/mitu-
skillologies



MITUSkillologies

Web Resources

<https://mitu.co.in>

<http://tusharkute.com>

contact@mitu.co.in

tushar@tusharkute.com