

Feature Reduction

Reducing number of features makes the model faster

Performance of Model may improve when only important features are provided to the model.

Makes data easier to visualize

Data is easier to store & obtain

Noise is removed

Reduces complexity of Model

Types of Dimensionality reduction

feature selection

select from existing dimensions

feature extraction

Make new variables of k dimensions

Curse of Dimensionality

Increasing number of features increases performance upto a point. Then the performance starts to degrade



Occam's razor: Remove all that is unnecessary

Unnecessary features \rightarrow Remove (Dimensionality Reduction)

Unnecessary mapping \rightarrow Remove (Regularization)

Feature extraction

Yields New data set from old dataset

Prevent redundancy

Prevent irrelevancy

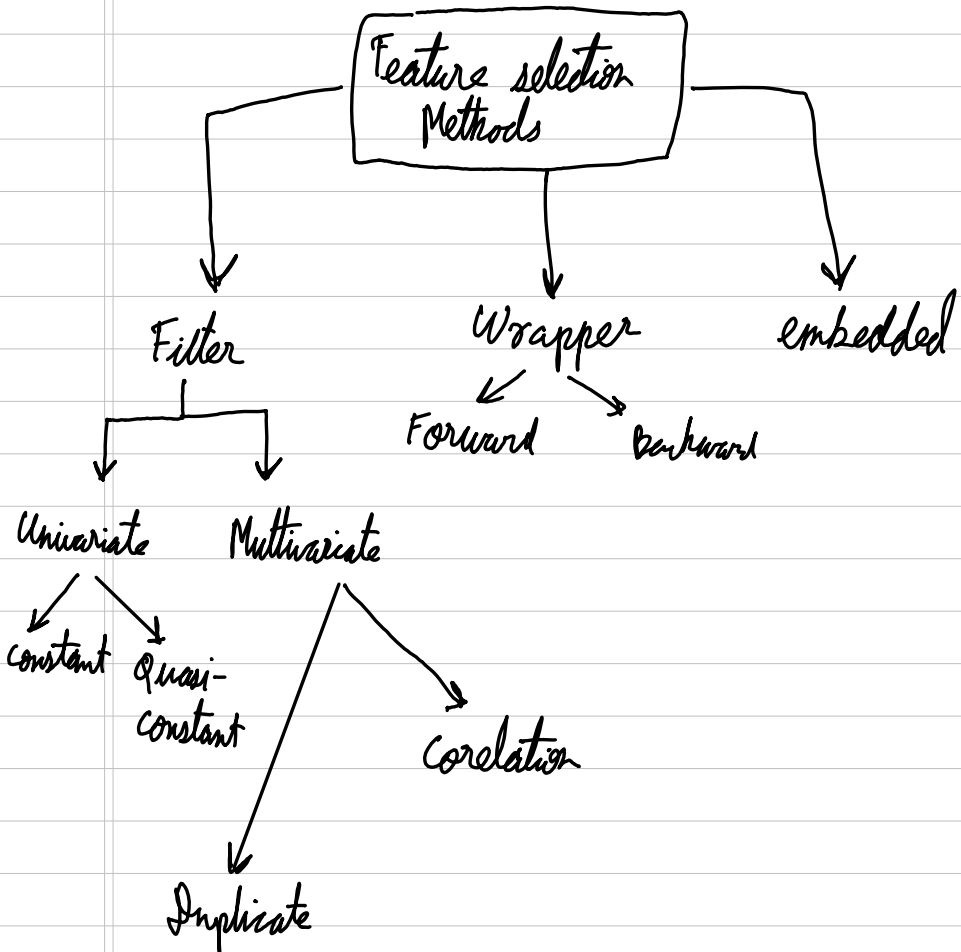
eg PCA

Feature selection

Features are selected & omitted

New space is not formed

More interpretable than Feature extraction



Filter Method

- select features independent of any model
- Rely only on characteristics of variables
- Inexpensive methods very fast
- Remove irrelevant & constant redundant features

Univariate filter → Treat every feature independently

A) Constant features

eg. gender of all patients is female, then no point keeping it in dataset

B) Quasi Constant

One value occupies majority records

eg. age > 50 for 98% patients then remove age.

Multivariate filter → Relation to other features

A) Duplicated features

B) Correlation filters

If two variables are highly correlated among themselves then they are redundant

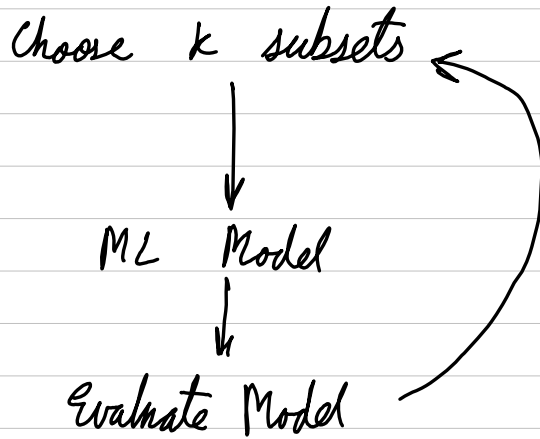
Use Pearson correlation coefficient to calculate it

Wrapper Methods

Greedy search approach

Evaluate all combinations by evaluation criterion like accuracy

Cross Validation based methods



High computational time required

High Overfitting

Feature selection : Forward selection

start with 0 parameters

take a feature & check

take another feature & check

⋮
all features

take 2 features & check

⋮

check for combinations

Backward selection

start with all features

Remove features & check

Embedded Methods

Perform Feature selection during model training

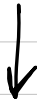
Learning algorithm takes advantage of its own variable selection process

More accurate like Wrapper

Fast like filter

Less prone to overfitting

Train Model



select features based on importance



Remove Non important features

example RIDGE & LASSO Regularization