

The Curse of Dimensionality refers to the various challenges and complications that arise when analysing and organising data in high-dimensional spaces which leads to other problems such as “over fitting”.

How to solve it?

1-Variance Threshold

Feature selector that removes all low-variance features.

This feature selection algorithm looks only at the features (X), not the desired outputs (y), and can thus be used for unsupervised learning.

2-missing value threshold

Finding columns with missing values then when the percent of missing value of column is above a certain ratio we apply mask

3-See how features are related to each other

-Pairwise correlation

We first see relation by corr function and see how each feature related to other features then we create mask where true values are related and false are not

When we gain this info we see that if 2 features are strongly related to the data we might wanna remove features close for 1 or -1 because these features will have same values

Feature Selection:- This module is used for feature selection/dimensionality reduction on given datasets. This is done either to improve estimators' accuracy scores or to boost their performance on very high-dimensional datasets.

Feature Extraction:- This module is used to extract features in a format supported by machine learning algorithms from the given datasets consisting of formats such as text and image.

The main difference:- Feature Extraction transforms an arbitrary data, such as text or images, into numerical features that is understood by machine learning algorithms. Feature Selection on the other hand is a machine learning technique applied on these (numerical) features.

We have another approach by selecting features

-we pick feature with small coeff and remove it since it does not have a huge effect accuracy by RFE in sklearn

Tree-based feature selection

We create mask for the importance of the features based on their contribution to the decision tree's performance.

Now with regression

-BY lasso regressor to reduce our features with alpha parameter

LassoCV to choose appropriate value of alpha by cross-validation where $\text{coff} \neq 0$