

- Heatmap
- boxplot
- countplot
- barplot

Pearson's correlation

Pandas Profiling

from pandas-profiling import ProfileReport

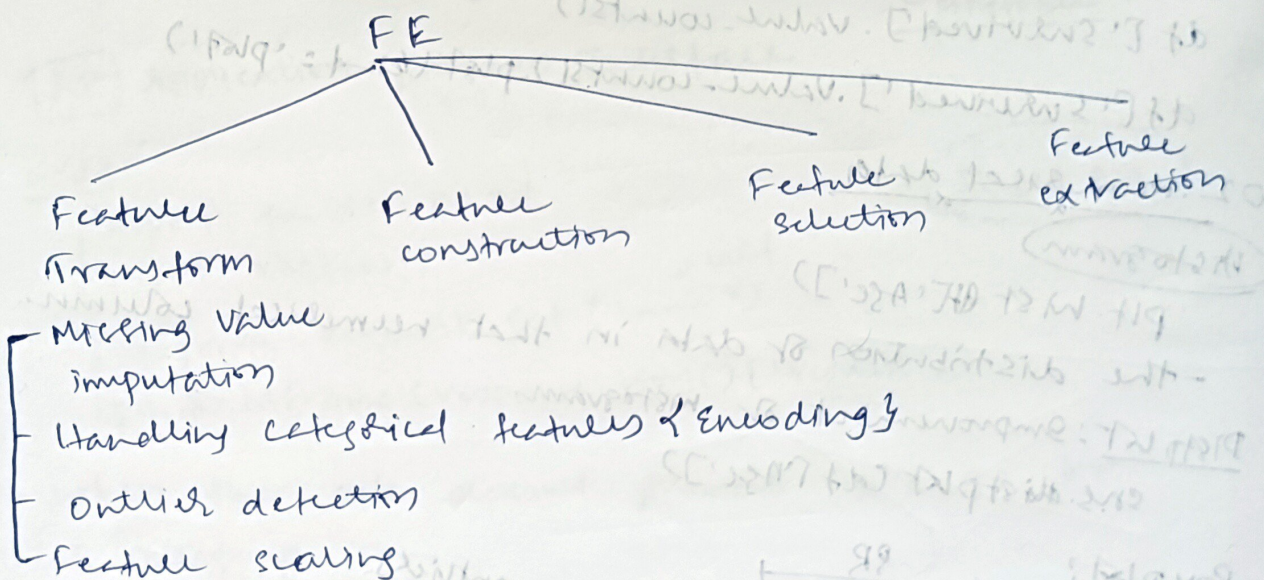
prof = ProfileReport(df)

prof.to_file(output_file="output.html")

- makes a webpage which contains the report of our data.

Feature Engineering

- The process of using domain knowledge to extract features from raw data.



Standardization

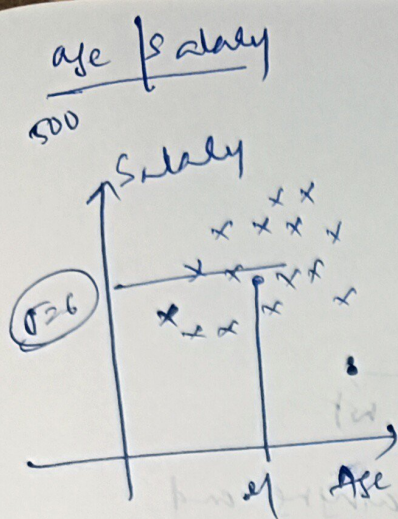
- also called 2-side normalization

Age	Salary
27	
15	
33	
63	
⋮	
50000	

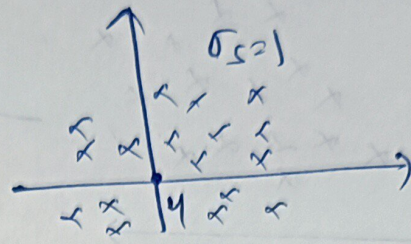
$$x' = \frac{x - \bar{x}}{\sigma}$$

$$\mu = 0 \quad \sigma = 1$$

after standardization.



Standardization



Standard scaler

from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()

scaler.fit(X_train)

X_train_scaled = scaler.transform(X_train)

X_test_scaled = scaler.transform(X_test)

scaler.mean

- scaling doesn't effect decision tree but effects logistic regression.

K-means

KNN

PLA

ANN

Gradient descent

use standardization.

Normalization

- To change the values of numeric columns in the dataset to use common scale.

MinMax scaling

Robust scaling

Mean Normalization

Max absolute

- sets the data in a single range

MinMax Scaling

weights

130

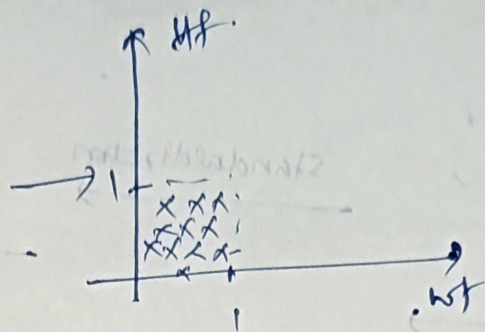
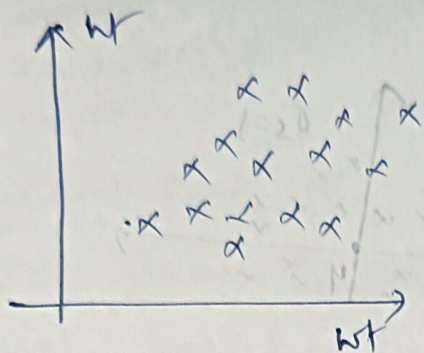
67

92

88

$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

wt, ht



- Before scaling make sure to split training and testing data.

from sklearn.preprocessing import MinMaxScaler

scaler = MinMaxScaler()

scaler.fit(X_train, y_train)

X_scaled_train = scaler.transform(X_train)

X_scaled_test = scaler.transform(X_test)

Mean Normalization:
$$X_i' = \frac{X_i - X_{\text{mean}}}{X_{\text{max}} - X_{\text{min}}}$$

Max Absolute scaling:
$$X_i' = \frac{X_i}{|X_{\text{max}}|}$$
 (for sparse matrix)

Robust scaling:
$$X_i' = \frac{X_i - X_{\text{median}}}{\text{IQR}}$$

IQR = 75th percentile - 25th percentile

Robust to outliers

$$X_i' = \frac{X_i - X_{\text{min}}}{X_{\text{max}} - X_{\text{min}}}$$