

# Feature engineering

→ There are a lot of features in any dataset

The dataset may not be pure and may contain a lot of noise

If such a raw dataset is fed into a ML model the performance of the model will suffer.

ML algorithms require the data to have a specific characteristic

Feature engineering is the process of transforming raw data into features that are more suitable for the model

It is the process of selecting, extracting & transforming the most relevant processes

These techniques highlight the most relevant patterns

## Goals of feature engineering

① Preparing proper input dataset compatible with the model

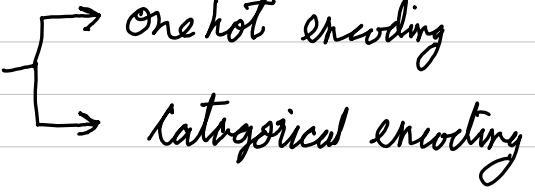
② Improving performance of ML model by cleaning the dataset

# Clean data Principles

These principles must be followed for good accuracy

- ① Each variable must have their own column → each column is a variable  
ie Name must be split into m name  
L name
- ② No cell must contain missing values
- ③ End result must be a numeric value

# List of Techniques

- ① Imputation
- ② Outliers
- ③ Binning
- ④ Data encoding 
  - one hot encoding
  - categorical encoding
- ⑤ Grouping
- ⑥ Feature split
- ⑦ Scaling
- ⑧ Log transforms

# Imputation

- ① Dataset contain missing values
- ② In order to remove the missing values, three things can be done

- ① Dropping
- ② Numerical Imputation
- ③ Categorical Imputation

① Dropping → If blank values are more in number, then drop the column or row  
threshold → 70%.

i) DROP columns where more than 70% data is missing

ii) DROP rows where more than 70% data is missing

② Numerical Imputation → If blank values are less then replace it with some values

Replacing can be done by

- i) Mean of non blank values
- ii) Median of non blank values

③ Categorical Imputation → Replace blank data by a category

If blank values make sense then replace them by a category like "other" "N.A." or "0"

# Missing value Patterns

In real life datasets, certain datapoints may have missing values due to human errors and unavailable respondents

## ① Missing Completely at Random

Data is missing at random without any relation to the characteristics of data.

eg. receptionist forgot to enter the age of a patient at random.

## ② Missing at Random.

Data is missing but missing because someone has not disclosed it.

The probability of the missing value depends on some observable data

eg. Young people may not disclose their name

name  $\rightarrow$  Missing data  
Age: young  $\rightarrow$  observed data

Missing data is due to an observed data.

## ③ Missing Not at random.

Data is missing because people want to hide it

eg. Rich people may not disclose their income.

The income is not disclosed because it is too high or too low

The missing is related to characteristic of missing unobserved data itself

# Outliers

- Ⓐ Outliers are exceptions
- Ⓑ Outliers are not to be treated like normal data points
- Ⓒ Example ages of students in a school

16, 18, 15, 22, 16, 14, 19, 55, 14, 12, 13



outlier : odd one out

- Ⓓ If we feed such data into model, then model will give wrong output
- Ⓔ even simple average of all students age will be a wrong parameter to predict student age
- Ⓕ If such outliers are very few in number then they must be removed
- Ⓖ Outliers can also be due to human error in data entry  
eg 15 was written as 55 by mistake

# Outlier detection

Outliers can be detected by statistical methods

## (A) Standard deviation

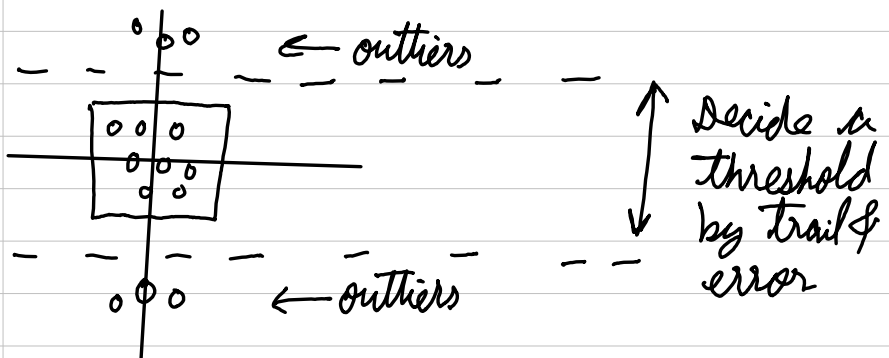
if a value has distance from average that is higher than  $\pm \times$  S.D

$$x - \bar{x} > t \times \text{S.D.}$$

↪ threshold

## (B) Thresholding

Depending on dataset, visualize the dataset



$x > t_1$  or  $x < t_2$  are outliers

## (C) Percentile

Decide a percentile for outliers

eg Top 5% are outliers

or Bottom 2.5% are outliers

But adjustment needs to be done on basis of graph visualization

# Outlier Handling

What to do after outliers are found?

① Remove outliers

→ If outliers are real values  
eg man of 55 does attend school

② Cap them to highest value / lowest value

55 becomes 23

③ Put realistic values for human error

eg a data point is "-17" we can understand that -ve is typo error



# Binning

- (a) Binning means placing data into bins or categories
- (b) Binning is used to prevent overfitting
- (c) Binning is used to reduce data
- (d) Binning leads to sacrificing the data & makes data more regularized

## (A) Numerical Binning

Value	Bin
0-30	→ Low
31-70	→ Mid
71-100	→ High

## (B) Categorical Binning

Value	Bin
Mumbai	→ Maharashtra
Pune	→ Maharashtra
Kolkata	→ Bengal

# One Hot encoding

① used to convert categorical data to numeric data

User	City
1	Mumbai
2	Pune
1	Pune
3	Kolkata
2	Kolkata
1	Kolkata
1	Mumbai

User	Mumbai	Pune
1	1	0
2	0	1
1	0	1
3	0	0
2	0	0
1	0	0
1	1	0

Kolkata  
city values  
can be inferred

← } Kolkata  
← }  
← }

②  $N$  distinct values converted to  $N-1$  columns  
3 values to 2 columns for city

# Catagorical Encoding

In this, simply replace the catagorical values by numbers

User	city
1	Mumbai
2	Pune
1	Pune
3	kolkata
2	kolkata
1	kolkata
1	Mumbai

User	city
1	1
2	2
1	2
3	3
2	3
1	3
1	1

# Grouping

Pivot tables can be used in similar way as one hot encoding but with non binary values

User	City	Visit
1	Mumbai	1
2	Pune	2
1	Pune	1
3	Kolkata	1
2	Kolkata	4
1	Kolkata	3
1	Mumbai	3

User	Mumbai	Pune	Kolkata
1	4	1	3
2	0	2	4
3	0	0	1

$N$  columns for  $N$  cities

Count the values (Unlike one hot encoding)

# Feature split

- ① sometimes, a feature might be represented in a different manner in dataset
- ② separation of relevant features need to be done

Name		F. Name	L. name
Harry Potter	→	Harry	Potter
Tom Riddle		Tom	Riddle

- ③ Extracting date can be done by various methods

- i) traditional date split into dd/mm/yy or other format & remove unnecessary columns
- ii) Extracting timeperiod eg time passed since
- iii) Extracting features eg weekday  
holiday  
old/new

# log transform

- ① sometimes data is skewed
- ② log transform helps to even such data
- ③ It also decreases the effect of outliers
- ④ eg age difference between ages 10-15 is more relevant than age difference in ages 65-70

5 years of difference is higher for small magnitude

$$x = \log(x+1)$$

only works for +ve values.

# Scaling

- ① Numerical features differ between themselves and don't have same range

eg age & income won't have the same range

- ② But ML model expects everything to be in the same range

- ③ Normalization based on range

$$X_{\text{norm}} = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$$

MinMax normalization (range will be 0-1)

- ④ Standardization (Z-score)

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

standard deviation is used to convert the values (range not 0-1)  
downscales all values, makes some -ve values

Mnemonic

F I S H B L O G

Feature extraction

Imputation

Scaling

Handling outliers

Binning

Log transform

One Hot encoding

Grouping