# Feature Engineering

## Types of Encoding

**1**    One hot Encoding

→ Assigning the values 1 or 0.

gt for example

| germany | France | Spain |
|---------|--------|-------|
| 1 | 0 | 0 |
| 0 | 1 | 0 |
| 0 | 0 | 1 |

↓

3 - 1 = 2

Columns coz
we Removed Spain

Dummy variable Trap.

↓

Performed with pandas and Sklearn

↓

pd.get - dummies

→ Skipping the Spain Column

   because we can derieve it by last Raw which is 0,0 so it is understood that last row will be 1 so it is Spain.

# Disadvantages of one hat Encoding :-

### Pincode Example :-

560011
560012 → If by previous Example
560013     i convert them
|          into something call
|          dummy then 99
and so on.   columns will get
             created.

→ If we perform one hat Encoding here. Then more columns will get Increased and it may cause. curse of dimension. ←

Leading us to this.

→ Many numbers of categories. So Don't need to Apply one hat Encoding.

## 2. Label Encoding →

### ORdinal category

Education
| | |
|---|---|
| BE | 2 |
| masters | 3 |
| PhD | 4 |
| Statistician | 1 |

## 3. One hat Encoding with multiple categories ↓

nominal category ↓

## 1. Target Guided ordinal categories.

Classification problem

$$\begin{array}{ll}
\text{A} & 1 \\
\text{B} & 1 \\
\text{C} & \\
\text{D} & \\
\text{A} & \\
\text{B} &
\end{array}$$

O/P
1 → mean → 0.73
1 → 0.6
0 → 0.4
1
0
0

→ where the value of A is 0 and 1 considering these features only.

→ when you consider the mean you are finding the number of values which is for A = 1.

→ Rearranging the values by there Rank } due to ordinal category.

LABEL } due to ordinal.

4
3
2
1
4
3
2
1

# Mean Encoding → Nominal

| f1 | | O/P | |
|----|----|----|----|
| A | ←→ | 1 | |
| B | | 0 | 0.73 |
| C | | 1 | 0.6 |
| D | | 1 | |
| A | | 1 | 0.5 |
| B | | 0 | 0.4 |
| C | | 1 | |
| D | | 1 | |

→ we will convert this into mean values.

| f1 | O/P |
|----|----|
| 56011 | 1 |
| 56022 | 0 |
| | 1 |
| | 1 |

↓

Finding out the mean and this values 56011 will be Replaced by the values of mean.

# Why feature Scaling

| Features | cm Height | kg weight | BMI |
|---|---|---|---|
| ⌐ magnitudes | 180 | 78 | |
| ⌐ units | 170 | 84 | |

→ Not perform feature Scaling.

① Decision Tree
② Random Forest
③ X-G Boost.

# Handle missing Values in Categorical Variables

1  delete the Rows

2  Replace with the most Frequent Values.

3  Apply classifier Algorithm to predict.

4  Apply unsupervised ml.

## Ordinal numbering Encoding or Label Encoding

Ordinal categorical Variables :-

→ ordinal data is categorical, Statistical data type where the variables have natural, ordered categories and the distances between the categories is not known.

# Categorical

| Nominal | Ordinal |
|---|---|
| Pen, Pencil, Eraser | Excellent, good, Bad |
| Cow, Dog, Cat | Fantastic, okay, don't Like |

# Life Cycle of a Data Science Projects

1. **Data Collection Strategy** :- From company side, 3rd party Api's, Surveys

2. **Feature Engineering** :- Handling missing values.

why are there missing values? Survey - depression Survey.

1. They hesitate to put down the into.
2. Survey information are not that valid.
3. men -- salary
4. Women --- age
5. people may have died --- NAN.

3. Data that will be missing

1) continuous data
2) categorical Data

What are the different types of missing data

1. missing completely at random (MCAR)

2. missing at random (MAR)

3. Not missing at Random (NMAR)

For getting the null values in a
Particular Row or column.
df [df ['Embarked'].isnull()]

1. MCAR :- It means that there is no
relationship between the
two particular thing in the
Dataset. ↓
                      No Relationship

2. NMAR :- There is absolutely a Relationship
between the data missing and
any other values. ↓

                      Having Relationship

3. MAR :- men :- hide their salary
         women :- hide their Age.

# Random Sample Imputation :-

→ It consists of taking Random observation from the dataset and we use this observation to Replace the nan values.

## When Should it be used?

→ It assumes that the data are missing completly at Random.

## Advantages

1. Easy to Implement

2. There is less distortion in variance.

## Dis-Advantages

1. Every Situation Randomness wont work.

# Capturing NAN Values with New Feature

You can use this by :- **MNAR** :-

**Advantages :-**
(i) Easy to Implement
(ii) captures the Importance of missing values.

**Dis Advantages:-** Creating Additional Features.
[ Curse of Dimensionality].

## End of Distribution Imputation

# End of Distribution Imputation

Advantages :- Easy to Implement

→ Captures the Importance of missingness if there is one.

Disadvantages :- Distorts the original distribution of the variable

→ If the number of NA is big, it will mask true outliers in the distribution

# Arbitrary Value Imputation

The Technique was derived from kaggle competition. It consist of the method to Replace the NAN Values by arbitrary Values.

Advantages :- (i) Easy to Implement

(ii) Captures the importance of missingness if there is one

Disadvantages :- (i) Distorts the original distribution

(ii) Hard to decide which value to use.

# Handling Categorical Missing Values

## 1 Frequent category Imputation

Advantages :- (i) Easy to implement.

Disadvantage :- (i) It may lead the distortion in the Relation of the most frequent label

(ii) Since we are using the more frequent labels, it may use them in an over represented way, if there are many nan's.

## 2 Adding a Variable to capture NAN