# One hot encoding

## Types of Encoding :-

Data Science Life cycle :-

1) Data Ingestion (collection of data)
2) EDA
3) Processing
4) Model
5) Evaluate & Validate Model

} Core ML pipelines .

<u>Statistics</u> — It is a science of collecting, organising and analysing data.

Collect, Organise, Interpretation, Analysis

↓

Insight .

# Types of Data

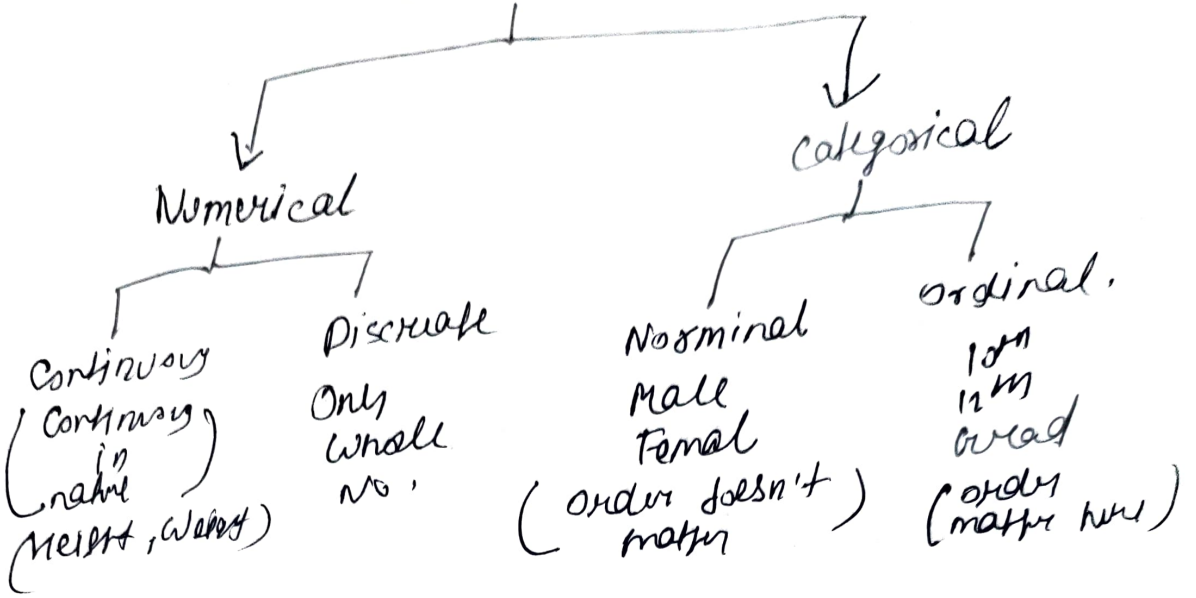Batch Data : Historical Data, Minibatch Data (Periodic)
Streaming Data : Continuous Data (Live Data)

① Structured Data : Table (Row × Column) → ML
② Unstructured " : Videos, Images, Voice, Sound, Text etc.
③ Semi-structured : JSON, XML. ↳ DL

Example of Structured Data :-

| Feature 1 | Feature 2 | Feature 3 |
|-----------|-----------|-----------|
| Weight | Height | BMZ |

# Structured Data

```
                    Structured Data
                    /              \
                   ↓                ↓
              Numerical         Categorical
              /      \           /        \
             ↓        ↓         ↓          ↓
       Continuous  Discrete  Norminal    ordinal.
       (Continuous  Only      Male        10th
        nature)     Whole     Femal       12th
       (Height,Weight) No.   (order doesn't  Grad
                              matter)     (order
                                          matter here)
```

Univariate : Single column.
Bivariate : Two Column.
Multi-variate: More than two column.

Independent & Dependent Variable :-

[Age (Height, Sex)]    Weight
         ↓               ↓
     Independent      Dependent

Q) First EDA is required. or (Preprocessing)
   ↳ EDA ⊙ is required.
   Order: EDA → Pre-processing ⟹ [Models]

EDA (Analysis):—

1) Profiling ⟶ All types of plots
2) Statistic Analysis:                    ← Variance
3) Graph Based Analysis                      Covariance
                                             stel
         Row                                 Correlation
         Column                              chisquare test
         Missing                             T-test
         Category                            2-test
         Numeric                             Annova test
         Duplicate
         Dtype                               Mean Median Mode
         Ram
```

# Pre-processing of Data :-

① Missing Value Handle
② Outlier Handle
③ Scaling of Data
④ Transformation -
⑤ Encoding
⑥ Imbalance Data

⑦ Feature Selection
⑧ Dimension Reduction
⑨ Duplicate value / Duplicate column
⑩ Split / Merge / Prop / Add.

There are ~~3~~ steps of feature engi'm : ——

① Missing Null Value → Missing Value (EDA) Handler (PP)

② Outlier → Handle.

③ Categorical (man, Women) → Encoding.

④ Skewed Range → Scale ( within a certain range )

⑤ Count of Feature → { Handle Imbalanced Data. Feature Selection Dimension Reduction ( PCA, tsne) }

Encoding → To change categorical data into numerical data is called encoding.

Types of Encoding :-
↳ We are discussing about categorical var.

Gender < Male / Female

we follow diff encoding technique to convert it into maths.

# Types :-

1. Nominal Encoding ── Nominal cat. var
2. Ordinal Encoding ── For ── Ordinal cat. var
   └→ Rank

$\begin{cases} BE \\ Bcom \\ Phd \\ Master \end{cases}$

①
- One hot encoding
- ② one hot encoding with many categorical
- mean③ encoding

②
- Label④ encoding
- Target⑤ guided Ordinal encoding

Dummy variable Trap.

① 

| State | India | Pak | China |
|-------|-------|-----|-------|
| India | 1 | 0 | 0 |
| Pak | 0 | 1 | 0 |
| China | 0 | 0 | 0 |

→ when it is 0 0 then China.

Disadvantage of one hot encoding

In place of country (pincode)
- So on that time we need to create lots of columns.

└→ For this ④ Label Encoding.

Education

BE    1     → we will give rank.
was   2
phd   3

② (One hot encoding with multiple cat)

### KDD Orange

Which top 10 categories repeated more
⮑ And create 9 columns.

③ Target Guided

| $f_1$ ⟷ O/P | | → Mean | Classify | |
|---|---|---|---|---|
| A | ! | | 0.73 | } Assign rank. |
| B | ! | | 0.6 | |
| C | 0 | | 0.4 | |
| D | 0 | | | |

③ Mean Encoding

| | O/P | Mean |
|---|---|---|
| A. | ! | 0.73 |
| A | 0 | 0.6 |
| B | ! | 0.5 |
| C | 0 | 0.4 |
| D | | |

Pincode

| | A | O/P | |
|---|---|---|---|
| | 56001 | ! | 0.73 |
| | 56001 | 2 | 0.6 |

Placed at pincode position

---

## Why Feature Scaling :-

|  | cm | kg | |
|---|---|---|---|
| Features | Height | Weight | BML |
| ⮑ magnitude | 187 | 78 | → Scale down this value |
| ⮑ Units | 170 | 84 | |

① Linear Regression ,

② K Means
③ KNN

More feature scaling used.

global Minima

① Decision tree ⎫ No need
② RF ⎬ of
③ xgboost ⎭ feature scaling.

## How to handle missing values :—

Cons

① Delete the Rows     Impt data may be deleted.
② Replace the most frequent → Imbalance data occur.
    value.
③ Apply classifier algorithm → $f_2, f_3$ & o/p used to predict $f_1$.
    to predict.
④ Apply unsupervised ML → Take $f_2$ & $f_3$ we start
    (clustering technique)   creating into 2 categories

$$
\begin{bmatrix}
 & f_1 & f_2 & f_3 & \%p \\
Male & 23 & .24 & Yes \\
Female & 24 & 25 & NO \\
\end{bmatrix}
$$

## Handle Categorical Features :—

— Replace each label of the categorical variable
by the count.     Cons :—
We are losing some key of info.

| | $X_1$ | count |
|---|---|---|
| 0 | w | 150 |
| 1 | v | 80 |
| 2 | x | 72 |
| 3 | y | 75 |
| 4 | | |

## Handling Ordinal Categories :—
    (Ordinal Encoding)

Cons :—
Does not add machine learning option.

create — ⬭ A   B   C   Fail
      4   3   2   1

| Sum | Math |
|-----|------|
| Mon | Bot . |
| Tue | (Holiday) |
| Wed | ( Phy |
| Tm | .Eng . |
| Fri | Chem |
| Sat | Zoology |

Chem $\underline{2}$ Phy

Sun × Holiday

Holiday
Phy ↑ Botany

M _ Phy = Zo _ Eng

Math × Mon
× ·Thes

Math ↓i
Not holiday

R i    X $\underline{3}$R
W

P    Q
Y    $\frac{3}{W}$
     ·

mn    ANTARCTICA

~~AANAARCLIA~~

·NANRATCCI

| 10 | U |   |
|----|---|---|
| 9 | · |   |
| 8 | T | T |
| 7 | 2 | 2 |
| 6 | Ø | U |
| S | ● | ● |
| 4 | · |   |
| 3 | · |   |
| 2 | · |   |
| 1 | S |   |