# Feature Engineering

Feature? $\longrightarrow$ Problem and data specific

The input variables of your model.

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix} \longrightarrow \boxed{model} \longrightarrow \hat{y}$$

Churn Prediction (Bank)

i/p variables : $\begin{cases} \text{Salary} \\ \text{Credit Score} \\ \text{Bank balance} \\ \text{\# products} \\ \text{Credit Card User or not} \end{cases}$

favourite color

NLP : Sentiment analysis

(+ve / -ve)

The direction and the cinematography of the movie is really good. $\longrightarrow$ (+ve)

## Feature Engineering

1) **Feature pre-processing** $\longrightarrow$ preprocessing the i/p variables so that it can be
   used in the model.

2) **Feature Extraction** $\longrightarrow$ i/p variables $\longrightarrow$ features.

$$\phi(x) \longrightarrow x^*$$

$$f^* : \underset{i/p}{\mathbb{R}^{n \times d}} \longrightarrow \underset{o/p}{\mathbb{R}}$$

$d \rightarrow$ dimension

$\underset{\text{dimensionality}}{\underline{d' < d}}$ reduction

$$\hat{y} = f^*(x)$$

$$\phi : \mathbb{R}^{n \times d} \longrightarrow \mathbb{R}^{n \times d'}$$

$$\hat{y}' = \hat{f}(\phi(x)) \nearrow \text{model}$$
$\longrightarrow$ visualization
$\searrow$ improve the performance

3) **Feature Selection :**

Feature Preprocessing :        Feature Extration

$x_1 \longrightarrow x_1^*$              $x_1 x_2$ , $x_1/x_2$ , $(1 + x_1 x_2)^3$

$x_2 \longrightarrow x_2^*$

$\vdots$                            $(x_1)$             $(x_2)$

$x_n \longrightarrow x_n^*$        Units Sold      price per unit      amount $(x_1 \cdot x_2)$

RFM : Recency  Frequency  Monetary          $111 \rightarrow 555$

$\mathcal{U}$        1            1          1

           $\vdots$        $\vdots$       $\vdots$

           5            5          5

GDP        # population      GDP / person $\mathcal{U}$

Raw features $\longrightarrow$ feature Extration $\longrightarrow$ preprocessing $\longrightarrow$ Feature Selection

$x_2$ , $x'_2$ , $x'_1$

$$(X_1, X_2) \xrightarrow{PCA} (X'_1, X'_2)$$

$X_1$

Sphere

$(x, y, z)$

$(\phi, \lambda)^2$

$X'_1 \rightarrow$ Low variance along $X'_2$

$a \cdots b \xrightarrow{\text{projection}} a^* \ b^*$

$x'_2$

$x'_1$  high variance along $x'_1$

Feature Selection

Filter Technique : Filter out less important features using some criteria

Wrapper Technique : (set of features) $\longrightarrow$ (subset)
$\downarrow$
model

(RFE)

# Filter Methods:

1) Variance threshold

$f1 \longrightarrow$ variance $= 0$

$f1$
1
1
1
1
.
.
.
.
1

$f2$
1
1
0
1
1
0
.
.
1

$(80\% = 1, \quad 20\% = 0)$

$p(1-p)$

$= 0.8 \times 0.2$

$= \underline{0.16}$

$-$

$f_3$
0.3
0.5
.
1.2
,
.
.

$\boxed{0.8}$.

2) **Univariate F-test**

feature (X)

|  |  |
|---|---|
|  | 2 |
| 3 | 3 |
|  | 4 |
|  | 5 |
| 6 | 6 |
|  | 7 |
|  | 8 |
| 1.5 / 0.5 / 9 |  |
| 9.5 / 0.5 / 10 |  |
| 1.5 / 11 |  |

Class(y)

A ⎫
A ⎬
A ⎭
B ⎫
B ⎬
B ⎭
C ⎫
C ⎪
C ⎬
C ⎭

**Step-1 :** Compute Group mean. of X

$$\bar{X}_A = \frac{2+3+4}{3} = 3$$

$$\bar{X}_B = 6 \quad , \quad \bar{X}_C = 9.5$$

**Step-2 :** Overall mean

$$\bar{X} = \frac{65}{10} = 6.5$$

**Step-3 :** Between Group Variance

$$SSB = 3(\bar{X}_A - \bar{x})^2 + 3(\bar{X}_B - \bar{x})^2 + 4(\bar{X}_C - \bar{x})^2$$

$$= 3 \cdot 3.5^2 + 3 \cdot (0.5)^2 + 4 \cdot 3^2$$

$$= 73.5$$

$$MSB = \frac{SSB}{df} = \frac{73.5}{2} = 36.75$$

**Step-4 :** Within Group Variance:

$$SSW = \sum_j \sum_i (x_{ij} - \overline{x}_j)^2 = 9 \checkmark$$

$$MSW = \frac{SSW}{dof} = \frac{SSW}{\underset{\text{\#data}}{N} - \underset{\text{\#class}}{k}} = \frac{SSW}{10 - 3} = \frac{9}{7} = 1.28$$

**Step-5 :** F-statistics:

$$F = \frac{MSB}{MSW} = \frac{36.75}{1.28} \simeq \underline{28.7}$$