# Features vs Targets

**What is a feature?**
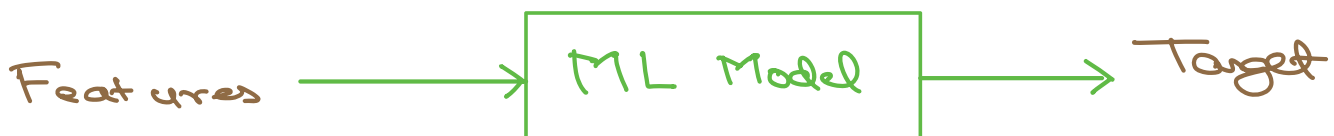
## Customer Dataset

*inputs/Features* · Target

| index | Education | Gender | Income | Fitness . . . . — — — — . . . Usage | Product |
|---|---|---|---|---|---|
| 0 | 10th | M | 20,000 | ------------------------- | P₁ |
| 1 | 12th | F | 30,000 | ------------------------- | P₁ |
| 2 | Masters | M | 60,000 | ------------------------- | P₃ |
| ,, | ,, | ,, | ,, | ------------------------- | ,, |
| ,, | ,, | ,, | ,, | ------------------------- | ,, |
| ,, | ,, | ,, | ,, | ------------------------- | ,, |
| ,, | ,, | ,, | ,, | ------------------------- | ,, |
| n | | | | | |

1 cust →

**Goal:** Given 'details' of a person build ML 'Model' to Recommend Best Product.

| Features | Targets |
|---|---|
| Education | Product |
| Gender | |
| Income | |
| ,, | |
| ,, | |

Features ⟶ [ ML Model ] ⟶ Target

# Feature Engineering

Processing features to improve ML Models

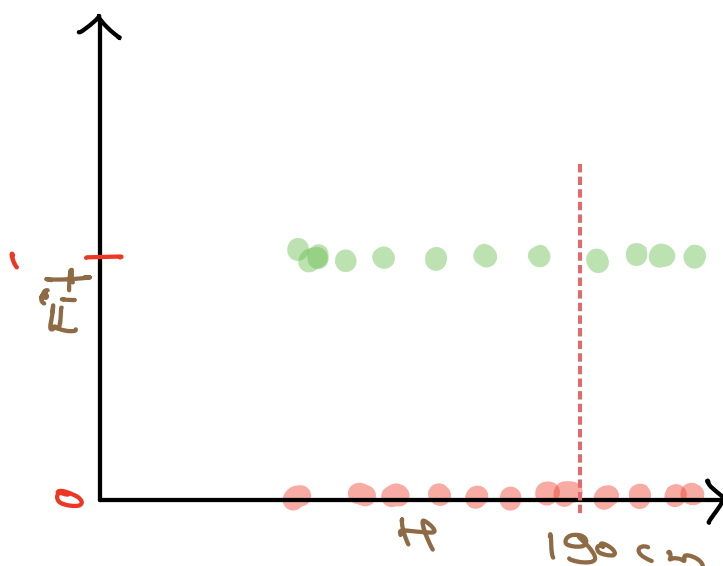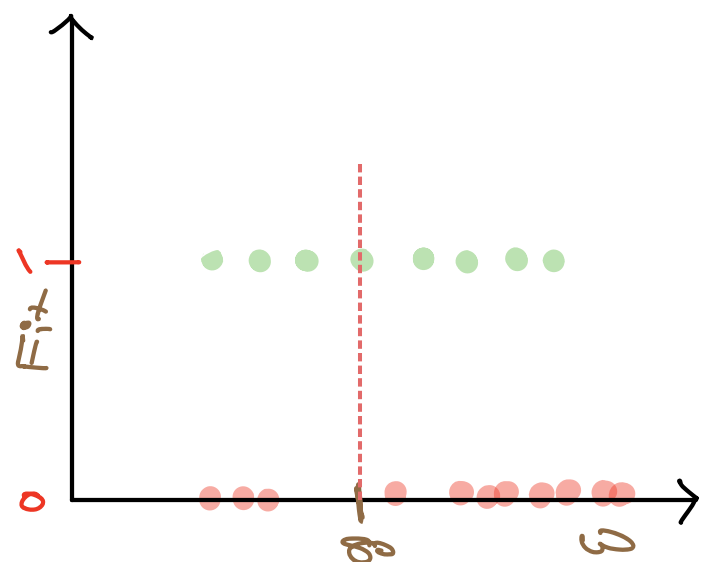→ create New Feature

→ Apply transformation on Existing Features

Ex: Given weight and Height of a person predict
if they are Fit or Not Fit

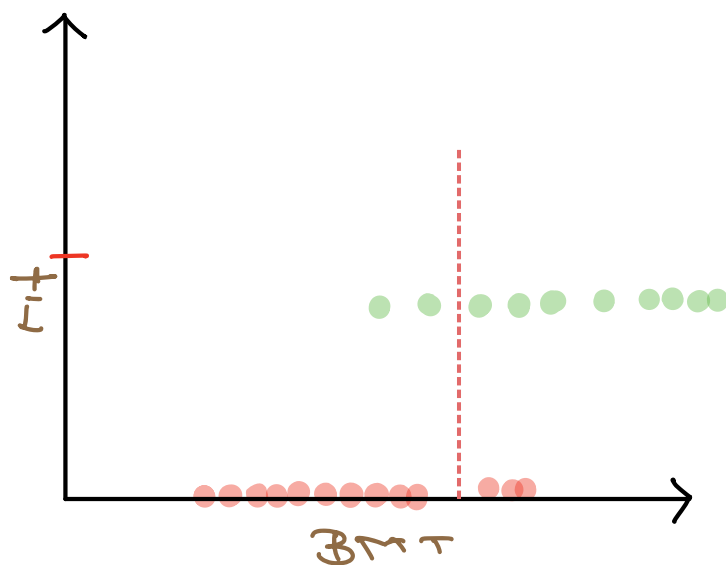| H | W | Fit |
|---|---|-----|
| 170 | 70 | Yes |
| 160 | 65 | Yes |
| 165 | 80 | No |
| 150 | 90 | No |
| ,, | ,, | ,, |
| ,, | ,, | ,, |

Features = H, W

Target = Fit

Can we plot the features against Targets



Q= Is there anything that we can do to Seperate Fit and Non-Fit people clearly?

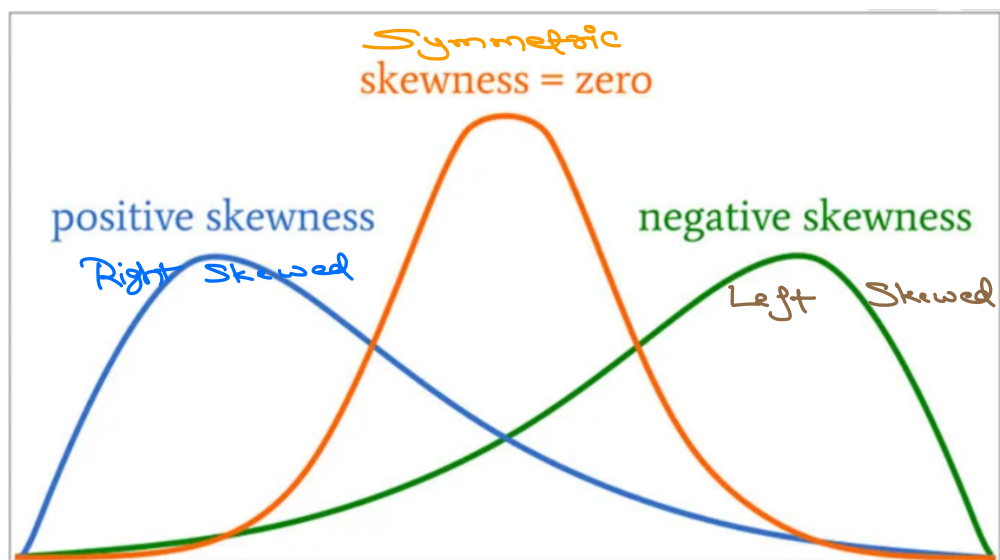| H | W | BMI | Fit |
|---|---|---|---|
| 170 | 70 | $N_1$ | Yes |
| 160 | 65 | $N_2$ | Yes |
| 165 | 80 | $N_3$ | No |
| 150 | 90 | ,, | No |
| ,, | ,, | ,, | ,, |
| ,, | ,, | ,, | ,, |

Feature Engineering can help build Better ML Models Efficiently

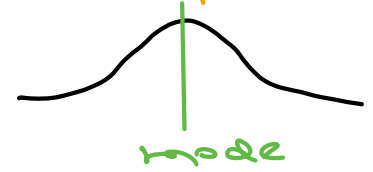Lets Dive Deeper with Loan Status Case study in Colab Notebook

## Skewness

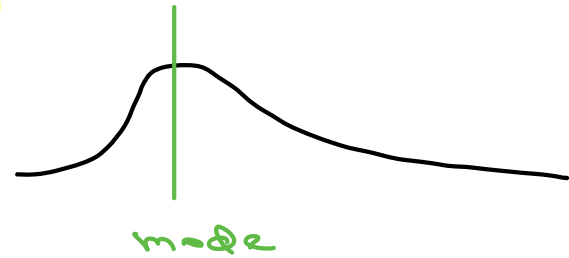Skewness is a measure of Asymmetry in the distribution of Data

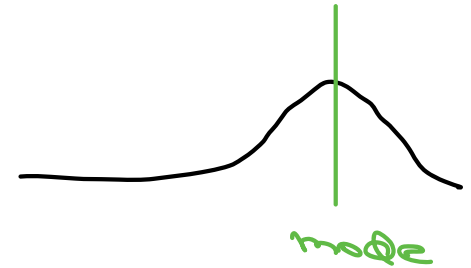⑤ Perfectly Balanced with even Spread on both sides

⑤ Mean ≈ mode ≈ Median



mode

⑤ Long Tail on Right Side

⑤ Bulk data on Left



mode

⑤ Long Tail on Left Side

⑤ Bulk data on Right



mode

---

The skewness ($g_1$) of a dataset can be calculated using the formula:

$$g_1 = \frac{n}{(n-1)(n-2)} \sum_{i=1}^{n} \left( \frac{x_i - \bar{x}}{s} \right)^3$$

Where:

- $n$ = Number of observations

- $x_i$ = Individual data point

- $\bar{x}$ = Mean of the data

- $s$ = Standard deviation of the data

Alternatively, **Pearson's moment coefficient of skewness** can be simplified as:

$$\text{Skewness} = \frac{\text{Mean} - \text{Median}}{\text{Standard Deviation}}$$

Skewness > 0     Right Skewed (Positive)
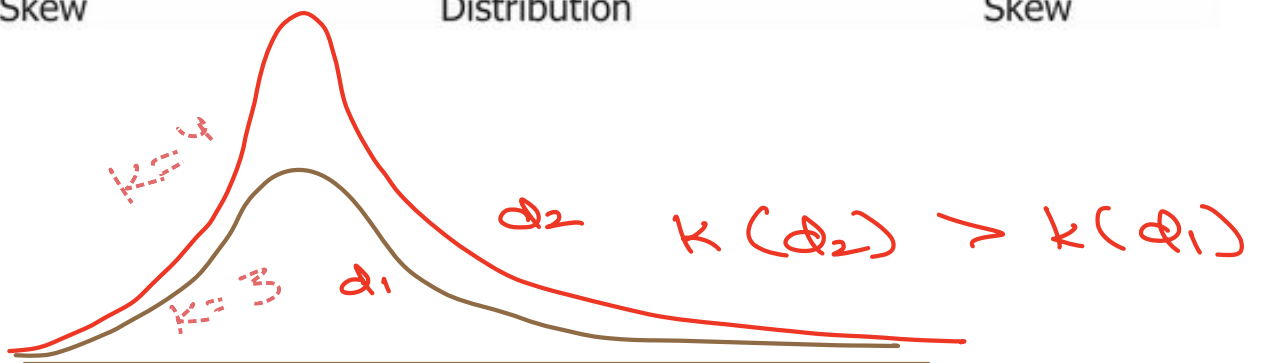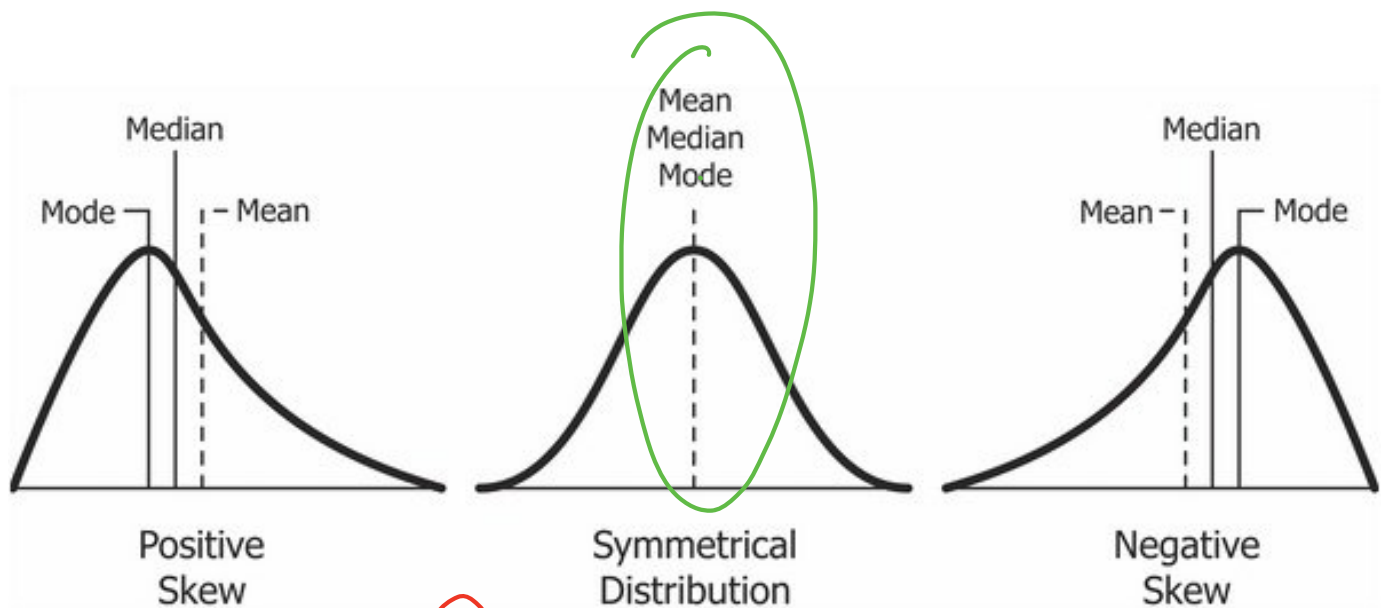
Skewness < 0     Left Skewed (Negative)

Skewness = 0     No Skew (Symmetric)

## Mean vs Median vs Mode



| Positive Skew | Symmetrical Distribution | Negative Skew |

$K = 4$     $K = 3$     $d_1$     $d_2$     $K(d_2) > K(d_1)$

# Kurtosis

Kurtosis measure "Sharpness of Peak" of Data Distribution

| High Kurtosis | Low Kurtosis |
|---|---|
| ⮌ High Peak | ⮌ Low Peak |
| ⮌ Heavy tails | ⮌ Light tails |
| ⮌ More Outliers | ⮌ Less Outliers |

**Excess Kurtosis** ⮌ Measures Kurtosis w.r.t Normal Distribution (kurt = 3)



$K > 3$
$E-K > 0$
Positive Kurtosis

Negative Kurtosis

Normal Distribution
$K = 3$

$K < 3$
$E-K < 0$

$$\text{Excess } k = \text{Kurt}_{dist} - 3$$

① Leptokurtic ($EK > 0$)

$$\boxed{\begin{array}{l} \partial_{i} \exists \ 3.1 \\ E\text{-}k \exists \ 0.1 \end{array}}$$

② Mesokurtic ($E\text{-}k \, \lessgtr 0$)

Approx Normal

③ Platykurtic ($E\text{-}k < 0$)

A dist with negative $E\text{-}k$

**What is the relationship between skewness and**
Excess **kurtosis in a normal distribution?**

4 options

**Active Duration**(Most preferred: 30 seconds)

| Appears for | 60 Secs | ⌄ |
|---|---|---|

Skew = 0
$E\text{-}k = 0$

$K = 3$
$E\text{-}k = K\text{-}3$
$3\text{-}3 = 0$

## New Feature

**Feature:** Able to Pay EMI

Consider Following two Loan applicants

① P1 earns 20L and applies for 50L Loan

② P2 earns 30L and applies for 5CR Loan

Who is more Likely to Get Loan?