# FEATURE ENGINEERING-2

# ANOVA

> 2 Categories

① Gaussian → QQ Plot
→ Shapiro
→ KS Test

Kolmogorov Smirnoff Test.
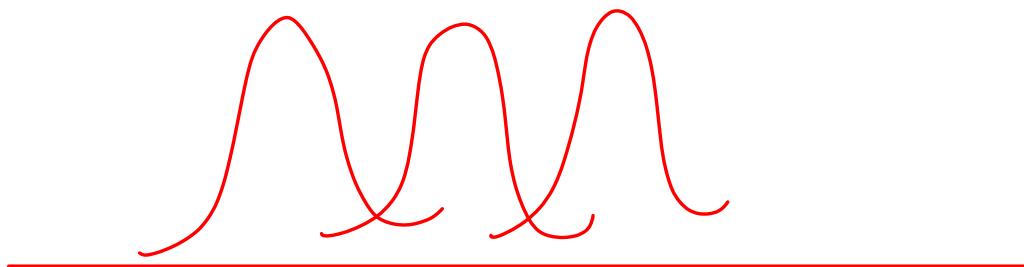
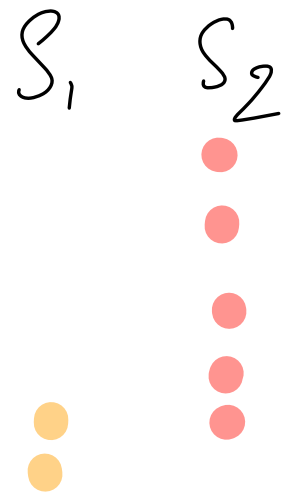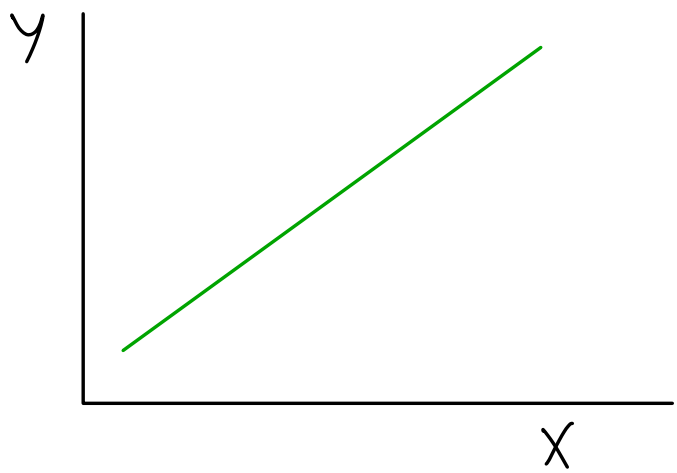② Independant

③ Equal Variances → Levene.

KRUSKAL'S TEST

$$f = \frac{MSB}{MSW}$$

Any distribution ↔ Gaussian distrib$^n$

% percentiles          % percentile

$1^{st} S_1$          $1^{st} S_2$

$2^n S_1$          $2^{nd} S_2$

$3^{rd} S_1$          $3^r S_3$

⋮          ⋮

$100^{th} S_1$          $100^{th} S_2$

# KS Test

## ttest $\longrightarrow$ $\mu_1$ $\mu_2$

$H_a = \mu_1 = \mu_2$

$H_a = \mu_1 \neq \mu_2$

pdf 1

$\mu_1$ $\mu_2$

pdf 2

$\mu_2$ $\mu_1$

cdf 1

cdf 2

$S_1$     $S_2$

$t_{tour}$   (Ha)

$x$   $M_1$     $M_2$    $x$

$x_2$

Random distribution
$\quad\quad\quad \hookrightarrow$ Gaussian $\quad \longrightarrow$ {QQ Plot

Gaussian

Not Gau

{ KS Test

{ Shapiro     pvalue

$H_0$: Gaussian

$H_a$: Not Gaussian.

Levene:

| | | $W/H^2$ | | |
|---|---|---|---|---|
| H | W | BMI | fitness | Risk of heart attack |
| ● | | $\overline{(X)}$ ● | | |
| ● | | | | |
| | ○ | | | |

Record / Data point

features

Target

feature Importance → feature Selection

mathematical func.

Garbage in

Garbage out

Prob

Credit History

| | |
|---|---|
| YES | 1 |
| NO | 0 |
| YES | 1 |
| NO | 0 |
| YES | 1 |
| YES | 1 |
| YES | 1 |

$$\frac{1 + 0 + 1 + 0 + 1 + 1 + 1}{7}$$

$$\frac{5}{7} \% \times 100$$

$0.84\%$ → Dataset

Credit History

M1　M2

$C_1$　$C_{12}$

$C_1$

$C_{12}$

$C_{10}$ $C_{11}$ $C_{12}$

1week

1week

$-3.1, -3, -2.9, -2.8 \cdots 0 \cdots$

$3$

$1^{st}$

$1^o L$

$1^o /.$ $\longleftarrow$ $-3$ /

\* Z test → $\mu, \sigma \checkmark$ (Numerical v/s Categories)

\* T test ( t test-1 samp, t test-rel, t test-ind, t test-ind-fromstats)
  $\sigma X$ → $n < 30$, (Numerical v/s Categories)

\* KS Test → $\mu_1 \approx \mu_2$, distributions are diff.
  (CATEGORICAL)

\* $X^2$ ⌐ Goodness of fit
       └ Test of Independance (CAT v/s CAT)

\* ANOVA → More than 2grps (Numerical)

\* KROSKAL → when assump^n of ANOVA FAIL

\* CORRELATION → (NUMERICAL v/s NUMERICAL)

Normality
 - Q-Q Plot
 - Shapiro
 - KS Test

 └ Variance
    └ Levene

# HYPOTHESIS TESTING CHEAT SHEET

## 1 CONCEPTS

### CENTRAL LIMIT THEOREM (CLT)

The **distribution of sample means is Gaussian**, no matter what the shape of the original distribution is.

**Assumptions:** population mean and standard deviation should be finite and sample size >=30

### HYPOTHESIS TESTING

- A method of statistical inference to decide whether the **data** at hand sufficiently **support** a **particular hypothesis**.
- A test statistic **directs us to either reject or not reject the null hypothesis.**

**Null Hypothesis (H₀)** represents the assumption that is made about the data sample whereas,
**Alternative Hypothesis (Hₐ)** represents a counterpoint.

### P-VALUE

Probability of observing the Test statistic as extreme or more than $T_{observed}$ considering the null hypothesis as true.

If **p-value < significance level; reject the null hypothesis,** else fail to reject the null hypothesis.
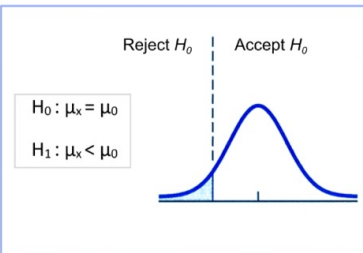
### CRITICAL VALUE

A **cut-off value** used to mark the **start of a region** where the **test statistic** is **unlikely** to fall in.
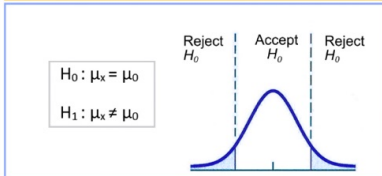
## 2 TYPES OF HYPOTHESIS TESTING

**Type I error (α) - Reject** a null hypothesis that is true.
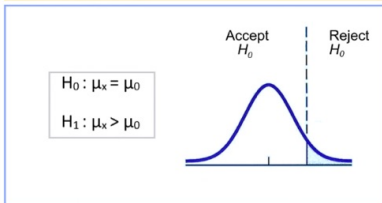**Type II error (β) - Not reject** a null hypothesis that is false.

### ONE TAILED - LEFT



$H_0 : \mu_x = \mu_0$

$H_1 : \mu_x < \mu_0$

### TWO TAILED



$H_0 : \mu_x = \mu_0$

$H_1 : \mu_x \neq \mu_0$

### ONE TAILED - RIGHT



$H_0 : \mu_x = \mu_0$

$H_1 : \mu_x > \mu_0$

### Framework for Hypothesis testing

1. Define the experiment and a sensible test statistic variable.
2. Define the null hypothesis and alternate hypothesis.
3. Decide a test statistic and a corresponding distribution.
4. Determine whether the test should be left-tailed, right-tailed, or two-tailed.
5. Determine the p-value.
6. Choose a significance level.
7. Accept or reject the null hypothesis by comparing the obtained p-value with the chosen significance level.

## 3 TESTS

### ONE TAILED - LEFT

- Used to determine whether the population mean is significantly different from an assumed value.
- It uses Standard normal distribution as the baseline.

**Assumptions:**
- either the **standard deviation of the population** should be **known** or,
- we should estimate them well when the **sample size is not too small (n>30).**

$$Test\ statistic = Z = \frac{\bar{x} - \mu_0}{\frac{\sigma}{\sqrt{n}}}$$

### Two sample Z-test

Used to **compare** the **means** of two populations.

**Assumptions:**
- either the **standard deviation (σ₁, σ₂)** of the populations should be **known**
- we should estimate them when the **sample sizes are not too small** ( n₁, n₂ ≥ 30).

$$Test\ statistic = t = \frac{(\bar{x}_1 - \bar{x}_2) - 0}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

### ONE SAMPLE T-TEST

The test statistic follows a **t - distribution** it is used when:
- the **sample size** is **too small (n < 30)** and/or,
- the **population standard deviation (σ)** is **unknown.**

$$Test\ statistic = z = \frac{\bar{x} - \mu_0}{\frac{s}{\sqrt{n}}}$$

### Two sample t-test

Used when
- the **sample sizes** are **too small ( n1 , n2 < 30)** and/or,
- the **population standard deviations (σ1, σ 2) are unknown.**

$$Test\ statistic = t = \frac{(\bar{x}_1 - \bar{x}_2) - 0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

### ANOVA (Analysis of variance)

- Used to determine if there is a **statistically significant difference** between two or more **categorical groups** by testing for **differences of means** using variance.
- The test statistic **f** follows the **f** distribution represented by two parameters **(k-1)** and **(n-k)**. k = no. of groups, n = total sample size.

$$Test\ statistic = \text{f} = \frac{MSB}{MSW}$$

where,
**MSB** = mean of the squared distances between the groups and **MSW** = the mean of the squared distances within the groups.

$$MSB = \frac{\sum_{i=1}^{k} n_i (\bar{X}_i - \bar{X})^2}{k-1} \qquad MSW = \frac{\sum_{i=1}^{k} \sum_{j=1}^{m} (X_i j - \bar{X}_i)^2}{n-k}$$

**Assumptions:**
- the variance of each group should be the same or close to each other.
- the total n observations should be independent of each other.

### KS (Kolmogorov - Smirnov) test

- A **non - parametric test** used for determining whether the **distributions** of two samples are the **same or not**
- The test statistic $T_{ks}$ follows a distribution called the kolmogorov distribution
- $T_{KS}$ = the maximum absolute value of the difference in the CDFs of the two samples X and Y

## 4 CORRELATION

Degree of the mutual relationship between two variables

### PEARSON CORRELATION COEFFICIENT(PCC)

$$\rho_{xy} = \frac{Cov(X,Y)}{\sigma_x . \sigma_y}$$

**Limitation of PCC** is that it only **captures the linear relationship** between the variables. It fails to capture the non-linear patterns.

### SPEARMAN RANK CORRELATION COEFFICIENT

A statistical measure of the strength of a **monotonic relationship** between paired data. it captures the monotonicity of the variables rather than the linearity.

$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

where,
**d** =difference between the two ranks of each observation and,
**n** = number of observations

281　　481　　781