

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 (H_0) represents the assumption that is made about the data sample whereas, **Alternative Hypothesis (H_a)** 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

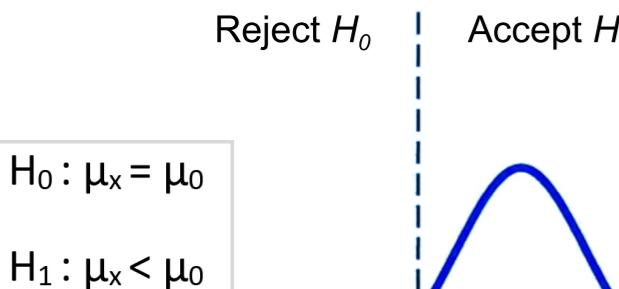
A cut-off value used to mark the start of a region where the test statistic is unlikely to fall in.

2 TYPES OF HYPOTHESIS TESTING

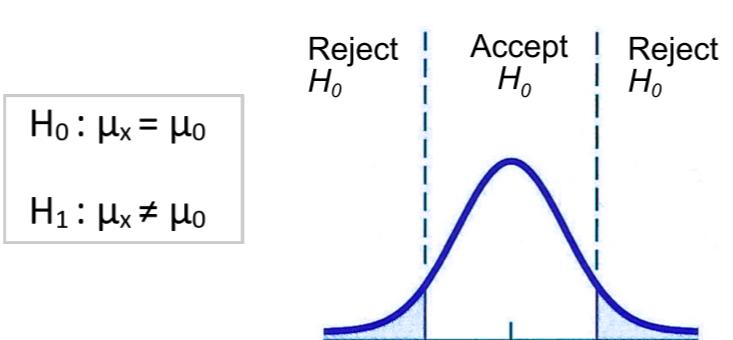
Type I error (α) - Reject a null hypothesis that is true.

Type II error (β) - Not reject a null hypothesis that is false.

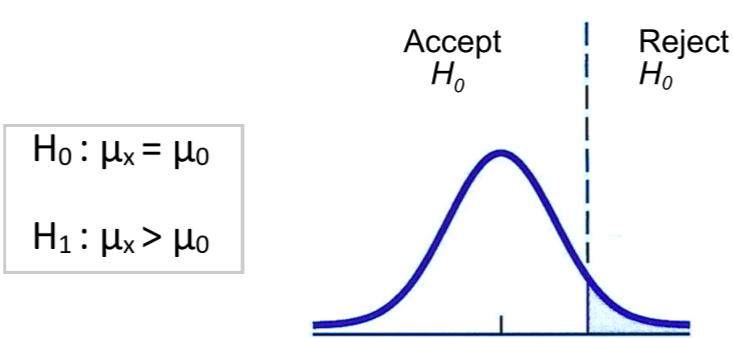
ONE TAILED - LEFT



TWO TAILED



ONE TAILED - RIGHT



FRAMEWORK FOR HYPOTHESIS TESTING

- Define the experiment and a sensible test statistic variable.
- Define the null hypothesis and alternate hypothesis.
- Decide a test statistic and a corresponding distribution.
- Determine whether the test should be left-tailed, right-tailed, or two-tailed.
- Determine the p-value.
- Choose a significance level.
- 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 ($n > 30$).

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

TWO SAMPLE Z-TEST

Used to compare the means of two populations.

Assumptions:

- either the standard deviation (σ_1, σ_2) of the populations should be known
- we should estimate them when the sample sizes are not too small ($n_1, n_2 \geq 30$).

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

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.

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

TWO SAMPLE T-TEST

Used when

- the sample sizes are too small ($n_1, n_2 < 30$) and/or,
- the population standard deviations (σ_1, σ_2) are unknown.

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

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.

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

where,

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

$$MSB = \frac{\sum_{i=1}^k n_i (\bar{X}_i - \bar{X})^2}{k-1}$$

$$MSW = \frac{\sum_{i=1}^k \sum_{j=1}^m (X_{ij} - \bar{X}_i)^2}{n-k}$$

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)

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

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