

STATS

- to get useful information from the data.

2 types

DESCRIPTIVE | INFERENTIAL

↑
- To describe data in hand

↑
To infer something from data for a larger population.

- Depending on no. of values resulted in mode → Unimodal, bimodal, uniform, multimodal.



* Variance vs SD. → Square root of variance.

↓
Average of squared differences from mean.

Squaring so that +ve -ve values won't cancel out.

$$\sigma^2 = \sum (x_i - \bar{x})^2 / N$$

Both SD and variance indicates how much data is spreaded across w.r.t Mean.

* SD's unit same as the underlying data.

Conditional Probability.

AND → *

OR → +

Eg:- Picking a house in Street1 given we already picked a house in same Street. ($P_1 * P_2$)

DISCRETE PROBABILITY DISTRIBUTIONS.

- Binomial Distributions → Eg:- Coin toss.
- Poisson Distributions → Eg:- A Particular cell within a large population will acquire a mutation.
- Geometric Distributions
- Negative Binomial Distributions

CONDITIONS:

- 1) Total. no. of trials fixed.
- 2) Each trial is binary
- 3) Prob. of success is same in all trials.

FORMULA

$$nC_r(P)^r(1-P)^{n-r}$$

Poisson - To describe distribution of rare events in a large population.
- In another way, it shows how many times an event is likely to occur within a specified period of time.

$$e^{-\lambda} \lambda^x / x!$$

Areas of Descriptive Stats

* Measures of central tendency

- mean (EV), median, mode, quantiles.

* Measures of dispersion/variation - SD, variation, range.

EV - Expected Value - same as average

(For Casinos, portfolio managing strategies for investment)

PMF vs PDF



Prob. Mass Function

Distribution Func.

Eg:- Bar charts / Pie charts / Tally charts

For discrete random variables

For continuous random variables

Normal Distribution / Gaussian

- distribution that is symmetric about the mean.

1-2-3 rule. about SD of mean.
68% 95% 99.7%

Sample - Subset of population.

Sampling Distribution - Distribution of means of different samples.

Central Limit Theorem (CLT)

- Sampling distribution mean = Population mean.
- For $n > 30$, Sampling distribution will be normally distributed.
- Confidence Interval - Range of possible values (CI) for an unknown parameter.
↓
An associated confidence level with this (that indicates certainty).

$$CI = \left(\bar{x} - \frac{Z^* S}{\sqrt{n}}, \bar{x} + \frac{Z^* S}{\sqrt{n}} \right)$$

\bar{x} → sample mean

S → sample SD

n → sample size

Z^* → Z-score associated with Confidence Level.

NOTE: Z-score is the value associated with a particular cumulative probability and Z^* is Z-score associated with confidence level.

Sampling Useful in:

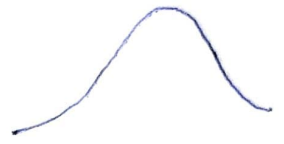
- Quality Control
- Pilot Testing
- Market Research
- Market Campaign Efficacy.

Skewness → measure of symmetry

Kurtosis → measure of whether data heavily-tailed or lightly-tailed.



+ve skew
(Right skewed Distribution)



-ve skew
(Left skewed Distribution)

STANDARD $-1.96 \leq S_K \leq +1.96$.

To find out skewness

- Univariate plots.
- Shapiro-Wilk's test
(from `scipy.stats` import `shapiro`)
- `.skew()` method from `pandas`.

Handling Skewness

- * Square-Root Transformation
- * Log Transformation
- * Reciprocal Transformation
- * Box-Cox transformation
- * Yeo-Johnson

Sampling Types.

Random Sampling (RS)

* Simple RS with replacement

* Simple RS without replacement

* Stratified RS

* Cluster Sampling

* Systematic Sampling

Non-Random

* Convenience

* Quota

* Judgemental

* Snowball.

HYPOTHESIS TESTING - method of statistical inference

Null Hypothesis (H_0): status-quo / default position / no relation b/n variables.

Alternate Hypothesis (H_1): alternate to Null.

So 2 cases:
- Rejection of Null Hypothesis
Failure to Reject Null Hypothesis.

NOTE: - H_0 assumed to be true and statistical evidence required to reject it in favour of Alternate Hypo.

Type of Errors:

	Truth about population		
	H_0 true	H_a true	
Decision based on sample.	Reject H_0	Type-1 Error	Correct Decision
	Accept H_0	Correct Decision	Type-2 Error

Type-1 error \rightarrow Rejection of True Null Hypothesis.

Type-2 error \rightarrow fail to reject a False Null Hypothesis.

In terms of Confusion Matrix,

Type-1 \equiv False Positives

Type-2 \equiv False Negatives.

Methods that supplement Critical Hypothesis:

- CRITICAL VALUE METHOD \rightarrow * Find the Z-critical with the significance level (α)
 - P-VALUE METHOD.
- \downarrow
- * Find upper & lower critical values with Z-critical with $\mu \pm (Z_c * \sigma_x)$

p-value is the probability of obtaining extreme - test result, assuming Null hypothesis is correct.

* If p-value less than α , then reject Null Hypothesis.

T-tests: \rightarrow If sample size less than 30 and population SD/variance not known

Different types:

- 1-sample t-test \rightarrow `scipy.stats.ttest_1samp()`
- 2-sample t-test \rightarrow `scipy.stats.ttest_ind()`
- Paired t-test \rightarrow `stats.ttest_rel()`

ANOVA

(Analysis of Variance)

- To check statistical similarity of 2 or more groups.
- measured using F-Ratio.

- Other Forms of ANOVA \rightarrow MANOVA, ANCOVA.

CHI-SQUARE TESTS.

- Used when we have single categorical variable.

- 2 types: Test of Independence
Goodness of Fit Test.