

Probability & Statistics

1) Random Variable → A variable that can take any value from the sample space.

or

Can take an outcome as its value from the set of all the possible outcomes available.

2) Sample Space → A set of all the possible outcomes of a random experiment.

3) Outlier → An extreme observation.

Population & Sample Size

↳ All the subjects for which an estimate need to be wanted for.

e.g. estimate the height of all the humans on earth.

μ → population mean.

→ It is often not possible to estimate the statistic over all the subjects in the population.

Sample → A subset of a population that can be used to measure a statistic.

e.g. Estimate the height of a few people in each state of a country.

\bar{x} → Sample mean.

→ As the sample size increases, sample mean comes closer to population mean.

Types of Events

1) Mutually Exclusive & Events that can't happen at the same time.

ex: Sachin scoring a century & Sachin getting duck.

2) Joint Events & Events that can happen together.

ex: Sachin scoring 150 & India winning the match.

3) Independent Events & Occurrence of one event don't affect the occurrence of other.

ex: outcome of Sachin's performance & weather conditions.

4) Exhaustive Events & Set of events that cover all the possible outcomes in the sample space.

ex: Sachin scoring 0 to 50.

Sachin scoring 50 to 100.

Sachin scored 150 & India won.

Events of Set operations & 1) Intersection & Both events happen.

fd. union(A, B)
how: union)

2) union & Either one of the events happen.

fd. union(Cevents, eventB)). drop-duplicates()

ex: scored 150 or India won.

3) Complement & Event don't happen.

Probability Rules

1) Addition Rule & $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

2) Multiplication Rule & $P(A \cap B) = P(A|B) P(B)$

3) Conditional Probability & $P(A|B) = \frac{P(A \cap B)}{P(B)}$

Gaussian Distribution & (normal distribution)

1) Distribution → Very simple model of nature.

Given → 1/3 Type of distribution → One can easily
by mean of variance model the probability
of a property frequency of occurrence
of each possible random variable of that property.

2) Why Gaussian Distribution is important?

→ Most of the naturally occurring phenomena follow gaussian distribution.

ex- Height & weight of people.

3) Mathematical representation →

$$P(X = x) = \frac{1}{\sqrt{2\pi}} \cdot e^{\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)}$$

↑ quadratic power.
↑ exponential fn.

Simplified → standard normal

$$\mu = 0, \sigma^2 = 1$$

$$y = P(x) = \frac{1}{\sqrt{2\pi}} \cdot e^{\left(-\frac{x^2}{2}\right)}$$

remove constants & shape of fn don't change when constants are removed.

$$y = e^{-x^2} = \frac{1}{e^{x^2}}$$

↑ exponential with symmetry

quadratic power.

→ Due to quadratic exponentiation if x increases or decreases $\frac{1}{e^{x^2}}$ always decreases.
if x changes twice \rightarrow y reduces by 100 times.

$$\text{ex- } P(x=1) = \frac{1}{e} = 0.3678$$

$$P(x=2) = \frac{1}{e^4} = 0.018$$

$$P(x=3) = \frac{1}{e^9} = 0.000123$$

μ & σ^2 \rightarrow Parameters of a gaussian dist.

\hookrightarrow only these two are required to get the prob. density fn of a distribution.

\rightarrow Maximum probability is always at mean of the distribution.
 \rightarrow Mean = mode.

CDF \rightarrow Cumulative density fn.

Pdf is used when we want $P(X = x)$.

CDF is used when we want $P(X \leq x)$ or $P(X \geq x)$ $\text{at a specific point}$

\hookrightarrow Data points in a range.

\rightarrow Cdf of normal distribution is S shaped.

lower σ^2 \rightarrow lighter shape.
higher σ^2 \rightarrow looser shape.

Properties related to shape of a distribution -

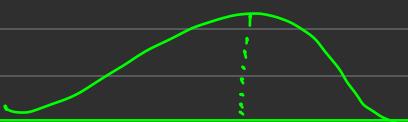
1) Symmetry \rightarrow A distribution is symmetric if $\&$ a point \bar{x} about which the # of points left of $\bar{x} =$ # of points right of \bar{x} .

mathematically $\rightarrow f(x_0 - x_i) = f(x_0 + x_i)$

\rightarrow Gaussian distribution is symmetric.

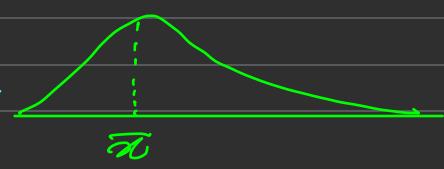
2) Skewness \rightarrow Tailsness of a distribution.

\hookrightarrow Left (-ve) skewed \rightarrow longer tail towards left side of mean.



\hookrightarrow Right (+ve) skewed \rightarrow longer tail towards right side of mean.

mathematically $\rightarrow \frac{\sqrt{n} \sum (x_i - \bar{x})^3}{(\sqrt{n} \sum (x_i - \bar{x})^2)^{3/2}}$



3) Kurtosis \rightarrow Another measure for flatness of a distribution.
 \rightarrow Used to detect outliers.

mathematically \rightarrow

$$\text{Kurtosis}(\alpha) = \frac{\sqrt{n} \sum (\alpha_i - \bar{\alpha})^4}{\left(\sqrt{n} \sum (\alpha_i - \bar{\alpha})^2 \right)^2}$$

$$\text{Gauss Kurtosis}(\alpha) = \text{Kurtosis}(\alpha) - 3$$

(\Rightarrow relative measure.

\Rightarrow Kurtosis of Gaussian distribution.

Standard Normal \rightarrow

$$x_{\text{std}} = \frac{\alpha_i - \mu}{\sigma} \quad \forall i \in \{1, \dots, n\}$$

Distribution where $\mu = 0$ & $\sigma = 1$.

- Description of the distribution is clearer than normal distribution.

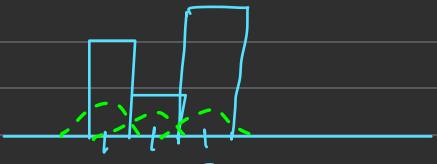
68% of data within (-1 + 1) std. deviation.

95% of data within (-2 + 2) std. deviation.

99% of data within (-3 + 3) std. deviation.

Kernel Density Estimate \rightarrow Convert histograms to probability density.

ex \rightarrow



\rightarrow Draw a normal dist. around each point.

\rightarrow Distributions overlap for rounded points & are mutually exclusive for non-rounded points.

2) Kde(α) = \sum values of dist. overlapping at α .

Value \nearrow higher probability \rightarrow Points are rounded.
 \searrow low probability \rightarrow Less rounded points.

Sampling Distribution & Central Limit Theorem

Given $\rightarrow X \sim N(\mu, \sigma^2)$ population
 (e.g. Variance
 population-mean)

1) Collect n samples of size n , S_1, S_2, \dots, S_m from the population.

S_1, S_2, \dots, S_m total collection = $m \times n \approx 1/30000$.

2) Get the mean of each sample.

$S_1 \rightarrow \bar{x}_1, S_2 \rightarrow \bar{x}_2, \dots, S_m \rightarrow \bar{x}_m$

3) Get the distribution of sample means \rightarrow

$Dist(\bar{x}_m) =$ Sampling distributions of sample means.

Central Limit Theorem $\rightarrow Dist(\bar{x}_m) \approx N(\mu, \sigma^2/m)$ or \approx

The Sampling distributions of sample means is normally distributed around the population mean with σ^2/m variance.

Population (mean, variance) $\xrightarrow{\text{can be computed by}} \text{Sampling distribution of sample means.}$

of samples from the population.

Quantile - Quantile Plot \rightarrow Q-Q plot

why - To check if a random variable is normally distributed or not.

How - Let X be a random variable we want to check if

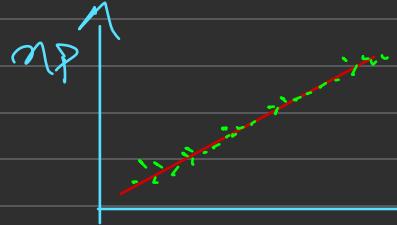
1) Take n observations of X .

2) Divide the n observations into percentiles. $\rightarrow x_{P_i}$

1) Take another r.v $N(0, 1)$ obtain 100 observations.

2) Divide them into percentiles. $\rightarrow Y_P$

plot the each percentile x_i in X_P & y_i in Y_P



If points lie on a straight line
 $X \sim Y$ Normally distributed.
else
 $X \not\sim$ Not Gaussian.