# Probability and Statistics

Rolling a dice — random experiment.

Its outcome as a random variable $X = \{1, 2, 3, 4, 5, 6\}$.

Similarly, tossing a coin. — r.e.

r.v. $Y = \{H, T\}$.

$$P(X=1) = \frac{1}{6} = P(X=2)$$

$$P(X = \text{even}) = \frac{1}{2} = P(X=2) + P(X=4) + P(X=6).$$

$P(X = x_i) \rightsquigarrow$ sometimes also written as $\underline{P(x_i)}$.
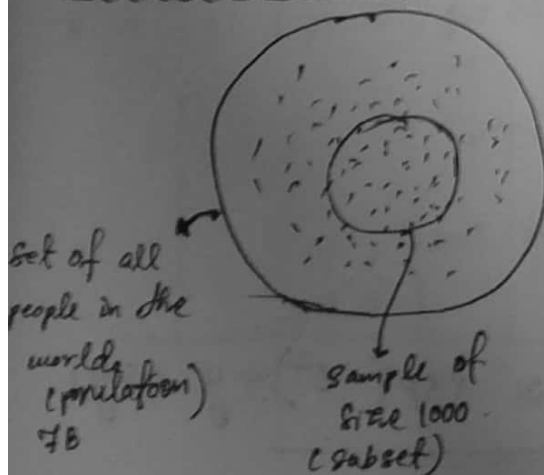
$X$ is a discrete random variable.

let $X$ is a random variable of storing outcome of random expt of measuring heights.

So, $X =$ continuous random variable.

$$X = \{122.2, \ 146.4, \ 132.5, \ \cdots \ \underset{\substack{\downarrow \\ \text{outlier (maybe error or} \\ \text{genuine)}}}{(12.26)}, \ 156.23 \}.$$

---

## Population & Sample



Set of all people in the world (population) ⇄ B

sample of size 1000 (subset)

The mean height of human $(\mu) = \dfrac{1}{7B} \displaystyle\sum_{i=1}^{7B} h_i$

denotes population mean

$\bar{h} = \dfrac{1}{1000} \displaystyle\sum_{i=1}^{1000} h_i$ (sample mean).
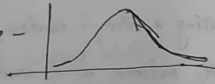
As sample size increases,

$$\boxed{\bar{x} = \mu}$$

Gaussian distribution (Normal distb$^n$).

PDF of a ~~random~~ random variable having :–
gaussian distribution

Graphical:

X : continuous random variable.

why bother learning about this particular kind of distribut$^n$?

→ They occur in nature. Ex→ the sepal length & petal length of Iris
flowers,

Similarly, heights & weights of peoples also have a gauss. distb$^n$

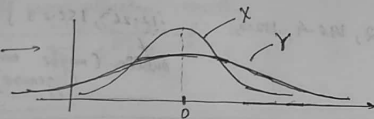wikis– normal–distb$^n$:
↳ graph with diff $\mu$ & $\sigma$.

($\mu$ & $\sigma^2$) are parameters of Gaussian distb$^n$.

Knowing only these 2 values enough. (No need of sample datas)

$X \sim N(\mu, \sigma) \Rightarrow$ X as a random variable that follows normal
distb$^n$ with mean $\mu$ & std deviat$^n$ $\sigma$.

Ex $X \sim N(0,2)$.

$Y \sim N(0,4)$

Mathematical:
$X \sim N(\mu, \sigma^2)$.

$$P(X=x) = \frac{1}{\sqrt{2\pi}\,\sigma} \exp\left\{ -\frac{(x-\mu)^2}{2\sigma^2} \right\}$$
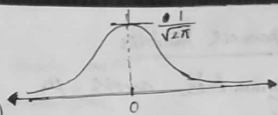
simplify to analyze, $\sigma=1$, $\mu=0$.

$$P(x) = \frac{1}{\sqrt{2\pi}} \exp\left\{ -\frac{1}{2}x^2 \right\} \qquad P(x) \approx \exp(-x^2) = y \;(say)$$
$\underset{c_1}{\underbrace{\phantom{xxx}}}$  $\underset{c_2}{\underbrace{\phantom{xxx}}}$
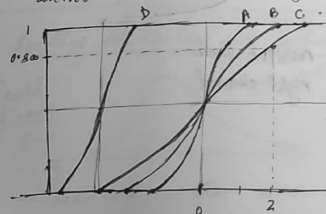
$y = \exp(-x^2)$

① symmetric
$\exp(-x^2) = \exp(-(-x)^2)$

② as $x$ moves away from $\mu$, $y$ reduces exponentially.

CDF of gaussian distb$^n$

wikis– normal distb$^n$ page.

A: $\mu=0$ $\sigma^2 = 0.5$
B: $\mu=0$ $\sigma^2 = 1$
C: $\mu=0$ $\sigma^2 = 2$
D: $\mu=-4$ $\sigma^2 = 0.4$.

$P(X \leq 2) = 0.8$ for C from the CDF.

× Since gaussian dist. symmetric about mean, in CDF, 0.5 for $\mu$.

$$P(X \leq \mu) = 0.5$$

68 – 95 – 99 rule
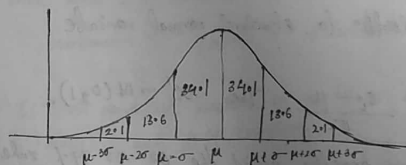
× 50% of pts lie on either side of mean

✱ In range $[\mu-\sigma, \mu+\sigma]$, 68% of pts lie.
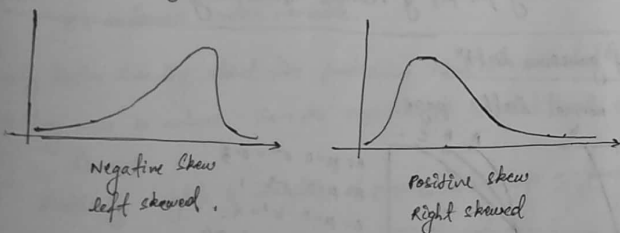
y " " $[\mu-2\sigma, \mu+2\sigma]$, 95% " " "

y " " $[\mu-3\sigma, \mu+3\sigma]$, 99% " " ".

34.1  34.1
2.1  13.6     13.6  2.1
$\mu-3\sigma$  $\mu-2\sigma$  $\mu-\sigma$  $\mu$  $\mu+\sigma$  $\mu+2\sigma$  $\mu+3\sigma$

## Symmetric, Skewness, Kurtosis

Gaussian is symmetric across $\mu$

Measure of asymmetricness — skewness



Negative Skew
left skewed.

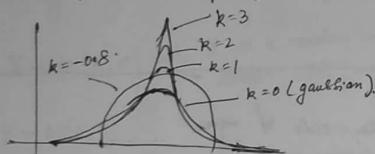Positive skew
Right skewed

$\mu$ — location where central value lies

$\sigma$ — spread.

skewness — how asymmetric is distrib$^n$.

kurtosis — peakness (sharpness). See wikipedia's kurtosis for graph.



$k = -0.8$
$k = 3$
$k = 2$
$k = 1$
$k = 0$ (gaussian).

## Standard Normal Variate

$Z \sim N(0,1)$  i.e. a random variable having gaussian distribution
with mean = 0 & std-dev$^n$ = 1.

Can convert any gaussian distrib$^n$ to standard normal variate
using standardizat$^n$.

$X \sim N(\mu, \sigma)$.   $x_i' = \dfrac{x_i - \mu}{\sigma}$   then  $x_i' \sim N(0,1)$.

Advantage of this :— know exactly about the $68-95-99.7$ rule.
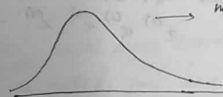& many more properties applicable.

## Kernel density estimation

↳ converts histogram to density curve.
wiki-pages: kernel-density-estimation.

for each point in the sample space, a gaussian kernel is drawn (with
& values at each point is summed up to get total height. the point
being mean)

variance — (bandwidth) — simulated & selected a suitable
value @ which gives a smooth curve.

## Sampling distribution and Central Limit theorem



distb
of incomes.
$X \to$
→ not necessarily gaussians

→ let's say we take $m$ random samples, each of size $n$ (say 30).
$s_1, s_2 \ldots \ldots \ldots s_m$ ($m$ samples).

$\bar{x}_1, \bar{x}_2, \ldots \ldots \bar{x}_m$ are the mean of the $m$ samples.

$\bar{x}_i \sim$ distb$^n$. called the "sampling distb$^n$ of sample mean".

CLT:  If $X$ (original population) has finite $\mu$ & $\sigma^2$, then

( Pareto distb$^n$ has
infinite mean,
var).

$$\bar{x}_i \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right) \quad \text{as } n \to \infty.$$

i.e. the $\bar{x}_i$ distb$^n$ of the means of the samples is a Gaussian distb$^n$
with mean $= \mu$ (mean of original population) &

$$\text{variance} = \frac{\sigma^2}{n} \quad \text{(where } \sigma^2 \text{ is variance of original pop}^n\text{).}$$

This CLT is powerful because it works on data having any kind of
distb$^n$, not just Gaussians

# Quantile-Quantile Plot

Given $X: x_1, x_2, x_3, \cdots \cdots x_{500}$

Q.? Is X Gaussian distb? → QQ plot (graphical)
                          → Statistical testing (KS, AD).

## QQ

① Sort $x_i$'s & computer percentiles.

$$x_1, x_2, \cdots \quad x_{500}$$

↓ sort (asc)

$x_1', x_2', x_3', \cdots \quad x_{500}' \xrightarrow{\text{percentiles}} \begin{array}{ccccc} 1 & 2 & 3 & 4 & 100 \\ x_5', x_{10}', x_{15}', x_{20}', \cdots x_{500}' \\ x^{(1)} \ x^{(2)} \ x^{(3)} \ x^{(4)} \ \cdots x^{(100)} \end{array}$
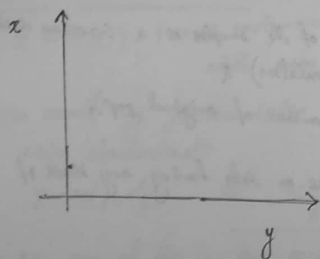
$x^{(i)}$ — $i^{th}$ percentile value of $x_i$'s

② $Y \sim N(0,1)$ — std normal distb".

$y_1, y_2, y_3 \cdots \quad y_{1000}$ — 1000 obsv" fm N(0,1).

↓ sort (asc)

$y_1', y_2', \cdots \quad y_{1000}' \xrightarrow{\text{percentiles}} y^{(1)}, y^{(2)}, \cdots y^{(100)}$.

③ plot Q-Q plot using $x^{(1)}, x^{(2)}, \cdots x^{(100)}$
                                $y^{(1)}, y^{(2)}, \cdots y^{(100)}$.

$\left( y^{(i)}, x^{(i)} \right)$ → $\forall i \in [1,100]$
                              ordered pairs
                              plotted

If all these points roughly lie on a straight line, then we can say x & y have the same distb". Hence $x \in N(\mu, \sigma)$.

---

but we can't conclude that X also has mean=0 & variance =1.

$\underset{x}{\text{Mean}(x)}$ = value corresponding to $y=0$ in plot.

Refer ipynb .

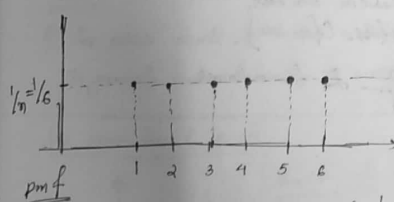Q-Q plot :generally answers—

Q. If X & Y come from same distribution?

---

## Uniform distribution
→ discrete.
→ continuous

① discrete (wiki-page: discrete-uniform-distb")
Pdf is called pmf here (probability mass function).

Ex— throwing a dice.
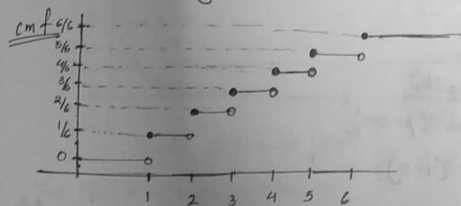All are equiprobable. That's what equiprobable. uniform distb" means.



pmf

Parameters — a & b.
$a \in \mathbb{Z}$
$b \in \mathbb{Z}$, $b \geq a$.
$n = b - a + 1$.

parameters completely describe about the distb".



cmf

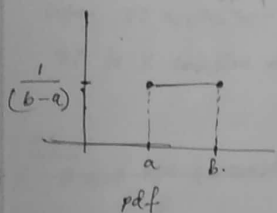See wiki page for

mean :— $\left( \frac{a+b}{2} \right)$

median : $\frac{a+b}{2}$

variance: $\frac{(b-a+1)^2 - 1}{12}$

skewness :— 0.

③ continuous

parameters:- $a$ & $b$
$$a, b \in R \quad \& \quad b \geq a.$$



pdf.

Application of uniform distb$^n$

 ↳ Random no. generator (uniform distb$^n$).
 
 see rough for random sampling example.
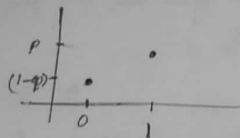
Bernoulli Distribution          wiki

 ↳ discrete distb$^n$.
 
 has 2 outcomes.  Probability → $P$ & $(1-P)$.
 
 $X$ → random variable of getting head in coin toss.
 
 $X \sim$ Bernoulli($p=0.5$) ––→ coin toss. (fair coin).

$$pmf \begin{cases} (1-P) & \text{for } k=0 \\ P & \text{for } k=1 \end{cases} \text{— the two outcomes.}$$



Binomial distb$^n$          wiki.

 coin-toss: $X \sim$ Bernoulli($p=0.05$).
 
 coin tossed n-times (n=10).
 
 e.g. probability of success
   Let $Y$ be a random variable denoting no of heads in 10 coin
      tosses.

$Y \sim Bin(n, p)$                    2 parameters.
   no. of coin — prob. of getting heads
   toss

In general the parameters of binomial distb$^n$ are
 ① no. of trials
 ② probability of success.

$$pmf: \binom{n}{k} p^k (1-p)^{n-k} \quad \text{— } P(Y=k)$$

Log Normal Distribution          wiki

 $X \sim$ log-normal ($\mu, \sigma$) if $\log(X) \sim$ normal distb$^n$.
 
   wiki—for graph of pdf. & cdf.

obsv$^n$: as $\sigma^2 \uparrow$, the pdf curves become more skewed.

Applications— ① length of comments posted in internet discussions.
  wiki      follow log normal distb$^n$.

② The users dwell on online articles follows a log-normal distb$^n$.
  In general, human behaviour is mostly log-normal distb$^n$.

See power law & pareto



Pareto distribution graph

This kind of rel$^n$ b/w variables is
called power law relationship in maths

$x_m = 5$ for all
  $(x_m \& \alpha)$ — parameters.

As $\alpha \to \infty$, distb$^n$ approaches
  $\delta(x - x_m)$ where $\delta$ is dirac
  delta function.

(as $\alpha \downarrow$, tails fatness $\uparrow$)

— dirac delta fun$^n$

Occurrence in nature

① file size distb^n in internet traffic.
② hard disk drive error rates.
③ value of oil reserves in oil fields.

Q)^p How to check if a distb^n is power pareto?

→ Draw log-log graph. — straight line

y (probability value)   x (feature value)   both's pb log taken & plotted

Of course we can always use Q-Q plot.

Pareto distribution to Gaussian (box-cox transform^n).

Pareto: ~ X [ $x_1, x_2, x_3, \cdots, x_n$].
Gaussian: Y [ $y_1, y_2, y_3, \cdots, y_n$].

① box-cox(X) = lambda ($\lambda$)
   └ all n observat's.

Some complex
math done.

② $y_i = \begin{cases} \dfrac{x_i^{\lambda} - 1}{\lambda}, & \text{if } \lambda \neq 0. \\ \ln(x_i), & \lambda = 0. \end{cases}$

$\forall i \in [1, n]$.

※ $\lambda = 0$ means X is a log-normal distb^n.

Y = scipy. stats. boxcox (X). ↙

Quantifying rel ship among features

→ Let say 2 features height and weight.
Are they co-related? i.e. h↑ ⟶ w↑/w↓.   & vice-versa?!

3 measures
① co-variance
② pearson    co-relation coeff
③ spearman      "       " .

① co-variance

variance(X) = $\dfrac{1}{n} \sum\limits_{i=1}^{n} (x_i - \mu)(x_i - \mu)$.

co-variance(X, Y) = $\dfrac{1}{n} \sum\limits_{i=1}^{n} (x_i - \mu_x)(y_i - \mu_y)$.

So, covariance(X, X) = variance of (X).

Now, co-variance(X, Y) = $\begin{cases} (+)\text{ve. implies } X\uparrow \text{ then } Y\uparrow \\ (-)\text{ve implies } X\uparrow \text{ then } Y\downarrow. \end{cases}$

Intuition



|  | $(x_i - \mu_x)$ | $(y_i - \mu_y)$ |  |
|---|---|---|---|
| $(x_1, y_1)$ | (+) | (+) | = (+) |
| $(x_2, y_2)$ | (−) | (−) | = (+) |

cov(X, Y) = (+)ve { summat^n over all such pts}

|  |  |  |  |
|---|---|---|---|
| $(x_1, y_1)$ | (+) | (−) | = (−) |
| $(x_2, y_2)$ | (−) | (+) | = (−). |

cov(X, Y) - (−)ve.

But we don't have any idea of increase or decrease.

One major drawback of covariances- co-var(X, Y) $\neq$ co-var(X, Y)
cm×kg                    ft×lbs

Pearson takes care of that.

② pearson correlation coeff [refer wiki].

$\sigma_x$ = std-devn of $x$.

$$\rho_{x,y} = \frac{covar(X,Y)}{\sigma_x \; \sigma_y}$$

See graphs of diff value of $\rho$ in wiki.

$$-1 \le \rho_0 \le +1$$



$\rho = 1$    $\rho = -1$    $0 < \rho < 1$

$-1 < \rho < 0$    $\rho = 0$

limitations

① $\leftarrow$ only +1 when linear rel^n ship b/w $x$ & $y$.
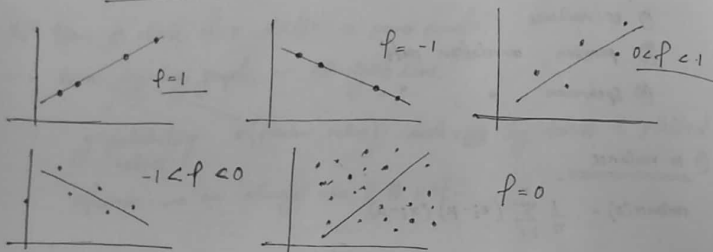   So, if $y = x^2$, $\rho < 1$. (even though monotonical increase).

② slope of st. line doesn't affect $\rho$.

③ complex relat^n ship not captured. (see wiki page). $\rightarrow$

$\rho = 0$ (beautiful sine rel^n neglected)

fix - using spearman coeff.

Spearman rank - coeff
   $\leftarrow$ sort the feature values & find the position in sorted order & find pearson coeff of these positions.

| | X | Y | $r_x$ | $r_y$ | |
|---|---|---|---|---|---|
| ex $\rightarrow$ | | | | | $\rho_{r_x, r_y}$ is Spearman coeff. |
| $S_1$ | 160 | 52 | 4 | 3 | |
| $S_2$ | 150 | 66 | 2 | 4 | So, if $X\uparrow$ then $Y\uparrow$, (doesn't matter |
| $S_3$ | 170 | 68 | 5 | 5 | linear rel^n) then Spearman coeff = 1. |
| $S_4$ | 140 | 46 | 1 | 1 | |
| $S_5$ | 158 | 51 | 3 | 2 | Similarly -1 if $X\uparrow \rightarrow Y\downarrow$ |

---

$\times$ Spearman more robust to outliers than pearson.

Co-relation Vs Causation

$\times$ Just because two random variables are co-related ($X\uparrow$ then $Y\uparrow$) doesn't mean X causes Y or vice versa.

$\times$ causal models - advanced statistics used for those purposes.

---

Confidence Interval

disb : X : heights.
(any disb^n) $\{x_1, x_2, \ldots x_{10}\}$ - random sample of size 10.

ex Estimate the populat^n mean of $X = \mu$.

$$\mu \approx \bar{x} \;(\text{sample mean}).$$

— This is point estimate.
Not bad, but we can do better.

If we say $\mu \in [162.1, 174.9]$ with 95%. — Interval with some confidence value. - Richer than previous in terms of informat^n.

Now, how to calculate C.I. ??

$\rightarrow$ ① Computing C.I. given the underlying distribution

say $X \sim N(\mu, \sigma)$. Let $\mu = 168$ cm
(heights)    $\sigma = 5$ cm

fm knowledge of gaussian distb^n,
$(\mu - 2\sigma, \mu + 2\sigma)$ contains 95% of my observat^n.



2.5%    2.5%
$\mu-2\sigma$  $\mu-\sigma$  168 cm $\mu+\sigma$  $\mu+2\sigma$
158 . 163    173  178

So, we can say heights of people lie b/w [158, 178] with 95% probability.
$\leftarrow$ C (confidence).
Similarly other values like 90%, 80% can be
found using Normal distb tables

**C.I. for mean($\mu$) of r.v.**

$X \sim$ some distb$^n$ with mean $= \mu$ & std-dev $= \sigma$. — population parameters

$\{x_1, x_2, \cdots x_{10}\}$ — we are given only this sample of X

$\S^{\circ}$ what's the C.I. of 95% C.I. of $\mu$? (given this sample)

**Case I :** — we are given the std-dev$^n$ ($\sigma$) of population.

for **CLT**, if $\bar{x}$ is the sample mean then (analogy to petal length on iris dataset)

say in mind for $\bar{x}$

then $\bar{x} \sim N\left(\mu, \dfrac{\sigma}{\sqrt{n}}\right)$.

So, we can say $\bar{x} \in \left[\mu - \dfrac{2\sigma}{\sqrt{n}}, \mu + \dfrac{2\sigma}{\sqrt{n}}\right]$ with 95% probability.

Or, i.e. $\bar{x}$'s value is b/w $\mu - \dfrac{2\sigma}{\sqrt{n}}$ & $\mu + \dfrac{2\sigma}{\sqrt{n}}$ with a 95% chance

Or $\mu - \dfrac{2\sigma}{\sqrt{n}} \leq \bar{x} \leq \mu + \dfrac{2\sigma}{\sqrt{n}}$

$-\dfrac{2\sigma}{\sqrt{n}} \leq \bar{x} - \mu \leq \dfrac{2\sigma}{\sqrt{n}}$

$-x - \dfrac{2\sigma}{\sqrt{n}} \leq -\mu \leq \dfrac{2\sigma}{\sqrt{n}} - x$.

$\bar{x} + \dfrac{2\sigma}{\sqrt{n}} \geq \mu \geq \bar{x} - \dfrac{2\sigma}{\sqrt{n}}$ Or $\mu \in \left[\bar{x} - \dfrac{2\sigma}{\sqrt{n}}, \bar{x} + \dfrac{2\sigma}{\sqrt{n}}\right]$

with 95% chance.

So, we have calculated C.I. for $\mu$.

---

**Case 2 :** we don't know $\sigma$ (pop$^n$-std-dev)

use student's t-distb$^n$

$\bar{x} \sim t(n-1)$ ← degrees of freedom.

$\bar{x}$ follows t-distb$^n$.

what about C.I. of other statistical measures like median or 90$^{th}$ percentile?

$\hookrightarrow$ **bootstrap C.I.** — using higher computational powers to do simulations.

**C.I. using empirical bootstrap**

task :- estimate 95% CI for median of X.

using only the given sample of X.

$\rightarrow$ Generate k samples, each of size m ($m \leq n$).
← (sampling with repetition)

find medians of each one of them.

So, now we'll have k-medians.

Say $k = 1000$.

Now, sort the 1000 medians.

$m_1, m_2, \ldots m_{25}, \ldots m_{975}, m_{\ldots}, m_{1000}$ ← sorted order.

$\boxed{\text{C.I. (95%) } \in [m_{25}, m_{975}]}$

$\dfrac{950}{1000} = 95\%$.

Similarly other params can also be calculated.

Also, the larger the value of n, the narrower will be interval.

# Hypothesis Testing

"As the name suggests, we test for the touchness of an assumed hypothesis based on the observations we've got while experimenting.

Ex 1:-

Task:- Given a coin determine if the coin is biased towards head/not.
$\begin{cases} \text{biased towards head:- } P(H) > 0.5 \\ \text{not biased towards head:- } P(H) = 0.5 \end{cases}$

design expt: flip coin 5 times & count no of heads = X. → Test statistic.

perform expt: f, f, f, f, f    $X = 5$ — observation fm performing expt.
$\quad\quad\quad\quad$ H, H, H, H, H.

$P(X=5 \mid \underbrace{\text{coin is not biased}}_{\text{towards head}}) = P(obs \mid H_0) = \frac{1}{2^5} \simeq 0.03 = 3\%.$
obsr$^n$ $\quad\quad$ an assumption.
$\quad\quad\quad\quad$ └ called Null Hypothesis (H_0)

$H_0$: coin is not biased towards head

So, $P(X=5 \mid H_0) = 3\%$.
└ There is a 3% chance of getting 5 heads in 5 flips if the coin is not biased towards head.

probability of observation given assumption is 3%, quite low.

Since the observation is done practically, i.e. the ground truth. Hence, our assumption may be wrong.

* $P(obs \mid assumption)$ — also called p-value.
$\quad\quad\quad$ Typically p-value < 5% is said to be small

Hence $H_0$ may be incorrect ⟹ we reject our null hypothesis.
$\quad\quad\quad\quad\quad\quad\quad$ the idea that coin isn't biased.

---

we accept the fact that coin is biased towards heads.

$H_0$: coin not biased. — Null hypothesis

$H_1$: coin biased towards head — Alternate hypothesis.

$\quad$ rejecting $H_0$ ⟹ accepting $H_1$
$\quad$ rejecting $H_1$ ⟹ accepting $H_0$

In the expt:- flip was done 5 times → sample size
$\quad$ if we had flipped 3 times →
$\quad\quad$ f, f, f $\quad$ $X=3$ $\quad$ $P(X=3 \mid H_0) = \frac{1}{2^3} = 12.5\% > 5\%$
$\quad\quad$ H, H, H

hence $H_0$ acan't be rejected ⟹ accepted.
$\quad$ So, sample size matters.

* p-value = 3%
$\quad$ ↳ implies $P(obs \mid H_0)$ ✓ not $P(H_0)$ ✗✗

---

Ex 2

Task: determine if the population mean of heights of people in these two cities is same or not.

→ Its impossible to calculate population mean so, we use sample mean.

expt: Measure height of 50 random people fm each city.
$\quad$ let $\mu_1$ & $\mu_2$ be sample means of both cities. say (162 & 167)
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ $\mu_1$ $\quad$ $\mu_2$

test statistic:- $\mu_2 - \mu_1 = 167 - 162 = 5cm.$ (X).

Null hypothesis:- there is no difference in population mean of both the cities.

compute: $P(X=5 \mid H_0)$
$\quad\quad\quad\quad$ ↳ diff is sample mean with sample size of 50

$P(X=5 \mid H_0)$

↳ probability of observing a diff of 5cm in sample mean height of sample size 50 between $c_1$ & $c_2$ if there is no population mean difference.

case 1: $P(X=5 \mid H_0) = 0.2 = 20\%$

There is a 20% chance of observing a diff of 5cm in sample mean of height & width sample size of 50 if there is no population mean difference (of $c_1$ & $c_2$) so we accept the $H_0$.

And vice-versa

Computing $P(X=5 \mid H_0)$ (Resampling).

① Take all heights of both cities & put them together to a new set. (S)

② Randomly select 50 pts fm S to $S_1$ & remaining 50 to $S_2$. This is called resampling.

* The idea of resampling is that now since $S_1$ & $S_2$ are coming fm same distrib" randomly, this will simulate the two cities having same population means or simulate the null hypothesis.

calcute $\mu_1$ & $\mu_2$ and $\mu_2 - \mu_1 = \delta$

③ repeat the 2nd step k no of times.

④ sort the $S_i$'s. in inc. order.

case 1: $\underbrace{S_1 \le S_2 \le \ldots \ldots \le S_k^*}_{\text{simulated differences.}}$

Observed difference = 5cm

$P(\text{diff} \gtrsim 5cm \mid H_0) = \underbrace{0.2}_{\text{significant}}$ ✓

& vice versa

{ Say k = 1000.
& out observed diff = $S_{801}$.
So, 20% of sim. diff greater than obs. diff. }

---

KS-Test for similarity of two distribution (refer wiki)

let $X_1$ & $X_2$ be the two samples. of size m & n.

let $D_{m,n}$ be the max$^m$ diff in their CDFs.

If $\boxed{D_{m,n} > c(\alpha) \sqrt{\frac{(m+n)}{mn}}}$ ✓ then the belong to diff distrib", else same distrib".

$\alpha$ & $c(\alpha)$ values are taken fm table.

scipy.stats.kstest()

---

## Dimensionality Reduction.

why? → to visualize high dimension data.

how? → PCA & t-SNE.

* By default, a vector is column vector.

$x_i \in \mathbb{R}^d$ — column vector $d \times 1$

> Representing dataset

$D = \{ x_i, y_i \}_{i=1}^{n} \leftarrow$ data pt.
$x_i \in \mathbb{R}^d$; $y_i \in \{$ setosa, virginica, versicolor $\}$

D is usually represented in 2 matrices — X & Y.

X — rows are data pts & cols are features.

Y — rows are data pts usually 1 col. only.

---

Data pre-processing: Column normalization

why? — so that the data becomes nicely structured so that data modelling algos will perform well.

normalizat" — for every feature, we get

$\boxed{a_i' = \frac{a_i - a_{min}}{a_{max} - a_{min}}}$    $a_i' \in [0, 1]$