

Random or stochastic process :-

→ in essence there processes whose outcomes are random within a fixed set of or range of values.

Probability:-

→ can be defined as the measure of certainty or uncertainty that a certain event or outcome will occur given a certain stochastic or random process.

population and sample :-

→ population refers to whole or complete data.
→ A collection of some randomly selected data from population that means a subset of population and which can serve as a good approximation of the population's characteristics and closely represent population is called sample.

pop probability distribution :-

→ refers to how the data is distributed among the set or range of values it contains.

Random variables :-

→ Random variables are variables that represent the collection of outcomes of a stochastic process.

Discrete & continuous random variables:-

(2)

→ Random variables that represent discrete & non continuous data are called Discrete random variable.

→ whereas those which represent continuous data are represent are called Continuous random variable.

Statistics:-

→ science of techniques & methodologies that are used for the collection, presentation and analysis of quantitative data for decision making.

→ Descriptive statistics (summarizing data)
→ Inferential statistics (drawing conclusions)

* Measures of centrality

→ is a number that summarises the typical value or "tendency" of a random variable. These measures give the general indication as to where most values in a distribution fall.

MEAN → average value of random variable

MEDIAN → Middle value or average of middle two values in "SORTED" Random variable.

MODE → Value with the Maximum frequency of occurrence in random variable.

Note:- data with no repetition have no mode.

MEASURES OF DISPERSION :-

→ Measure of dispersion is a number or a range that summarises the "SCATTER" or "SPREAD" of a random variable with "RESPECT TO ITS MEAN".

MEAN ABSOLUTE DEVIATION:-

→ average of the absolute deviation of all the values of a random variable with respect to mean.

$$\text{Mean Absolute Deviation} = \frac{1}{n} \sum_{i=1}^n |x_i - \text{mean}(x)|$$

VARIANCE:-

→ average squared distance of the values of a R.V with respect to its mean. (squared to eliminate negative values)

$$\text{Variance} = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

standard deviation:-

→ square root of the variance. it modifies the variance to provide measure of dispersion that is of same units as that of R.V.

Median absolute deviation:-

$$\text{Median} (\text{abs}(\mu - x_i)) \text{ for all } i$$

$$= \text{Median} (|x_i - \text{mean}|)_{i=1}^n$$

IQR - Inter quartile range

$\rightarrow Q_3 - Q_1 \rightarrow$ gives mid 50% of dataset.

$\rightarrow Q_2$ - median

parameters:-

\rightarrow The statistical concepts of centrality, dispersion and shape when used in context to a population are called parameters.

statistics:-

\rightarrow concepts of centrality, dispersion & shape when used in context of sample are called statistics.

probability functions:-

\rightarrow mathematical characterizations of a probability distribution is called probability function. gives distribution type and its parameters we can completely define the probability distribution of a r.v.

probability Mass functions:-

\rightarrow probability functions of discrete random variables are called PMF. They mathematically define the probability of occurrence of each value or state within a discrete random variable.

Example: Binomial distribution.

5

Probability Density function:-

→ PDF used to express the frequency distributions of continuous r.v. maps the values of r.v to its probabilities.

- probability densities can have values greater than 1.
- area under the curve defined by PDF represents the probability.
- probabilities of continuous r.v expressed w.r.t ranges $P(x_1 \leq X \leq x_2)$.
- PDF is non-negative everywhere and its integral over the entire space is equal to 1.

CUMULATIVE DENSITY FUNCTION:-

→ gives us the probability of any value in the random variable X being less than or equal to a particular value x .

$$CDF(x) = P(X \leq x)$$

→ CDF for both continuous & discrete r.v.

KERNEL DENSITY ESTIMATION:-

→ KDE is a statistical tool that allows us to create a smooth curve that approximates the PDF of a continuous r.v.

(i) replace each datapoint with PDF of normal distribution whose mean is equal to value of data point & standard deviation as function of bandwidth chosen.

- Those PDF are called Kernels.
- Band width is the range of values on either side of data point which we want to consider.
- PDF of the dataset as a whole is approximated by summing up the values of individual kernels at suitably chosen intervals.
- The smoothness depends on size of the bandwidths chosen. Larger bandwidths - less responsive to the sparsity of the data.

PROBABILITY DISTRIBUTIONS:-

THE BERNoulli DISTRIBUTION:-

- Discrete distribution, used to represent or model stochastic processes that can have only two possible outcomes, success or failure
- P - probability of success
- 1-P - probability of failure
- Bernoulli is the basic building block for other discrete distributions like binomial, hypergeometric & poisson.

$$\rightarrow \underline{\text{Bernoulli}(P) \sim X}$$

mean $\rightarrow P$

Variance $\rightarrow P(1-P)$

$$P \text{ estimation} = \frac{n(1)}{n(0) + n(1)}$$

The Binomial distribution:-

→ Binomial distribution is extension of the Bernoulli distribution. That is a binomial distribution or trial is the sequence of Bernoulli events.

→ $X \sim B(n, p)$

- ↑ number of trials or events per outcome
- ↓ probability of success.

$$\text{PMF} \rightarrow P(X=k) = P(k) = {}^n C_k p^k (1-p)^{n-k}$$

$$\text{CDF} \rightarrow P(X \leq k) = \sum_{i=0}^k {}^n C_k p^i (1-p)^{n-i}$$

$$E(X) = \text{mean} = np$$

$$\text{Var}(X) = npq = np(1-p)$$

THE HYPERGEOMETRIC DISTRIBUTIONS:-

→ similar to binomial distribution, but not same, it is picking without replacement, where binomial is picking with replacement. In balls experiment if the # of balls is large relative to # of draws the distribution seem very similar to binomial distribution because the chance of success changes less.

→ $X \sim H(N, K, n)$

- ↑ population size
- ↑ # of trials per outcome
- ↓ number of success states in population

The poission distribution :-

- deals with frequency with which an event occurs in a specific interval of time, given that the probability of the event occurring is constant.
- $\lambda \rightarrow$ rate parameter
- $x \sim po(\lambda)$ expected or average frequency of success
- can be used to model situations where the general rate of occurrence of event is known and when we want to know the probability of same event occurring at the different rate of occurrence.

$$PMF = P(X=k) = P(k) = \frac{\lambda^k e^{-\lambda}}{k!}$$

$$CDF = P(X \leq k) = F(k) = e^{-\lambda} \sum_{i=0}^{[k]} \frac{\lambda^i}{i!}$$

$$E(X) = \text{mean} = \lambda$$

$$\text{Var}(X) = \lambda$$

THE EXPONENTIAL DISTRIBUTION:-

- exponential closely related to poission, while poission dist represents or models # of successful events per unit time, then exponential dist models the time period b/w 2 consecutive successful events., that follows a poission dist.

$$x \sim E(\beta) \rightarrow \beta = \lambda$$

$$\text{PDF} = f(n) = \begin{cases} \lambda e^{-\lambda n} & \text{if } n \geq 0 \\ 0 & \text{if } n < 0 \end{cases}$$

$E(n) = \text{mean}$

$$= \beta = \lambda$$

$$\text{CDF} = F(n) = \begin{cases} 1 - e^{-\lambda n} & \text{if } n \geq 0 \\ 0 & \text{if } n < 0 \end{cases}$$

$$\text{Var}(n) = \beta^2 = \lambda^2$$

→ memoryless distribution $\rightarrow P(\text{event occur in 40 min} | \text{waited 30 min}) = P(\text{event occur in 10 min})$.

NORMAL OR GAUSSIAN DISTRIBUTION:-

→ used to represent r.v whose outcomes are continuous and exhibit 'centrality'. The concentration of the frequency distribution is maximum near the central or mean value, and reduces symmetrically as we move away from it.

$\xrightarrow{\mu}$ mean - location parameter
 $\xrightarrow{\sigma}$ standard deviation - scale parameter

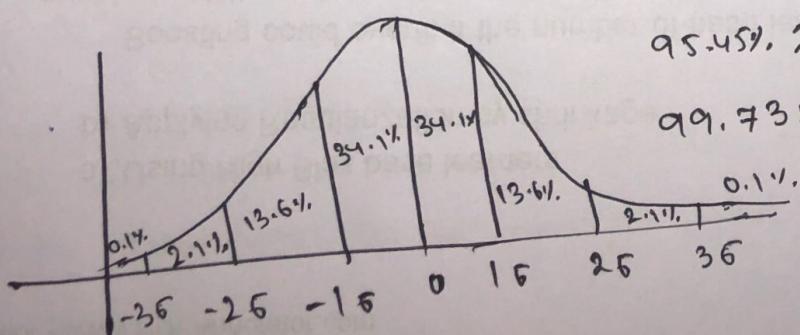
$$\text{PDF} = f(n) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{n-\mu}{\sigma} \right)^2}$$

$$E(n) = \text{mean} = \mu$$

$$\text{Var}(n) \approx \sigma^2$$

$$\text{CDF} = F(n) = P(X \leq n) = \int_{-\infty}^n f(n) dn. \quad S.D = \sigma$$

68 - 95 - 99.7 rule:-



$$68.27\% \approx \Pr(\mu - \sigma \leq X \leq \mu + \sigma)$$

$$95.45\% \approx \Pr(\mu - 2\sigma \leq X \leq \mu + 2\sigma)$$

$$99.73\% \approx \Pr(\mu - 3\sigma \leq X \leq \mu + 3\sigma)$$

The log normal distribution:-

- If $y'x'$ is said to be lognormally distributed when log of 'x' is normally distributed.
- Log normal distribution can only contain the real values, and asymmetric around mean and have heavy tails.

$$\ln(x) \sim N(\mu, \sigma)$$

Transformation of sampled data to gaussian:-

- Normal distribution is regarded as ideal and it is often assumed by many statistical methods. Hence it is often the practice that data is transformed to normal whenever possible.

→ Commonly used transformation

(i) The Box-Cox transform

(ii) The Yeo-Johnson transform

(iii) The Quantile transform.

EXPLORATORY DATA ANALYSIS:-

- ↳ UNIVARIATE → single variable analysis
- ↳ BI-VARIATE → two variable analysis
- ↳ MULTIVARIATE → multiple variable analysis.

$$\textcircled{1} \text{ MEAN} = \frac{1}{n} \sum_{i=1}^n m_i$$

$$\textcircled{2} \text{ Variance} = \frac{1}{n} \sum_{i=1}^n (m_i - \bar{m})^2 \quad \text{std-dev} = \sqrt{\text{Variance}}$$

$$\textcircled{3} \text{ MEDIAN} = \begin{cases} \text{ODD} & \frac{n+1}{2} \\ \text{EVEN} & \text{Mean } \left(\frac{n}{2} \text{ & } \frac{n+1}{2} \right) \end{cases}$$

$$\textcircled{4} \text{ PERCENTILE} = \frac{\text{rank of } n - \# \text{ values below } n}{n} \times 100$$

Quantiles $\rightarrow 25^{\text{th}}, 50^{\text{th}}, 75^{\text{th}}, 100^{\text{th}}$

$$\downarrow \quad \text{percentile value} = \frac{\text{median}}{100}$$

$$IQR = Q_3 - Q_1$$

$$\downarrow \quad 5.25 \rightarrow 5^{\text{th}} \text{ val} + (0.25)(6^{\text{th}} - 5^{\text{th}})$$

$[Q_1 - (1.5 * IQR), Q_3 + (1.5 * IQR)] \rightarrow \text{points lying outside are called outliers.}$

$$\textcircled{5} \text{ MAD} = \text{MEDIAN} \left(|m_i - \text{median}| \right)_{i=1}^n$$

$$\textcircled{6} \text{ AXIOMS} - \begin{aligned} & (i) 0 \leq P(E) \leq 1 \quad (ii) P(S) = 1 \quad (iii) E_1, E_2, E_3 \dots E_n \text{ are} \\ & \text{mutually exclusive} \quad P(E_1 \cup E_2 \cup E_3 \dots E_n) = \sum_{i=1}^n P(E_i) \end{aligned}$$

$$* P(E^c) = 1 - P(E), * P(E \cup E^c) = P(S) = 1$$

$\textcircled{7} \text{ PRINCIPAL OF INCLUSION \& EXCLUSION:}$

$$P(P_1 \cup P_2 \cup P_3 \dots P_n) = \sum_{i=1}^n P_i - \sum_{i=1}^n (P_i \cap P_j) + \sum_{j=2}^n (P_i \cap P_j \cap P_k) \dots (-1)^{n-1} \prod_{i=1}^n P_i$$

* if $E \subseteq F$ then $P(E) \leq P(F)$

$$\dots (-1)^{n-1} \prod_{i=1}^n P_i$$

$$\text{Conditional Probability} = \frac{P(E|F) = P(E \cap F)}{P(F)} \quad P(F) \neq 0 \quad (2)$$

$${}^n C_r = \frac{n!}{r!(n-r)!}, \quad {}^n P_r = \frac{n!}{(n-r)!}$$

MULTIPLICATION THEOREM:

$$P(E_1, E_2, \dots, E_n) = P(E_1) P(E_2 | E_1) P(E_3 | E_2 \cap E_1) \dots P(E_n | E_1 \cap E_2 \cap \dots \cap E_{n-1})$$

MATCHING PROBLEM - NO ONE picked hat - 0.3678

$$\text{Prob of exactly } k \text{ from } n \text{ picked} = {}^n C_k \sum_{i=0}^{n-k} \frac{(-1)^i}{i!} \frac{(n-k)!}{(n-i)!}$$

INDEPENDENT EVENTS - $P(E|F) = P(E) \quad \& \quad P(E \cap F) = P(E)P(F)$

MUTUALLY EXCLUSIVE - $E \cap F = \emptyset, \quad P(E \cap F) = 0$

$E \cup F$ independent - $E \cap F^c$ also independent

* infinite sequence trials, INDEPENDENT TO EACH OTHER

'p' success PROBABILITY

ATLEAST ONE SUCCESS - $1 - (1-p)^n$

EXACTLY k SUCCESS OUT OF n - ${}^n C_k p^k (1-p)^{n-k}$

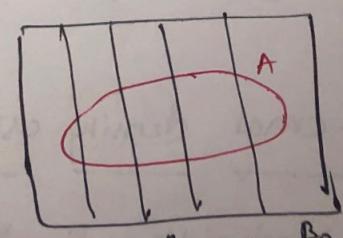
ALL n SUCCESS - p^n

* LAW OF TOTAL PROBABILITY :-

Mutually exclusive $B_1 \cap B_2 \cap \dots \cap B_n = \emptyset$

Mutually exhaustive $B_1 \cup B_2 \cup \dots \cup B_n = S$

$$P(A) = \sum_{i=1}^n P(A|B_i) P(B_i)$$



$$\text{BAYES THEOREM} - P(A|B) = \frac{P(B|A) P(A)}{P(B)} = \frac{P(B|A) P(A)}{P(B|A) P(A) + P(B|A^c) P(A^c)} \quad (3)$$

ODD IN FAVOUR OF A - $\frac{P(A)}{P(A^c)} = \frac{P(A)}{1-P(A)}$

PROBABILITY MASS FUNCTION - $P(X=i) = C \cdot \frac{1}{i!}, C = e^{-1}$

CDF - $F_X(a) = \sum_{n \leq a} P(X)$

(i) $0 \leq P(X=a) \leq 1$ (ii) $F(-\infty) = 0$ & $F(\infty) = 1$
 $a \geq b$ thus $F(a) \geq F(b)$

EXPECTATION - $E(n) = \sum_n n \cdot P(n)$ ↑ PMF $\left| \begin{array}{l} \int n \cdot f(n) dn \\ \downarrow \text{PDF} \end{array} \right.$

VARIANCE - $\text{Var}(n) = E[(X-\mu)^2] = E[X^2] - \mu^2$

$E[X+Y] = E[X] + E[Y], E[X \cdot Y] = E[X] E[Y]$

BERNOULLI DISTRIBUTION:- $X \in \{0, 1\}$ Discrete distribution

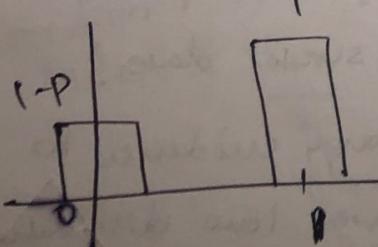
p - success prob

$q = (1-p)$ - failure prob

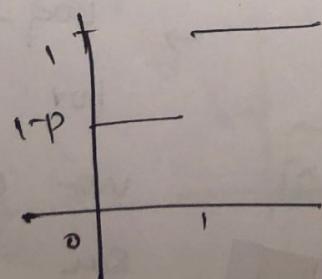
$X \sim \text{Bernoulli}(p)$

MEAN - p

PMF:-



CDF:-



Variance - pq

P estimation - $\frac{n(1)}{n(0)+n(1)}$

$n(0)+n(1)$

BINOMIAL DISTRIBUTION:- $X \sim (n, p)$ → probability of success \uparrow
 ↳ # of trials

$$\text{PMF} = {}^n C_k p^k q^{(n-k)}$$

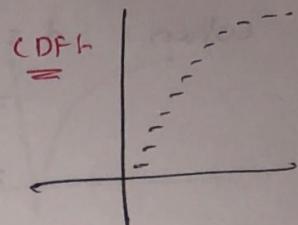
$$\text{CDF} = \sum_{i=0}^k {}^n C_i p^i (1-p)^{(n-i)}$$

$$\text{MEAN} = np$$

$$\text{Variance} = npq = np(1-p)$$

$$\text{Pestimation} = \frac{\sum_{i=1}^m n_i}{n \times m}$$

discrete curves



POISSON DISTRIBUTION:- $X \sim \text{poisson}(\lambda) \rightarrow$ rate → Discrete

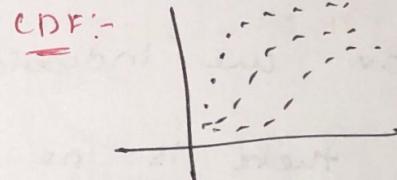
$$\text{PMF} = \frac{\lambda^K e^{-\lambda}}{K!} = P(X=k)$$

$$\text{CDF} = P(X \leq k) = \sum_{n=0}^k e^{-\lambda} \cdot \frac{\lambda^n}{n!}$$

$$\text{MEAN} = \lambda$$

$$\text{PMF} =$$

$$\text{Variance} = \lambda$$



UNIFORM DISTRIBUTION:-

$$\text{PDF} = \begin{cases} \frac{1}{b-a} & \text{if } x \in [a, b] \\ 0 & \text{otherwise} \end{cases}$$

$$\text{CDF} = \begin{cases} 0 & \text{if } x < a \\ \frac{x-a}{b-a} & \text{if } x \in [a, b] \\ 1 & \text{if } x \geq b \end{cases}$$

median

$$\text{Mean} = \frac{a+b}{2}$$

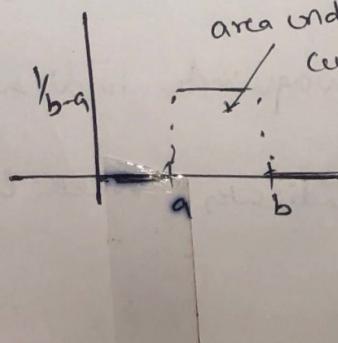
$$\text{Variance} = \frac{(b-a)^2}{12}$$

ESTIMATION of $a \& b$:-

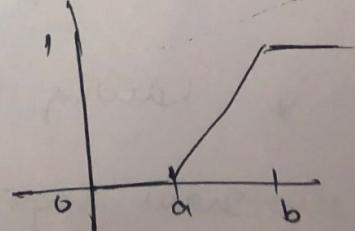
$a \rightarrow$ minimum of all the values

$b \rightarrow$ maximum of all the values

PDF



CDF



EXponential DISTRIBUTION:- $X \sim \text{EXPO}(\lambda)$

(3)

$$\text{PDF} - f(n) = \begin{cases} \lambda e^{-\lambda n} & n \geq 0 \\ 0 & n < 0 \end{cases} \quad \underline{\text{PDF}}$$

$$\text{CDF} - F(n) = \begin{cases} 1 - e^{-\lambda n} & n \geq 0 \\ 0 & n < 0 \end{cases}$$

MEAN - γ_1

Variance - γ_{12}

ESTIMATION - $\frac{n}{\sum_{i=1}^n n_i}$

* MEMORYLESS DISTRIBUTION

→ Exponential

→ Geometrical

Normal distribution: $X \sim N(\mu, \sigma^2)$ * Binomial is normal for large n , and for P not too close to 0 or 1

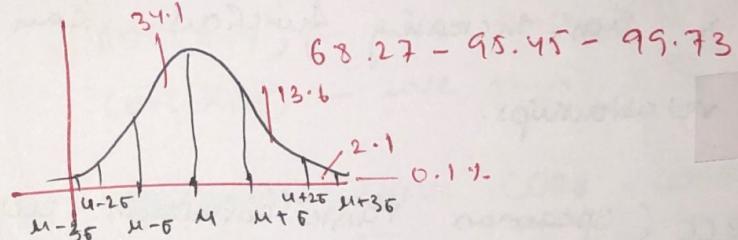
$$\text{PDF} - \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left\{-\frac{(n-\mu)^2}{2\sigma^2}\right\}, \sigma \neq 0$$

$$\text{CDF} - \gamma_2 \left[1 + \operatorname{erf}\left(\frac{n-\mu}{\sigma\sqrt{2}}\right) \right]$$

$$\text{Mean} - \frac{\sum_{i=1}^n n_i}{n}$$

$$\text{Variance} - \frac{\sum_{i=1}^n (n_i - \bar{n})^2}{n}$$

$$X \sim N(\mu, \sigma^2) \rightarrow Z = \frac{n-\mu}{\sigma} \sim N(0, 1)$$



SKEWNESS:-

$$b_1 = \frac{\gamma_3}{\gamma_2^{3/2}} = \frac{\frac{1}{n} \sum_{i=1}^n (n_i - \bar{n})^3}{\left[\frac{1}{n} \sum_{i=1}^n (n_i - \bar{n})^2 \right]^{3/2}}$$

$$= \frac{E[X^3] - 3\mu\sigma^2 - \mu^3}{\sigma^3}$$

* normal distribution's skewness = 0.

KURTOSIS:- outliers, univariate normal distribution = 3

$$\text{Excess kurtosis} = \frac{\frac{1}{n} \sum_{i=1}^n (n_i - \bar{n})^4}{\left[\frac{1}{n} \sum_{i=1}^n (n_i - \bar{n})^2 \right]^2} - 3$$

LOWER EXCESS KURTOSIS : LOWER OUTLIER | Lighter tails

HIGHER EXCESS KURTOSIS : higher outliers | heavier tails

CLT — $X \sim N(\mu, \sigma^2) \rightarrow \bar{X} \sim N(\mu, \sigma^2/n)$ as $n \rightarrow \infty$ (6)

CHEBYSHEV'S INEQUALITY — $P(|X-\mu| \geq k\sigma) \leq 1/k^2$

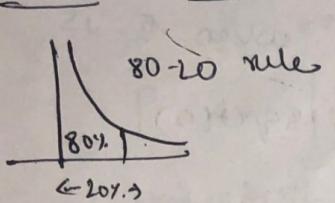
$$P[\mu - k\sigma < X < \mu + k\sigma] \geq 1 - 1/k^2$$

* valid only for finite mean & finite variance

LOG NORMAL DISTRIBUTION:- $X \sim \text{log-normal}(\mu, \sigma^2)$

R.V takes only +ve values $Y = \log_e(X) \sim N(\mu', \sigma'^2)$

POWER LAW (PARETO) DIST:- $X \sim \text{pareto}(\alpha_m, \alpha)$



scale like μ , shape like σ^{-2}
 α^2 -tails become even fatter.

BOX-COX TRANSFORMATION:-

$$y_i = \begin{cases} \frac{n_i^\lambda - 1}{\lambda} & \text{if } \lambda \neq 0 \\ \ln(n_i) & \text{if } \lambda = 0 \end{cases}$$

COVARIANCE:-

$$\text{cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$\text{cov}(x, y)$ — +ve then $x \uparrow y \uparrow$
 $\text{cov}(x, y)$ — -ve then $x \uparrow y \downarrow$

PEARSON CORRELATION COEF:-

$$r_{x,y} = \frac{\text{cov}(x, y)}{\sigma_x \cdot \sigma_y} \quad -1 \leq r \leq 1$$

SPEARMAN RANKS CORR-COFF:-

$$r = \rho_{r_x, r_y}$$

LINER BIASED REMOVED

K-S TEST:-

$$D_{n,m} = \sup_n |F_{1,n}(n) - F_{2,m}(n)|$$

significance level

$$D_{n,m} > C(\alpha) \sqrt{\frac{n+m}{nm}}, \quad C(\alpha) = \sqrt{-1/2 \ln(\alpha/2)}$$

$$P\text{-value} = 2 \cdot e^{-2D^2 \left(\frac{nm}{nm} \right)}$$

		truth	
		H ₀ true	H ₀ False
Decision based on sample	Fail to reject H ₀		Type II error
	Rejecting H ₀	Type I error	