



Exactly one empty box when n balls are randomly placed in n boxes:

$$\frac{\binom{n}{2} n!}{n^n}$$

$$P(\bar{A}|A) = 1 - P(A|A)$$

$$P(\bar{A}|B) = 1 - P(A|B)$$

GROUP I : SETS, COUNTING

SETS

Union :

$$A \cup B$$

Intersection:

$$A \cap B$$

Disjoint events: $A \cap B = AB = \emptyset$

PERMUTATIONS

permutations of n distinct objects taken k at a time

$$A_n^k = \frac{n!}{(n-k)!} = n(n-1)\dots(n-k+1)$$

→ ORDERED, NO REPEAT

n -permutations of n objects

$$P_n = n!$$

→ ORDERED, NO REPEAT

COMBINATION

n -combinations, taken k

$$C_n^k = \frac{n!}{(n-k)!k!} = \binom{n}{k}$$

→ NO ORDER, NO REPEAT

ways to choose k from n with replacement

$$\bar{A}^k = n^k$$

The number of distinct permutations of n things, n_1 of one kind, n_2 , ...

$$\frac{n!}{n_1! n_2! n_3! \dots n_k!}$$

Circles: $(n-1)!$

GROUP 2: PROBABILITY

DEFINITIONS

Sample space Ω Set of all possible outcomes

Event E $E \subseteq \Omega$

\Rightarrow Assume Ω have equally likely elements $P(E) = \frac{|E|}{|\Omega|}$

PROPERTIES

ADDITION RULES

$$P(A \cup B) = P(A) + P(B) - P(A \cap B) \rightarrow \text{disjoint} = 0$$

$$\begin{aligned} P(A \cup B \cup C) &= P(A) + P(B) + P(C) - P(AB) - P(AC) - P(BC) \\ &\quad + P(ABC) \end{aligned}$$

CONDITIONAL PROBABILITY

$$P(A|B) = \frac{P(AB)}{P(B)}$$

$\rightarrow A, B$ are independent if: $P(B|A) = P(B)$ $\rightarrow \neq 0$

$$P(A|B) = P(A) \rightarrow P(A \cap B) = P(A) \cdot P(B)$$

PRODUCT RULE

$$\begin{aligned} P(AB) &= P(A) P(B|A) \\ P(BA) &= P(B) P(A|B) \end{aligned} \quad \left. \begin{array}{l} \text{BOTH CAN} \\ \text{OCUR} \end{array} \right\}$$

$= P(A)$ if
independent

NOTE: Sum of r -digit numbers, taken from n numbers, without repetition:

$$\left({}^{n-1}P_{r-1} \right) (\text{sum of } n \text{ numbers}) \left(\underbrace{111\dots1}_{r \text{-times}} \right)$$

TOTAL PROBABILITY

$$P(A) = \sum_n P(A|B_n) \cdot P(B_n)$$

BAYES' THEOREM

$$P(B|A) = \frac{P(A|B) \cdot P(B)}{P(A)}$$

BERNOULLI'S TRIAL

A $\begin{cases} \rightarrow k \text{ successes} \\ \rightarrow n-k \text{ failures} \end{cases}$ $\rightarrow p$
 $\rightarrow (1-p)$

$$P(A) = \binom{n}{k} p^k (1-p)^{n-k}$$

MUTUALLY EXCLUSIVE / DISJOINT

$$\left\{ \begin{array}{l} P(AB) = 0 \\ P(A \cup B) = P(A) + P(B) \end{array} \right.$$

INDEPENDENT

$$\left\{ \begin{array}{l} P(AB) = P(A)P(B) \\ P(A \cup B) = P(A) + P(B) - P(A)P(B) \end{array} \right.$$

Note: $P(A) \Rightarrow P(B)$ but $P(B) \not\Rightarrow P(A)$ then $P(B) \geq P(A)$

GROUP 3: DISCRETE RANDOM VARIABLES

RANDOM VARIABLES

Suppose we conduct an exp with sample space Ω

A RV is a numeric function of the outcome : $X : \Omega \rightarrow \mathbb{R}$

$$\omega \rightarrow X(\omega)$$

$\rightarrow \Omega_X$ is the set of all variables of X

\equiv range / support

$\hookrightarrow \Omega_X$ finite or countably infinite : DISCRETE

$\hookrightarrow \Omega_X$ is uncountably large : CONTINUOUS

PMF

- DISCRETE

$$P_X : \Omega_K \rightarrow [0,1]$$

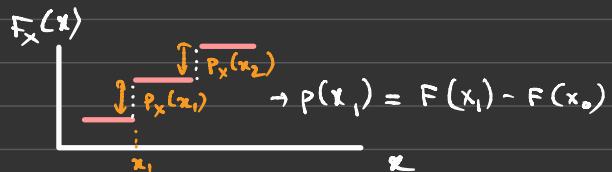
$$k \rightarrow P(X=k)$$

$$\rightarrow \left\{ \begin{array}{l} \sum_{i \in \Omega_X} P_X(i) = 1 \\ p(x) = \text{Prob}(X=x) \end{array} \right.$$

NOTE: Given $f(x)$ \rightarrow $\begin{cases} x_{\text{med.}}: f_x(x) = \frac{1}{2} \\ x_{\text{mod.}}: f_x(x) = \max \text{ (consider } f_x) \end{cases}$

CDF
— DISCRETE

$$F_X(x) = P[X \leq x] \rightarrow$$



$$P(a < X \leq b) = F_X(b) - F_X(a)$$

$$0 \leq F_X(a) \leq 1$$

$$P(X < b) = P(X \leq a) + P(a \leq X < b)$$

$$F(-\infty) = 0 ; F(+\infty) = 1$$

EXPECTATION
— DISCRETE

$$E[X] = \sum_{k \in \Omega_X} k \cdot p_X(k)$$

Average of possible values, weighted by their prob

$$E[X^2] = \sum x^2 p_X(x)$$

$$\text{Var}(X) = E[X^2] - E[X]^2 \geq 0$$

How "spread out" the distribution is

VARIANCE
— DISCRETE

$$\sigma_X = \sqrt{\text{Var}(X)} \rightarrow$$

$$\text{Var}[aX + b] = a^2 \text{Var}(X) \quad \forall a, b \in \mathbb{R}$$

STANDARD DEVIATION
— DISCRETE

**LINEARITY OF
EXPECTATION**
— DISCRETE

$$\begin{aligned} E[X+Y] &= E[X] + E[Y] \\ E[aX+b] &= aE[X] + b \\ \Rightarrow E[aX+bY+c] &= aE[X] + bE[Y] + c \end{aligned}$$

a, b, c are scalars
 X, Y are RVs

LOTUS
— DISCRETE

X : RV, $g: D \rightarrow \mathbb{R}$ is a function defined at least over $\Omega_X \subseteq D$:

$$E[g(X)] = \sum_{b \in \Omega_X} g(b) p_X(b)$$

Note: $E[g(X)] \neq g(E[X])$

$$\Rightarrow \sigma^2_{g(X)} = E[g(X) - \mu_{g(X)}]^2 = \sum_{x \in S_X} [g(x) - \mu_{g(X)}]^2 p_X(x)$$

DISCRETE RV DISTRIBUTIONS

BERNOULLI

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

$\Leftrightarrow \Omega_X = \{0,1\}$ and:

$$\text{PMF} \equiv p_X(k) = \begin{cases} p & , k=1 \\ 1-p & , k=0 \end{cases}$$

$$\rightarrow E[X] = p$$

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

YES/NO QUESTIONS

BINOMIAL

$$X \sim \text{B}(n, p)$$

$\Leftrightarrow \forall k \in \Omega_X = \{0, 1, 2, \dots, n\}$,

$$\text{PMF} \equiv p_X(k) = \binom{n}{k} p^k (1-p)^{n-k}$$

$$\rightarrow E[X] = np$$

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

SUCCESSES/FAILURES ?

GEOMETRIC

$$X \sim \text{Geo}(p) \Leftrightarrow \Omega_X = \{1, 2, 3, \dots\}$$

$$\text{PMF} \equiv p_X(k) = (1-p)^k p \quad k = 1, 2, \dots$$

$$\rightarrow E[X] = \frac{1}{p}$$

$$\text{and } \text{Var}(X) = \frac{1-p}{p^2}$$

SAO LÂU NỮA ĐÍNH ĐƯỢC
1 TẶNG GÓI MÈ ?

How many flips to have a heads for the first time?

How many ^{independent} flips to get the r -th heads?

NEGATIVE BINOMIAL

$$X \sim \text{NegBin}(r, p) \Leftrightarrow \Omega_X = \{r, r+1, \dots\}$$

$$\text{PMF} = p_X(k) = \binom{k-1}{r-1} p^r (1-p)^{k-r}$$

$$\rightarrow E[X] = \frac{r}{p} \quad \text{and} \quad \text{Var}(X) = r \frac{(1-p)}{p^2}$$

$k = r, r+1, \dots$

sum of Geo(p) vars

UNIFORM

- DISCRETE

$$X \sim \text{Unif}(a, b), a < b \in \mathbb{Z} \Leftrightarrow \text{PMF}:$$

$$p_X(k) = \begin{cases} \frac{1}{b-a+1} & , k \in \{a, a+1, \dots, b\} \\ 0 & , \text{ow} \end{cases}$$

$$\rightarrow E[X] = \frac{a+b}{2} \quad \text{and} \quad \text{Var}(X) = \frac{(b-a)(b-a+2)}{12}$$

KNOWN, FINITE
OUTCOMES EQUALLY
LIKELY TO HAPPEN

POISSON

$$X \sim P(\lambda) \Leftrightarrow \Omega_X = \{0, 1, 2, \dots\}$$

$$\text{PMF} \equiv p_X(k) = e^{-\lambda} \frac{\lambda^k}{k!}, k=0,1,2,\dots$$

$\lambda > 0$ is the RATE OF OCCURRENCE

$$\Rightarrow E[X] = \text{Var}(X) = \lambda$$

OUTCOMES OCCURRING IN A TIME INTERVAL

BINOMIAL APPROXIMATION By POISSON DISTRIBUTION

$$B(n, p) \rightarrow P(\lambda) \text{ as } n \rightarrow \infty, p \rightarrow 0$$

$$\text{Ideal: } \lambda = np < >$$

GROUP 4: CONTINUOUS RVs

PDF
—CONTINUOUS

PDF: $f_x : \mathbb{R} \rightarrow \mathbb{R}$:

$$f_x(x) \geq 0, \forall x \in \mathbb{R} \quad \left. \int_{-\infty}^{+\infty} f(x) dx = 1 \right\} \text{condition}$$

$$P(a < x < b) = \int_a^b f_x(x) dx$$

CDF
—CONTINUOUS

CDF: $F_x : \mathbb{R} \rightarrow \mathbb{R}$

$$F_x(x) = P(X \leq x) = \int_{-\infty}^x f_x(t) dt, t \in \mathbb{R}$$

$$\frac{dF_x(x)}{dx} = f_x(x)$$

$$P(a \leq X \leq b) = F_x(b) - F_x(a)$$

$$F_x(-\infty) = 0, \quad F_x(+\infty) = 1$$

$$F_x(c) \leq F_x(d), \quad c < d$$

EXPECTATION

- CONTINUOUS

$$\mathbb{E}[X] = \int_{-\infty}^{+\infty} x f_X(x) dx$$

LOTUS

- CONTINUOUS

$$\mathbb{E}[g(X)] = \int_{-\infty}^{+\infty} g(x) f_X(x) dx$$

$$\rightarrow \sigma_y^2 = \mathbb{E}[g(x) - \mu_{g(x)}]^2 = \int_{-\infty}^{+\infty} [g(x) - \mu_{g(x)}]^2 f_X(x) dx$$

CONTINUOUS RV DISTRIBUTIONS

UNIFORM

- CONTINUOUS

$X \sim \text{Unif}(a,b)$, $a < b \in \mathbb{R} \Leftrightarrow$ PDF:

< x can take ~~#~~ value in $[a,b]$. >

$$f_X(x) = \begin{cases} \frac{1}{b-a}, & x \in [a,b] \\ 0 & \text{ow} \end{cases}$$

$$\rightarrow E[X] = \frac{a+b}{2} \text{ and } \text{Var}(X) = \frac{(b-a)^2}{12}$$

$$P(a \leq X \leq b) = \int_a^b f_X(x) dx$$

$$F_X(x) = \begin{cases} 0, & x < a \\ \frac{x-a}{b-a}, & a \leq x \leq b \\ 1, & x > b \end{cases}$$

NORMAL

- CONTINUOUS

$X \sim N(\mu, \sigma^2) \Leftrightarrow \Omega_X = (-\infty, +\infty)$:

$$\text{PDF} \equiv f_X(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

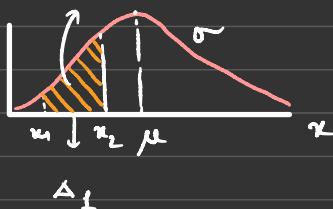
$$\rightarrow E[X] = \mu, \quad \text{Var}(X) = \sigma^2$$

STANDARD NORMAL DISTRIBUTION

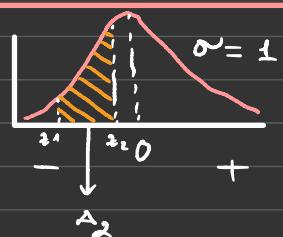
$$Z \sim \mathcal{N}(0, 1), z = \frac{x - \mu}{\sigma}$$

$$\text{CDF}(z) = \varphi_z(a) = F_z(a) = P(Z \leq a)$$

$$P(x_1 < X < x_2)$$



$$\varphi(-a) = 1 - \varphi(a)$$



$$\rightarrow A_1 = A_2$$

$$z = \frac{x - \mu}{\sigma}$$

NORMAL APPROXIMATION TO THE BINOMIAL DISTRIBUTION

$$P(x_1 < X < x_2) = \Phi\left(\frac{x_2 + 0.5 - \mu}{\sigma}\right) - \Phi\left(\frac{x_1 - 0.5 - \mu}{\sigma}\right)$$

Ideal: $np > 5$ and $n(1-p) > 5$

$$\sqrt{npq}$$

EXPONENTIAL

$$X \sim \text{Exp}(\lambda) \Leftrightarrow \Omega_X = [0, +\infty)$$

$$\text{PDF} \equiv f_X(x) = \begin{cases} \lambda e^{-\lambda x}, & x \geq 0 \\ 0, & \text{ow} \end{cases}$$

X is the waiting time until the first occurrence of an event in a Poisson process with para. $\lambda \rightarrow$ POISSON DISTRIBUTION: # events occurring in a time interval

EXPONENTIAL DISTRIBUTION:
Time taken between 2 events

$$\rightarrow E[X] = \frac{1}{\lambda},$$

$$\text{Var}(X) = \frac{1}{\lambda^2}$$

$$\text{CDF} \equiv F_X(x) = \begin{cases} 1 - e^{-\lambda x}, & x \geq 0 \\ 0, & \text{ow} \end{cases}$$

JOINT CDF OF TWO RVs

PROPERTIES

$$F_{x,y}(x,y) = P[X < x, Y < y]$$

$$0 \leq F_{x,y}(x,y) \leq 1$$

$$\begin{aligned} F_x(x) &= F_{x,y}(x, \infty) \\ F_y(y) &= F_{x,y}(\infty, y) \end{aligned} \quad \left\{ \text{MARGINAL CDF} \right.$$

$$F_{x,y}(-\infty, y) = F_{x,y}(x, -\infty) = 0$$

$$F_{x,y}(\infty, \infty) = 1$$

If $x_1 < x_2, y_1 < y_2$ then:

$$\begin{aligned} P[x_1 \leq X \leq x_2, y_1 \leq Y \leq y_2] &= F_{x,y}(x_2, y_2) - F_{x,y}(x_2, y_1) \\ &\quad - F_{x,y}(x_1, y_2) + F_{x,y}(x_1, y_1) \end{aligned}$$

GROUP 5 : JOINT DISCRETE RVs

JOINT PMF

let X, Y be discrete random variables

$$\rightarrow \text{Joint PMF: } P_{X,Y}(a,b) = P(X=a, Y=b)$$

Joint range: $\Omega_{X,Y} = \{(c,d) : P_{X,Y}(c,d) > 0\} \subseteq \Omega_X \times \Omega_Y$



$$\sum_{(s,t) \in \Omega_{X,Y}} P_{X,Y}(s,t) = 1$$

$X \backslash Y$	y_1	...	y_j	...	y_n	\sum_j
x_1	p_{11}	...	p_{1j}	...	p_{1n}	$P[X = x_1]$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
x_i	p_{i1}	...	p_{ij}	...	p_{in}	$P[X = x_i]$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
x_m	p_{m1}	...	p_{mj}	...	p_{mn}	$P[X = x_m]$
\sum_i	$P[Y = y_1]$...	$P[Y = y_j]$...	$P[Y = y_n]$	$\sum_i \sum_j = 1$

THEOREM

$$0 \leq p_{ij} \leq 1, \quad \forall i = \overline{1, m}, \quad \forall j = \overline{1, n}$$

$$\sum_{i=1}^m \sum_{j=1}^n p_{ij} = 1$$

$$\rightarrow F_{X,Y}(x,y) = \sum_{x_i < x} \sum_{y_j < y} p_{ij}$$

MARGINAL PMF

$$P_X(a) = \sum_{b \in \Omega_Y} P_{x,y}(a, b)$$

$$P_Y(b) = \sum_{a \in \Omega_X} P_{x,y}(a, b)$$

Let's $g : \mathbb{R}^2 \rightarrow \mathbb{R}^2$

$$E[g(x,y)] = \sum_{x \in \Omega_X} \sum_{y \in \Omega_Y} g(x,y) P_{x,y}(x,y)$$

We display $P_X(x)$ and $P_Y(y)$ by rewriting the matrix in Example 1 and placing the row sums and column sums in the margins.

$P_{X,Y}(x,y)$	$y = 0$	$y = 1$	$y = 2$	$P_X(x)$
$x = 0$	0.01	0	0	0.01
$x = 1$	0.09	0.09	0	0.18
$x = 2$	0	0	0.81	0.81
$P_Y(y)$	0.10	0.09	0.81	

$$P_X(x) = \begin{cases} 0.01, & x = 0, \\ 0.18, & x = 1, \\ 0.81, & x = 2, \\ 0, & \text{otherwise.} \end{cases} \quad P_Y(y) = \begin{cases} 0.1, & y = 0, \\ 0.09, & y = 1, \\ 0.81, & y = 2, \\ 0, & \text{otherwise.} \end{cases}$$

GROUP C: JOINT CONTINUOUS RVs

JOINT PDF

$$\text{Joint PDF: } f_{x,y}(a,b) > 0 = \frac{\partial^2 F_{x,y}(x,y)}{\partial x \partial y}$$

$$\Rightarrow \text{Joint range: } \Omega_{x,y} = \{(c,d) : f_{x,y}(c,d) > 0\} \subset \Omega_x \times \Omega_y$$

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{x,y}(x,y) dx dy = 1$$

MARGINAL PDF

$$f_x(x) = \int_{-\infty}^{\infty} f_{x,y}(x,y) dy \quad \rightarrow \text{domain of } y$$

$$f_y(y) = \int_{-\infty}^{+\infty} f_{x,y}(x,y) dx \quad \rightarrow \text{domain of } x$$

LOTUS $g: \mathbb{R}^2 \rightarrow \mathbb{R}^2$

$$\mathbb{E}[g(x,y)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x,y) f_{x,y}(x,y) dy dx$$

$$\Rightarrow P(a \leq x \leq b, c \leq y \leq d) = \int_a^b \int_c^d f_{x,y}(x,y) dy dx$$

EXPECTED VALUE OF TWO RVs

For R.Vs : X and Y , the expected value of $W = g(X, Y)$ is :

DISCRETE

$$\begin{aligned} E[W] &= \sum_{x \in S_X} \sum_{y \in S_Y} g(x, y) \cdot P_{X,Y}(x, y) \\ &= \sum_i \sum_j g(x_i, y_j) P[X=x_i, Y=y_j] \end{aligned}$$

CONTINUOUS

$$E[W] = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} g(x, y) f_{X,Y}(x, y) dx dy$$

THEOREMS

$$E[X+Y] = E[X] + E[Y]$$

$$\text{Var}[X+Y] = \text{Var}[X] + \text{Var}[Y] + 2 \cdot E[(X - \mu_X)(Y - \mu_Y)]$$

GROUP 7: COVARIANCE, CORRELATION

DEFINITION

Covariance of X and Y :

$$\begin{aligned}\text{Cov}(X, Y) &= E[(X - E[X])(Y - E[Y])] \\ &= E[XY] - E[X]E[Y]\end{aligned}$$

$$\begin{aligned}&\sum_i \sum_j x_i y_j p_{ij} \\ &\sum xy f(x,y)\end{aligned}$$

PROPERTIES

$$\text{Cov}(X, X) = \text{Var}(X)$$

$$\text{Cov}(X, Y) = \text{Cov}(Y, X)$$

$$\text{Cov}(ax + b, cy + d) = a.c. \text{Cov}(X, Y)$$

$$\text{Cov}(X_1 + X_2, Y) = \text{Cov}(X_1, Y) + \text{Cov}(X_2, Y)$$

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y)$$

COVARIANCE MATRIX

$$\Gamma = \begin{bmatrix} \text{Cov}(X, X) & \text{Cov}(X, Y) \\ \text{Cov}(Y, X) & \text{Cov}(Y, Y) \end{bmatrix} = \begin{bmatrix} \text{Var}(X) & \text{Cov}(X, Y) \\ \text{Cov}(Y, X) & \text{Var}(Y) \end{bmatrix}$$

CORRELATION

$$r_{x,y} = E[xy] \rightarrow \begin{cases} x, y \text{ are orthogonal: } r_{x,y} = 0 \\ x, y \text{ are uncorrelated: } \text{Cov}[x, y] = 0 \end{cases}$$

CORRELATION COEFFICIENT

$$\rho_{x,y} = \frac{\text{Cov}[x, y]}{\sqrt{\text{Var}(x) \text{Var}(y)}} = \frac{\text{Cov}[x, y]}{\sigma_x \sigma_y}$$

$$-1 \leq \rho_{x,y} \leq 1$$

$$y = ax + b : \rho_{x,y} = \begin{cases} -1, a < 0 \\ 0, a = 0 \end{cases}$$

CONDITIONING BY AN EVENT

CONDITIONAL JOINT PMF

$\rho_{x,y|B}(x,y) = P[X=x, Y=y | B]$

B is a region

$$= \begin{cases} \frac{\rho_{x,y}(x,y)}{P(B)} & , (x,y) \in B \\ 0 & \text{else} \end{cases}$$

CONDITIONAL JOINT PDF

$$f_{x,y|B}(x,y) = \begin{cases} \frac{f_{x,y}(x,y)}{P(B)} & , (x,y) \in B \\ 0 & \text{else} \end{cases}$$

CONDITIONAL EXPECTED VALUE

$$E[w|B] = \sum_{x \in S_x} \sum_{y \in S_y} g(x, y) P_{x,y|B}(x, y)$$

$$E[w|B] = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} g(x, y) f_{x,y|B}(x, y) dx dy$$

CONDITIONAL VARIANCE

$$\text{Var}[w|B] = E[w^2|B] - (E[w|B])^2$$

CONDITIONING
BY A RV

CONDITIONAL PMF

$$\begin{aligned} P_{x|y}(x|y) &= P[x=x | y=y] \\ \Rightarrow P_{x,y}(x,y) &= P_{x|y}(x|y) P_y(y) \\ &= P_{y|x}(y|x) P_x(x) \end{aligned}$$

CONDITIONAL EXPECTED VALUE

$$E[x|y=y] = \sum_{x \in S_x} x P_{x|y}(x|y)$$

CONDITIONAL PDF

$$f_{y|x}(y|x) = \frac{f_{x,y}(x,y)}{f_x(x)}$$

$$f_{x|y}(x|y) = \frac{f_{x,y}(x,y)}{f_y(y)}$$

GROUP 8 : INDEPENDENT RVs

INDEPENDENCE

$$\text{Discrete: } p_{x,y}(x,y) = p_x(x) p_y(y)$$

$$\text{Continuous: } f_{x,y}(x,y) = f_x(x) f_y(y)$$

$$\rightarrow \begin{cases} p_{x|y}(x,y) = p_x(x), & p_{y|x}(y|x) = p_y(y) \\ f_{x|y}(x,y) = f_x(x), & f_{y|x}(y|x) = f_y(y) \end{cases}$$

PROPERTIES

$$E[g(X) h(Y)] = E[g(X)] E[h(Y)]$$

$$r_{x,y} = E[XY] = E[X] E[Y]$$

$$\text{Cov}[X,Y] = \rho_{x,y} = 0$$

$$\text{Var}[X+Y] = \text{Var}[X] + \text{Var}[Y]$$

$$E[X|y=y] = E[X]$$

$$E[Y|X=x] = E[Y]$$

GROUP 9: LIMIT THEOREMS

SAMPLE MEAN

Let $X_{1,n}$ be independent identically distributed RVs, μ and σ^2
Sample mean:

$$\bar{X}_n = \frac{x_1 + x_2 + \dots + x_n}{n}$$

$$\rightarrow E[\bar{X}_n] = \mu, \quad \text{Var}(\bar{X}_n) = \frac{\sigma^2}{n}$$

LIMIT THEOREMS

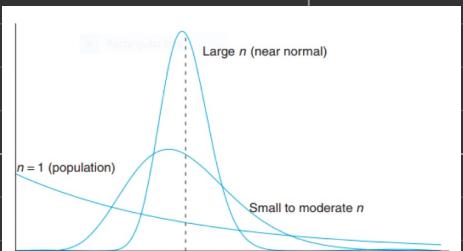
LAW OF LARGE NUMBERS

As $n \rightarrow \infty$, $\underbrace{\bar{X}_n}_{\text{prob}}$ converges to the true mean μ

$$P\left(\lim_{n \rightarrow \infty} \bar{X}_n = \mu\right) = 1$$

CENTRAL LIMIT THEOREM

The standardised mean approaches the standard Normal distribution



$$\bar{Z}_n = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \rightarrow N(0,1) \text{ as } n \rightarrow \infty$$

ideal: $n \geq 30$, not skewed

$$\text{Sum: } Z_n = \frac{T - n\mu}{\sqrt{n}\sigma} \sim N(n\mu, \sqrt{n}\sigma)$$

GROUP 10 : RANDOM SAMPLES - STATISTICS

SAMPLE MEAN,
MEDIAN, MOD

SAMPLE MEAN

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

SAMPLE MEDIAN

$$\hat{x} = \begin{cases} x_{(n+1)/2}, & n \text{ odd} \\ \frac{1}{2}(x_{n/2} + x_{n/2+1}), & n \text{ even} \end{cases}$$

Sorted in increasing
order

SAMPLE VARIANCE,
S.D., RANGE

SAMPLE VARIANCE

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$= \frac{1}{n(n-1)} \left[n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2 \right]$$

SAMPLE S.D.

$$s = \sqrt{s^2}$$

RANGE

$$R = x_{\max} - x_{\min}$$

Note

Sample	Population
n : number of measurements in the sample	N : number of measurements in the population
\bar{x} : sample mean	μ : population mean
s^2 : sample variance	σ^2 : population variance
s : sample standard deviation	σ : population standard deviation

GROUP II: SAMPLING DISTRIBUTION

S.D OF MEANS

$$\bar{x} = \frac{1}{n} (x_1 + x_2 + \dots + x_n)$$

$$\begin{cases} \mu_{\bar{x}} = \mu \\ \sigma_{\bar{x}}^2 = \frac{\sigma^2}{n} \end{cases}$$

↳ has a normal distribution

S.D OF DIFFERENCE OF MEANS

Theorem 3

If independent samples of size n_1 and n_2 are drawn at random from two populations, discrete or continuous, with means μ_1 and μ_2 and variances σ_1^2 and σ_2^2 , respectively, then the sampling distribution of the differences of means, $\bar{X}_1 - \bar{X}_2$, is approximately normally distributed with mean and variance given by

$$\mu_{\bar{X}_1 - \bar{X}_2} = \mu_1 - \mu_2 \quad \text{and} \quad \sigma_{\bar{X}_1 - \bar{X}_2}^2 = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}.$$

Hence,

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{(\sigma_1^2/n_1) + (\sigma_2^2/n_2)}}$$

is approximately a standard normal variable $\mathcal{N}(0, 1)$.

S.D OF s^2

$$s^2 \sim \sigma^2 : \text{random sample of size } n \quad \chi^2 = \frac{(n-1)s^2}{\sigma^2} = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{\sigma^2}$$

SPECIAL: STUDENT t- DISTRIBUTION

DEFINITIONS

Theorem 5

Let Z be a standard normal random variable and V a chi-squared random variable with ν degrees of freedom. If Z and V are independent, then the distribution of the random variable T , where

$$T = \frac{Z}{\sqrt{V/(n-1)}}$$

is given by the density function

$$h(t) = \frac{\Gamma(\nu+1)/2}{\Gamma(\nu/2)\sqrt{\pi\nu}} \left(1 + \frac{t^2}{\nu}\right)^{-(\nu+1)/2}, \quad -\infty < t < +\infty.$$

This is known as the **t-distribution** with ν degrees of freedom.

$$h(t) = \frac{\Gamma(\nu+1)/2}{\Gamma(\nu/2)\sqrt{\pi\nu}} \left(1 + \frac{t^2}{\nu}\right)^{-(\nu+1)/2}$$

↳ let X_1, \dots, X_n be independent RVs, normal with μ, σ^2 .

$$\bar{X} = \frac{x_1 + \dots + x_n}{n} \quad \text{and} \quad S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{X})^2$$

Then RV: $T = \frac{\bar{X} - \mu}{S/\sqrt{n}}$ has a t-distribution, $\nu = n-1$
 ↳ deg. of freedom

t-dis $\xrightarrow{n \rightarrow \infty}$ Normal

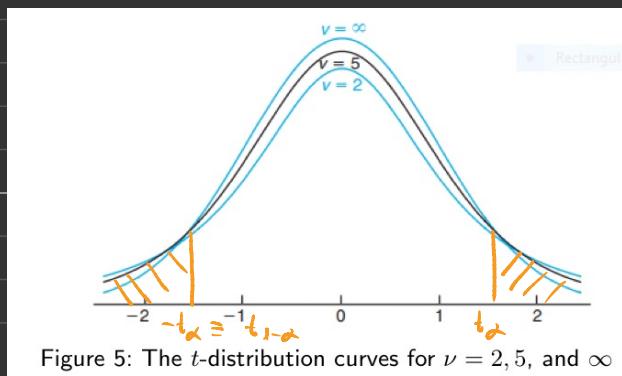


Figure 5: The t-distribution curves for $\nu = 2, 5$, and ∞

GROUP 12 : ESTIMATIONS

POINT ESTIMATE

UNBIASED ESTIMATOR

A statistic $\hat{\theta}$ is said to be an unbiased estimator of para.
 θ if: $E[\hat{\theta}] = \theta$
↳ Distribution centred at θ

→ Smallest variance = most efficient estimator of θ .

INTERVAL ESTIMATE

Interval $\hat{\theta}_L < \theta < \hat{\theta}_U$, $\hat{\theta}_{L,U}$ depend on the value of $\hat{\theta}$ for a sample
Sampling dist. of $\hat{\theta}$.

$$\rightarrow P(\underbrace{\hat{\theta}_L < \theta < \hat{\theta}_U}_{}) = 1-\alpha, 0 < \alpha < 1$$

confidence interval \rightarrow 100 $(1-\alpha)\%$ confidence interval

→ Shorter interval, with higher degree of confidence is WIDER!

MEAN ESTIMATING

KNOWN σ

TWO SIDED INTERVALS

$$z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

$$P(-z_{\alpha/2} < z < z_{\alpha/2}) = 1-\alpha$$

$$\underbrace{z}_{\text{~} \rightarrow} P\left(\bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} < \mu < \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right) = 1-\alpha$$

CONFIDENCE FOR μ

$$\bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} < \mu < \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

ERROR $\approx 90(1-\alpha)\%$ confident that the error not exceed:

$$e = \frac{z_{\alpha/2} \sigma}{\sqrt{n}}$$

SAMPLE SIZE Confident that error $\leq e$ when:

$$n = \left(\frac{z_{\alpha/2} \sigma}{e} \right)^2$$

ONE-SIDED BOUNDS

$$\xrightarrow{\text{CLT}} P\left(\mu < \bar{X} + \frac{z_{\alpha} \sigma}{\sqrt{n}}\right) = 1 - \alpha \quad \left\{ \begin{array}{l} \text{upper bound: } \bar{X} + z_{\alpha} \frac{\sigma}{\sqrt{n}} \\ \text{lower bound: } \bar{X} - z_{\alpha} \frac{\sigma}{\sqrt{n}} \end{array} \right.$$

UNKNOWN σ

$$P\left(\bar{X} - t_{\alpha/2} \frac{s}{\sqrt{n}} < \mu < \bar{X} + t_{\alpha/2} \frac{s}{\sqrt{n}}\right) = 1 - \alpha$$

< similar to known σ , but $\sigma \rightarrow s$ and $N \rightarrow t$ -dist. >

CONFIDENCE INTERVAL

BOUNDS

$$\bar{X} - t_{\alpha/2} \frac{s}{\sqrt{n}} < \mu < \bar{X} + t_{\alpha/2} \frac{s}{\sqrt{n}}$$

$\left\{ \begin{array}{l} \text{upper: } \bar{X} + t_{\alpha/2} \frac{s}{\sqrt{n}} \\ \text{lower: } \bar{X} - t_{\alpha/2} \frac{s}{\sqrt{n}} \end{array} \right.$	$\bar{X} + t_{\alpha/2} \frac{s}{\sqrt{n}}$
---	---

$$\overbrace{\sigma \text{ unknown}}^{n \geq 30} \Rightarrow s \leftarrow \sigma : \bar{X} \pm t_{\alpha/2} \frac{s}{\sqrt{n}}$$

PROPORTION ESTIMATING

A point estimator of the proportion p in a binomial exp is given by stat. $\hat{P} = \frac{x}{n}$, $x = \# \text{ successes}$ in n trials. $\rightarrow \hat{P} = \frac{x}{n}$: point est.

CLT $\Rightarrow \mu_{\hat{P}} = p, \sigma^2_{\hat{P}} = \frac{pq}{n}$

$$P(-z_{\alpha/2} < Z < z_{\alpha/2}) = 1-\alpha$$

$$, Z = \frac{\hat{P} - P}{\sqrt{pq/n}}$$

CONFIDENCE INTERVAL

$$\hat{q} = 1 - \hat{P}$$

$$\hat{P} - z_{\alpha/2} \sqrt{\frac{\hat{P}\hat{q}}{n}} < P < \hat{P} + z_{\alpha/2} \sqrt{\frac{\hat{P}\hat{q}}{n}}$$

$\hat{n}\hat{P}$ and $\hat{n}\hat{q} \geq 5$

ERROR

$$\epsilon = z_{\alpha/2} \sqrt{\frac{\hat{P}\hat{q}}{n}}$$

SAMPLE SIZE

$$n = \frac{z^2 \alpha/2 \hat{P}\hat{q}}{\epsilon^2} \quad // \text{at least } 90\% \text{ confident}$$

$$n = \frac{z^2 \alpha/2}{4\epsilon^2}$$

// at least 90% confident

GROUP 13: HYPOTHESIS TESTING

DEFINITIONS

Alternative hypothesis H_1 : which we will wish to support : $\mu_1 = \mu_2$
 Null hypothesis H_0 : contradiction

$\begin{cases} \mu_1 < \mu_2 : \text{left-tailed test} \\ \mu_1 > \mu_2 : \text{right-tailed test} \\ \mu_1 \neq \mu_2 : \text{two-tailed test} \end{cases}$

- Assume H_0 true standardised test statistic \rightarrow

$\begin{cases} \notin \text{rejection region} : \text{fail to reject } H_0 \\ \in \text{rejection region} : \text{reject } H_0 \end{cases}$

at $\alpha\%$ level of significance

	H_0 true	H_0 false	
- Error	Not reject H_0	Correct	Type II
	Reject H_0	Type I	Correct
	\downarrow		$\rightarrow \alpha \sim \frac{1}{\beta}$
		$P[\cdot] = \alpha$	

Approach to Hypothesis Testing with Fixed Probability of Type I Error

- 1 State the null and alternative hypotheses.
- 2 Choose a fixed significance level α .
- 3 Choose an appropriate test statistic and establish the critical region/rejection region based on α .
- 4 Reject H_0 if the computed test statistic is in the critical region/rejection region. Otherwise, do not reject.
- 5 Draw scientific or engineering conclusions.

SINGLE SAMPLE HYPOTHESIS TESTING

SINGLE MEAN,
VARIANCE KNOWN

TWO-TAILED

$$H_0: \mu = \mu_0$$

$$2) z = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{n}}$$

$$H_1: \mu \neq \mu_0$$

3) $\left[z < -z_{\alpha/2} \text{ or } z > z_{\alpha/2} \right] \rightarrow \text{Reject } H_0$

$-z_{\alpha/2} < z < z_{\alpha/2} \rightarrow \text{Fail to reject } H_0$

ONE-TAILED

$$H_0: \mu = \mu_0$$

$$\begin{aligned} H_1: \mu &> \mu_0 && (\text{Right}) \\ &< \mu_0 && (\text{Left}) \end{aligned}$$

$$2) z = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{n}}$$

3) $\left\{ \begin{array}{l} z > z_\alpha \rightarrow \text{Reject} \\ z < -z_\alpha \rightarrow \text{Fail} \end{array} \right. \quad \left\{ \begin{array}{l} \text{right} \\ \text{left} \end{array} \right.$

$\left\{ \begin{array}{l} z < -z_\alpha \rightarrow \text{Reject} \\ z > z_\alpha \rightarrow \text{Fail} \end{array} \right. \quad \left\{ \begin{array}{l} \text{right} \\ \text{left} \end{array} \right.$

SINGLE SAMPLE,
VARIANCE UNKNOWN

TWO-TAILED

$$H_0: \mu = \mu_0$$

$$H_1: \mu \neq \mu_0$$

$$t = \frac{\bar{x} - \mu_0}{s} \sqrt{n}$$

$$3, \quad t < -t_{\alpha/2}^{(n-1)} \text{ or } t > t_{\alpha/2} \rightarrow \text{reject}$$

$$-t_{\alpha/2} < t < t_{\alpha/2} \rightarrow \text{Fail}$$

ONE-TAILED

$$H_0: \mu = \mu_0$$

$$H_1: \mu > \mu_0 \text{ (L)}$$

$$\mu < \mu_0 \text{ (R)}$$

$$2, \quad t = \frac{\bar{x} - \mu_0}{s} \sqrt{n}$$

$$3, \quad \begin{array}{ll} t > t_\alpha & \rightarrow \text{Reject} \\ t < -t_\alpha & \rightarrow \text{Fail} \end{array} \quad \left. \begin{array}{l} \text{right} \\ \text{left} \end{array} \right\}$$

$$\begin{array}{ll} t < -t_\alpha & \rightarrow \text{Reject} \\ t > t_\alpha & \rightarrow \text{Fail} \end{array} \quad \left. \begin{array}{l} \text{left} \\ \text{right} \end{array} \right\}$$

$$n \geq 30 \rightarrow T \sim N(0,1)$$

$$t = \frac{\bar{x} - \mu_0}{s} \sqrt{n}$$

PROPORTIONS

TWO-TAILED

NOTE: $\begin{cases} np_0 \geq 5 \\ n(1-p_0) \geq 5 \end{cases}$

1) $H_0 : p = p_0$

$H_1: p \neq p_0$

2)

$$z = \frac{\hat{p} - p_0}{\sqrt{p_0(1-p_0)}} \quad \sqrt{n}$$

3) $z < -z_{\alpha/2} \text{ or } z > z_{\alpha/2} \rightarrow \text{reject}$

$-z_{\alpha/2} < z < z_{\alpha/2} \rightarrow \text{fail}$

Similarly for one-tailed ~

TWO SAMPLES HYPOTHESIS TESTING

CONDITIONS:

- 1) Samples are randomly selected
- 2) Samples are independent
- 3) Each sample size > 30 ; or each must have normal distribution, known standard deviation

σ_1^2, σ_2^2 KNOWN

$$H_0: \mu_1 - \mu_2 = D_0$$

$$\begin{cases} H_1: \mu_1 - \mu_2 > D_0 \text{ or } \mu_1 - \mu_2 < D_0 \\ H_1: \mu_1 - \mu_2 \neq D_0 \end{cases}$$

$$z = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

Rejection region: similar

σ_1^2, σ_2^2 UNKNOWN
 $n_1, n_2 \geq 30$

$$z = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

σ_1^2, σ_2^2 UNKNOWN
 $n_1, n_2 < 30$

CONDITIONS

- 1) Samples are randomly selected
- 2) Independent
- 3) Each sample has a normal distribution,
population variances are equal.

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\left(\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2} \right) \cdot \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

Rejection region

$$\left\{ \begin{array}{l} \text{One-tailed : } \begin{array}{ll} t > t_{\alpha}^{(n_1+n_2-2)} & \rightarrow \text{Reject (Right)} \\ t < -t_{\alpha}^{(n_1+n_2-2)} & \rightarrow \text{Reject (Left)} \end{array} \\ \text{Two-tailed : } t > t_{\alpha/2}^{(n_1+n_2-2)} \text{ or } t < -t_{\alpha/2}^{(n_1+n_2-2)} \end{array} \right.$$

DIFFERENCE BETWEEN
PROPORTIONS

CONDITIONS

- 1) Samples are randomly selected
- 2) Independent

3) Large enough for normal sampling distributions

$$\begin{cases} n_1 p_1 \geq 5, n_1(1-p_1) \geq 5 \\ n_2 p_2 \geq 5, n_2(1-p_2) \geq 5 \end{cases}$$

$$z = \frac{(\hat{p}_1 - \hat{p}_2) - (p_1 - p_2)}{\sqrt{\bar{p}(1-\bar{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}, \quad \bar{p} = \frac{n_1 \hat{p}_1 + n_2 \hat{p}_2}{n_1 + n_2}$$

H_0	H_1	Rejection region W_α
$p_1 = p_2$	$p_1 \neq p_2$	$(-\infty; -z_{\alpha/2}) \cup (z_{\alpha/2}; +\infty)$
$p_1 = p_2$	$p_1 > p_2$	$(z_\alpha; +\infty)$
$p_1 = p_2$	$p_1 < p_2$	$(-\infty; -z_\alpha)$