

Pop vs Sample

Pop	Sample
N	n
all items	subset of pop
params	statistics
hard to gather data	cheap & fast
	needs random sample

Types of Data

- Categorical
 - numerical
 - continuous
 - can be infinitely divided into smaller measurements
 - discrete

barplots vs histograms

- barplots: categorical
- histograms:
 - numerical
 - intervals as categories
 - frequency table in intervals

percentile

30th percentile 70th percentile

30 40 50 60 70

trim 10% means < 10% & > 90% were removed

Mean - $\frac{\text{Sum of all values}}{\# \text{ of values}}$

Median - The middle number

Mode - most common number

Experiment - process w/ randomness

Outcome - results of Experiment

Ex: Weather

S - Sample Space in $\{ \}$ braces

Event - Collection of outcomes

Ex: $A = \{1, 2, 5, 1, 2\}$

Nullset - no outcomes (\emptyset)

IQR

Ex: 7, 9, 9, 10, 10, 10, 11, 12, 12, 14

lower quartile median upper quartile

min max

box plot

low Q_1 median Q_2 max

IQR

$1.5 \times \text{IQR}$

$\text{IQR} = (\text{Upper median} - \text{Lower median})$

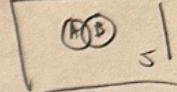
or $12 - 9 = 3$

$\text{IQR} = Q_2 - Q_1$

mid outlier: $1.5 \times \text{IQR}$

Extreme outlier: $3 \times \text{IQR}$

Venn Diagram

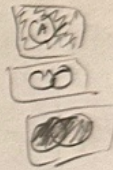


3 operations: Set operations

Complement (A^c) - opposite

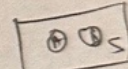
Intersection ($A \cap B$) - $A \text{ and } B$

Union ($A \cup B$) - $A \text{ or } B$



IF A & B disjoint, then they are mutually exclusive

$A \cap B = \emptyset$



1) $P(\emptyset) = 0$

2) $P(A^c) = 1 - P(A)$

3) for any event A , $P(A) \leq 1$

4) $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

5) $P(A \cup B \cup C) = P(A) + P(B) + P(C) - P(A \cap B) - P(A \cap C) - P(B \cap C) + P(A \cap B \cap C)$

3 probability Axioms

1) non-negativity: for all events A , $P(A) \geq 0$

2) unitarity: $P(S) = 1$

3) Sigma additivity: for disjoint events A_1, A_2, A_3, \dots

$P(A_1 \cup A_2 \cup A_3 \cup \dots) = P(A_1) + P(A_2) + \dots$

$P\left(\bigcup_{i=1}^{\infty} A_i\right) = P\left(\sum_{i=1}^{\infty} P(A_i)\right)$

$\bigcup_{i=1}^{\infty}, \bigcap_{i=1}^{\infty}, \sum_{i=1}^{\infty}$

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

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

Conditional Probability

$P(A|B)$ = Probability A will happen given B

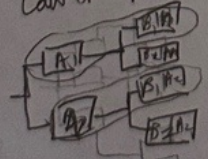
$P(A|B) = \frac{P(A \cap B)}{P(B)}$

$P(A \cap B) = P(A)P(B|A)$ multiplication rule

Bayes Theorem

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

Law of total probability



make a tree diagram, then add all probabilities that fulfill condition

$P(B) = P(A_1)P(B|A_1) + \dots + P(A_n)P(B|A_n)$

Independent Events

All 3 must be true

1) $P(A \cap B) = P(A)P(B)$

2) $P(A|B) = P(A)$

3) $P(B|A) = P(B)$

= A & B are independent

IF $P(A \cap B) = P(A)P(B)$

$\Rightarrow P(A, B) = \frac{P(A, B)}{P(B)} = \frac{P(B|A_i)P(A_i)}{\sum_{i=1}^k P(B|A_i) \cdot P(A_i)}$ $i = 1, \dots, k$

UNO replacement

order doesn't matter
(combination)

$$\binom{n}{r} = \frac{n!}{r!(n-r)!} = nCr$$

we choose r things from n values

order does matter
(permutation)

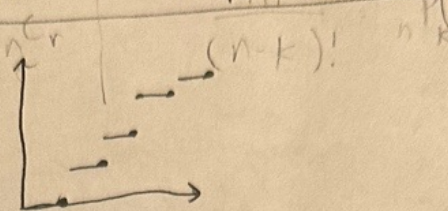
$$\frac{n!}{(n-k)!} = nPk$$

of spots remaining (k)

Addition rule

- Two separate tasks $A \nmid B$
 $A = m$ ways, $B = n$ ways
 then $A \nmid B$ can be done $m+n$ ways

PMF - Probability of a certain value (discrete only)
 Cdf - Describes distribution of RVs (discrete & continuous)



Product rule

- Subsequent tasks $A \nmid B$
 $A = m$ ways, $B = n$ ways
 then B can be done mn ways

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

a largest X less than a

$$P(a \leq X \leq b) = F(b) - F(a-1)$$

$$P(X \geq a) = 1 - P(X \leq a-1)$$

$$P(x) \leq P(X=x) \Rightarrow \text{PMF}$$

$$F(x) = P(X \leq x) = \sum_{y: y \leq x} P(y)$$

\hookrightarrow cdf

x	3	5	6	8
F(x)	0.1	0.2	0.4	0.5
P(X=x)	0.1	0.3	0.4	0.9

add F(x) and prev value

x	4	7	9	11
P(X=x)	0.3	0.4	0.9	1
F(x)	0.3	0.7	0.9	1

add to find pmf

Expected value & Variance of RV

$$E(X) = \sum_{all x} xP(x), \quad E[g(x)] = \sum_{all x} g(x)P(x)$$

$$E[(X-\mu)^2] = \sum_{all x} (x-\mu)^2 P(x) = V(x) = \sigma^2$$

$$E[(X-\mu)^2] = E(X^2) - [E(X)]^2$$

Linearity properties

$$E(ax+b) = aE(X) + b$$

$$V(ax+b) = \sigma_{ax+b}^2 = a^2 \cdot \sigma_X^2 \text{ and } \sigma_{ax+b} = |a| \cdot \sigma_X$$

$$\sigma_{ax} = |a| \cdot \sigma_X, \quad \sigma_{x+b} = \sigma_X$$

$$SD(x+b) = SD(x)$$

Binomial RV

$$b(x; n, p) = \begin{cases} \binom{n}{x} p^x (1-p)^{n-x} & x=0, 1, \dots, n \\ 0 & \text{otherwise} \end{cases}$$

In R: `dbinom(x, size, prob)`

$$X \sim \text{Bin}(n, p) \Rightarrow B(x; n, p) = P(X \leq x) = \sum_{j=0}^x b(j; n, p)$$

$x=0, 1, \dots, n$

$$E(x) = np, \quad V(x) = np(1-p)$$

- Fixed trials, X = # of successes in n trials

Negative Binomial RV

$$nb(x; r, p) = \binom{x+r-1}{r-1} p^r (1-p)^x \quad x=0, 1, \dots$$

$$E(x) = \frac{r(1-p)}{p}, \quad V(x) = \frac{r(1-p)}{p^2}$$

- x = # of failures before r successes
- In R: `pnbinom(x, size, prob)`

Poisson RV

$$P(x; \mu) = \frac{e^{-\mu} \cdot \mu^x}{x!} \quad x=0, 1, \dots, \infty$$

$\mu > 0$

$$E(x) = V(x) = \mu$$

avg # of events in a time frame (μ)

- No fixed trials
- X = # of events in a given time