1. Statistics
    — samples to learn about population

observed subset, size n    complete set of items that interest
                                    investigation

Parameter
    — numerical measurement that describes characteristic of population

Statistic
    — numerical measure that describes a specific characteristic of a sample

Sampling Errors
    — random differences between sample/population
    — cancel out on average
    — decrease as sample size grows

Non Sampling errors
    — Systematic Differences between sample/population
    — Don't necessarily cancel out on average
    — Don't necessarily decrease as sample size grows

Margin of Error

$$ ME \approx 1.96 \cdot \sqrt{\frac{p(1-p)}{n}} $$

Usually Report    $p \pm ME$
— Parameter should relate to what you are interested in

2. Types of Variables
        — Nominal ($=$ or $\neq$)                    } Categorical, qualitative differences
        — Ordinal ($<$ or $>$)

        — Interval ($+$ or $-$) (No natural 0)
        — Ratio ($\times, \div, \%, \log x \dots$ etc) (Absolute 0)    } Numerical, quantitative differences

        — discrete    vs    continuous
    takes values from a discrete  ⌐ any value in a range
        set

# Relative Frequency

- frequency / sample size

# Histogram

- plotting frequency
- on very bar widths ← even have 2 bins w/ different widths

# Measures of Central Tendency

- Mean

Sample Mean $\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$   } Sensitive to outliers

Outlier is an unusually small/large value

- Median

n is odd     $med(x) = x_{\frac{n+1}{2}}$

n is even    $med(x) = \dfrac{x_{\frac{n}{2}} + x_{\frac{n+1}{2}}}{2}$   } Mid point, less sensitive to outliers

- Quantile

  - generalization of median
  - a number $0 < \alpha < 1$ where $\alpha$-quantile $= \alpha \cdot 100\%$

Percentiles, Deciles, Quintiles, Quartiles
1%         10%      20%        25%

- Range

  - measure of variability / spread
  - 100% quantile − 0% quantile, max − min   } Very sensitive to outliers

- Inter Quartile Range (IQR)

  $IQR = 3^{rd}$ Quartile − $1^{st}$ Quartile

  $= 75\% - 25\%$

- Variance

- Variance

$$\sigma^2 \triangleq \frac{1}{n-1} \sum_{i=1}^{n} \left(x_i - \bar{x}\right)^2$$
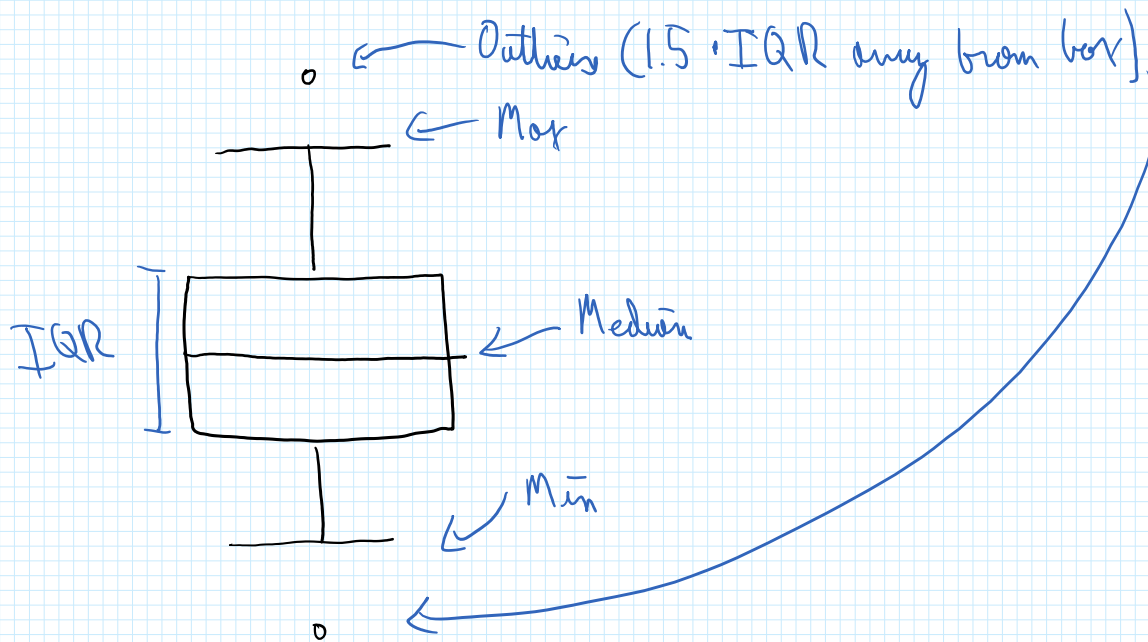
- Standard Deviation

$$\sigma \triangleq \sqrt{\sigma^2}$$

- Mean Absolute Deviation

$$MAD \triangleq \frac{1}{n} \sum_{i=1}^{n} |x_i - \bar{x}|$$

3. - Box Plot
   - visualize all 5 quartiles

o ⟵ ⟵ Outliers (1.5 · IQR away from box)

⟵ Max

IQR

⟵ Median

Min

o

Measures of Symmetry

$$\text{Skewness} := \frac{\frac{1}{n} \sum_{i=1}^{n} \left(x_i - \bar{x}\right)^3}{\sigma^3}$$

{
Negative - Left skewed (usually mean < median)

0 → Symmetry

Positive → right skewed (usually mean > median)
}

- Covariance

$$S_{xy} = \frac{1}{n-1} \sum_{i=1}^{n} \left(x_i - \bar{x}\right)\left(y_i - \bar{y}\right)$$

{
− ⟹ negative linear dependence

0 ⟹ No linear dependence

+ ⟹ positive linear dependence
}

~ Correlation

$$r_{xy} = \frac{s_{xy}}{s_x \cdot s_y}$$

$-1 \leq r \leq 1$   Some rules

— standardized and unitless

# Linear Transformations

$$W_i = b \cdot x_i + a$$

$$\bar{W} = b\bar{x} + a$$

$$s_w^2 = b^2 s_x^2$$

$$s_w = |b| s_x$$

$$s_{wy} = b \cdot s_{xy}$$

$$r_{wy} = \pm r_{xy}$$  } sign in $\pm$ depends on sign of $b$

|  | Sample Statistic | Population Parameter |
|---|---|---|
| Mean | $\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$ | $\mu_y := E(X_i)$ |
| Var | $s_x^2 = \frac{1}{n-1} \sum_{i=1}^{n} (x_i - \bar{x})^2$ | $\sigma^2 = E\left[(X_i - \mu_y)^2\right]$ |
| S.D. | $s_x = \sqrt{s_x^2}$ | $\sigma = \sqrt{\sigma^2}$ |
| Cov. | $s_{xy} = \frac{1}{n-1} \sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})$ | $\sigma_{xy} = E\left[(X_i - \mu_x)(y_i - \mu_y)\right]$ |
| Corr | $r_{xy} = \frac{s_{xy}}{s_x \cdot s_y}$ | $\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \cdot \sigma_y}$ |

# Two Pillars of Statistics

— LLN (Law of Large #'s)
  — as sample size ↑, statistics approximate population parameter better and better

— CLT (Central Limit Theorem)

— if sample size is sufficiently large, difference between sample statistics and population parameter follows a Gaussian distribution

# 4. Random Experiment
— one whose outcomes are random

## Basic Outcome
— finest grained relevant outcomes of a random experiment

## Sample Space ($\Omega$)
— set of all possible basic outcomes

## Event (E)
— subset of Sample space

## Occurrence
— After experiment, only one basic outcome will happen
— lets call is $\omega_{realized}$
— E has "occurred" if $\omega_{realized} \in E$

## Probability
— defined on events
— $P(E)$ is the probability of E
— Must Satisfy the 3 Axioms

① $0 \leq P(E) \leq 1$

② $P(\Omega) = 1$

③ if $E_1 \rightarrow E_N$ are mutually exclusive (empty intersect)
$$P(E_1 \cup \rightarrow E_N) = P(E_1) + \dots + P(E_N)$$

## Classical Probability
— all outcomes are equally likely
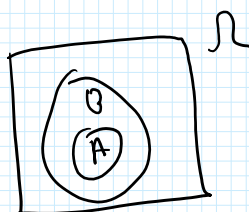
$- \ \mathbb{P}(E) = \dfrac{\#(E)}{N}$

# 5.

Complement

$$\mathbb{P}(E^c) = 1 - \mathbb{P}(E)$$

The Inclusion Rule
- event A logically implies event B
- $A \subseteq B$

If $A \subseteq B$ then $\mathbb{P}(A) \leq \mathbb{P}(B)$



The Union (Logical Addition) Rule

$$\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B)$$

Conditional Probability

$$\mathbb{P}(B/A) := \dfrac{\mathbb{P}(A \cap B)}{\mathbb{P}(A)}$$

$\Big\}$ Conditional probability of B given A

- can treat $\mathbb{P}(\cdot | A)$ as a restricted sample space
- must follow 3 axioms

Multiplication Rule

$$\mathbb{P}(A \cap B) = \mathbb{P}(B/A) \cdot \mathbb{P}(A)$$

comes from above

# 6. Statistical Independence
- A and B are statistically independent if

$$\mathbb{P}(A \cap B) = \mathbb{P}(A) \cdot \mathbb{P}(B)$$

$$\therefore \ \mathbb{P}(B) = \mathbb{P}(B | A)$$

# Pairwise Independence

$$E_1, E_2 \cdots E_n \quad \} \text{ For a list of events}$$

$$P(E_i \cap E_j) = P(E_i) \cdot P(E_j) \quad \} \begin{array}{l} \text{If all combos of pairs} \\ \text{are statistically independent} \\ \text{the list is pairwise independent} \end{array}$$

# Mutual Independence

For a similar list of events as ↑

$$P(E_1 \cap E_2 \cap \cdots \cap E_K) = P(E_1) \cdot P(E_2) \cdots P(E_K) \quad \} \begin{array}{l} \text{for only sublist of} \\ \text{events} \end{array}$$

# Law of Total Probability

$$E_1, E_2 \cdots E_K$$

Mutually Exclusive if $E_i \cap E_j = \emptyset \quad \} \text{For all pairs in the list}$

Exhaustive if $E_1 \cup E_2 \cdots \cup E_K = \Lambda$

∴ If $E_1 \to E_K$ are mutually exclusive and exhaustive

$$P(A) = P(A/E_1) \cdot P(E_1) + P(A/E_2) P(E_2) \cdots + P(A/E_K) \cdot P(E_K)$$

# Bayes Rule

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

Permutations: $\dfrac{n!}{(n-K)!}$

(order matters)

Combinations: $\dfrac{n!}{K! (n-K)!}$

(order doesn't matter)

$\} \begin{array}{l} N \text{ possibilities} \\ K \text{ slots} \end{array}$

# 2. Random Variable (X)

$\Omega$ is a sample space

An RV, $X$, is a fixed function that maps each basic outcome $\omega \in \Omega$ to a real #, $X(\omega)$

(can be discrete, continuous, or mixed)

## Realization (x)

- a particular numerical value in $\mathbb{R}$ that a RV takes on

$$\{X = x\} := \{\omega : X(\omega) = x\}$$

## Support

- the set of all possible realizations of a RV

$$Supp(X) := \{X(\omega) : \omega \in \mathbb{R}\}$$

(Set of all realizations)

## Probability Mass Function

Given $\Omega$, $\mathbb{P}$, $X$

$$p(x) = \mathbb{P}(X = x)$$

← Get out the probability it owns

↑ Plug in a realization

$$0 \leq p(x) \leq 1$$

$$\sum_{i=1}^{\hat{} } p(x_i) = 1$$ } — sum of pmf for all realizations in the support = 1

## Cumulative Distribution (CDF)

$$F(x_0) := \mathbb{P}(X \leq x_0)$$

Properties of CDF
- always increasing

$-F(-\infty) = 0$

$-F(\infty) = 1$

## Expectation

$$\mu_x = \mathbb{E}(X) := \sum_{x \in Supp(x)} x \cdot p(x)$$

$\Big\{$ The sum of all realizations in the support of $x$ multiplied by their probability

**8.** RV function on an RV?

$$g(X) \Rightarrow g(X(\omega))$$

Say $Y := g(X)$

$$Supp(Y) = \{g(x_1), g(x_2) \cdots g(x_n)\}$$

$$P_y = \mathbb{P}(g(x) = y) = \sum_{x: g(x) = y} P_x(x)$$

$$\mathbb{E}[Y] = \sum_{y \in supp(y)} y P_y(y) \quad or \quad \mathbb{E}[Y] = \sum_{x \in supp(x)} g(x) P_x(x)$$

$$\boxed{\mathbb{E}[g(x)] \neq g(\mathbb{E}(x))}$$

only true if $g$ is linear

$$\mathbb{E}(bx + a) = b \cdot \mathbb{E}(x) + a$$

Variance

$$var(x) = \sigma^2 = \mathbb{E}\left[(x - \mu)^2\right] = \mathbb{E}\left\{(x - \mathbb{E}(x))^2\right\}$$

$$\sigma^2 = \mathbb{E}(x^2) - (\mathbb{E}(x))^2$$

Standard Deviation

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

Linear Transformations?

$$Var(ax+b) = a^2 Var(X)$$

\# Try to use definitions whenever you can

Bernoulli Trial

– RV w/ parameter $P$ where

$$P(1) = P(X=1) = P$$
$$P(0) = P(X=0) = 1-P$$

$\}$ 2 outcomes, where "success" is 1 and has probability $P$, "failure" is 0 and has probability $1-P$

$$X \sim Ber(p)$$
$$E(Ber(p)) = p$$
$$Var(Ber(p)) = p(1-p)$$

Counts Multiple Bernoulli Trials ($n$ trials)

$$X \sim Binomial(n, p)$$

$$P(X) = \binom{n}{x} p^x (1-p)^{n-x}$$

← all of combinations $\frac{n!}{x!(n-x)!}$

↑ Probability of $x$ "successes" in $n$ trials

$$\frac{4!}{2!(4-2)!} = 6 \cdot \left(\frac{1}{3}\right)^x \left(\frac{2}{3}\right)^2 \quad \frac{8}{27}$$

$$\mathbb{E}(X) = n \cdot p$$

RV

– defined over $\Omega$, depend on $\omega$ realized

Constants

– don't depend on $\omega$

– $\mathbb{E}(X)$ is a constant