

Business Statistics - MS6107E

Arjun Anil Kumar

April 8, 2025



## 2 / 106

# Statistical Thinking [1]

- Statistics is the Science of the Data.
- Statistics helps to describe and understand **variability**.
- By variability, we mean that successive observations of a system or phenomenon do not produce exactly the same result.
- For example : The observed variability in gasoline mileage depends on many factors, such as the type of driving that has occurred most recently (city versus highway), the changes in condition of the vehicle over time, the brand and/or octane number of the gasoline used. These factors represent potential sources of **variability** in the system.



# Random Experiment

## Definition

An experiment that can result in different outcomes, even though it is repeated in the same manner every time, is called a **random** experiment.



# Sample Space & Event

## Definition

The set of all possible outcomes of a random experiment is called the sample space of the experiment. The sample space is denoted as  $S$ .

## Definition

A sample space is discrete if it consists of a finite or countable infinite set of outcomes. A sample space is continuous if it contains an interval (either finite or infinite) of real numbers.

## Definition

An event is a subset of the sample space of a random experiment.



# Basic Set Operations

## Definition

- The union of two events is the event that consists of all outcomes that are contained in either of the two events. We denote the union as  $E_1 \cup E_2$ .
- The intersection of two events is the event that consists of all outcomes that are contained in both of the two events. We denote the intersection as  $E_1 \cap E_2$ .
- The complement of an event in a sample space is the set of outcomes in the sample space that are not in the event. We denote the complement of the event  $E$  as  $E^c$ .



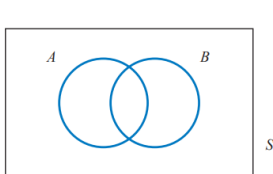
# Mutually Exclusive

## Definition

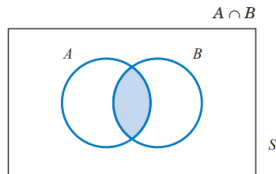
Two events, denoted as  $E_1$  and  $E_2$ , such that are said to be mutually exclusive if  $E_1 \cap E_2 = \emptyset$



# Venn Diagrams



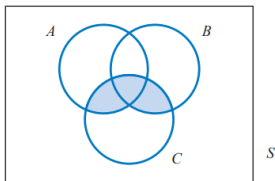
(a)



(b)

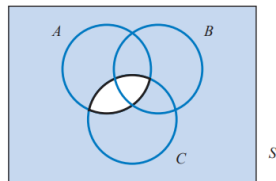
Sample space  $S$  with events  $A$  and  $B$

$$(A \cup B) \cap C$$



(c)

$$(A \cap C)'$$



(d)





# Counting Techniques

## Definition

Assume an operation can be described as a sequence of  $k$  steps, and the number of ways of completing step 1 is  $n_1$ , and the number of ways of completing step 2 is  $n_2$  for each way of completing step 1, and the number of ways of completing step 3 is  $n_3$  for each way of completing step 2, and so forth. The total number of ways of completing the operation is  $n_1 \cdot n_2 \cdot \dots \cdot n_k$ .



# Counting Techniques

## Permutation

A permutation of the elements is an ordered sequence of the elements.

### Definition

The number of permutations of  $n$  different elements is  $n!$  where  $n! = n * (n-1) * (n-2) * (n-3) .. 2 * 1$





# Counting Techniques

## Permutation of Similar Objects

### Definition

The number of permutations of  $n = n_1 + n_2 + ..n_r$  objects of which  $n_1$  are of one type,  $n_2$  are of a second type and  $n_r$  are of an  $r^{th}$  type is  $\frac{n!}{n_1!n_2!..n_r!}$ .



# Interpretations and Axioms of Probability

- Probability is used to quantify the likelihood, or chance, that an outcome of a random experiment will occur.
- The likelihood of an outcome is quantified by assigning a number from the interval  $[0, 1]$  to the outcome (or a percentage from 0 to 100%).
- Higher numbers indicate that the outcome is more likely than lower numbers.
- The probability of an outcome is interpreted as the limiting value of the proportion of times the outcome occurs in  $n$  repetitions of the random experiment as  $n$  increases beyond all bounds - Relative Frequency Interpretation of Probability.
- Whenever a sample space consists of  $N$  possible outcomes that are equally likely, the probability of each outcome is  $\frac{1}{N}$ .





## Addition Rules

## Definition

$$P(E_1 \cup E_2) = P(E_1) + P(E_2) - P(E_1 \cap E_2)$$

## Definition

$$P(E_1 \cup E_2 \cup E_3) = P(E_1) + P(E_2) + P(E_3) - P(E_1 \cap E_2) - P(E_2 \cap E_3) - P(E_1 \cap E_3) + P(E_1 \cap E_2 \cap E_3)$$

## Definition

A collection of events,  $E_1, E_2, E_k$ , is said to be mutually exclusive if for all pairs, then  $P(E_1 \cup E_2 \cup \dots E_k) = P(E_1) + P(E_2) + \dots P(E_k)$



# Conditional Probability

## Definition

The probability of an event B under the knowledge that the outcome will be in event A is denoted as  $P(\frac{B}{A})$ .

## Definition

The conditional probability of an event B given an event A, denoted as  $P(\frac{B}{A})$ , is

$$P(\frac{B}{A}) = \frac{P(A \cap B)}{P(A)} \quad (1)$$

where  $P(A) > 0$









# Concept of Independence

## Definition

In some cases, the conditional probability of  $P(B/A)$  might equal  $P(B)$ . In this special case, knowledge that the outcome of the experiment is in event  $A$  does not affect the probability that the outcome is in event  $B$ .

Two events A and B are independent when one of the three conditions are satisfied.

$$P(A/B) = P(A)$$

$$P(B/A) = P(B)$$

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





# Random Variables

## Definition

A **random variable** is a function that assigns a real number to each outcome in the sample space of a random experiment.

## Definition

A **discrete** random variable is a random variable with a finite (or countably infinite) range. A **continuous** random variable is a random variable with an interval (either finite or infinite) of real numbers for its range.



# Probability Distributions and Probability Mass Functions

## Definition

The probability distribution of a random variable  $X$  is a description of the probabilities associated with the possible values of  $X$ .

## Definition

For a discrete random variable  $X$  with possible values  $x_1, x_2, \dots, x_n$ , a probability mass function is a function such that

$$f(x_i) \geq 0$$

$$\sum_{i=1}^n f(x_i) = 1$$

$$f(x_i) = P(X = x_i)$$



# Cumulative Distributive Function

## Definition

The cumulative distribution function of a discrete random variable  $X$ , denoted as  $F(x)$ , is

$$F(x) = P(X \leq x) = \sum_{x_i \leq x} f(x_i)$$

$$F(x) \leq 1$$

$$\text{If } (x < y), \text{ then } F(x) \leq F(y)$$



# Mean and Variance of Discrete Random Variable

## Definition

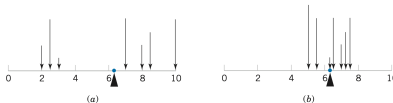
The mean or expected value of the discrete random variable  $X$ , denoted as  $E[X]$ , is

$$E[X] = \sum_x xf(x)$$

## Definition

The variance of  $X$  denoted as  $\sigma^2$  or  $V(X)$  is

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







# Discrete Distributions

## Discrete Uniform Distribution

### Definition

A random variable  $X$  has a discrete uniform distribution if each of the  $n$  values in its range, say,  $x_1, x_2, \dots, x_n$ , has equal probability. Then,  $f(x_i) = \frac{1}{n}$



# Discrete Distribution

## Binomial Distribution

### Examples of Binomial Random Variables

- Flip a coin 10 times. Let  $X$  be the number of heads obtained.
- A multiple-choice test contains 10 questions, each with four choices, and you guess at each question. Let  $X$  be the number of questions answered correctly.
- In the next 20 births at a hospital, let  $X$  be the number of female births.
- Of all patients suffering a particular illness, 35% experience improvement from a particular medication. In the next 100 patients administered the medication, let  $X$  be the number of patients who experience improvement.



# Discrete Distribution

## Binomial Distribution

### Definition

A random experiment consists of  $n$  Bernoulli trials such that The trials are independent. Each trial results in only two possible outcomes, labeled as “success” and “failure”. The probability of a success in each trial, denoted as  $p$ , remains constant.

The random variable  $X$  that equals the number of trials that result in a success has a binomial random variable with parameters  $0 \leq p \leq 1$  and  $n = 1, 2, \dots$ . The probability mass function of  $X$  is

$$f(x) = \binom{n}{x} p^x \cdot (1 - p)^{n-x}$$

where  $x = 0, 1, 2, 3, \dots, n$



## Discrete Distribution

## Binomial Distribution

- $\binom{n}{x}$  equals the total number of different sequences of trials that contain  $x$  successes and  $n - x$  failures.
- The total number of different sequences that contain  $x$  successes and  $n - x$  failures times the probability of each sequence equals  $P(X=x)$ .
- The sum of probability mass function of binomial distribution is 1.
- For a fixed  $n$ , the distribution becomes more symmetric as  $p$  increases from 0 to 0.5 or decreases from 1 to 0.5.
- For a fixed  $p$ , the distribution becomes more symmetric as  $n$  increases.



# Discrete Distribution

## Binomial Distribution

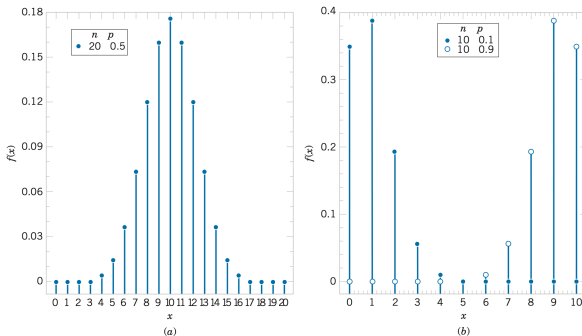


Figure: Binomial Distribution for selected  $p$  and  $n$



# Discrete Distribution

## Binomial Distribution

### Definition

If  $X$  is a binomial random variable with parameters  $p$  and  $n$ ,

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

$$\sigma^2 = V[X] = np.(1 - p)$$

















# Discrete Distribution

## Hypergeometric Distribution

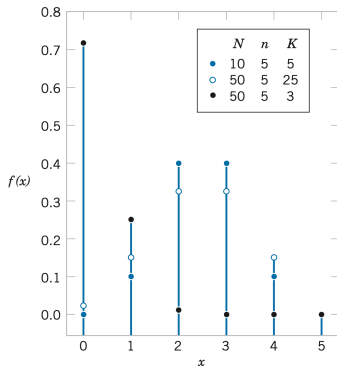


Figure: Hypergeometric Distributions for different  $N, K$  and  $n$























# Mean and Variance of Continuous Random Variable

## Definition

Suppose  $X$  is a continuous random variable with probability density function  $f(x)$ . The mean or expected value of  $X$ , denoted as  $\mu$  or  $E(X)$ , is  $\mu = E[X] = \int_{-\infty}^{\infty} x \cdot f(x) dx$

## Definition

The variance of  $X$  is denoted by  $V(X)$  or  $\sigma^2$   
 $\sigma^2 = V(X) = \int_{-\infty}^{\infty} (x - \mu)^2 \cdot f(x) dx$

$$\sigma^2 = E[X^2] - E[X]^2$$

The square root of variance ( $\sigma^2$ ) is the standard deviation ( $\sigma$ ).







## Continuous Uniform Distribution

## Definition

A continuous random variable  $X$  with probability density function  $f(x) = \frac{1}{b-a}$ , where  $a \leq x \leq b$  is a continuous uniform random variable

## Definition

If  $X$  is a continuous uniform random variable over  $a \leq x \leq b$

$$\mu = \frac{a+b}{2}$$

$$\sigma^2 = \frac{(b-a)^2}{2}$$











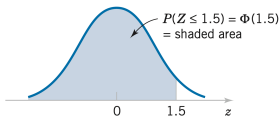
# Standard Normal Variable

## Definition

A normal random variable with  $\mu = 0$  and  $\sigma = 1$  is called a standard normal variable and is denoted as  $Z$ .

## Definition

The cumulative distribution function of a standard normal random variable is denoted as  $\phi(z) = P(Z \leq z)$

































# Descriptive Statistics

```
data <- c(1, 1, 1, 1, 4, 4, 4, 4, 8, 8, 8, 8, 12, 12, 12, 12)
summary(data)
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
1.00	3.25	6.00	6.25	9.00	12.00

Figure: Descriptive Statistics

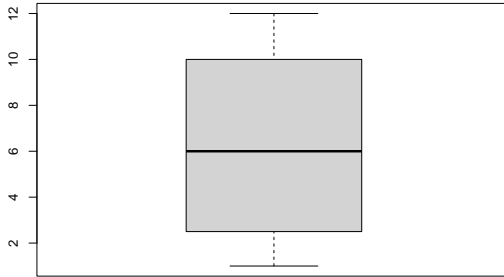


# Descriptive Statistics

- Descriptive Statistics summarizes the data set.
- The central tendency of any data distribution can be estimated using the **mean, mode and median**.
- The variation of any data distribution around the mean can be estimated using the **variance** or **inter-quartile range (IQR)**  
 $= (Q_{.75} - Q_{.25})$ .
- - Mean - Average of a data set.
  - Median - The data point that divides the data into equal halves.
  - Maximum/Minimum - The maximum/minimum of a data set.
  - Quantile -  $Q_\lambda$  separates the data set into two halves where the left half would contain around  $100.\lambda\%$  of the data points and the right half would contain around  $100.(1 - \lambda)\%$  of the data points. ( $0 \leq \lambda \leq 1$ )



# Boxplot



Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
1.00	3.25	6.00	6.25	9.00	12.00



# Boxplot

- Box-Whisker plot - Box + Whiskers
- Boxplot - Whiskers are drawn at  $Q_{.75} + 1.5 \cdot \text{IQR}$  and  $Q_{.25} - 1.5 \cdot \text{IQR}$ , where IQR is the inter-quartile range ( $Q_{.75} - Q_{.25}$ ).
- Higher IQR indicates higher variance in the data
- Boxplot does not have mean in the plot
- The 5 main points of a box plot are whisker 1, whisker 2, median  $Q_{.5}$  and the two quantiles ( $Q_{.25}$  &  $Q_{.75}$ ).





# Boxplot and Outliers

$c(100, 110, 110, 110, 120, 120, 130, 140, 140, 150, 170, 220)$   
 $Q_{.75} = 142.5, Q_{.25} = 110, [L, U] = [61.25, 191.25]$

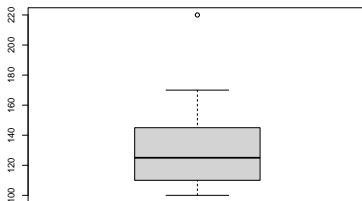


Figure: Interquartile Range  $Q_{.75} - Q_{.25} = 32.5$



# Symmetric Distributions

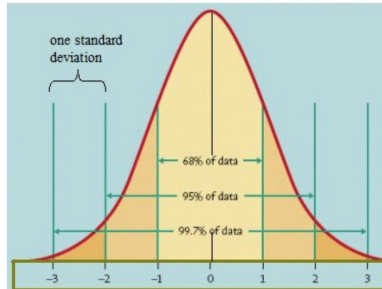


Figure: Symmetric Distribution

A **symmetric** distribution is a type of distribution where the left side of the distribution mirrors the right side.



# Symmetric Distributions

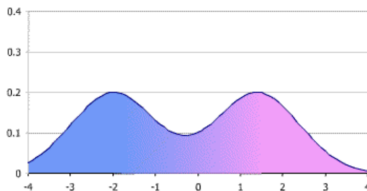
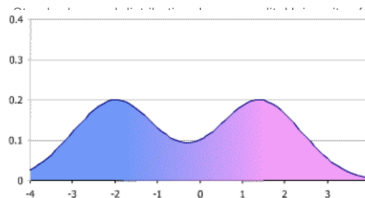
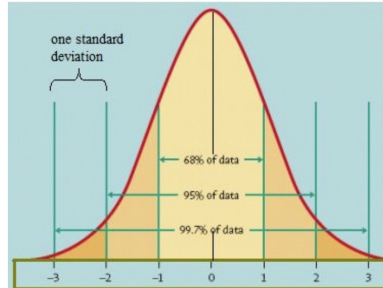


Figure: Symmetric Distributions can be Bi-Modal too!



# Symmetric Distributions - Mean, Median, Mode



# Symmetric Distributions

## Definition

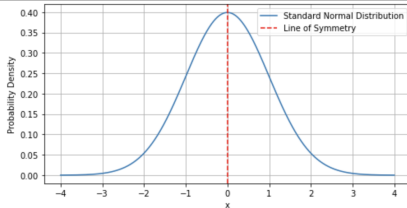
A probability distribution is said to be symmetric if and only if there exists a value  $x_0$ , such that  $f(x_0 + \delta) = f(x_0 - \delta)$ , for  $\forall \mathbb{R}$ , where where  $f$  is the probability density function if the distribution is continuous or the probability mass function if the distribution is discrete.

Symmetric Distributions can be either discrete or continuous.



# Symmetric Distributions

## Continuous



## Discrete

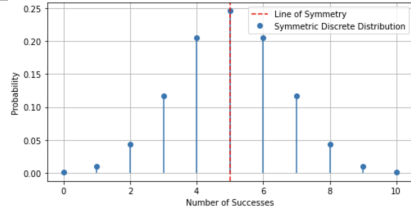


Table: Symmetric Distributions









# Bar Plot

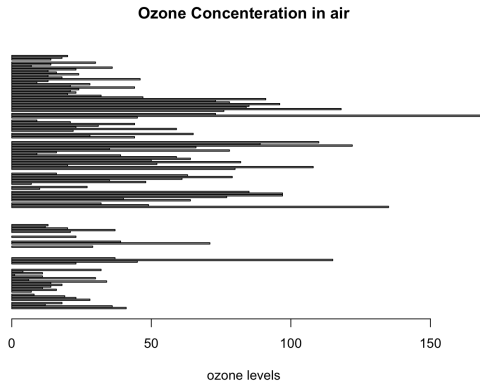


Figure: Bar Plot - Ozone Concentration in Air



# Histogram

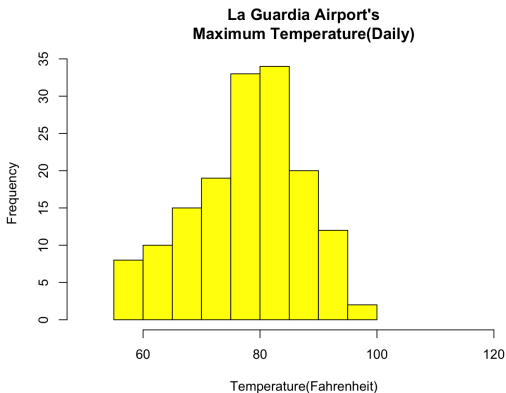
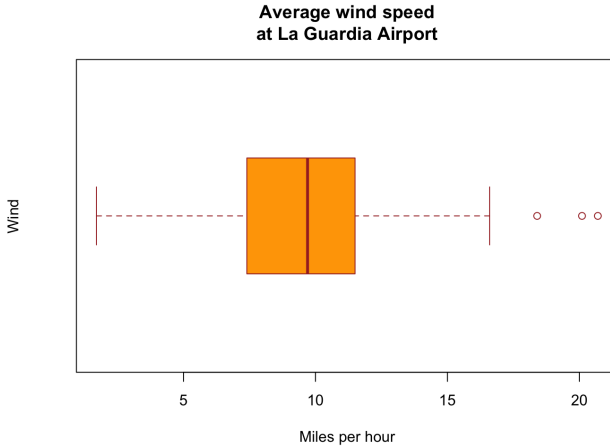


Figure: Histogram - Airport Maximum Temperature

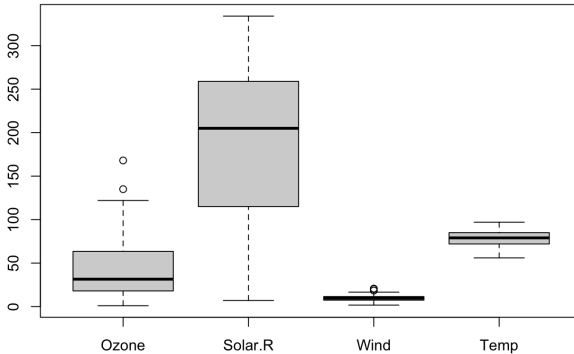


# Boxplot

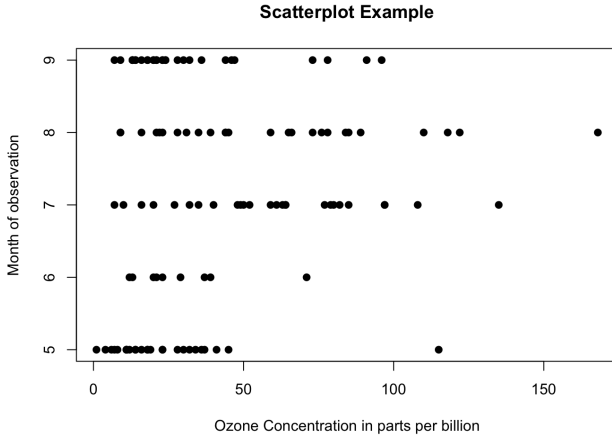


# Multiple Boxplot

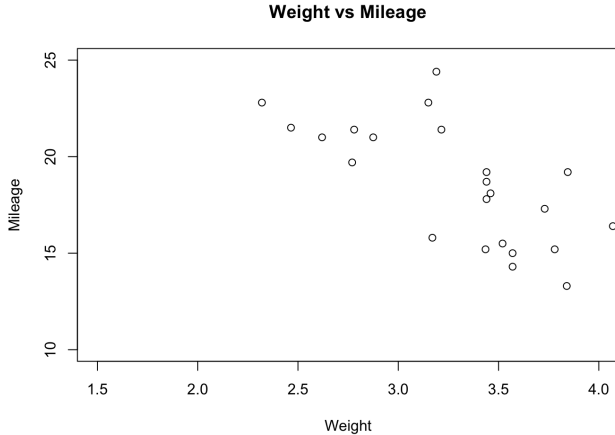
Box Plots for Air Quality Parameters



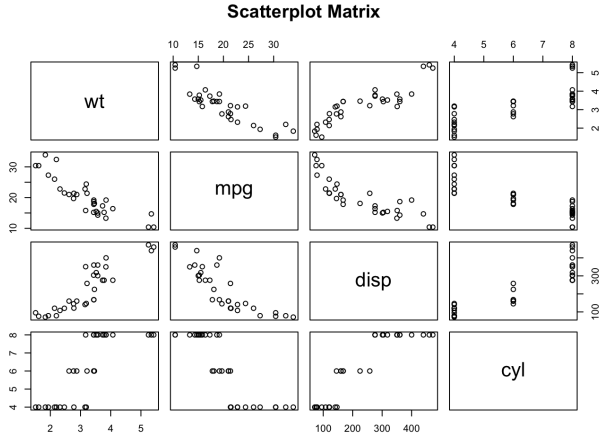
# Scatter Plot



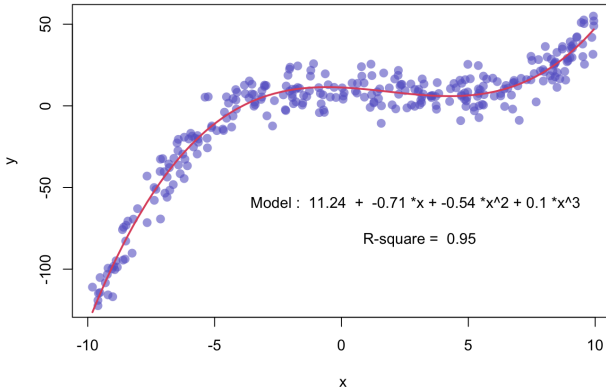
# Scatter Plot



# Scatter Plot Matrix

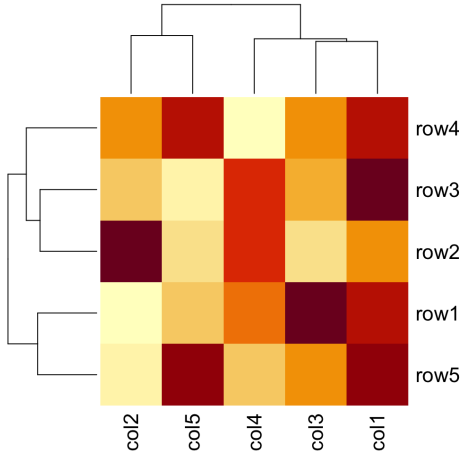


# Polynomial Scatter Plot



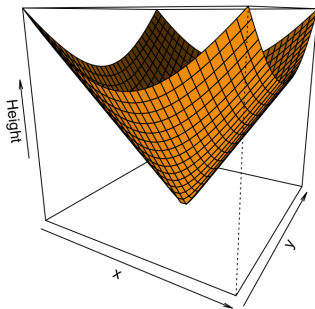


# Heat Map



# 3-D Plot

3d plot



# Univariate Graphs

## Simple Bar Chart

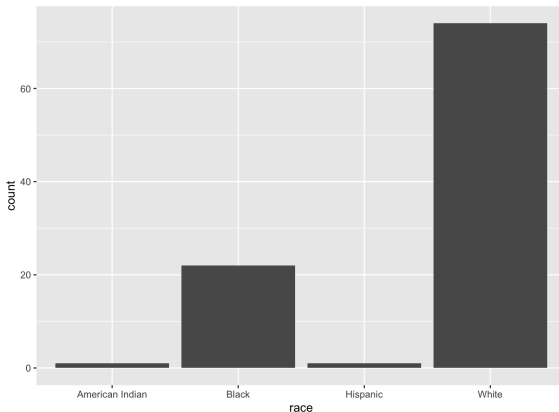


Figure: Barplot - Variable - Race



# Univariate Graphs

## Distribution - Categorical

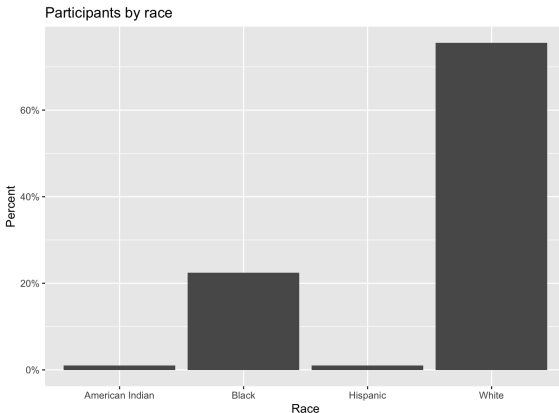


Figure: Distribution - Variable - Race



# Univariate Graphs

## Pie Chart - Categorical

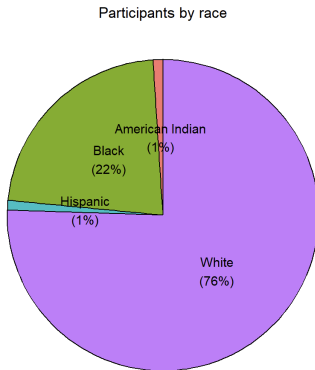


Figure: Pie Chart - Variable - Race



# Univariate Graphs

## Tree Map - Categorical

Marriages by officiate

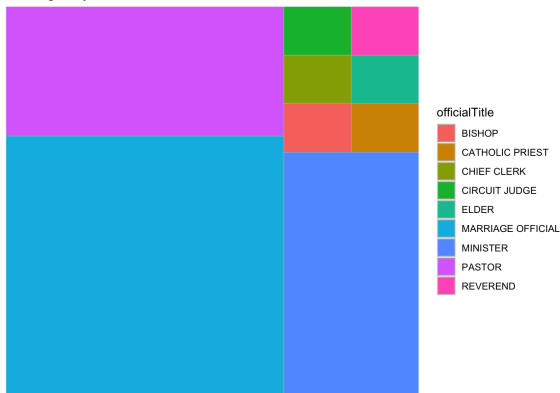


Figure: Tree Map - Marriages by Officiate



# Univariate Graphs

## Waffle Chart - Categorical

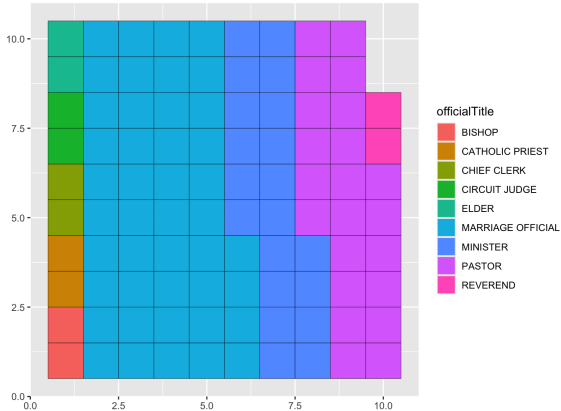


Figure: Waffle Chart - Marriages by Officiate



# Univariate Graphs

## Histogram - Continuous

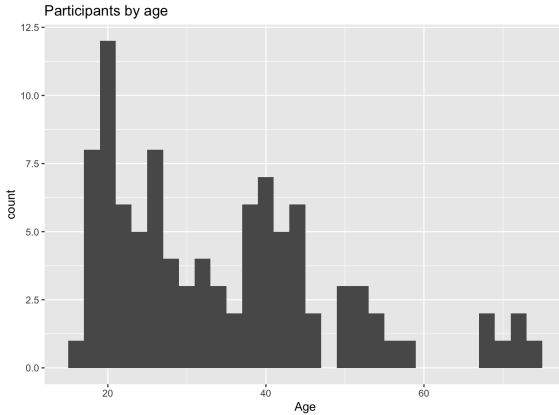


Figure: Histogram - Age





# Univariate Graphs

## Kernel Density - Continuous

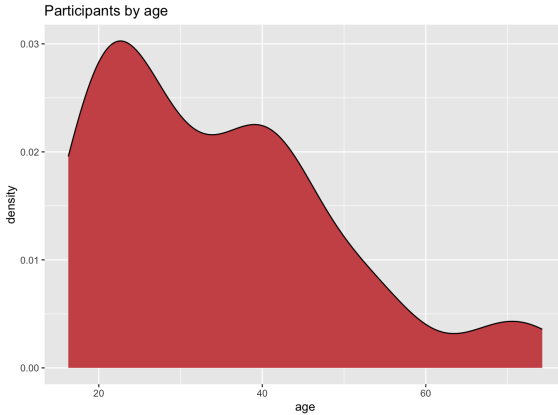


Figure: Kernel Density - Age



# Normal Distribution - 1

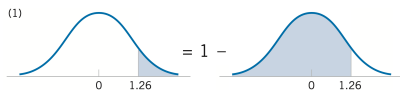


Figure:  $P(Z > 1.26)$

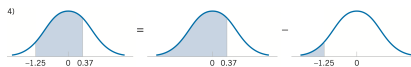


Figure:  $P(1.25 < Z < 0.372)$





# Symmetric Distributions

Case 1 - "A newly started IIM advertisement says "Mean salary is 18 lakh per annum". A candidate joins the institute and the faculty advisor tells him that the median salary is 6 lakhs. "Is there a cause of grave concern?"

Case 2 - "A premier national engineering college advertisement says "Mean salary is 18 lakh per annum". A candidate joins the institute and the faculty advisor tells him that the median salary is 25 lakh. "Is there a cause of grave concern?"



## References I



Douglas C Montgomery and George C Runger.  
*Applied statistics and probability for engineers.*  
John wiley & sons, 2010.

