

STAT 305: Chapter 3

Elementary Descriptive Statistics

Amin Shirazi

Course page:
ashirazist.github.io/stat305_s2020.github.io

Section 3.1

Elementary Graphical and Tabular Treatment of Quantitative Data

Summarizing Univariate Data

Introduction: Creative Writing Workshops

Two methods of teaching a creative writing workshop are being studied for their effectiveness of improving writing skills. First, two groups of creative writing students who were randomly assigned to one of two different 3-hour workshops. At the end of the workshop, the students were given a standard creative writing test and their score on the test was recorded.

Exam Scores for Two Groups of Students Following Different Courses

Group 1	Group 2
74 79 77 81	65 77 78 74
68 79 81 76	76 73 71 71
81 80 80 78	86 81 76 89
88 83 79 91	79 78 77 76
79 75 74 73	72 76 75 79

Exam Scores for Two Groups of Students Following Different Courses

Group 1	Group 2
74 79 77 81	65 77 78 74
68 79 81 76	76 73 71 71
81 80 80 78	86 81 76 89
88 83 79 91	79 78 77 76
79 75 74 73	72 76 75 79

We may have several questions we are interested in answering using this data. For instance,

- Which group did better on average?
- Which group has the most consistent scores?
- Were there any unusually low or high scores in either group?
- If we ignore unusual scores, which group is better?
- Which group had the most scores over 80?
- ...

However, none of these are immediately clear looking at the raw recorded data.

The Purpose of Summaries

Certain questions can and should be asked across many types of experiments.

But seeing data in this kind of *flat* format presents lots of problems for a person trying to understand the relationship between the two groups.

Summaries are tools (mainly mathematical or graphical) which help researchers quickly understand the data they have collected.

The purpose of a summary is to faithfully present aspects of the data in such a way that

- we are capable of answering the types of core questions about the data asked on the previous page,
- we are able to identify more complicated aspects of the data that we may want to investigate further.

Key Idea: Good summaries should be quickly interpreted, provide clear insight into the data, and be widely applicable.

Descriptive statistics

Descriptive statistics

Engineering data are always variable. Given precise enough measurement, even constant process conditions produce different responses. Thus, it is not the individual data values that are important, but their **distribution**. We will discuss simple methods that describe important distributional characteristics of data.

Descriptive statistics is the use of plots and numerical summaries to describe data without drawing any formal conclusions.

Descriptive statistics

Through the use of *descriptive statistics*, we seek to find the following features of data sets:

- **Center** : The point that the data are closest on average
- **Spread**: how wide the data look, how varied the points are
- **Shape** : common patterns/ trends that are present in the data
- **Outliers**: points that lie way beyond the rest of the data (wierd points)

Graphical and tabular displays of quantitative data

Almost always, the place to start a data analysis is with appropriate graphical and tabular displays. When only a few samples are involved, a good plot can tell most of the story about data and drive an analysis.

Dot diagrams and stem-and-leaf plots

When a study produces a small or moderate amount of **univariate quantitative data**, a *dot diagram* can be useful.

A **dot diagram** shows each observation as a dot placed at the position corresponding to its numerical value along a number line.

Summarizing

Intro

Purpose

Descriptive statistics

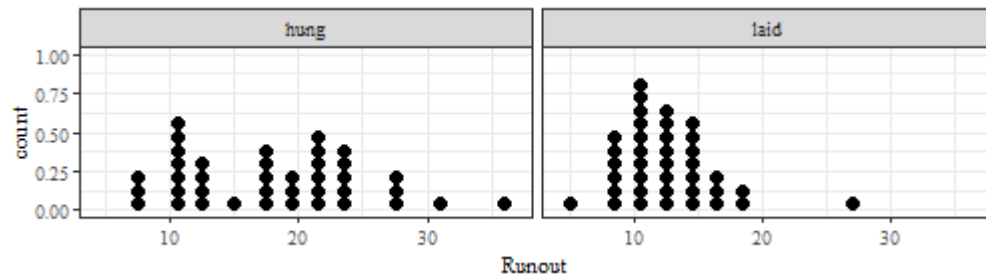
Plots

Dot diagrams and stem-and-leaf plots

Example:[Heat treating gears, cont'd]

Recall the example from Chapter 1. A process engineer is faced with the question, "How should gears be loaded into a continuous carburizing furnace in order to minimize distortion during heat treating?"

The engineer conducts a well-thought-out study and obtains the runout values for 38 gears laid and 39 gears hung.



Summarizing

Intro

Purpose

Descriptive statistics

Plots

Simple Graphical Summaries

Exaple:

Group 1	Group 2
74 79 77 81	65 77 78 74
68 79 81 76	76 73 71 71
81 80 80 78	86 81 76 89
88 83 79 91	79 78 77 76
79 75 74 73	72 76 75 79

Simple graphical summaries aim to provide a better view of the entire set of data. The best graphs are able to make important points clearly and give valuable insights with closer study. Producing good graphs is an **art**.

Two common graphical summaries

- Dot Diagrams
- Stem and Leaf Diagrams

Carries much the same visual information as a dot diagram while preserving the original values exactly

.left-column[layout:false

Summarizing

Simple Graphical Summaries

Intro

Purpose

Descriptive
statistics

Plots

Group 1

74 79 77 81

68 79 81 76

81 80 80 78

88 83 79 91

79 75 74 73

Group 2

65 77 78 74

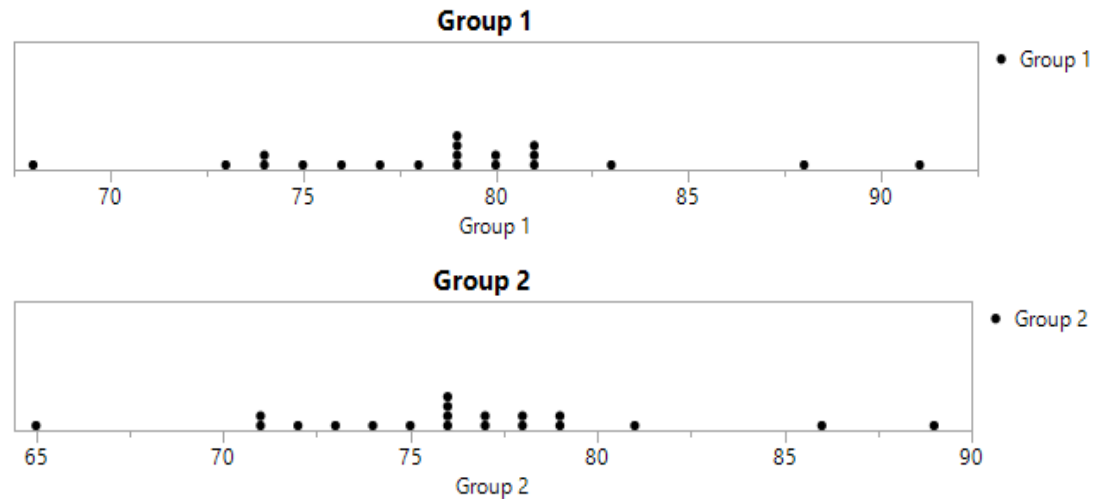
76 73 71 71

86 81 76 89

79 78 77 76

72 76 75 79

Dot Diagrams



Summarizing

Intro

Purpose

Descriptive
statistics

Plots

Dot diagrams are good for getting a general feel for the data (and can be done with pencil and paper), but do not allow the recovery of the exact values used to make them.

A **stem-and-leaf plot** is made by using the last few digits of each data point to indicate where it falls.

Summarizing

Intro

Purpose

Descriptive statistics

Plots

Example:

Group 1	Group 2
74 79 77 81	65 77 78 74
68 79 81 76	76 73 71 71
81 80 80 78	86 81 76 89
88 83 79 91	79 78 77 76
79 75 74 73	72 76 75 79

Stem and Leaf Diagrams

Stem and Leaf		
Stem	Leaf	Count
9	1	1
8	8	1
8		
8		
8	3	1
8	00111	5
7	89999	5
7	67	2
7	445	3
7	3	1
7		
6	8	1

6|8 represents 68

Stem and Leaf		
Stem	Leaf	Count
8	9	1
8	6	1
8		
8		
8	1	1
7	8899	4
7	666677	6
7	45	2
7	23	2
7	11	2
6		
6		
6	5	1

6|5 represents 65

Summarizing

Intro

Purpose

Descriptive statistics

Plots

Freq Tables

Frequency tables and histograms

Dot diagrams and stem-and-leaf plots are useful for getting to know a data set, but they are not commonly used in papers and presentations.

A **frequency table** is made by first breaking an interval containing all the data into an appropriate number of smaller intervals of equal length. Then tally marks can be recorded to indicate the number of data points falling into each interval. Finally, frequencies, relative frequencies, and cumulative relative frequencies can be added.

Summarizing

Intro

Purpose

Descriptive statistics

Plots

Freq Tables

Frequency tables and histograms

A frequency table is made by

- First breaking an interval containing all the data into an appropriate number of smaller intervals of **equal length**.
- Then tally marks can be recorded to indicate the number of data points falling into each interval.
- Finally, add frequency, relative frequency and cumulative relative frequency can be added.

Table 3.2

Frequency Table for Laid Gear Thrust Face Runouts

Runout (.0001 in.)	Tally	Frequency	Relative Frequency	Cumulative Relative Frequency
5-8		3	.079	.079
9-12		18	.474	.553
13-16		12	.316	.868
17-20		4	.105	.974
21-24		0	0	.974
25-28		1	.026	1.000
		38	1.000	

Summarizing

Intro

Purpose

Descriptive statistics

Plots

Freq Tables

Histogram

Histogram

After making a frequency table, it is common to use the organization provided by the table to create a histogram.

A **(frequency or relative frequency) histogram** is a kind of bar chart used to portray the shape of a distribution of data points.

Guidelines for making histograms:

- Use intervals of equal length
- Show the entire vertical axis starting at *zero*
- Avoid breaking either axes
- keep a uniform scale for axes (tick marks)
- Center bars of appropriate heights at midpoint of the intervals

Summarizing

Intro

Purpose

Descriptive
statistics

Plots

Freq Tables

Histogram

Histogram

Example:[Bullet penetration depth, pg. 67]

Sale and Thom compared penetration depths for several types of .45 caliber bullets fired into oak wood from a distance of 15 feet. They recorded the penetration depths (in mm from the target surface to the back of the bullets) for two bullet types.

200 grain jacketed bullets	230 grain jacketed bullets
63.8, 64.65, 59.5, 60.7, 61.3, 61.5, 59.8, 59.1, 62.95, 63.55, 58.65, 71.7, 63.3, 62.65, 67.75, 62.3, 70.4, 64.05, 65, 58	40.5, 38.35, 56, 42.55, 38.35, 27.75, 49.85, 43.6, 38.75, 51.25, 47.9, 48.15, 42.9, 43.85, 37.35, 47.3, 41.15, 51.6, 39.75, 41

Summarizing

Intro

Purpose

Descriptive
statistics

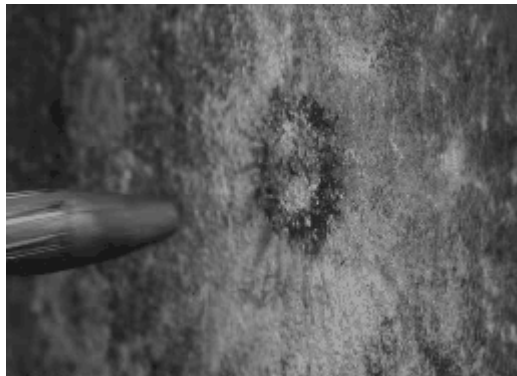
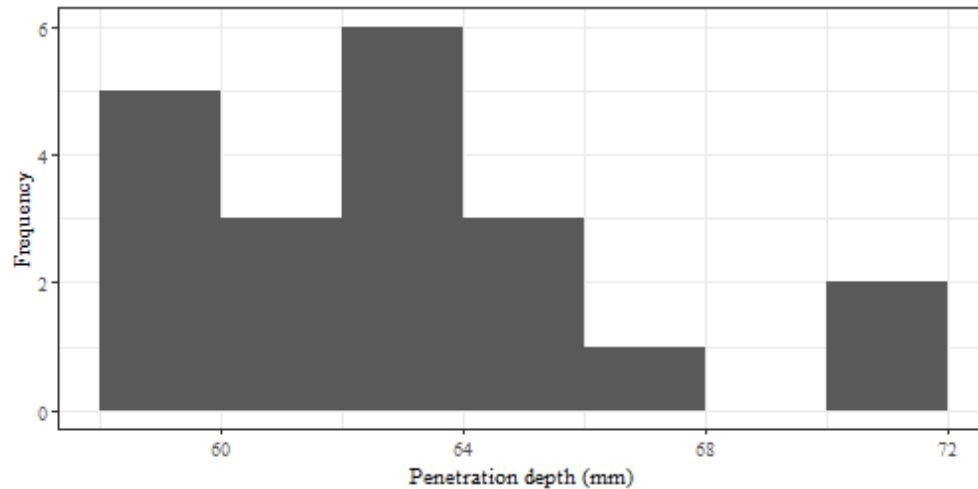
Plots

Freq Tables

Histogram

Histogram

Example:[Bullet penetration depth, pg. 67]



Summarizing

Intro

Purpose

Descriptive
statistics

Plots

Freq Tables

Histogram

Example:[Histogram]

Suppose you have the following data:

74, 79, 77, 81, 68, 79, 81, 76, 81, 80, 80

, 78, 88, 83, 79, 91, 79, 75, 74, 73

Create the corresponding *frequency table* and *frequency histogram*}

Why plotting data?

Summarizing

Intro

Purpose

Descriptive statistics

Plots

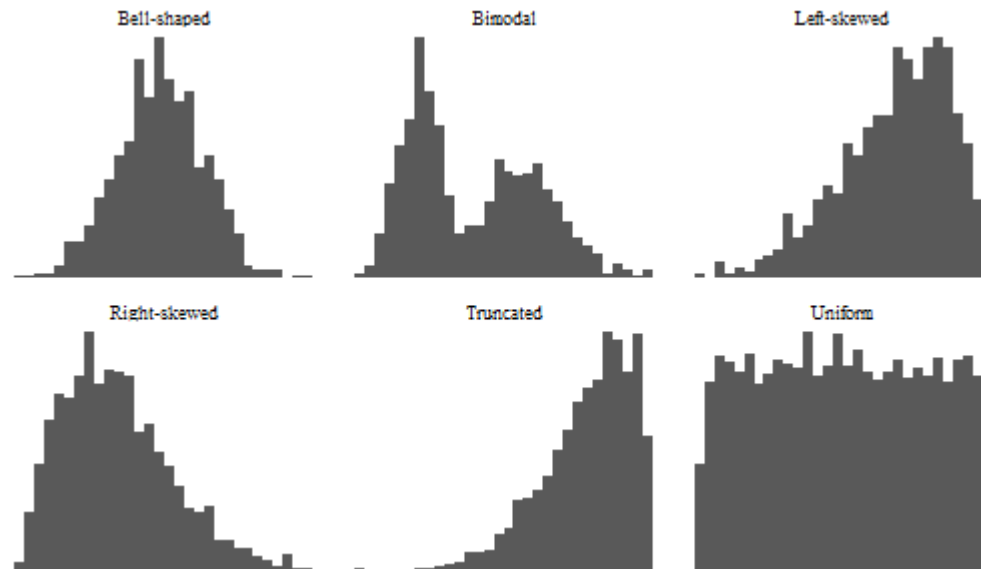
Freq Tables

Histogram

Why plotting data?

Why do we plot data?

Information on **location**, **spread**, and **shape** is portrayed clearly in a histogram and can give hints as to the functioning of the physical process that is generating the data.



Summarizing

Intro

Purpose

Descriptive
statistics

Plots

Freq Tables

Histogram

Why plotting data?

If data on the diameters of machined metal cylinders purchased from a vendor produce a histogram that is decidedly **bimodal**, this suggests

- the machining was done on 2 machines or by two operators or at two different times, etc. ...

If the histogram is **truncated**, this might suggest

- the cylinders have been 100% inspected

Summarizing

Intro

Purpose

Descriptive statistics

Plots

Freq Tables

Histogram

Scatter plots

Scatter plots

Dot-diagrams, stem-and-leaf plots, frequency tables, and histograms are univariate tools. But engineering questions often concern multivariate data and *relationships between the quantitative variables*.

A **scatterplot** is a simple and effective way of displaying potential relationships between two quantitative variable by assigning each variable to either the x or y axis and plotting the resulting coordinate points.

Summarizing

Intro

Purpose

Descriptive statistics

Plots

Freq Tables

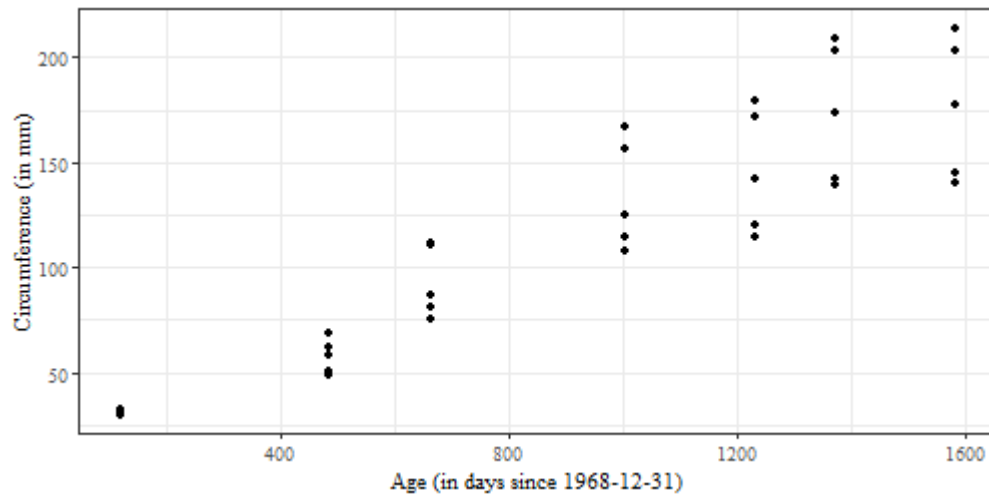
Histogram

Scatter plots

Scatter plots

Example:[Orange trees]

Jim and Jane want to know the relationship between an orange tree's age (in days since 1968-12-31) and its circumference (in mm). They recorded the data for 35 orange trees.



Summarizing

Intro

Purpose

Descriptive statistics

Plots

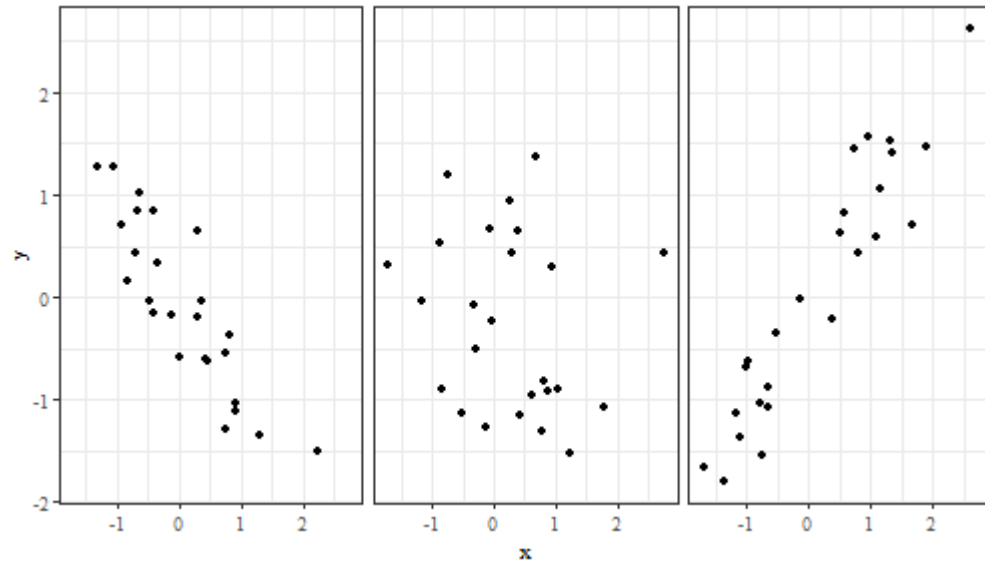
Freq Tables

Histogram

Scatter plots

Scatter plots

There are three typical association/relationship between two variables:



Summaries of Location and Central Tendency

Summarizing

Intro

Purpose

Descriptive statistics

Plots

Freq Tables

Histogram

Scatter plots

Center Stats

Summaries of Location and Central Tendency

Motivated by asking what is *normal/common/expected* for this data. There are three main types used:

Mean: A "fair" center value. The symbol used differs depending on whether we are dealing with a sample or population:

		Mean
Population		$\mu = \frac{1}{N} \sum_1^N x_i$
Sample		$\bar{x} = \frac{1}{n} \sum_1^n x_i$

N is the population size and **n** is the sample size.

Mode: The most commonly occurring data value in set.

Summarizing

Intro

Purpose

Descriptive statistics

Plots

Freq Tables

Histogram

Scatter plots

Center Stats

Summaries of Location and Central Tendency

Quantiles: The number that divides our data values so that the proportion, p , of the data values are below the number and the proportion $1 - p$ are above the number.

Median: The value dividing the data values in half (the middle of the values). The median is just the 50th quantile.

Range: The difference between the highest and lowest values (Range = max - min)

IQR: The Interquartile Range, how spread out is the middle 50% (IQR = $Q3 - Q1$)

Summarizing

Intro

Purpose

Descriptive statistics

Plots

Freq Tables

Histogram

Scatter plots

Center Stats

Summaries of Location and Central Tendency

Group 1	Group 2
74 79 77 81	65 77 78 74
68 79 81 76	76 73 71 71
81 80 80 78	86 81 76 89
88 83 79 91	79 78 77 76
79 75 74 73	72 76 75 79

Calculating Mean Think of it as an equal division of the total

- each value in the data is an " x_i " (i is a **subscript**)
- Group 1: $x_1 = 74, x_2 = 79, \dots, x_{20} = 73$
- The sum: $x_1 + x_2 + x_3 + \dots + x_{20}$
- divides : $(x_1 + x_2 + x_3 + \dots + x_{20})/20$
- Or using summation notation: $\frac{1}{20} \sum_{i=1}^{20} x_i$

Quantiles

Summarizing

Intro

Purpose

Descriptive statistics

Plots

Freq Tables

Histogram

Scatter plots

Center Stats

Quantiles

Summaries of Location and Central Tendency

The Quantile Function

Two useful pieces of notation:

floor: $\lfloor x \rfloor$ is the largest integer smaller than or equal to x

ceiling: $\lceil x \rceil$ is the smallest integer larger than or equal to x

Examples

- $\lfloor 55.2 \rfloor = 55$
- $\lceil 55.2 \rceil = 56$
- $\lfloor 19 \rfloor = 19$
- $\lceil 19 \rceil = 19$
- $\lceil -3.2 \rceil = -3$
- $\lfloor -3.2 \rfloor = -4$

Summarizing

Intro

Purpose

Descriptive
statistics

Plots

Freq Tables

Histogram

Scatter plots

Center Stats

Quantiles

Summaries of Location and Central Tendency

Quantiles

- Already familiar with the concept of "percentile".
 - e.g in the context of reporting scores on exams:

If a person has scored at the 80th percentile, roughly 80% of those taking the exam had worse scores, and roughly 20% had better scores.

- It is more convenient to work in terms of fractions between 0 and 1 rather than percentages between 0 and 100. We then use terminology **Quantiles** rather than percentiles.
- For a number **p** between 0 and 1, the **p quantile** of a distribution is a number such that a fraction p of the distribution lies to the left of that value, and a fraction 1-p of the distribution lies to the right.

Summarizing

Intro

Purpose

Descriptive
statistics

Plots

Freq Tables

Histogram

Scatter plots

Center Stats

Quantiles

Summaries of Location and Central Tendency

The Quantile Function

For a data set consisting of n values that when ordered are $x_1 \leq x_2 \leq \dots \leq x_n$ and $0 \leq p \leq 1$.

We define the **quantile function** $Q(p)$ as:

$$Q(p) = \begin{cases} x_i & \lfloor n \cdot p + .5 \rfloor = n \cdot p + .5 \\ x_i + (np - i + .5) (x_{i+1} - x_i) & \lfloor n \cdot p + .5 \rfloor \neq n \cdot p + .5 \end{cases}$$

(note: this is the definition used in the book - it's just written using *floor* and *ceiling* instead of in words)

Summarizing

Simple Graphs

Freq Tables

Histograms

Center Stats

Quantiles

Summaries of Location and Central Tendency

Example: Find the median, first quartile, 17th quantile and 65th quantile for the following set of data values:

58, 76, 66, 61, 50, 77, 67, 64, 41, 61

First notice that $n = 10$. It is possible helpful to set up the following table:

- **Step 1: sort the data**

data	41	50	58	61	61	64	66	67	76	77
i	1	2	3	4	5	6	7	8	9	10

Summarizing

Simple Graphs

Freq Tables

Histograms

Center Stats

Quantiles

Summaries of Location and Central Tendency

Example: Find the median, first quartile, 17th quantile and 65th quantile for the following set of data values:

58, 76, 66, 61, 50, 77, 67, 64, 41, 61

- **Step 2: find $\frac{i-.5}{n}$**

data	41	50	58	61	61	64	66	67	76	77
i	1	2	3	4	5	6	7	8	9	10
$\frac{i-.5}{n}$	0.05	0.15	0.25	0.35	0.45	0.55	0.65	0.75	0.85	0.95

Summarizing

Simple Graphs

Freq Tables

Histograms

Center Stats

Quantiles

Summaries of Location and Central Tendency

- **Step 3: find $Q(p)$**

data	41	50	58	61	61	64	66	67	76	77
i	1	2	3	4	5	6	7	8	9	10
$\frac{i-.5}{n}$	0.05	0.15	0.25	0.35	0.45	0.55	0.65	0.75	0.85	0.95

$$Q(p) = \begin{cases} x_i & \lfloor n \cdot p + .5 \rfloor = n \cdot p + .5 \\ x_i + (np - i + .5)(x_{i+1} - x_i) & \lfloor n \cdot p + .5 \rfloor \neq n \cdot p + .5 \end{cases}$$

Finding the first **quantile** ($Q(.25)$):

- $np + .5 = 10 \cdot .25 + .5 = 3.$
- since $\lfloor 3 \rfloor = 3$

then $i = 3$ and

$$Q(.25) = x_3 = 58$$

Your turn

Find the median

Summarizing

Simple Graphs

Freq Tables

Histograms

Center Stats

Quantiles

Summaries of Location and Central Tendency

data	41	50	58	61	61	64	66	67	76	77
i	1	2	3	4	5	6	7	8	9	10
$\frac{i-.5}{n}$	0.05	0.15	0.25	0.35	0.45	0.55	0.65	0.75	0.85	0.95

$$Q(p) = \begin{cases} x_i & \lfloor n \cdot p + .5 \rfloor = n \cdot p + .5 \\ x_i + (np - i + .5)(x_{i+1} - x_i) & \lfloor n \cdot p + .5 \rfloor \neq n \cdot p + .5 \end{cases}$$

- $np + .5 = 10 \cdot 0.5 + 0.5 = 5.5$.
- since $\lfloor 5.5 \rfloor = 5$ then $i = 5$ and

$$\begin{aligned} Q(.5) &= x_i + (n \cdot p - i + .5) \cdot (x_{i+1} - x_i) \\ &= x_5 + (10 \cdot 0.5 - 5 + .5) \cdot (x_{5+1} - x_5) \\ &= x_5 + (.5) \cdot (x_6 - x_5) \\ &= 61 + (.5) \cdot (64 - 61) \\ &= 62.5 \end{aligned}$$

Summarizing

Simple Graphs

Freq Tables

Histograms

Center Stats

Quantiles

Summaries of Location and Central Tendency

Finding $Q(.17)$

- $np + .5 = 10 \cdot 0.17 + 0.5 = 2.2.$

- since $\lfloor 2.2 \rfloor = 2$ then $i = 2$ and

$$Q(.17) = x_i + (n \cdot p - i + .5) \cdot (x_{i+1} - x_i)$$

$$= x_2 + (10 \cdot 0.17 - 2 + .5) \cdot (x_{2+1} - x_2)$$

$$= x_9 + (.2) \cdot (x_3 - x_2)$$

$$= 50 + (.2) \cdot (58 - 50)$$

$$= 51.6$$

Summarizing

Simple Graphs

Freq Tables

Histograms

Center Stats

Quantiles

Summaries of Location and Central Tendency

Finding $Q(.65)$

- $np + .5 = 10 \cdot 0.65 + 0.5 = 7.$

- since $\lfloor 7 \rfloor = 7$ then $i = 7$ and

$$Q(.65) = x_i + (n \cdot p - i + .5) \cdot (x_{i+1} - x_i)$$

$$= x_7 + (10 \cdot 0.65 - 7 + .5) \cdot (x_{7+1} - x_7)$$

$$= x_7 + (0) \cdot (x_8 - x_7)$$

$$= x_7 + 0$$

$$= 66$$

Section 3.2: Plots Using Quantiles

Plots and Quantiles

Quantile Plots

Quantile Plots:

Scatterplots using quantiles and their corresponding values

For each x_i in the data set, we plot $\left(\frac{i-.5}{n}, x_i\right)$ - meaning we are plotting $(p, Q(p))$. We connect the points with a straight line, which follows the values of $Q(p)$ exactly.

Consider the sample: 13, 15, 18, 19, 21, 34, 35, 35, 36, 39.
The following table which helps create the plot:

	1	2	3	4	5	6	7	8	9	10
p										
$Q(p)$										

##

Plots and Quantiles

Quantile Plots

Quantile Plots:

Scatterplots using quantiles and their corresponding values

For each x_i in the data set, we plot $\left(\frac{i-.5}{n}, x_i\right)$ - meaning we are plotting $(p, Q(p))$. We connect the points with a straight line, which follows the values of $Q(p)$ exactly.

Consider the sample: 13, 15, 18, 19, 21, 34, 35, 35, 36, 39.

Notice that we have $n = 10$ observations which means that $Q(0.05) = x_1 = 13$. We can get the quantile for each of our observations and create this table:

	1	2	3	4	5	6	7	8	9	10
p	0.05	0.15	0.25	0.35	0.45	0.55	0.65	0.75	0.85	0.95
$Q(p)$	13	15	18	19	21	34	35	35	36	39

##

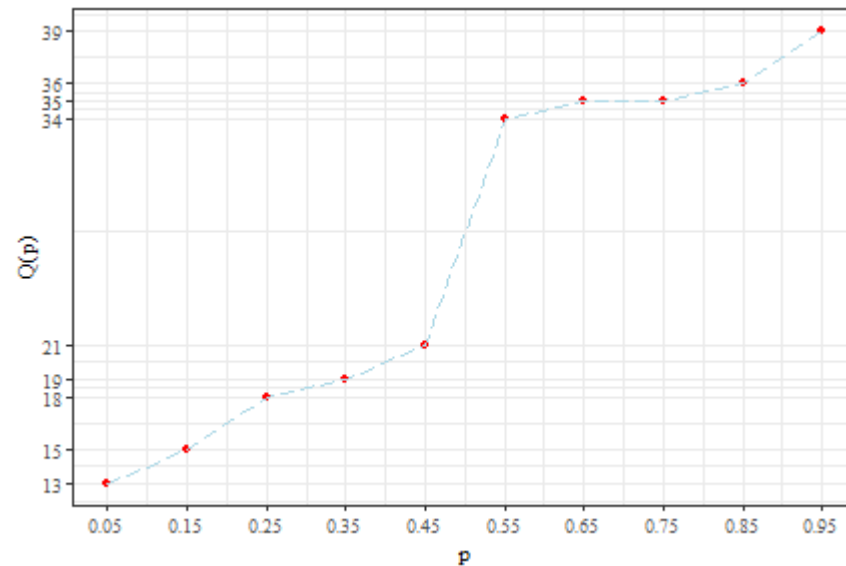
Plots and Quantiles

Quantile Plots

Quantile Plots:

Quantile plots

	1	2	3	4	5	6	7	8	9	10
p	0.05	0.15	0.25	0.35	0.45	0.55	0.65	0.75	0.85	0.95
$Q(p)$	13	15	18	19	21	34	35	35	36	39



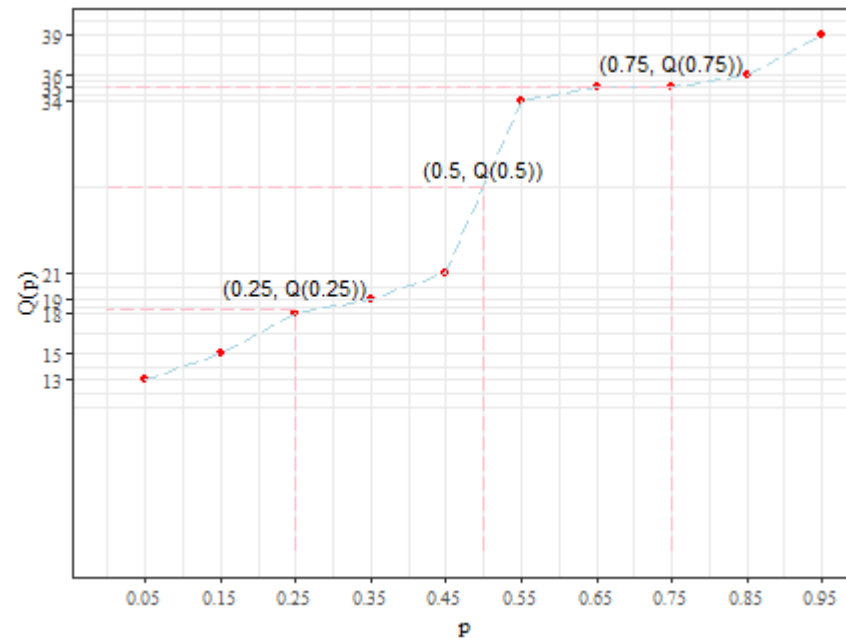
Plots and Quantiles

Quantile Plots

Quantile Plots:

Quantile plots

	1	2	3	4	5	6	7	8	9	10
p	0.05	0.15	0.25	0.35	0.45	0.55	0.65	0.75	0.85	0.95
$Q(p)$	13	15	18	19	21	34	35	35	36	39



Plots and Quantiles

Quantile Plots

QQ Plots

Quantile-Quantile Plots:

QQ plots are created by plotting the values of $Q(p)$ for a data set against values of $Q(p)$ coming from some other source.

- Compare the shape of two data sets (distributions).
- Two data sets having "equal shape" is equivalent to say their quantile functions are "**linearly related**".
- If the two data sets have different sizes, the size of smaller set is used for both.
- A **QQ plot** that is linear indicates the two distributions have similar shape.
- If there are significant departures from linearity, the character of those departures reveals the ways in which the shapes differ.

Plots and Quantiles

Quantile Plots

QQ Plots

Quantile-Quantile Plots:

Example: How similar the two data sets are?

- Set 1: 36, 15, 35, 34, 18, 13, 19, 21, 39, 35
- Set 2: 37, 39, 79, 31, 69, 71, 43, 27, 73, 71

	1	2	3	4	5	6	7	8	9	10
p										
Set 1 $Q(p)$										
Set 2 $Q(p)$										

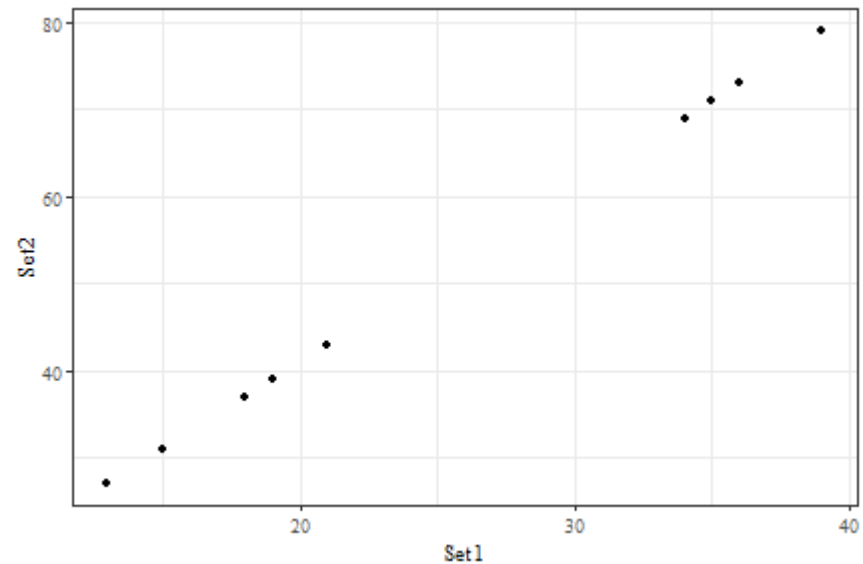
Plots and Quantiles

Quantile Plots

QQ Plots

Quantile-Quantile Plots:

	1	2	3	4	5	6	7	8	9	10
p	0.05	0.15	0.25	0.35	0.45	0.55	0.65	0.75	0.85	0.95
$Q_1(p)$	13	15	18	19	21	34	35	35	36	39
$Q_2(p)$	27	31	37	39	43	69	71	71	73	79



Plots and Quantiles

Quantile Plots

QQ Plots

Quantile-Quantile Plots:

Interpretation

The resulting plot shows some kind of linear pattern

This means that the quantiles increase at the same rate, even if the sizes of the values themselves are very different.

Plots and Quantiles

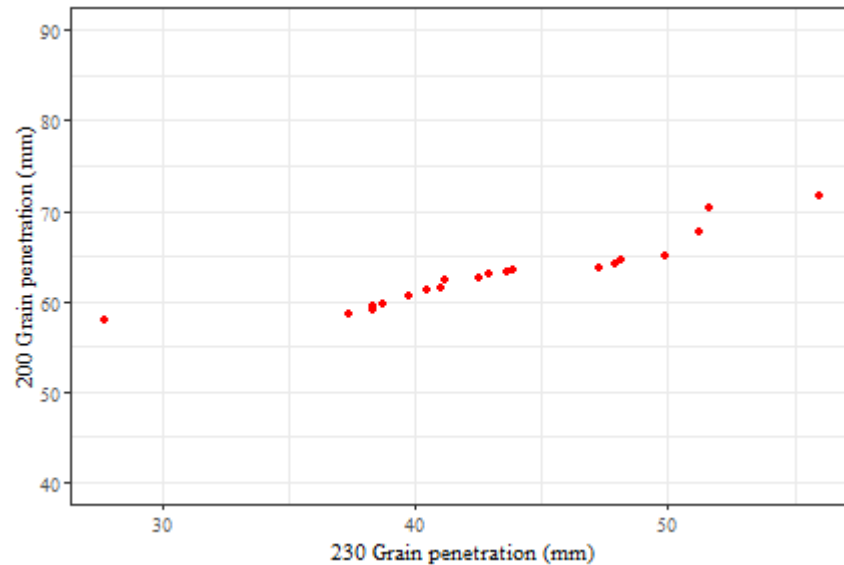
Quantile Plots

QQ Plots

Quantile-Quantile Plots:

Example 6 of chapter 3: Bullet penetration depth

- 230 Grain penetration (mm)
- 200 Grain penetration (mm)



Except extreme lower values, it **seems** the two distributions have similar shapes; however, it still needs more attention to make a rough decision (consider boxplots). Might want to figure out what has caused the extreme value

Plots and Quantiles

Quantile Plots

QQ Plots

Quantile-Quantile Plots:

The idea of QQ plots is most useful when applied to one quantile function that represents data and a second that represents a **theoretical distribution**

- Empirical QQ plots: the other source are quantiles from another actual data set.
- Theoretical QQ plots: the other source are quantiles from a theoretical set - we know the quantiles without having any data.

This allows to ask "Does the data set have a shape similar to the theoretical distribution?"

Boxplots

Plots and Quantiles

Quantile Plots

QQ Plots

Boxplots

Boxplots

A simple plot making use of the first, second and third quartiles (i.e., $Q(.25)$, $Q(.5)$ and $Q(.75)$).

1. A box is drawn so that it covers the range from $Q(.25)$ up to $Q(.75)$ with a vertical line at the median.
2. Whiskers extend from the sides of the box to the furthest points within 1.5 IQR of the box edges
3. Any points beyond the whiskers are plotted on their own.

Plots and Quantiles

Example: Draw boxplots for the groups using quantile function

Quantile Plots

QQ Plots

Boxplots

Group 1				Group 2			
74	79	77	81	65	77	78	74
68	79	81	76	76	73	71	71
81	80	80	78	86	81	76	89
88	83	79	91	79	78	77	76
79	75	74	73	72	76	75	79

solution: First we need the quartile values:

	$Q(.25)$	$Q(.5)$	$Q(.75)$
Group 1	75.5	79	81
Group 2	73.5	76	78.5

This means that Group 1 has $IQR = 5.5$ and

- $1.5 * IQR = 8.25$

while Group 2 has $IQR = 5$ and

- $1.5 * IQR = 7.5$

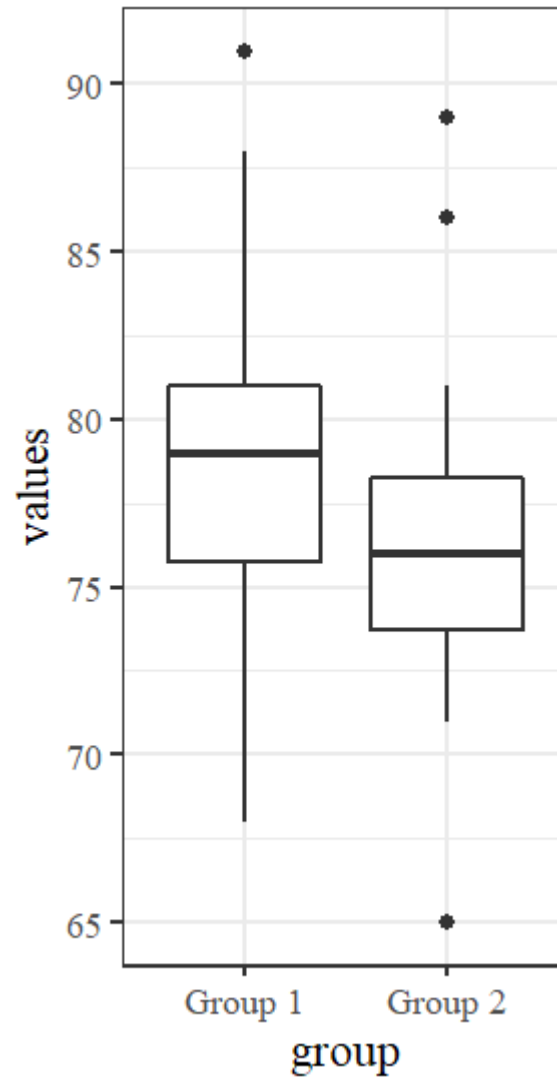
Plots and Quantiles

Quantile Plots

QQ Plots

Boxplots

Example:



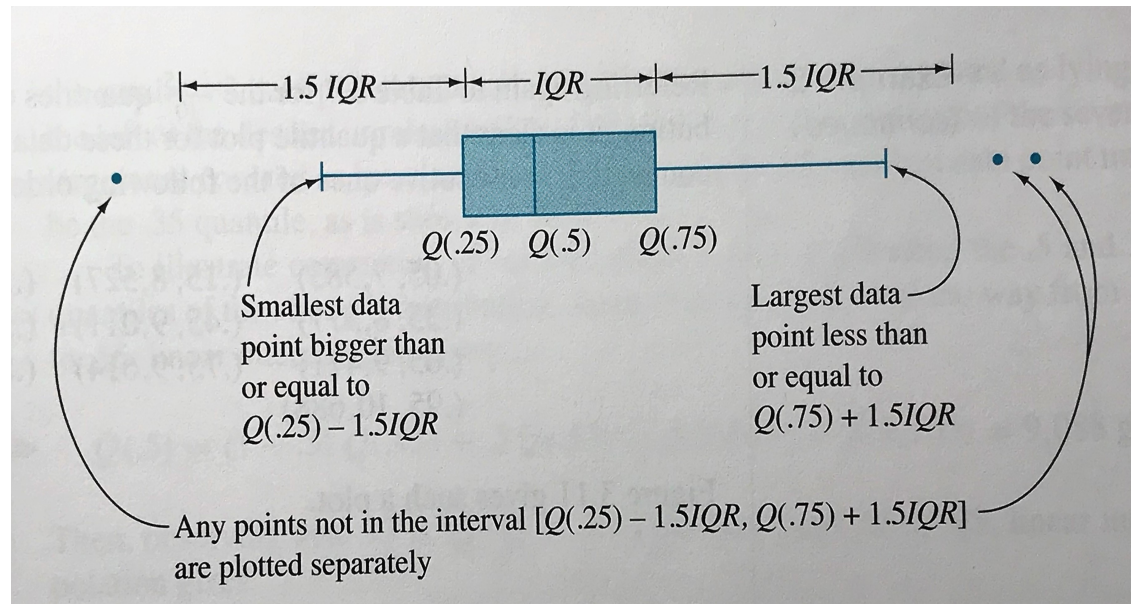
Plots and Quantiles

Quantile Plots

QQ Plots

Boxplots

Anatomy of a Boxplot



Plots and Quantiles

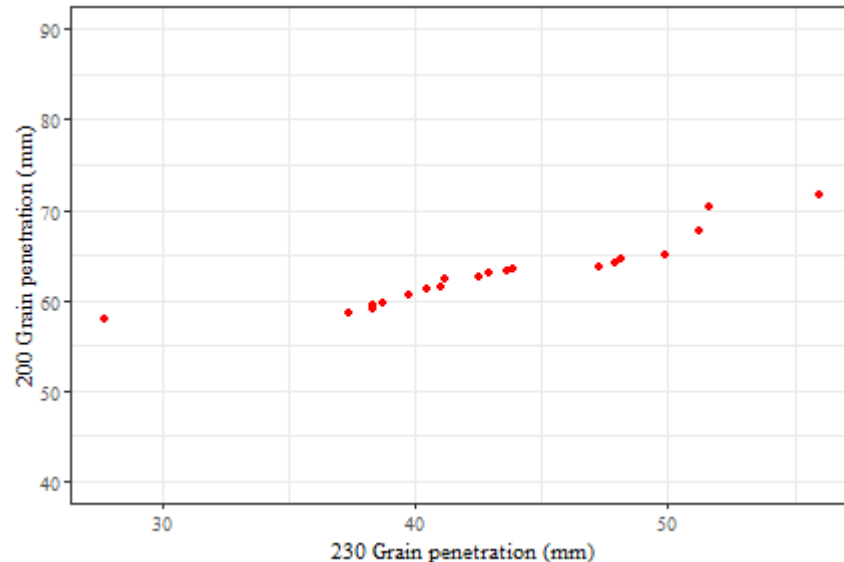
Quantile Plots

QQ Plots

Boxplots

Recap: Example 6 of chapter 3: Bullet penetration depth

- 230 Grain penetration (mm)
- 200 Grain penetration (mm)



Except extreme lower values, it **seems** the two distributions have similar shapes; however, it still needs more attention to make a rough decision (consider boxplots).

Plots and Quantiles

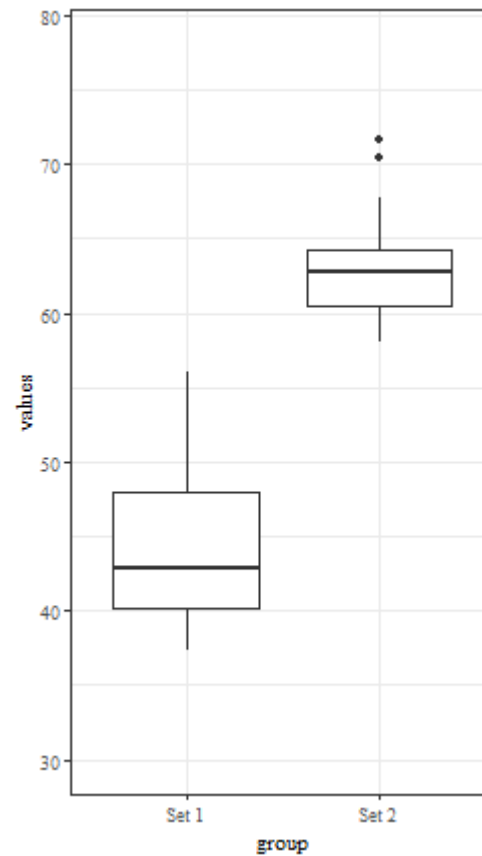
Quantile Plots

QQ Plots

Boxplots

Recap: Example 6 of chapter 3: Bullet penetration depth

Boxplot



Summarizing data Numerically

Location and central tendency

Measures of Spread

Recap

Location

Recap: Location and central tendency

Motivated by asking what is *normal/common/expected* for this data

Mean: A "fair" center value - $\frac{1}{n} \sum_{i=1}^n x_i$

Mode: The most commonly occurring value in set

Median: The value dividing the set in half (the middle of the values).

Group 1	Group 2
74 79 77 81	65 77 78 74
68 79 81 76	76 73 71 71
81 80 80 78	86 81 76 89
88 83 79 91	79 78 77 76
79 75 74 73	72 76 75 79

For group 1, the mean is 78.8, the median is 79, and the mode is 79.

For group 2, the mean is 76.45, the median is 76, and the mode is 76.

Recap

Location

Spread

Summaries of Variability (Measures of Spread)

Motivated by asking what kind of *variability is seen in the data* or *how spread out* the data is.

Range: The difference between the highest and lowest values (Range = max - min)

IQR: The Interquartile Range, how spread out is the middle 50% (IQR = Q3 - Q1)

Variance/Standard Deviation: Uses squared distance from the mean.

	Variance	Standard Deviation
Population	$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$	$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$
Sample	$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$	$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$

Recap

Location

Spread

Summarizing Data Numerically

Example: Taking a sample of size 5 from a population we record the following values:

65, 67, 57, 69, 58

Find the variance and standard deviation of this sample.

Example: Finding the Variance

Since we are told it is a sample, we need to use **sample variance**. The mean of 65, 67, 57, 69, 58 is 63.2

Example: Finding the Standard Deviation

With s^2 known, finding s is simple:

$$\begin{aligned}s &= \sqrt{s^2} \\ &= \sqrt{29.2} \\ &= 5.4037024\end{aligned}$$