



WELLCOME



STATISTICS WITH PYTHON

EMBARKING ON A JOURNEY
INTO DATA SCIENCE
YA MANON



You can have data without information but
you cannot have information without data.

-Daniel Keys Maran

COURSE OUTLINE

1

Why Statistics?

Discuss the role of statistics in the context of **business intelligence** and **decision-making**, and **introduce the statistics workflow**

2

Descriptive Statistics

Understand data using descriptive statistics, including **frequency distributions** and **measures of central tendency**, **variability**, **data visualization**

3

Probability Distributions

Model data with probability distributions, and use the **normal distribution** to calculate probabilities and make value estimates

4

Central Limit Theorem

Introduce the **Central Limit Theorem**, which leverages the normal distribution to make inferences on populations with any distribution

5

Confidence Intervals

Make estimates with **confidence intervals**, which use sample statistics to define a range where an unknown population parameter likely lies

6

Hypothesis Tests

Draw **conclusions** with **hypothesis tests**, which let you evaluate assumptions about population parameters using sample statistics

7

Regression Analysis

Make **predictions** with **regression analysis**, and estimate the values of a dependent variable via its relationship with independent variables

COURSE STRUCTURE



This is a **hands-on, project-based** course designed to help you apply statistical methods & techniques to real-world data analysis cases

Course resources include:

- ★ **Downloadable PDF ebook** to serve as a helpful reference when you're offline or on the go (*or just need a refresher!*)
- ★ **Quizzes** and **Project** to test and reinforce key concepts covered throughout the course, with detailed step-by-step solutions
- ★ **Interactive demos** to keep you engaged, with downloadable Excel files, code that you can use to follow along from home

THE COURSE PROJECT

THE OBJECTIVES

- Understand the data with descriptive statistics
- Model the data with probability distributions
- Make estimates with confidence intervals
- Draw conclusions with hypothesis tests
- Make predictions with regression analysis

SETTING EXPECTATIONS

-  The course simplifies essential statistics concepts
-  No math or statistics background is required to take this course
-  Focuses on the real-world application of statistical concepts
-  Uses Microsoft Excel and Jupyter notebooks (Python) for hands-on demos and assignments
-  Includes projects to test knowledge and apply it to different real-world scenarios

HELPFUL RESOURCES

Learn

Github

E-learning

- *Sunrise E-Learning*
- *scribbr.com/category/statistics/*

YouTub

Practice

Data Playground

Online Datasets

- *kaggle.com/datasets*
- *data.world/datasets/open-data*
- *vincentarelbundock.github.io/Rdatasets/articles/data*
- *https://archive.ics.uci.edu/*

DATA SCIENCE

INTRO TO DATA SCIENCE



In this section we'll **introduce the field of data science**, discuss how it compares to other data fields, and walk through each phase of the data science workflow

TOPICS WE'LL COVER:

What is Data
Science?

Machine Learning

Essential Skills

Data Science
Workflow

GOALS FOR THIS SECTION:

- Compare data science and machine learning with other common data analytics fields
- Introduce supervised and unsupervised learning, and examples of each technique
- Review the machine learning landscape and commonly used algorithms
- Discuss essential skills, and review each phase of the data science workflow



WHAT IS DATA SCIENCE?

What is Data
Science

Essential Skills

Machine Learning

Data Science
Workflow

Data science is about *using data to make smart decisions*.



Wait, isn't that **Business Interligence** ?

Yes! The differences lie in the **types of problems** you solve, and **tools and techniques** you use to solve them:

Business Interligence

What happened?

- Descriptive Analytics
- Data Analysis

Data Science

What's going to happen?

- Predictive Analytics

DATA SCIENCE SKILL SET



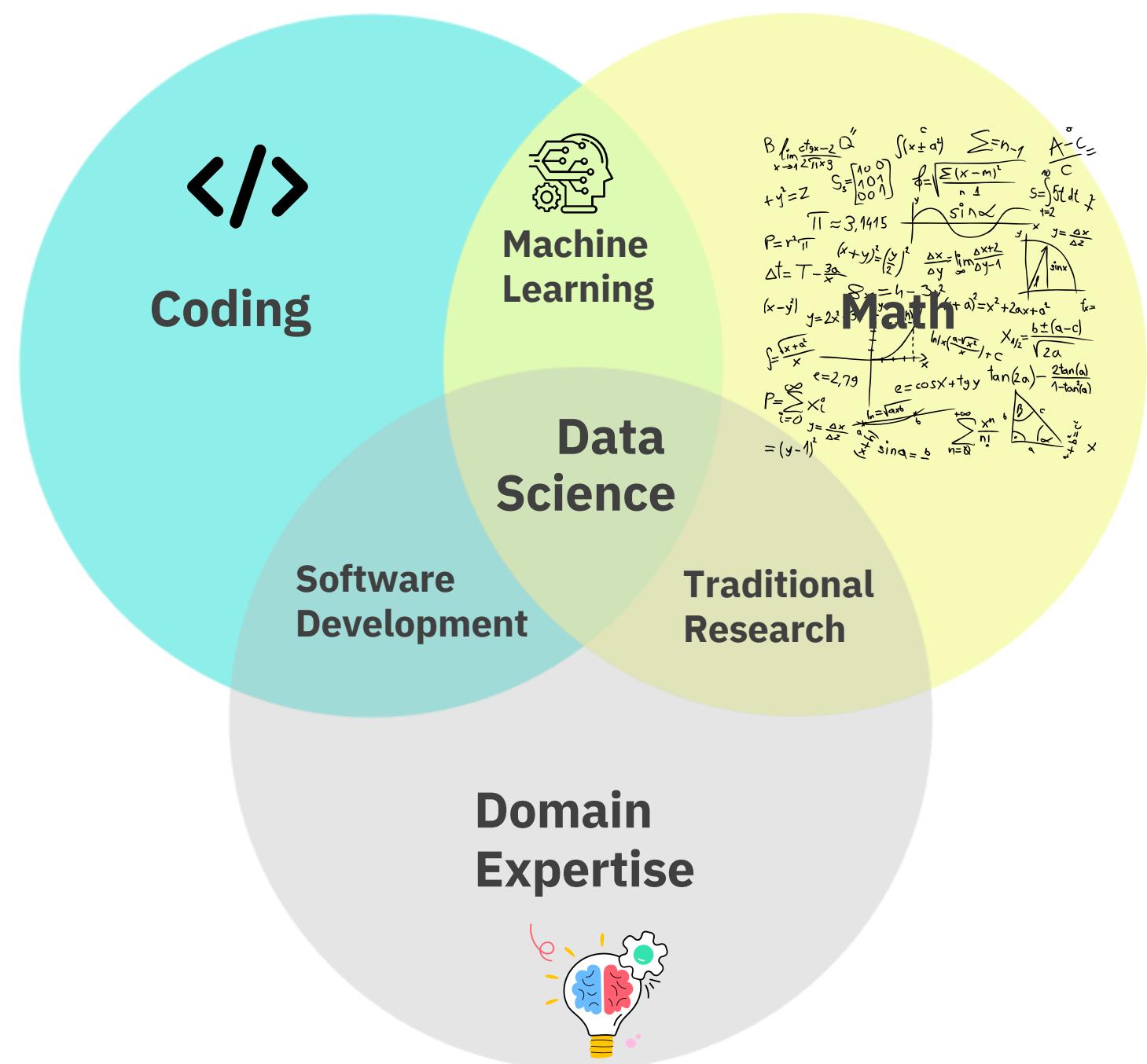
Data science requires a blend of **coding**, **math**, and **domain expertise**

What is Data Science

Essential Skills

Machine Learning

Data Science Workflow



The key is in applying these along with soft skills like:

- Communication
- Problem solving
- Curiosity & creativity
- Googling prowess



Data scientists & analysts approach problem solving in similar ways, but data scientists will often work with larger, more complex data sets and utilize advanced algorithms



WHAT IS MACHINE LEARNING?

What is Data
Science

Essential Skills

Machine Learning

Data Science
Workflow

Machine learning uses algorithms applied by data scientists to enable computers to learn and make decisions from data.

Machine learning algorithms fall into two broad categories:

Supervised Learning

Using historical data to predict the future



*What will house prices look like
for the next 12 months?*



*How can I flag suspicious emails
as spam?*

Unsupervised Learning

Finding patterns and relationships in data



*How can I segment my
customers?*



*Which TV shows should I
recommend to each user?*



DATA SCIENCE WORKFLOW

What is Data Science

Essential Skills

Machine Learning

Data Science Workflow

The **data science workflow** consists of scoping the project, gathering, cleaning and exploring the data, applying models, and sharing insights with end users



This is not a linear process! You'll likely go back to further gather, clean and explore your data



WELLCOME



STATISTIC WITH PYTHON

EMBARKING ON A JOURNEY
INTO DATA SCIENCE

YA MANON



WHY STATISTICS?

POPULATION & SAMPLES

A **population** contains all the data you're interested in to make your decision

- It's the data you wish you had, but are unlikely to get
- Any figure that summarizes a population is called a **parameter**

Populations

Statistics Workflow

A **sample** contains some of the data from the population

- It's the data you have (which should ideally represent the population)
- Any figure that summarizes a sample is called a **statistic**



Statistics lets you make reasonable estimates about **parameters** using **statistics**

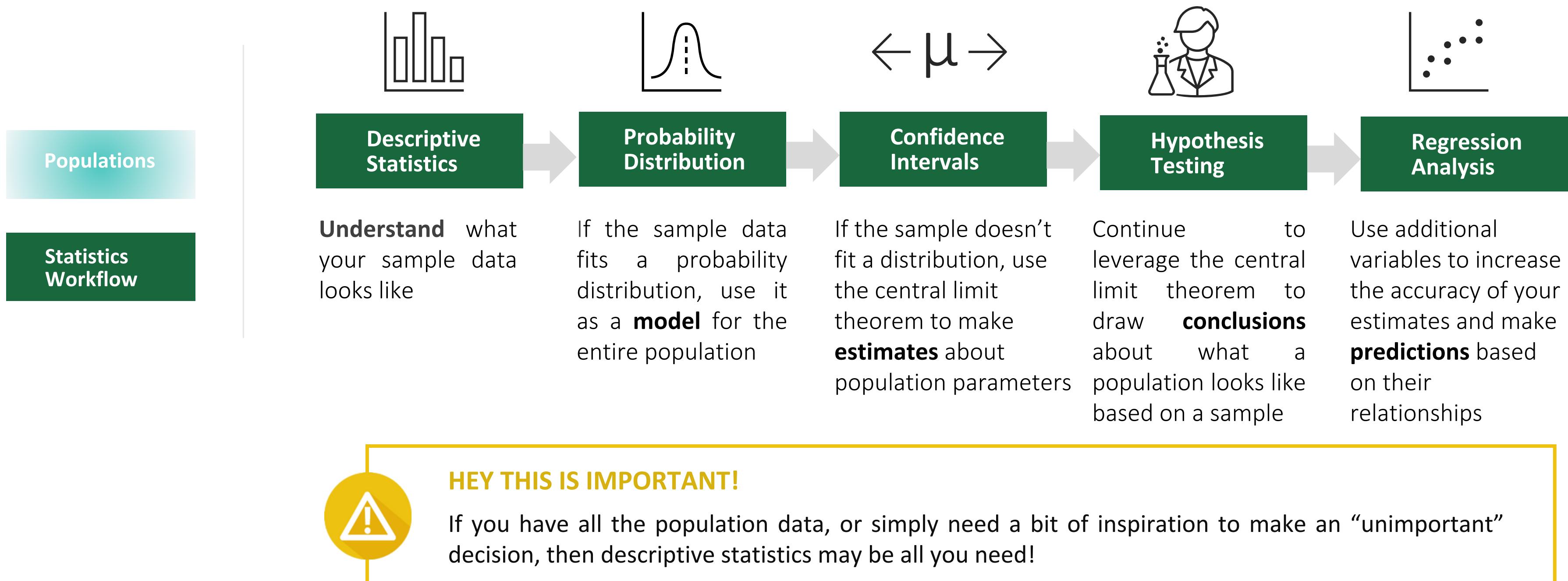


HEY THIS IS IMPORTANT!

Statistics can't create certainty out of uncertainty, it just helps you make controlled decisions under it!



THE STATISTICS WORKFLOW



DESCRIPTIVE STATISTICS

DESCRIPTIVE STATISTICS



In this section we'll cover understanding data with **descriptive statistics**, including frequency distributions, measures of central tendency, and measures of variability

TOPICS WE'LL COVER:

Statistics Basics

Central
Tendency

Distributions

Variability

GOALS FOR THIS SECTION:

- *Identify the different types of variables in a dataset, along with their use cases*
- *Create frequency tables and plot the distributions of numerical variables using histograms*
- *Calculate the mean, median, mode, and standard deviation of a numerical variable*
- *Visualize the key descriptive statistics of a numerical variable using a box plot*

DESCRIPTIVE STATISTICS

Descriptive statistics consists of the collection, organization, summarization and presentation of data.

They reduce a large array of numbers into a handful of figures that describe it accurately

Statistics Basics

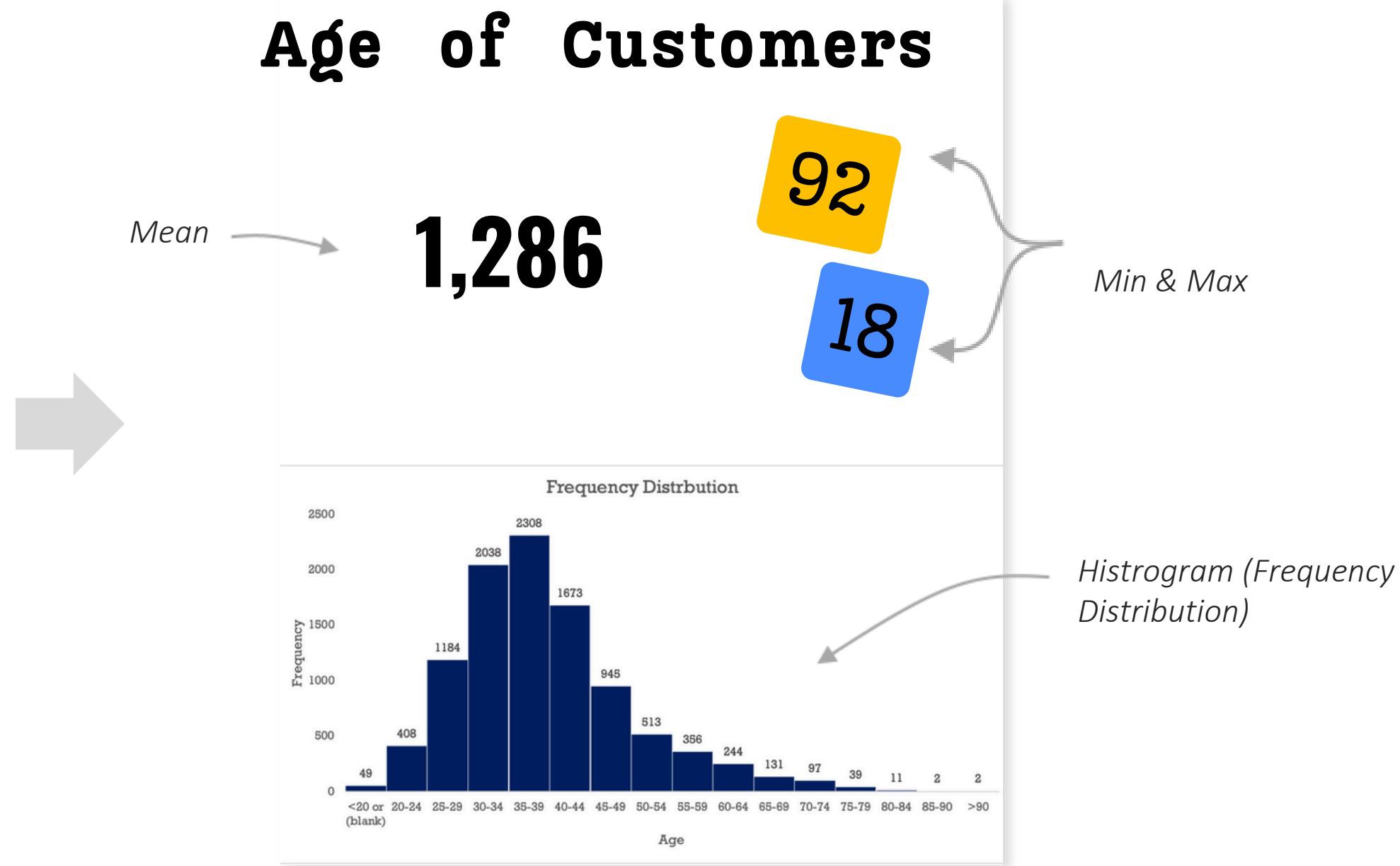
Distributions

Central Tendency

Variability

customer_id	age
15634602	42
15647311	41
15619304	42
15701354	39
15737888	43
15574012	44
15592531	50
15656148	29
15792365	44
15592389	27
15767821	31
15737173	24
15632264	34
15691483	25
15600882	35
15643966	45
15737452	58
15788218	24
15661507	45
15568982	24

n=10000



TYPES OF VARIABLES

There are two main types of variables in a dataset: Numerical & Categorical

- Numerical or Quantitative variables
- Categorical or Qualitative variables

Statistics Basics

Distributions

Central Tendency

Variability

NUMERICAL:

customer_id	age	tenure	balance	products_nur	credit_card	active_memk	estimated_sa	churn
15634602	42	2	0	1	1	1	101348.88	1
15647311	41	1	83807.86	1	0	1	112542.58	0
15619304	42	8	159660.8	3	1	0	113931.57	1
15701354	39	1	0	2	0	0	93826.63	0
15737888	43	2	125510.82	1	1	1	79084.1	0
15574012	44	8	113755.78	2	1	0	149756.71	1
15592531	50	7	0	2	1	1	10062.8	0
15656148	29	4	115046.74	4	1	0	119346.88	1
15792365	44	4	142051.07	2	0	1	74940.5	0
15592389	27	2	134603.88	1	1	1	71725.73	0
15767821	31	6	102016.72	2	0	0	80181.12	0
15737173	24	3	0	2	1	0	76390.01	0
15632264	34	10	0	2	1	0	26260.98	0
15691483	25	5	0	2	0	0	190857.79	0
15600882	35	7	0	2	1	1	65951.65	0
15643966	45	3	143129.41	2	0	1	64327.26	0
15737452	58	1	132602.88	1	1	0	5097.67	1
15788218	24	9	0	2	1	1	14406.41	0
15661507	45	6	0	1	0	0	158684.81	0
15568982	24	6	0	2	1	1	54724.03	0
15577657	41	8	0	2	1	1	170886.17	0
15597945	32	8	0	2	1	0	138855.46	0
15699309	38	4	0	1	1	0	118913.53	1
15725737	46	3	0	2	0	1	8487.75	0
15625017	28	5	0	1	1	1	197616.16	0

CATEGORICAL:

country	gender
France	Female
Spain	Female
France	Female
France	Female
Spain	Female
Spain	Male
France	Male
Germany	Female
France	Male
France	Male
France	Male
Spain	Male
France	Female
France	Female
Spain	Female

TYPES OF DESCRIPTIVE STATISTICS

There are 3 main **types of descriptive statistics** that can be applied to a variable:

Statistics Basics

Distributions

Central Tendency

Variability

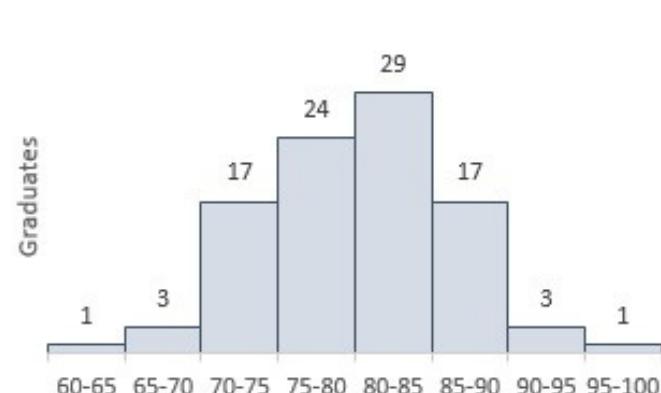
Distribution

Represents the **frequency** of each value

Examples:

- Frequency Tables
- Histograms

Grade Distribution



Central Tendency

Represents the **middle** of the values

Examples:

- Mean, Median, and Mode
- Skew

Class Average

80.17

Variability

Represents the **dispersion** of the values

Examples:

- Min, Max, and Range
- Quartiles & Interquartile Range
- Box & Whisker Plots
- Variance & Standard Deviation

Highest Grade

96.1

Lowest Grade

62.6

HEY THIS IS IMPORTANT!

Most measures of central tendency and variability can only be applied to numerical variables

FREQUENCY DISTRIBUTIONS

A **frequency distribution** counts the observations of each possible value in a variable

They are commonly depicted using frequency tables

Statistics Basics

Distributions

Central Tendency

Variability

Undergrad Degree	Undergrad Grade
Business	78.9
Business	74
Business	74.6
Engineering	79.3
Engineering	70.1
Business	88.8
Business	66
Art	82.9
Business	93.6
Business	75.6
Finance	67.5
Computer Science	68.7
Business	76
Computer Science	67.7
Engineering	75.3
Engineering	68.1
Finance	63.3

FREQUENCY TABLE:

Undergrad Degree	Frequency	Relative Frequency
Art	1	6%
Business	8	47%
Computer Science	2	12%
Engineering	4	24%
Finance	2	12%

The relative frequency shows the count of each value as a % of the total



TIP: Use a PivotTable or the COUNTIFS() function to calculate frequencies for categorical variables in Excel

FREQUENCY DISTRIBUTIONS

For numerical variables, a frequency distribution typically counts the number of observations that fall into defined ranges or “bins” (1-5, 6-10, etc.)

Statistics Basics

Distributions

Central Tendency

Variability

Undergrad Degree	Undergrad Grade
Business	78.9
Business	74
Business	74.6
Engineering	79.3
Engineering	70.1
Business	88.8
Business	66
Art	82.9
Business	93.6
Business	75.6
Finance	67.5
Computer Science	68.7
Business	76
Computer Science	67.7
Engineering	75.3
Engineering	68.1
Finance	63.3



FREQUENCY TABLE:

Undergrad Grade	Frequency
63.3	1
66	1
67.5	1
67.7	1
68.1	1
68.7	1
70.1	1
74	1
74.6	1
75.3	1
75.6	1
76	1
78.9	1
79.3	1
82.9	1
88.8	1
93.6	1

This isn't a meaningful representation of the distribution of the data

They are commonly depicted using grouped frequency tables or histograms

FREQUENCY DISTRIBUTIONS

Statistics Basics

Distributions

Central Tendency

Variability

Undergrad Degree	Undergrad Grade
Business	78.9
Business	74
Business	74.6
Engineering	79.3
Engineering	70.1
Business	88.8
Business	66
Art	82.9
Business	93.6
Business	75.6
Finance	67.5
Computer Science	68.7
Business	76
Computer Science	67.7
Engineering	75.3
Engineering	68.1
Finance	63.3



GROUPED FREQUENCY TABLE:

Undergrad Grade	Frequency	Cumulative Relative Frequency
60-65	1	6%
65-70	5	35%
70-75	3	53%
75-80	5	82%
80-85	1	88%
85-90	1	94%
90-95	1	100%
Grand Total	17	

The cumulative relative frequency shows the running total of the relative frequencies



TIP: Group the numerical values in a PivotTable or use the FREQUENCY() function with the upper limits to calculate frequencies for each bin in Excel

HISTOGRAMS

Histograms are used to visualize the distribution of a numerical variable

They also provide a glimpse of the variable's central tendency and variability

Statistics Basics

Distributions

Central Tendency

Variability

Undergrad Degree	Undergrad Grade
Business	78.9
Business	74
Business	74.6
Engineering	79.3
Engineering	70.1
Business	88.8
Business	66
Art	82.9
Business	93.6
Business	75.6
Finance	67.5
Computer Science	68.7
Business	76
Computer Science	67.7
Engineering	75.3
Engineering	68.1
Finance	63.3

n=17



Histogram of Undergrad Grades for MBA Graduates



TIP: Create a histogram by using a column chart to plot the variable's frequency table, instead of using Excel's native histogram chart type (not as customizable)

HISTOGRAMS

Histograms are used to visualize the distribution of a numerical variable

- They also provide a glimpse of the variable's central tendency and variability

Statistics Basics

Distributions

Central Tendency

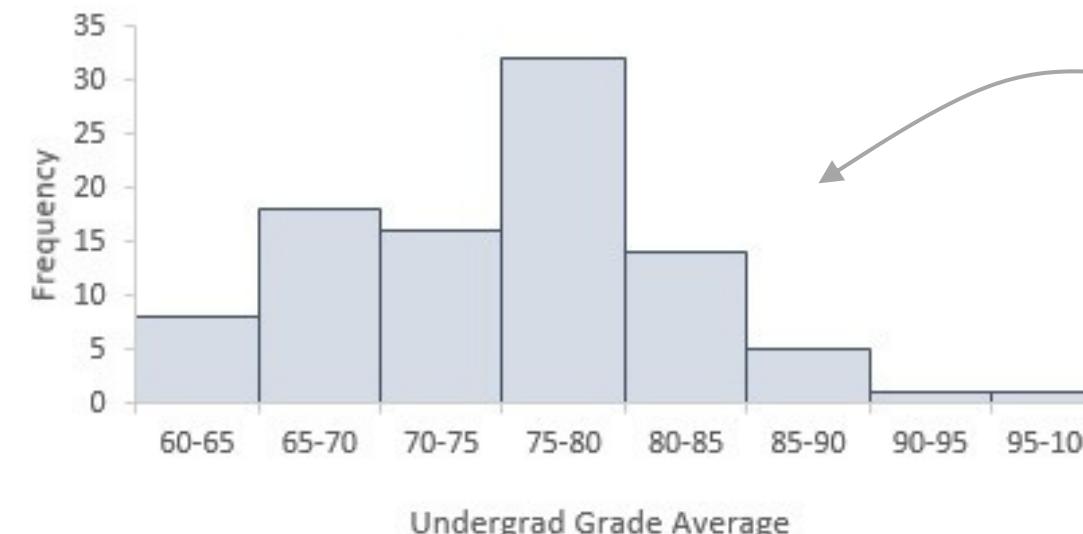
Variability

Undergrad Degree	Undergrad Grade
Business	78.9
Business	74
Business	74.6
Engineering	79.3
Engineering	70.1
Business	88.8
Business	66
Art	82.9
Business	93.6
Business	75.6
Finance	67.5
Computer Science	68.7
Business	76
Computer Science	67.7
Engineering	75.3
Engineering	68.1
Finance	63.3

n=95



Histogram of Undergrad Grades for MBA Graduates



Histograms are best suited for variables with many observations, to reflect the true population distribution

TIP: Bin size can significantly change the shape and “smoothness” of a histogram, so select a bin width that accurately shows the data distribution



MEAN

The **mean** is the calculated “average” value in a set on numbers

- It is calculated by dividing the sum of all values by the count of all observations
- It can only be applied to numerical variables (*not categorical*)

Statistics Basics

Distributions

Central
Tendency

Variability

customer_id	credit_score	age
15634602	619	42
15647311	608	41
15619304	502	42
15701354	699	39
15737888	850	43
15574012	645	44
15592531	822	50
15656148	376	29
15792365	501	44
15592389	684	27
15767821	528	31
15737173	497	24
15632264	476	34
15691483	549	25
15600882	635	35
15643966	616	45
15737452	653	58
15788218	549	24
15661507	587	45
15568982	726	24
15577657	732	41
15597945	636	32
15699309	510	38
15725737	669	46
15625047	846	38
15738191	577	25
15736816	756	36
15700772	571	44
15728693	574	43



Mean = 37.92

(average **37.9**

TIP: Use the **AVERAGEIFS()** function if you want

 to calculate the mean for values that meet a specified criteria (i.e., Mean by Undergrad Degree)

LIMITATIONS OF THE MEAN

The main **limitation of the mean** is that it is sensitive to outliers (*extreme values*)

“The average income in America is not the income of the average American”

Statistics Basics

Distributions

Central
Tendency

Variability



HEY THIS IS IMPORTANT!

While the mean is typically great for making a “best-guess” estimate of a value, it’s important to complement this value with other descriptive statistics like the distribution, median, and mode to see if the mean value is being distorted by outliers

MEDIAN

The **median** is the “middle value” in a sorted set of numbers

- Unlike the mean, the median is NOT sensitive to outliers
- When there are two middle-ranked values, the median is the average of the two

Statistics Basics

Distributions

Central Tendency

Variability

Undergrad Degree	Undergrad Grade
Business	78.9
Business	74
Business	74.6
Engineering	79.3
Engineering	70.1
Business	88.8
Business	66
Art	82.9
Business	93.6
Business	75.6
Finance	67.5
Computer Science	68.7
Business	76
Computer Science	67.7
Engineering	75.3
Engineering	68.1
Finance	63.3



Undergrad Grade
63.3
66
67.5
67.7
68.1
68.7
70.1
74
74.6
75.3
75.6
76
78.9
79.3
82.9
88.8
93.6

n=17

Median = **74.6**

MEDIAN

The **median** is the “middle value” in a sorted set of numbers

- Unlike the mean, the median is NOT sensitive to outliers
- When there are two middle-ranked values, the median is the average of the two

Statistics Basics

Distributions

Central Tendency

Variability

Undergrad Degree	Undergrad Grade
Business	78.9
Business	74
Business	74.6
Engineering	79.3
Engineering	70.1
Business	88.8
Business	66
Art	82.9
Business	75.6
Finance	67.5
Computer Science	68.7
Business	76
Computer Science	67.7
Engineering	75.3
Engineering	68.1
Finance	63.3



Undergrad Grade
63.3
66
67.5
67.7
68.1
68.7
70.1
74
74.6
75.3
75.6
76
78.9
79.3
82.9
88.8



Median = **74.3**

(average of **74** and **74.6**)

MODE

The **mode** is the “most frequent” value in a variable

- It can be applied to both numerical and categorical variables

Statistics Basics

Distributions

Central Tendency

Variability

Mode = “Business”

Undergrad Degree	Undergrad Grade
Business	78.9
Business	74
Business	74.6
Engineering	79.3
Engineering	70.1
Business	88.8
Business	66
Art	82.9
Business	93.6
Business	75.6
Finance	67.5
Computer Science	68.7
Business	76
Computer Science	67.7
Engineering	75.3
Engineering	68.1
Finance	63.3

Mode = N/A

MODE

The **modal class** is the group with the highest frequency

Statistics Basics

Distributions

Central Tendency

Variability

Undergrad Degree	Undergrad Grade
Business	78.9
Business	74
Business	74.6
Engineering	79.3
Engineering	70.1
Business	88.8
Business	66
Art	82.9
Business	93.6
Business	75.6
Finance	67.5
Computer Science	68.7
Business	76
Computer Science	67.7
Engineering	75.3
Engineering	68.1
Finance	63.3

GROUPED FREQUENCY TABLE:

Undergrad Grade	Frequency
60-65	1
65-70	5
70-75	3
75-80	5
80-85	1
85-90	1
90-95	1
Grand Total	17

Mode = **65-70, 75-80**

This is a **multi-modal** distribution, which indicates that there may be another variable impacting the undergrad grades

SKEW

Statistics Basics

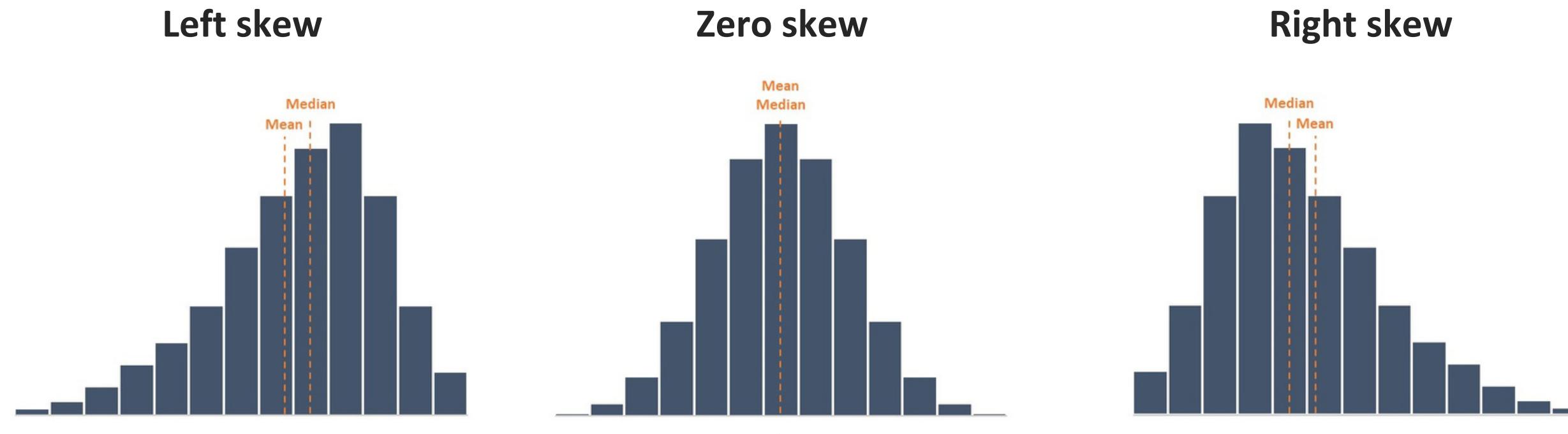
Distributions

Central Tendency

Variability

The **skew** represents the asymmetry of a distribution around its mean

- In a **zero-skewed** distribution, the mean and median are equal
- In a **right-skewed (or positive)** distribution, the mean is typically greater than the median
- In a **left-skewed (or negative)** distribution, the mean is typically smaller than the median



*This is one of the properties
of a **normal distribution**
(more on that later!)*

RANGE

The **range** is the spread from the lowest (*min*) to the highest (*max*) value in a variable

Statistics Basics

Distributions

Central Tendency

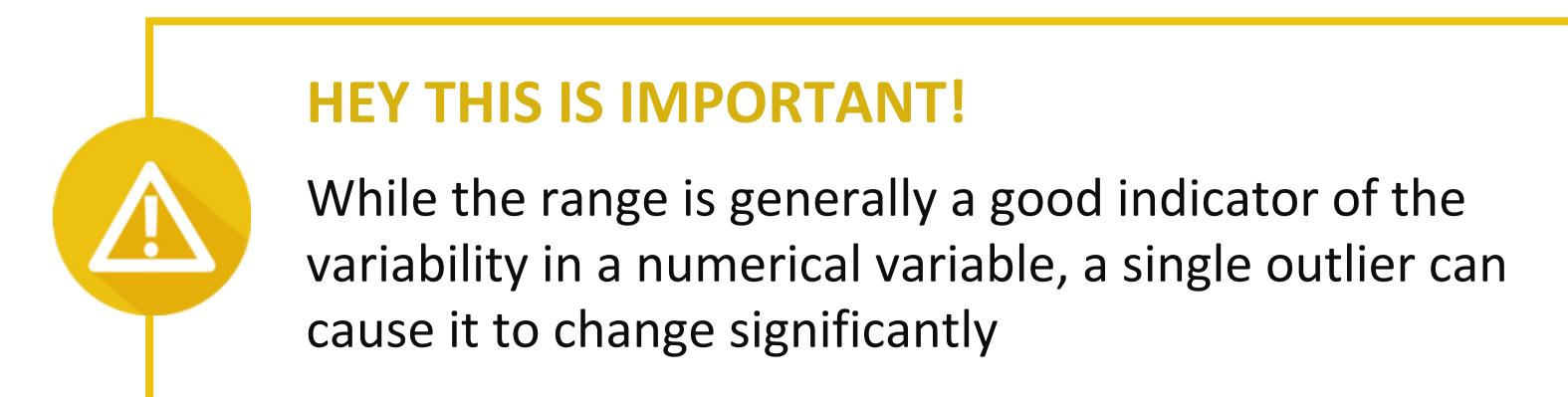
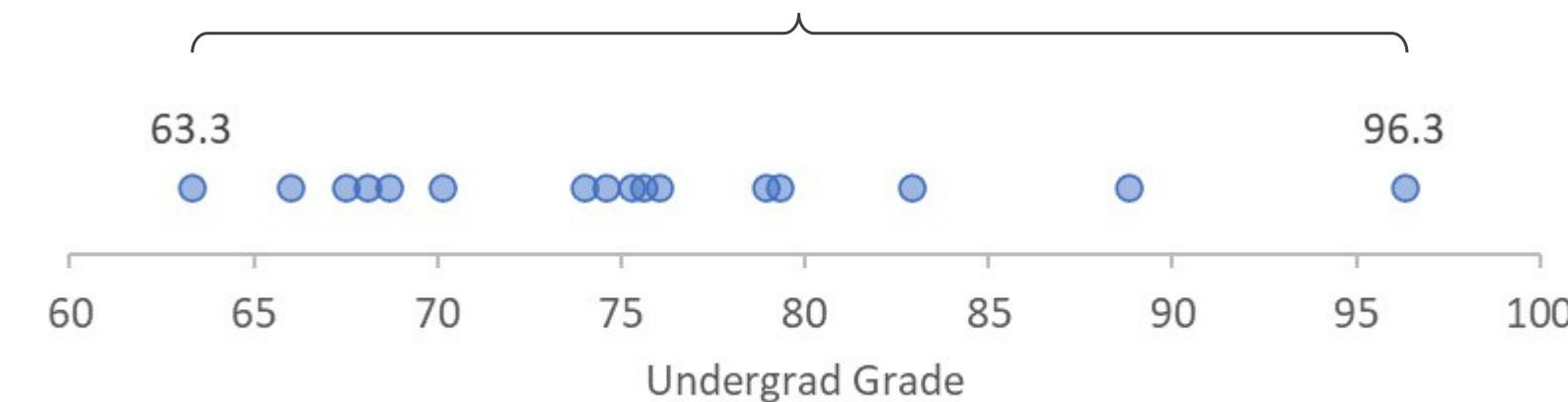
Variability

Undergrad Degree	Undergrad Grade
Business	78.9
Business	74
Business	74.6
Engineering	79.3
Engineering	70.1
Business	88.8
Business	66
Art	82.9
Business	96.3
Business	75.6
Finance	67.5
Computer Science	68.7
Business	76
Engineering	75.3
Engineering	68.1
Finance	63.3

n=16



$$\text{range} = 96.3 - 63.3 = 33$$



INTERQUARTILE RANGE

The **interquartile range** is the spread of the *middle half* of the values in a variable

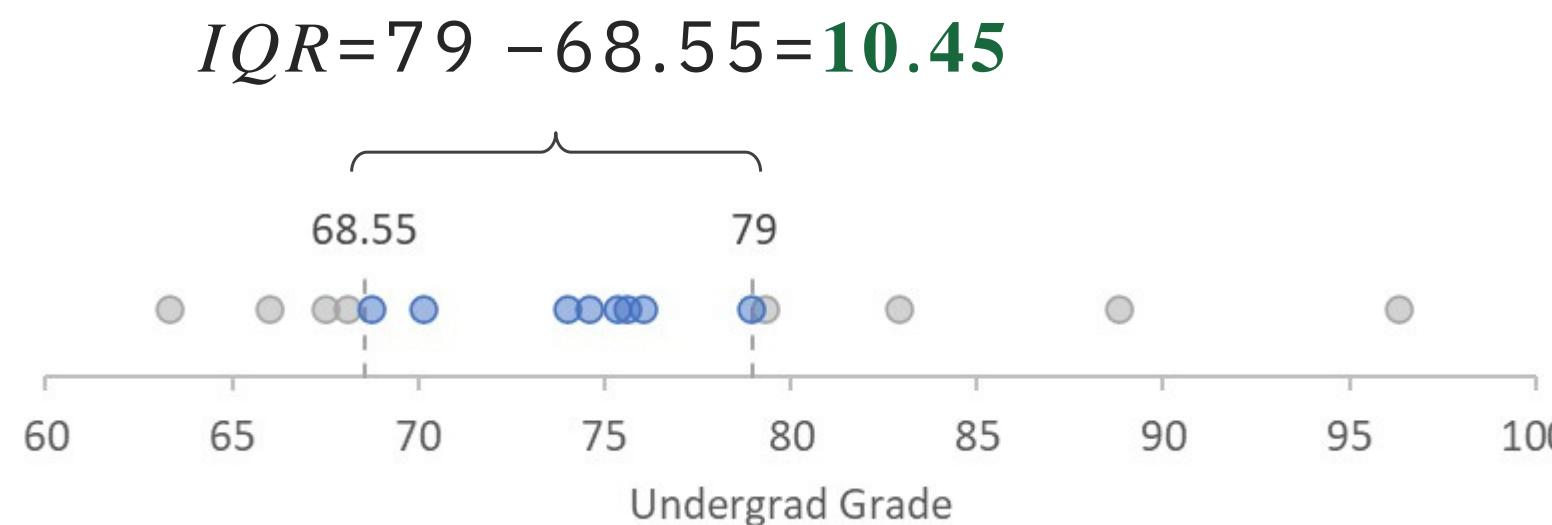
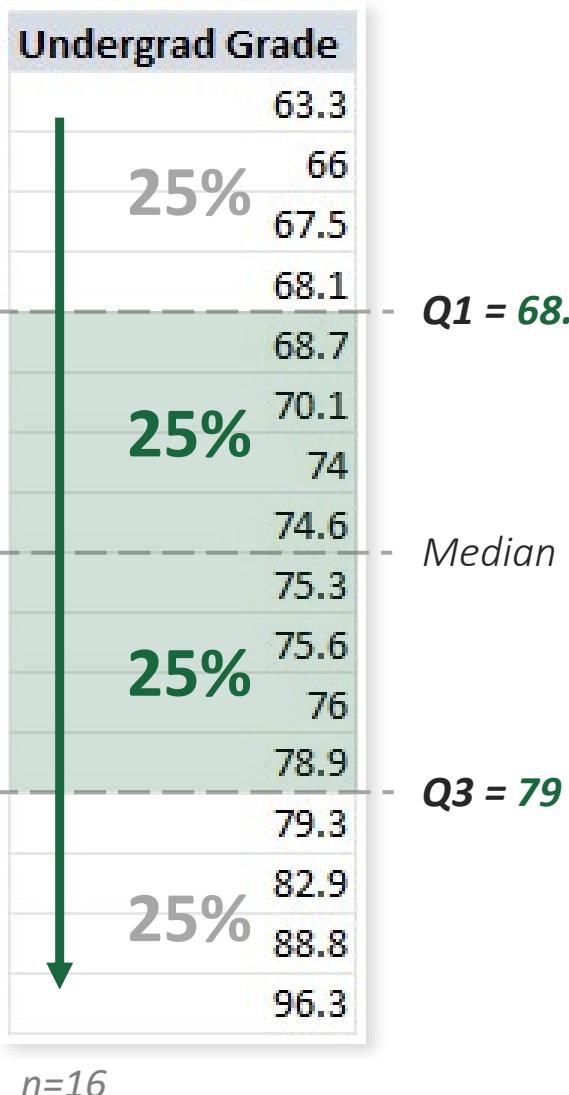
- In other words, it's the spread from the **first quartile** to the **third quartile**

Statistics Basics

Distributions

Central Tendency

Variability



BOX & WHISKER PLOTS

Box & whisker plots are used to visualize key descriptive statistics

Statistics Basics

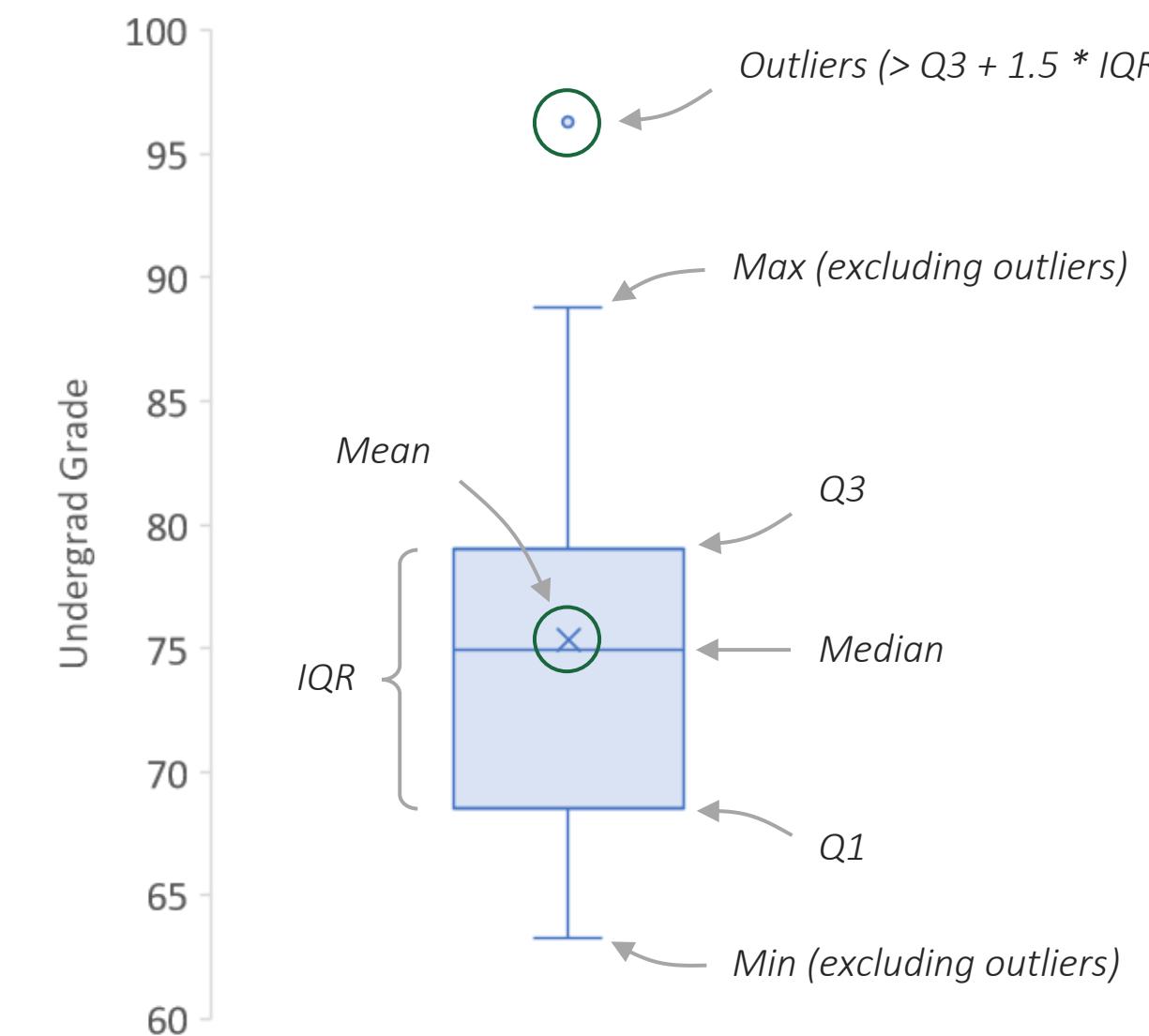
Distributions

Central Tendency

Variability

Undergrad Degree	Undergrad Grade
Business	78.9
Business	74
Business	74.6
Engineering	79.3
Engineering	70.1
Business	88.8
Business	66
Art	82.9
Business	96.3
Business	75.6
Finance	67.5
Computer Science	68.7
Business	76
Engineering	75.3
Engineering	68.1
Finance	63.3

n=16



BOX & WHISKER PLOTS

Box & whisker plots are used to visualize key descriptive statistics

- They can be used to quickly compare statistical characteristics between categories

Statistics Basics

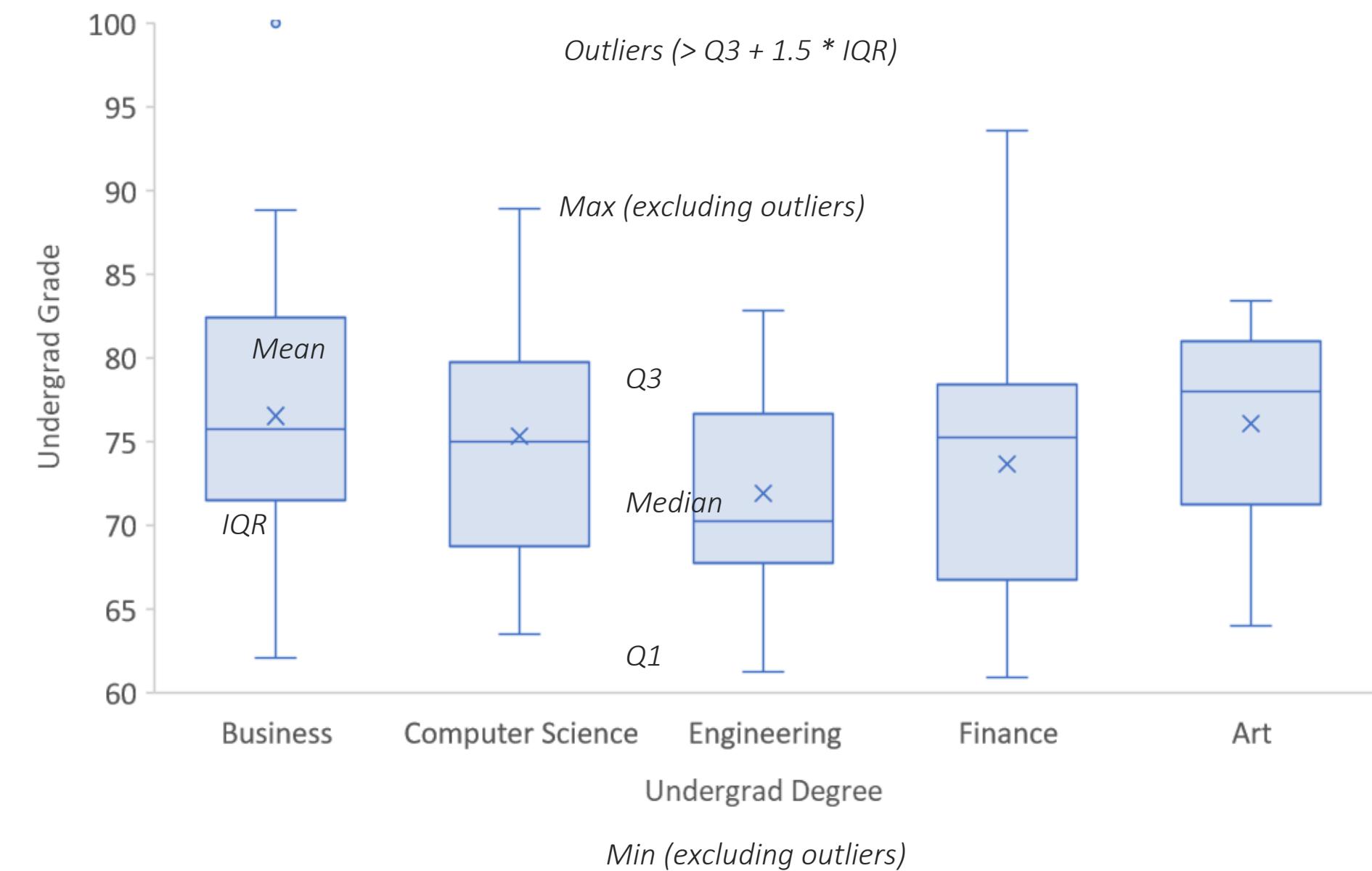
Distributions

Central Tendency

Variability

Undergrad Degree	Undergrad Grade
Business	78.9
Business	74
Business	74.6
Engineering	79.3
Engineering	70.1
Business	88.8
Business	66
Art	82.9
Business	96.3
Business	75.6
Finance	67.5
Computer Science	68.7
Business	76
Engineering	75.3
Engineering	68.1
Finance	63.3

n=1965



STANDARD DEVIATION

The **standard deviation** measures, on average, how far each value lies from the mean

The *higher* the standard deviation, the *wider* a distribution is (*and vice versa*)

Statistics Basics

Distributions

Central Tendency

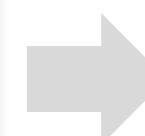
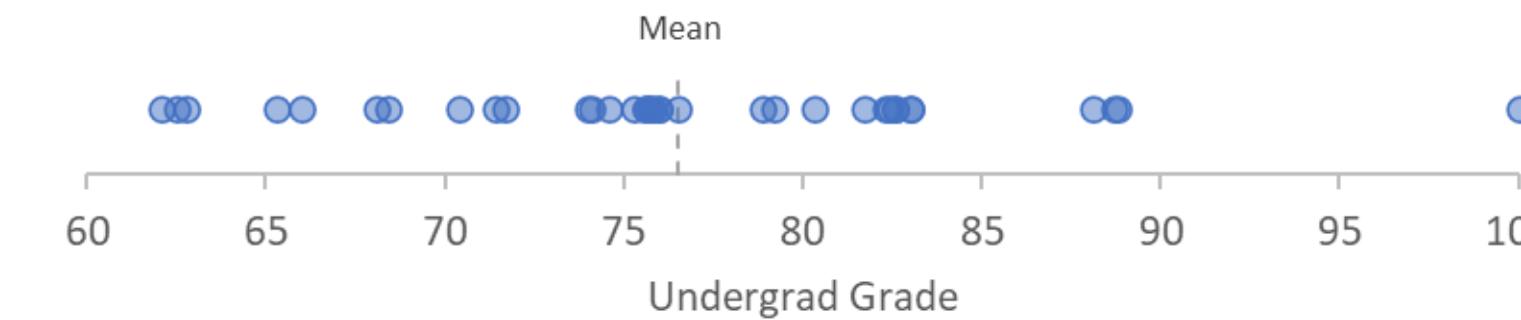
Variability

Undergrad Degree	Undergrad Grade
Business	78.9
Business	74
Business	74.6
Engineering	79.3
Engineering	70.1
Business	88.8
Business	66
Art	82.9
Business	96.3
Business	75.6
Finance	67.5
Computer Science	68.7
Business	76
Engineering	75.3
Engineering	68.1
Finance	63.3

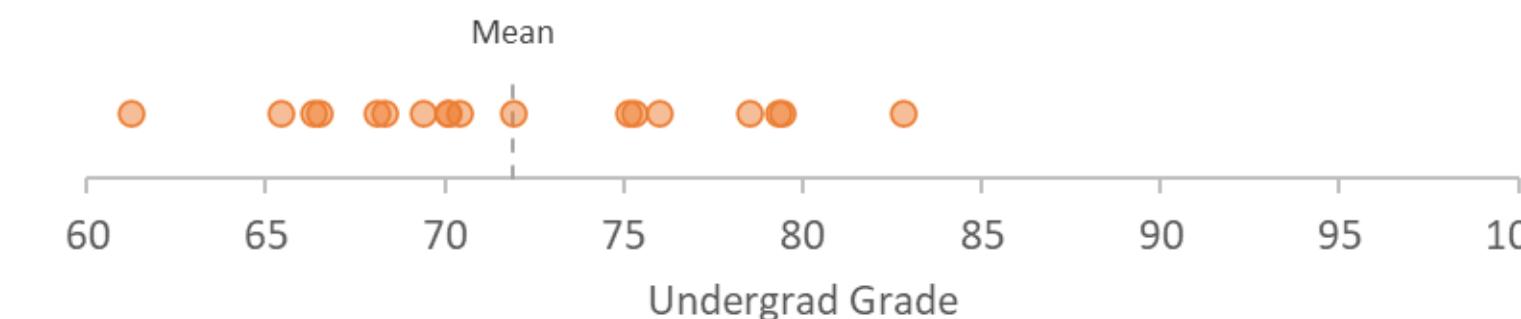
n=95



Business Undergrads



Engineering Undergrads



VARIANCE

The **variance** is the square of the standard deviation

Statistics Basics

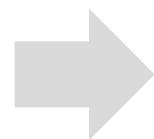
Distributions

Central Tendency

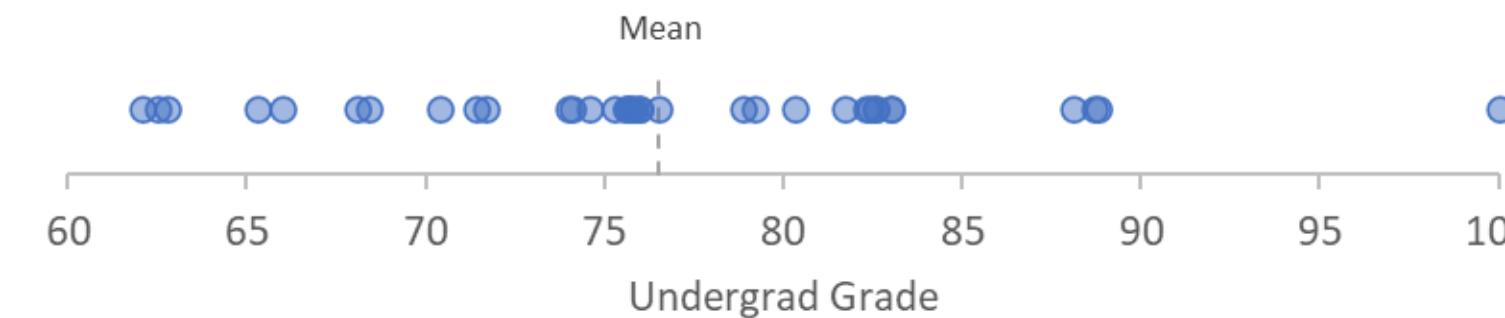
Variability

Undergrad Degree	Undergrad Grade
Business	78.9
Business	74
Business	74.6
Engineering	79.3
Engineering	70.1
Business	88.8
Business	66
Art	82.9
Business	96.3
Business	75.6
Finance	67.5
Computer Science	68.7
Business	76
Engineering	75.3
Engineering	68.1
Finance	63.3

n=95



Business Undergrads



HEY THIS IS IMPORTANT!

The variance does have its place in some statistical tests, so it shouldn't be discarded, but as a single numerical measure of a variable's dispersion the standard deviation is more effective

KEY TAKEAWAYS: DESCRIPTIVE STATISTICS



There are two main types of variables: numerical & categorical

- Numerical variables are meant to be aggregated, and categorical variables are used to create groups*



The distribution represents the “shape” of a variable

- Histograms are a great way to visualize this “shape” by plotting the frequency of each value (or class)*



The mean & median locate the “center” of a distribution

- Don’t focus on using one instead of the other, rather on using both to complement each other*



The standard deviation measures the dispersion around the mean

- Use a box plot alongside the standard deviation to provide additional context on the variability and center*