

04 | Exploratory Data Analysis

Ivan Corneillet

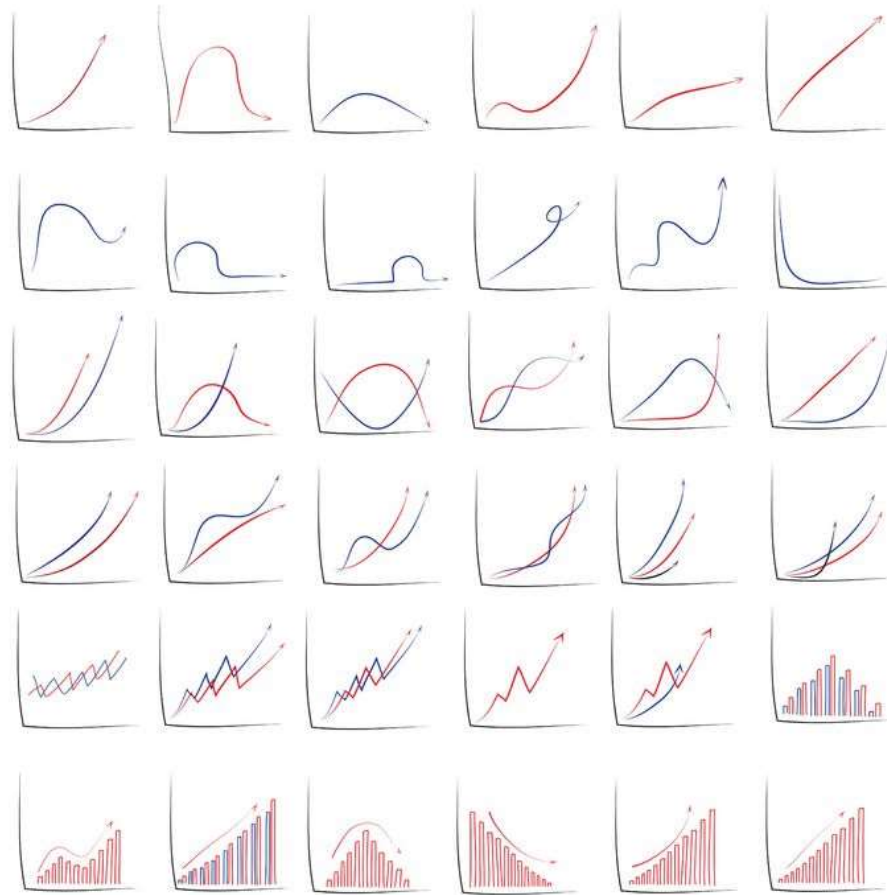
Data Scientist

Learning Objectives

After this lesson, you should be able to:

- Identify variable types
- Use the *pandas* (and *NumPy*) libraries to analyze datasets using basic summary statistics: mean, median, mode, max, min, quartile, inter-quartile range, variance, standard deviation, and correlation
- Create data visualizations – including boxplots, histograms, and scatter plots – to discern characteristics and trends in a dataset

Today, our main goal is to gain enough descriptive statistics knowledge to perform exploratory data analysis



Napat Polchoke © 123RF.com

- Descriptive vs. inferential statistics; populations vs. samples
- Types of data and types of measurement scales
- Measures of central tendency and measures of dispersion
- Boxplots
- Outliers
- Histograms
- Correlation



DS

Review

The pandas Library

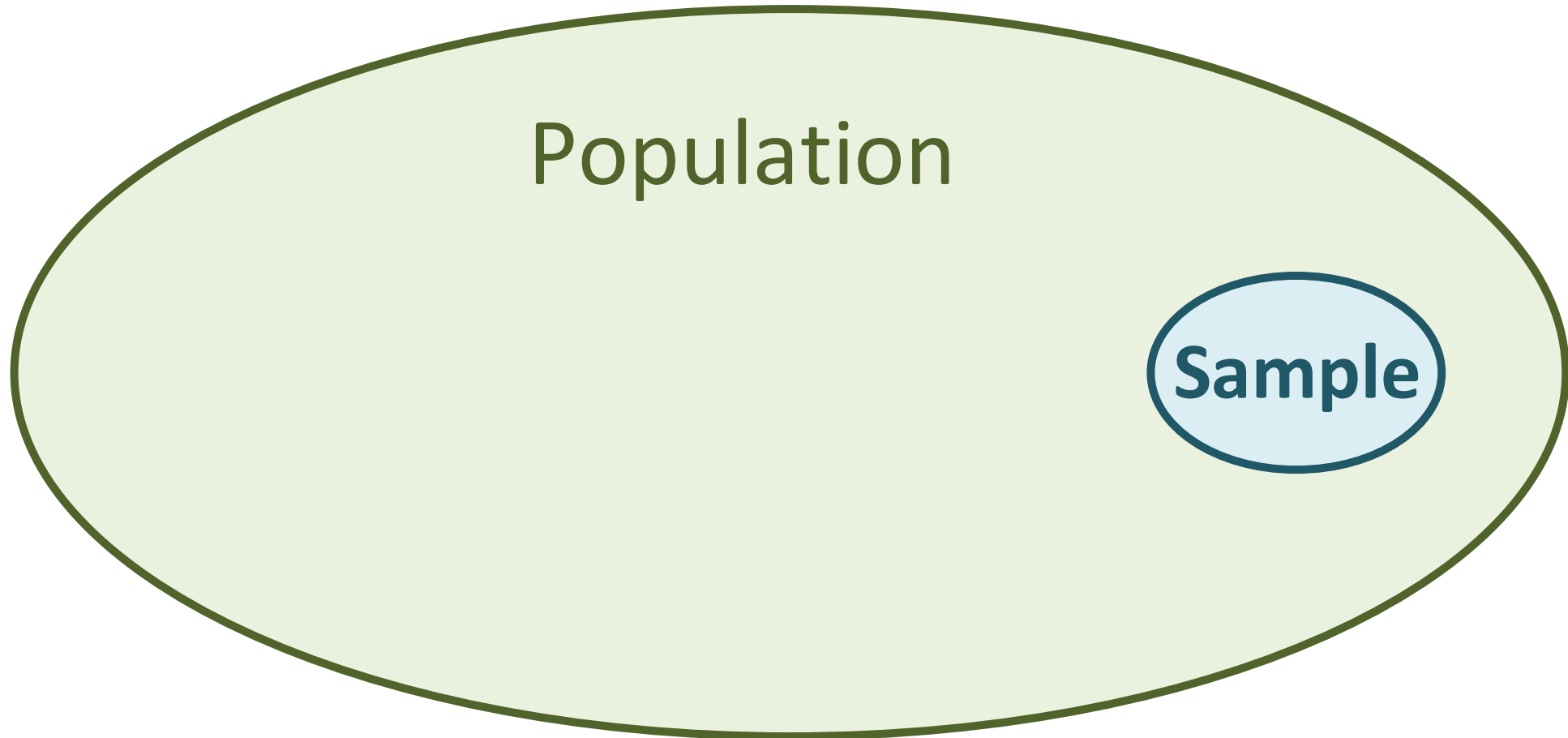
	DataFrame	Series
Column subsetting		
by name (Columns names are stored in <code>df.columns</code>) (<code>df.columns.get_loc('X1')</code> returns X1's column index)	# New DataFrame with column named X1 <code>df[['X1']]</code> # 2+ columns (in the order listed) <code>df[['X1', 'X2', ...]]</code>	<code>df['X1']</code> <code>df.X1</code>
by location	# New DataFrame with column at location i (numbering starts at 0) <code>df[[column_i]]</code> # 2+ columns (in the order listed) <code>df[[column_i, column_j, ...]]</code>	
Row subsetting		
by index label	<code>df.loc[[index_label_i]]</code> <code>df.loc[[index_label_i, index_label_j, ...]]</code> # Can use a range if the index is made of numbers (rows "a" to "b" included) <code>df.loc[index_label_a : index_label_b]</code>	<code>df.loc[index_label_i]</code>
by location	<code>df.iloc[[row_i]]</code> <code>df.iloc[[row_i, row_j, ...]]</code> # (rows "a" to "b" excluded) <code>df.iloc[row_a : row_b]</code> or <code>df[row_a : row_b]</code>	<code>df.iloc[location_i]</code>
Cell/scalar lookup		
by index label/column name	<code>df.at[index_label, 'X1']</code>	
by location	<code>df.iat[row_i, column_j]</code>	



DS

Populations and Samples

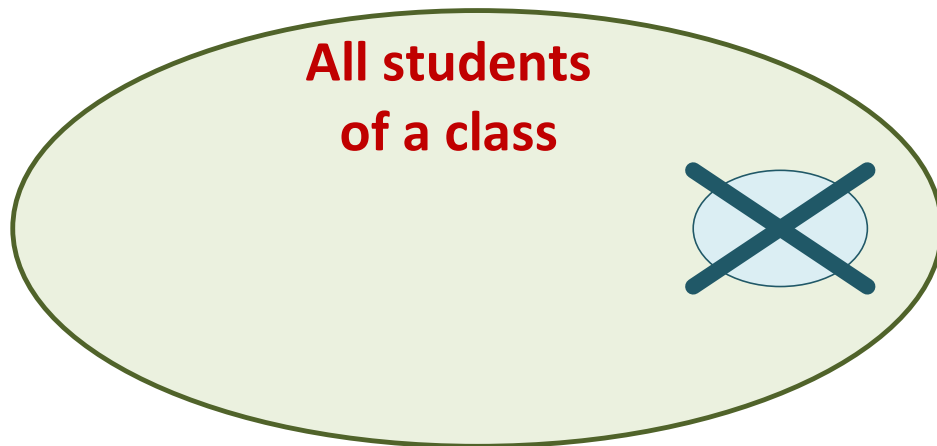
Populations and Samples



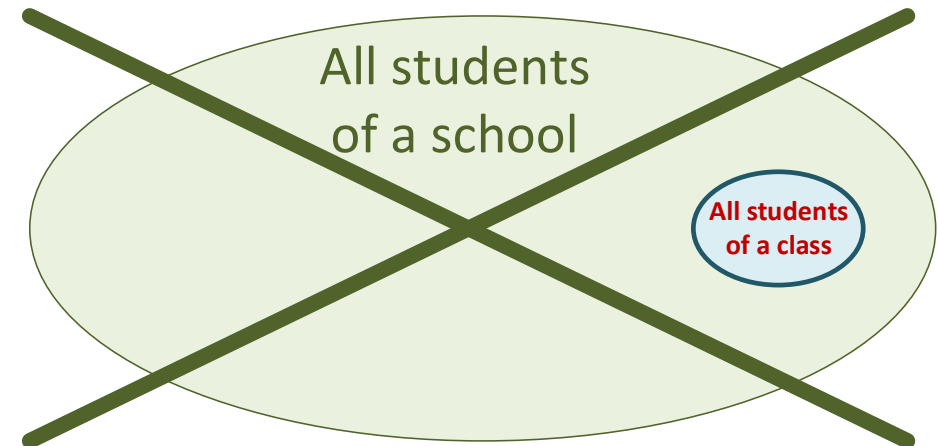
A dataset may be considered either as a population or a sample, depending on the reason for its collection and analysis

- Students of a class are a population if the analysis describes the distribution of scores in that class
- But they are a sample if the analysis infers from their scores the scores of other students (e.g., all students from that school)

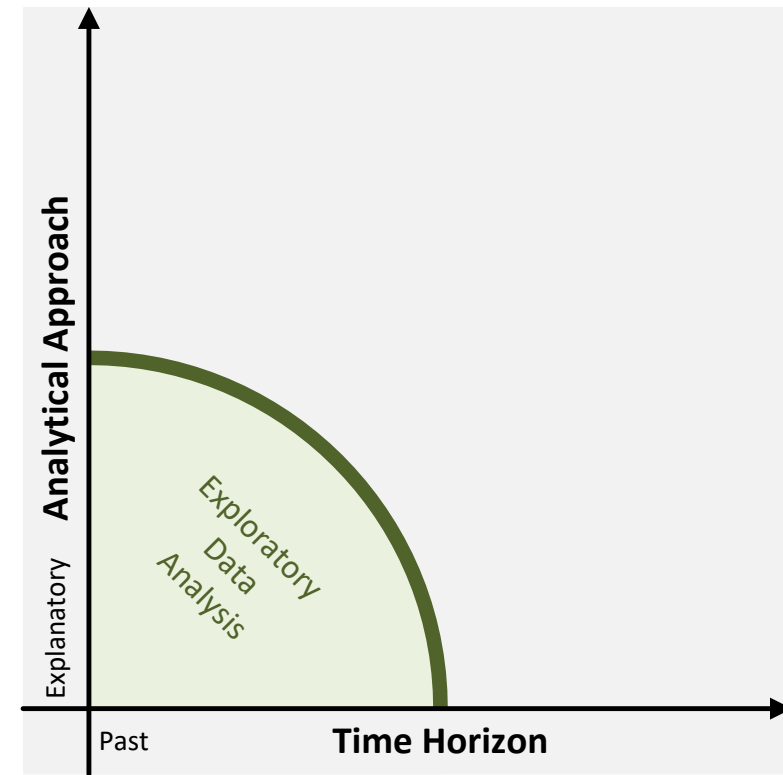
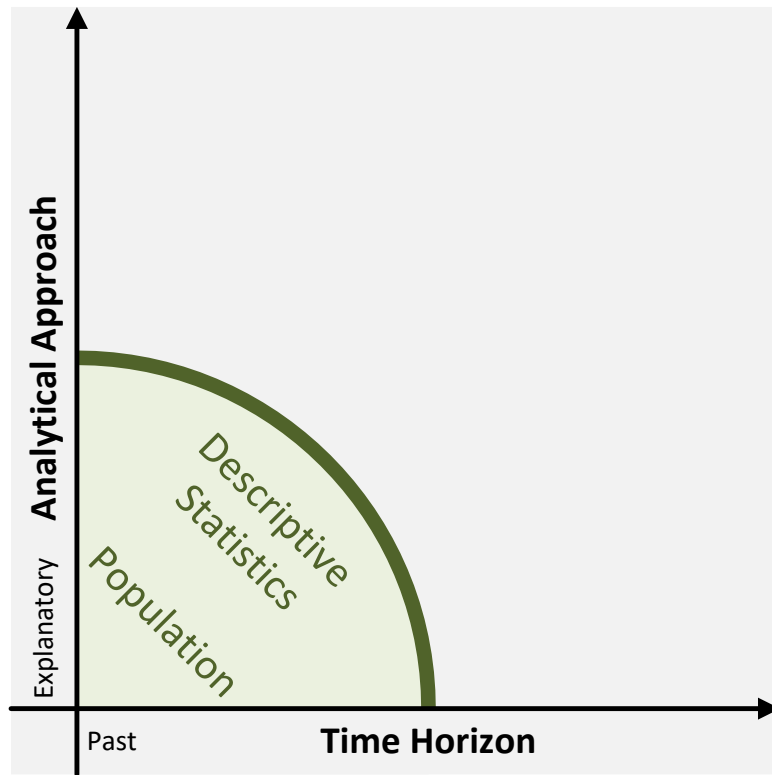
Descriptive Statistics



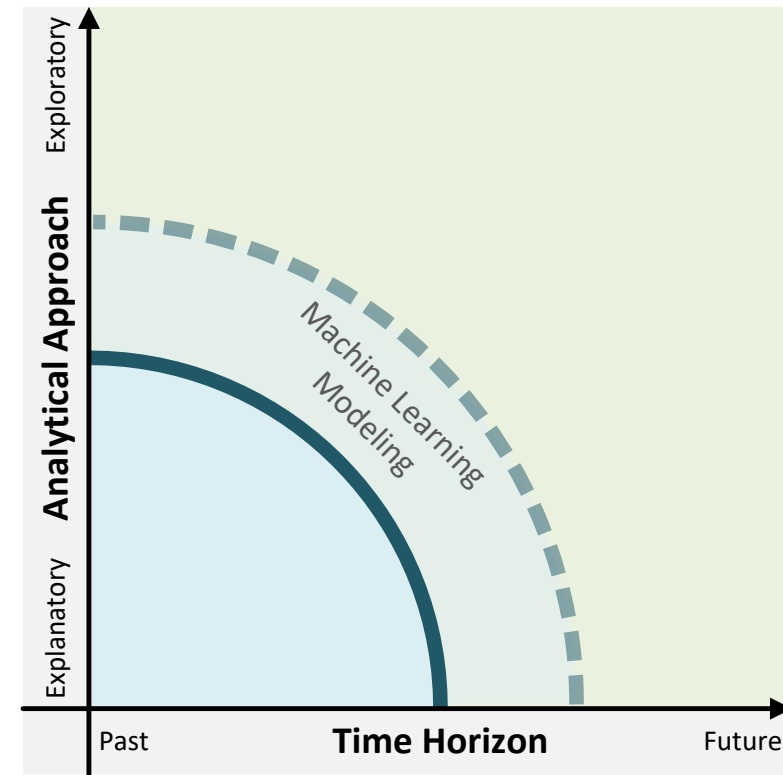
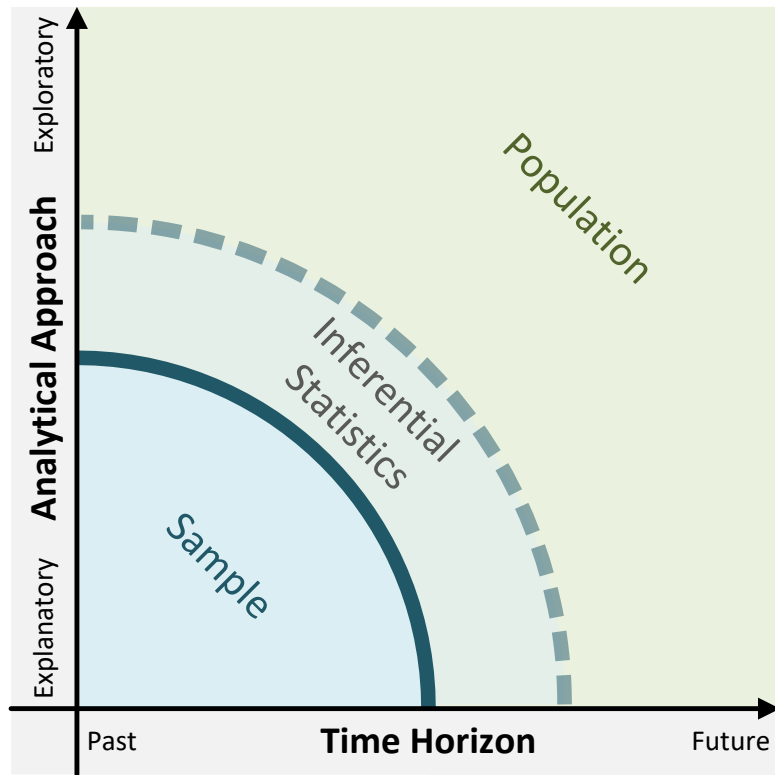
Inferential Statistics



Exploratory data analysis is concerned with descriptive statistics (e.g., “what happened last quarter?” and “how many units were sold?”)



Machine learning modeling concerns itself with inferential statistics (e.g., “what if ...?”, “what will happen next?”, and “what if these trends continue?”)



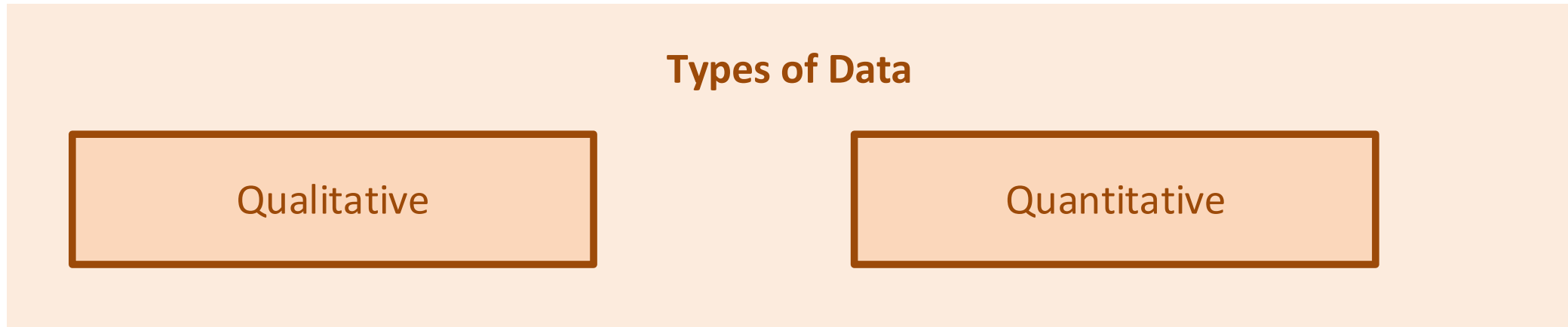


DS

Types of Data

Types of Measurement Scales

Types of Data



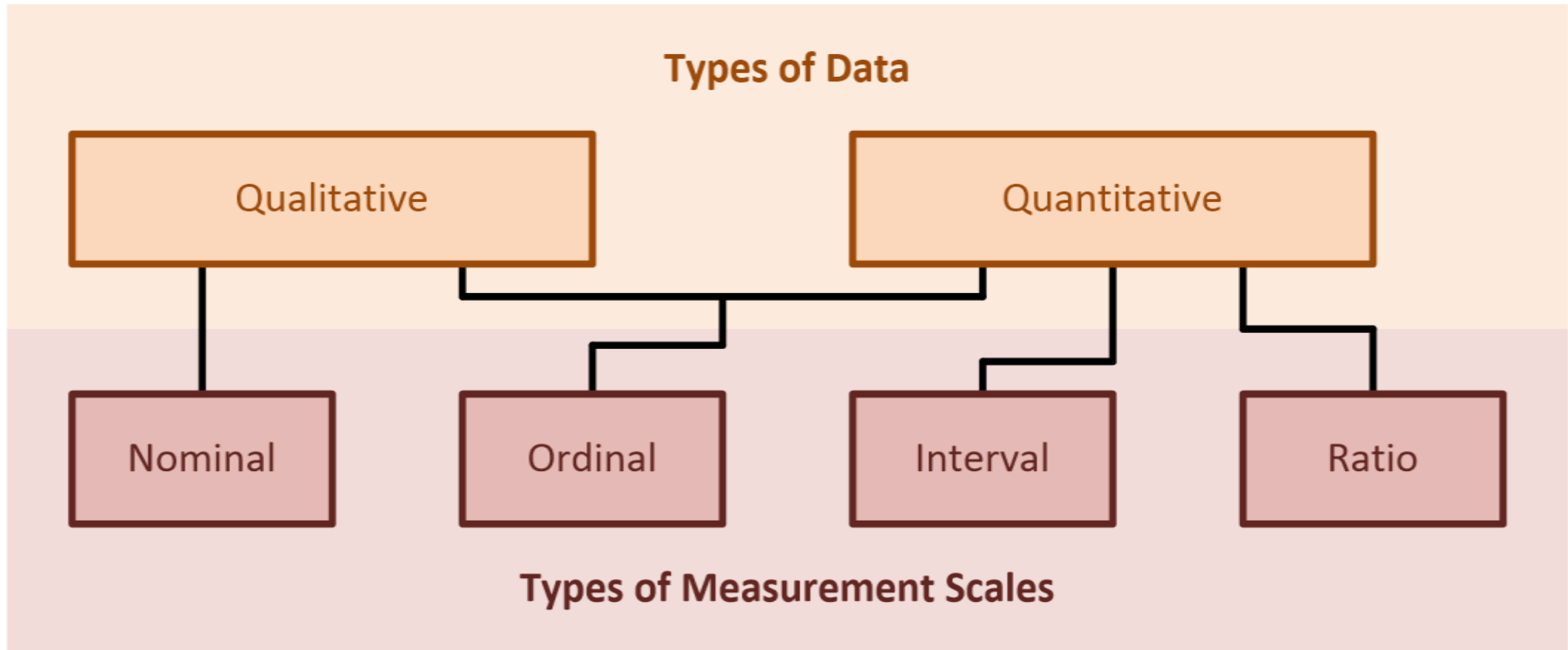
- Qualitative Data

- Uses descriptive terms to measure or classify something of interest, e.g., education level

- Quantitative Data

- Uses numerical values to describe something of interest, e.g., age

Types of Measurement Scales


















Types of Measurement Scales (cont.)

	Nominal	Ordinal	Interval	Ratio
<i>e.g.</i>	<i>Gender</i>	<i>Movie ratings</i>	<i>Temperature</i>	<i>Salary</i>
Categorize?	✓ (male, female)	✓	✓	✓
Rank-order?	✗	✓ (★ < 2★ < 3★ < 4★)	✓	✓
Add and subtract?	✗	✗ (4★ - 3★ ≠ ★)	✓ (75°C is 50°C warmer than 25°C)	✓
Multiply and divide?	✗	✗ (4★ not 4× better than 1★)	✗ (75°C not 3× as warm as 25°C) (0°C doesn't mean no temperature!)	✓ (Salary of \$200K is 2× that of \$100K) (\$0 means no salary ☹)

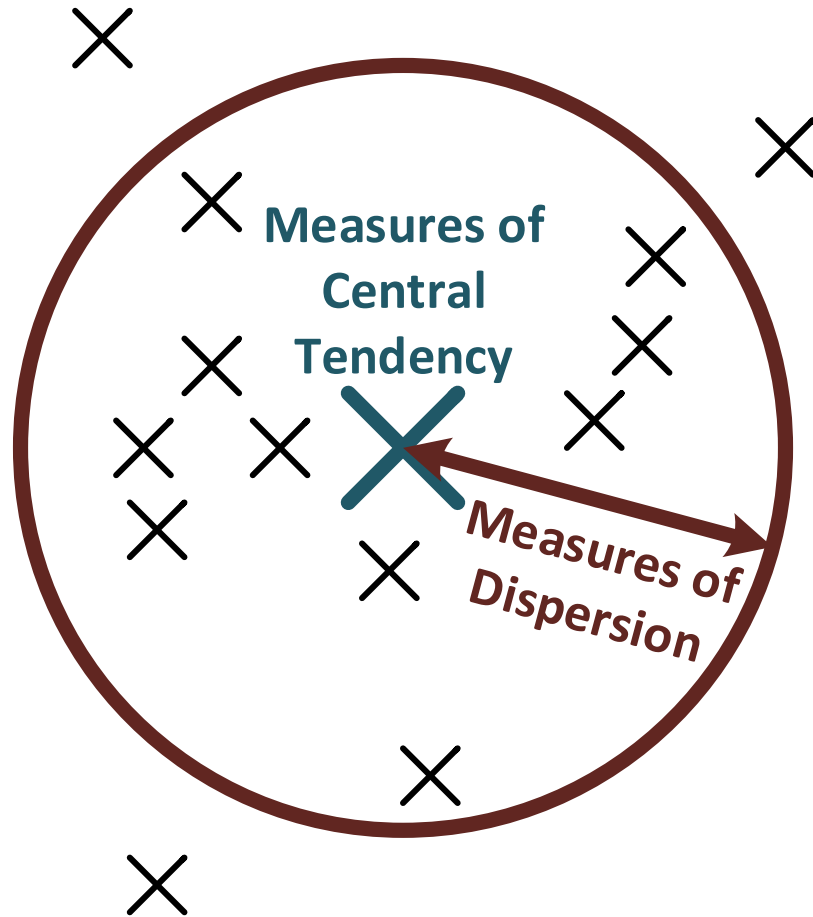
DS

Measures of Central Tendency Measures of Dispersion

Mean, Median, and Mode | Trade-offs

	Value is in the dataset	Value is easy to compute	Value is resistant to outliers	Corresponding measure of Dispersion	Used extensively by mathematical models
Mean	 (Unlikely)			 (Variance, standard deviation)	
Median	 (50% chance)	 (need to rank the values)		 (Interquartile Range)	
Mode	 (Always)	 (Need to count and rank the count)		 (Not really)	 (Mode might not be defined or you might have multiple values)

Measures of Central Tendency and Measures of Dispersion



- Measures of Central Tendency
 - (Or measures of location)
 - Answer the question: “What’s the typical or common value for a variable?”
 - Mean, Median, Mode
- Measures of Dispersion
 - (Or measures of variability/spread)
 - Answer the question: “How far do values stray from the typical value?”
 - Variance, Standard Deviation, Range, Interquartile Range (IQR)

(Arithmetic) Mean, Variance, and Standard Deviation

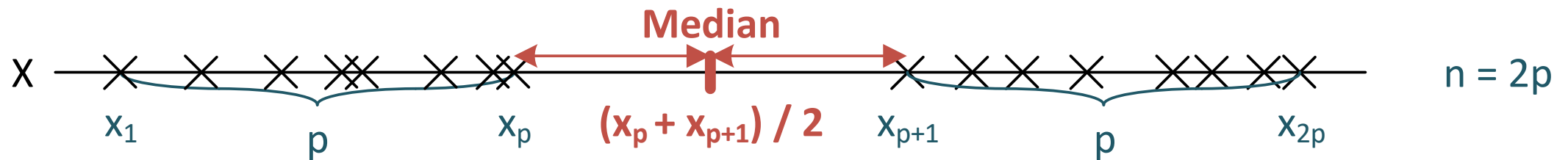
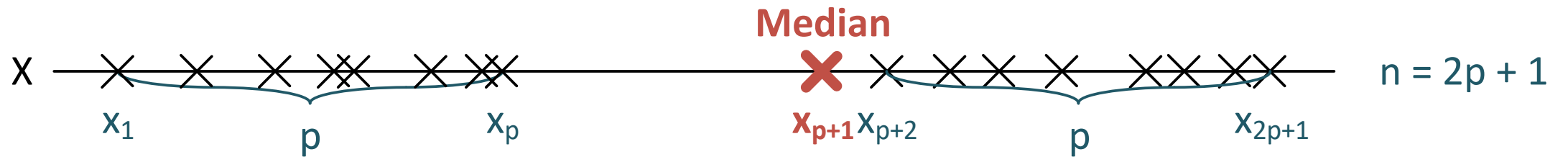
	Ordinal ✖	Nominal ✖	Interval ✔	Ratio ✔
	Population		Sample	
(Arithmetic) Mean <i>(a.k.a., the first moment)</i> (Mean has unit of $X:[X]$)	$\mu = \frac{1}{N} \sum_{i=1}^N x_i = E[X^1]$ (mu)		$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ (x-bar)	
Variance <i>(a.k.a., the second moment)</i> $[X^2]$	$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$ $= E[(X - \mu)^2]$ (sigma-squared)		$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$	
Standard Deviation $[X]$	$\sigma = \sqrt{\sigma^2}$ (sigma)		$s = \sqrt{s^2}$	

(mean, variance, and standard deviations are based on the values of x_i)

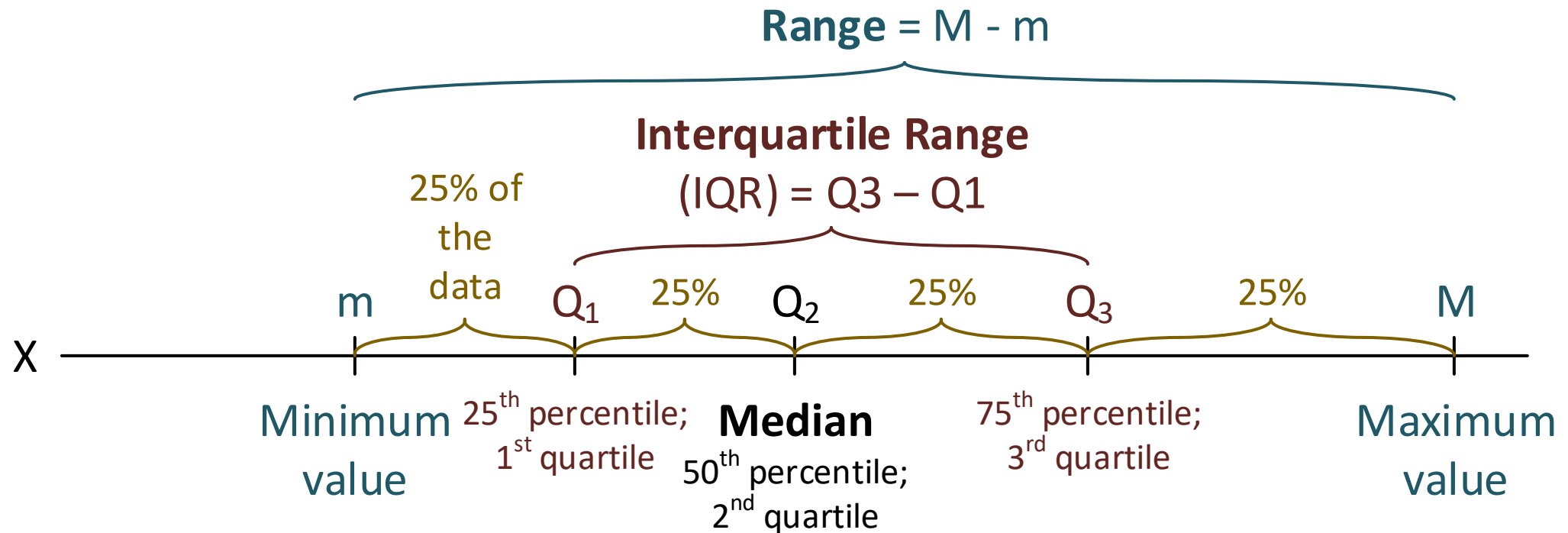
DS

Median, Range, and Interquartile Range

Median



Median, Range, and Interquartile Range



Median, Range, and Interquartile Range (cont.)

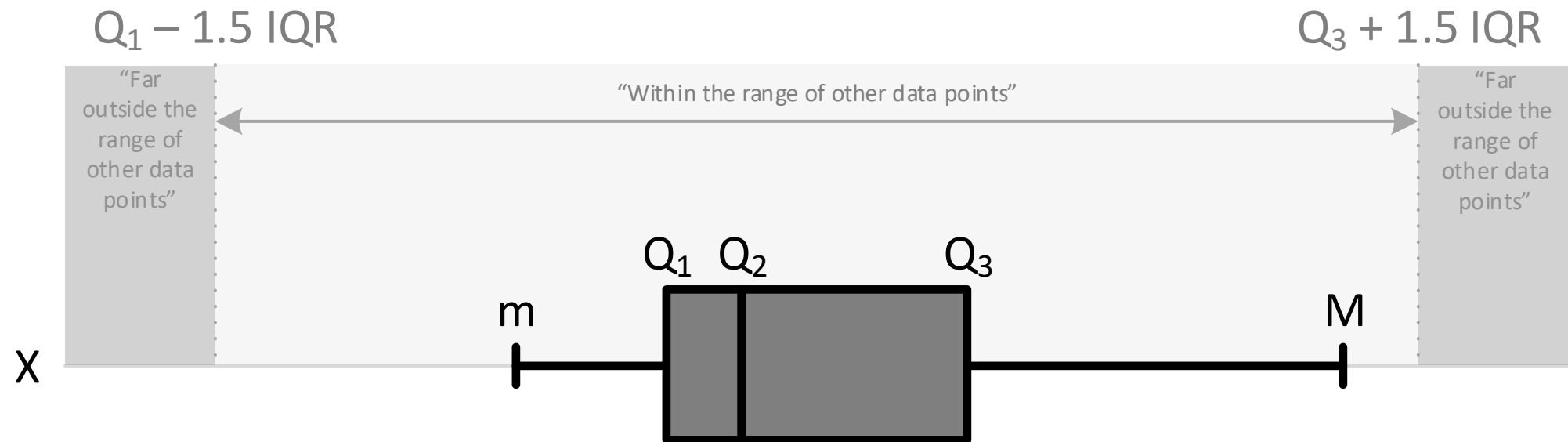
Nominal ✖		Ordinal ✖		Interval ✓		Ratio ✓	
Median		$median = \begin{cases} x_{p+1} & \text{if } n = 2p + 1 \\ \frac{x_p + x_{p+1}}{2} & \text{if } n = 2p \end{cases}$					
Range		$range = x_n - x_1$					
Percentile		$q_k = \begin{cases} x_{[p]} & \text{if } p = \frac{nk}{100} \text{ not integer} \\ \frac{x_p + x_{p+1}}{2} & \text{otherwise} \end{cases}$					
Quartile		$Q_1 = q_{25}; Q_3 = q_{75}$					
Interquartile Range		$IQR = Q_3 - Q_1$					

(median, range, and interquartile range are based on the ranks of x_i ; x_i ranked from smallest to largest)

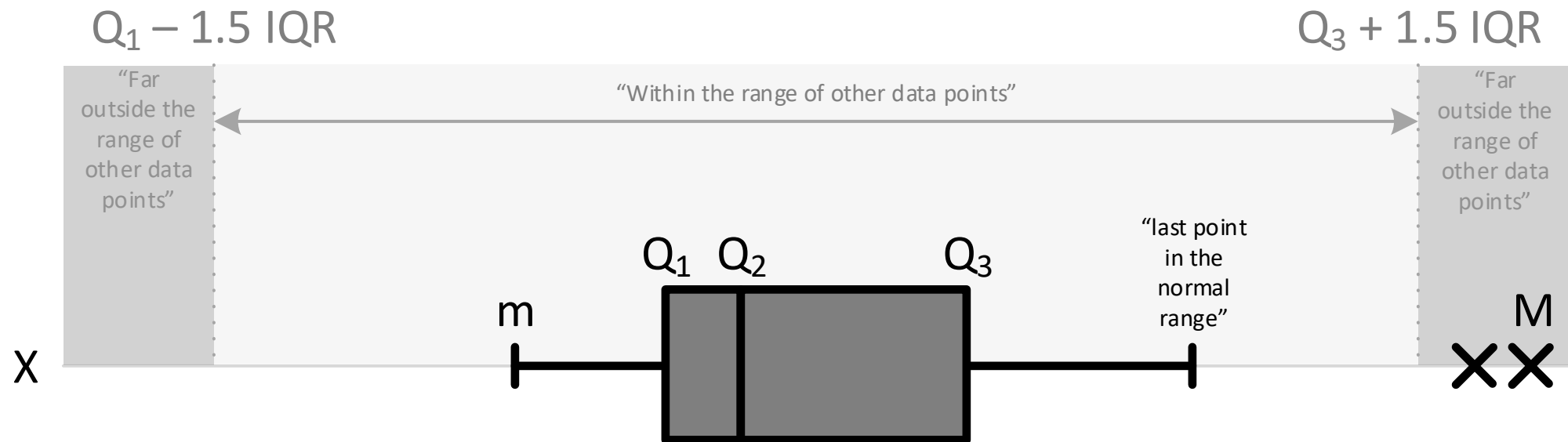
DS

Boxplots

Boxplot #1 | Median, Range, Interquartile Range; no Outliers



Boxplot #2 | Median, Range, Interquartile Range; with Outliers

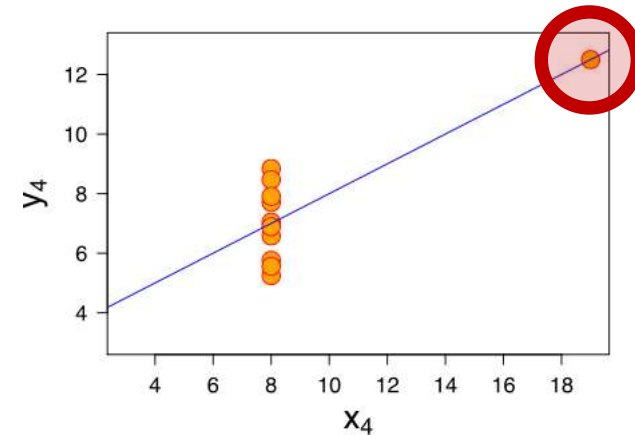
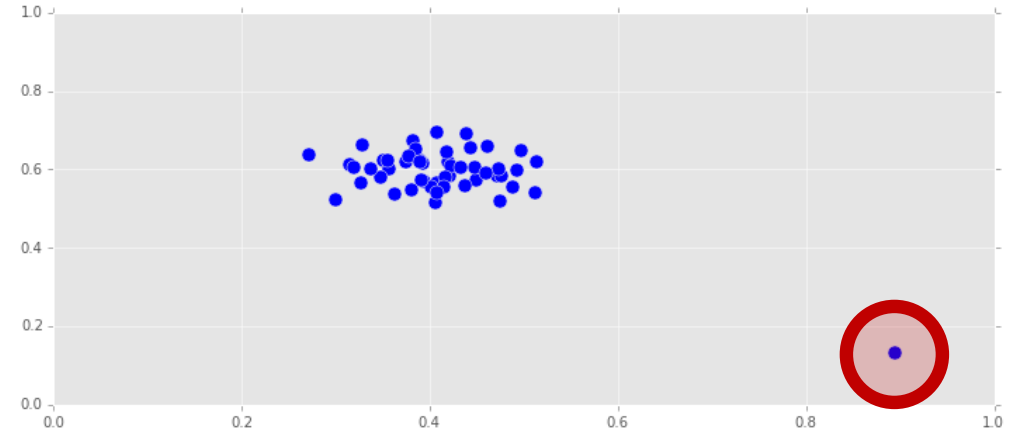


DS

Outliers

Think twice before discarding outliers; they might be the most important points of your dataset

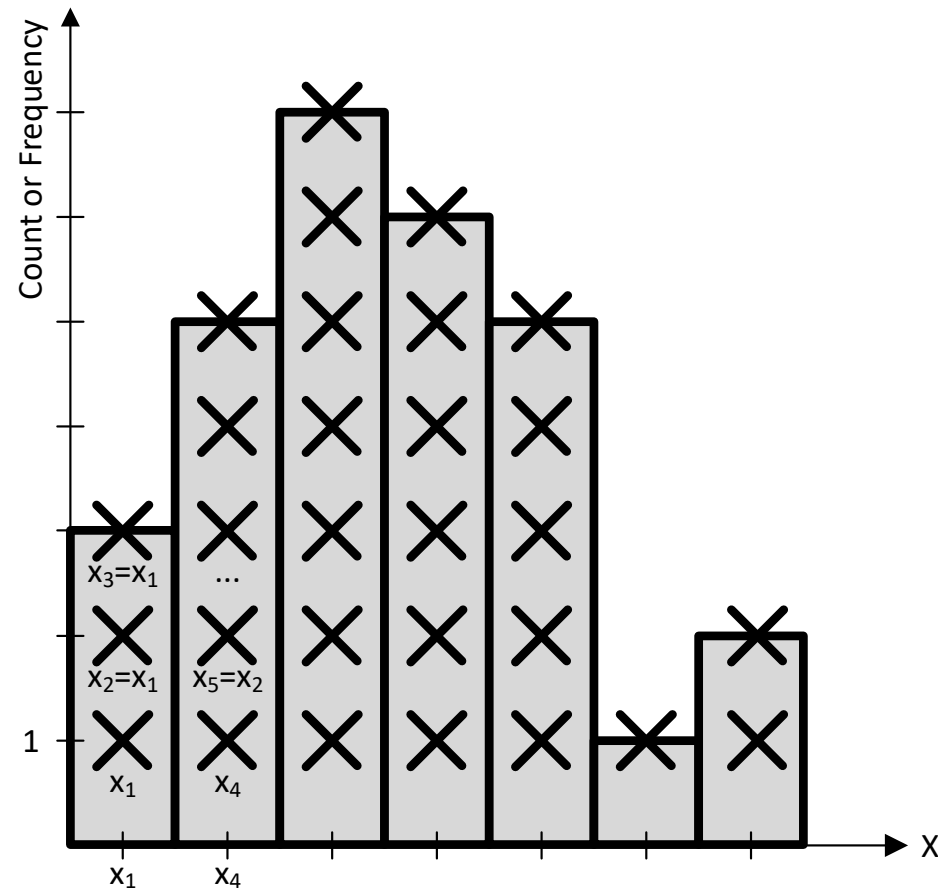
- Outliers are values that are “far” from the central tendency
- No formal definition among statisticians on how to define outliers (how do you define “far”?)
- However, general agreement that they be identified and dealt with appropriately (e.g., keep or discard)
 - They might be the most important points of your dataset



DS

Histograms

Histograms. $x_1 = x_2 = x_3 < x_4 = x_5 \dots$





Mode

Modes and Histograms

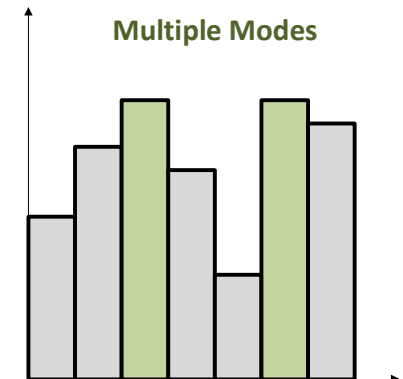
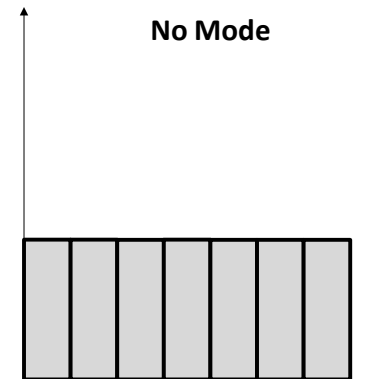
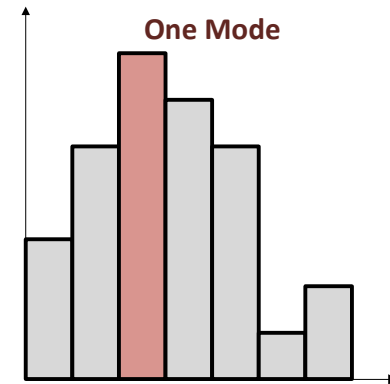
Nominal ✓

Ordinal ✓

Interval ✓

Ratio ✓

- The Mode is the value(s) that occur(s) most often



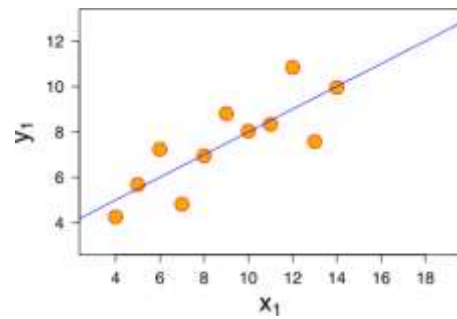


DS

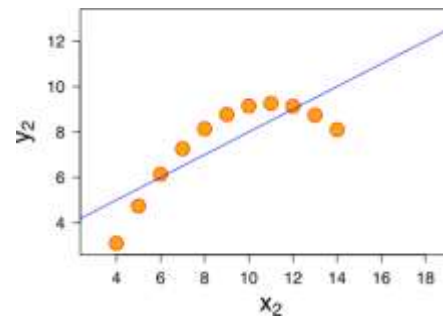
Plot the Data!

Don't rely on basic statistic properties and **plot the data!** 4 datasets (Anscombe's quartet) that have nearly identical simple statistical properties, yet are very different

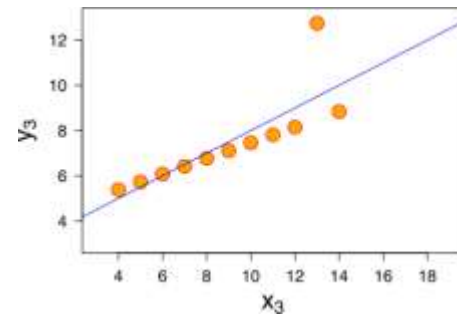
Scatter plot appears to be a simple linear relationship, corresponding to two variables correlated and following the assumption of normality.



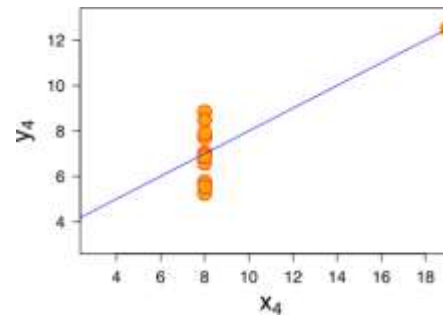
Not distributed normally; while an obvious relationship between the two variables can be observed, it is not linear, and the linear correlation is not relevant.



Distribution is linear, but with a different regression line, which is offset by the one outlier which exerts enough influence to alter the regression line.



Example when one outlier is enough to produce a high correlation coefficient, even though the relationship between the two variables is not linear.



Property	Value
Mean of x_i	9
Sample variance of x_i	11
Mean of y_i	7.50
Sample variance of y_i	4.122 or 4.127
Correlation between x_i and y_i	0.816
Linear regression line in each case	$y_i = 3.00 + 0.500 x_i$

DS

(Linear) Correlation

Correlation

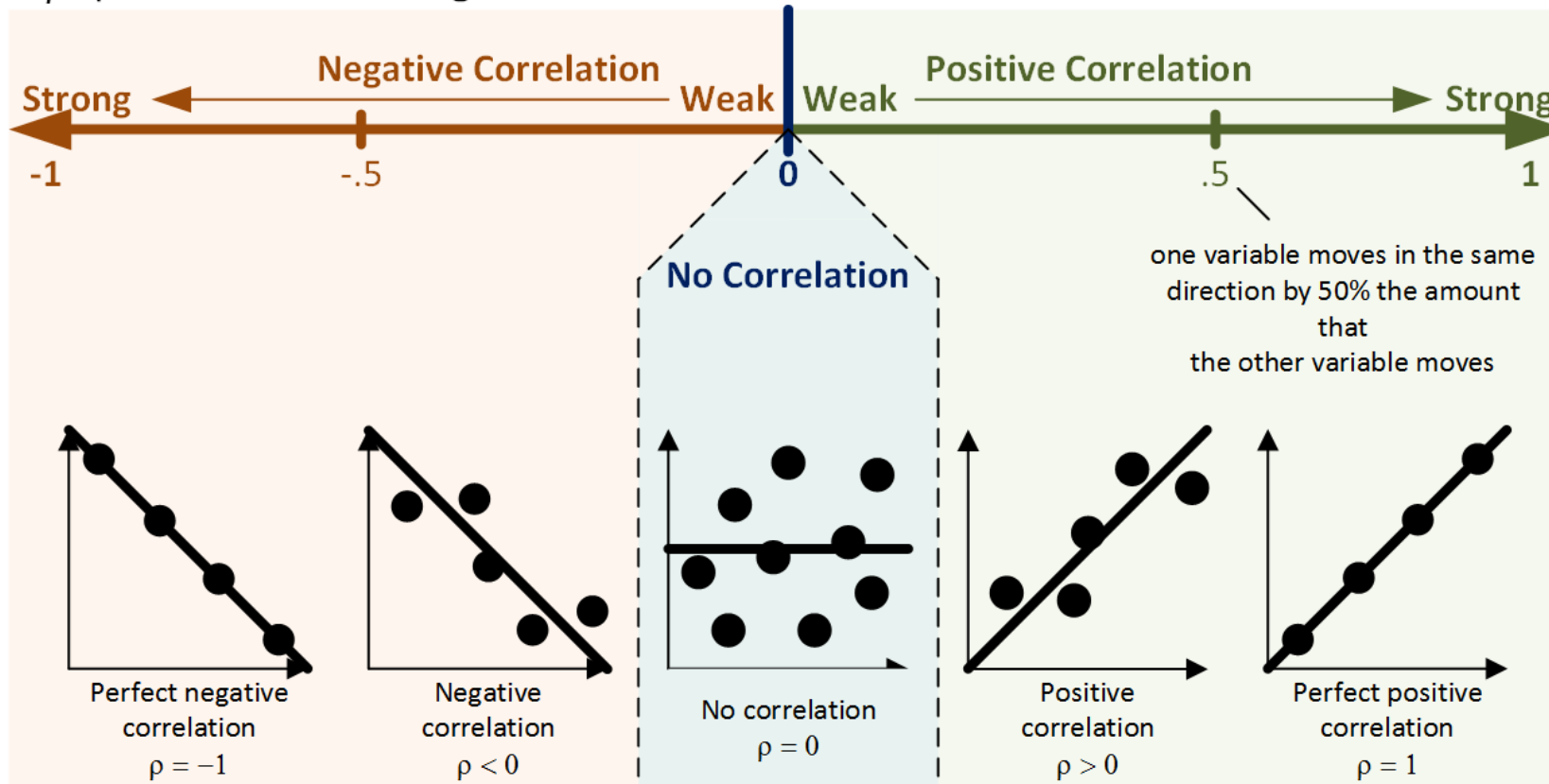
- A measure of strength and direction for a **linear association** between two random variables

$$\rho_{X,Y} = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y}$$

- $\rho = 0$ means that the two variables don't have a linear association
 - It doesn't imply that they are independent!

Correlation (cont.)

ρ quantifies the strength and direction of movements of two random variables



Slides © 2017 Ivan Corneillet Where Applicable
Do Not Reproduce Without Permission