

# Descriptive Statistics for 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



DS

# Announcements and Exit Tickets

**DS**

# Review

A black circle containing the white text 'DS' in a bold, sans-serif font.

DS

# Review

① *IDENTIFY the problem*

*The SMART Framework for Data Science*

# ① IDENTIFY the problem | The SMART Framework for Data Science:

<b>S</b> <sub>PECIFIC</sub>	The dataset and key variables are clearly defined
<b>M</b> <sub>EASURABLE</sub>	The type of analysis and major assumptions are articulated
<b>A</b> <sub>TTAINABLE</sub>	The question you are asking is feasible for your dataset and is not likely to be biased
<b>R</b> <sub>EPRODUCIBLE</sub>	Another person (or you in 6 months!) can read your state and understand exactly how your analysis is performed
<b>T</b> <sub>IME-BOUND</sub>	You clearly state the time period and population for which this analysis will pertain

Trends often change over time and vary by the population of source of your data. It is important to clearly define who/what you included in your analysis as well as the time period for the analysis

A black circle containing the white text 'DS' in a bold, sans-serif font.

DS

# Review

③ *Parse the Data*

*Tidy Data and pandas*

### ③ PARSE the Data | Tidy Data: a tabular format suitable for *pandas* and machine learning algorithms

- ▶ The three rules of tidy data:
  - ▶ Each observation is placed in its own row
  - ▶ Each variable in the dataset is placed in its own column
  - ▶ Each value is placed in its own cell

The screenshot shows an Excel spreadsheet with the following data:

ID	Address	Latitude	Longitude	DateOfSale	SalePrice	SalePriceUnit	IsAStudio	BedCount	BathCount	Size	SizeUnit	Lo
1500000000	1000000000	37804392	-122406590	12/11/2015	1.5	\$M	FALSE	1	1	1060 sqft	N/	
1500000000	1000000000	37804240	-122405509	1/15/2016	970000	\$	FALSE	2	2	1299 sqft	N/	
1500000000	1000000000	37804240	-122405509	12/17/2015	940000	\$	FALSE	2	2	1033 sqft	N/	
1500000000	1000000000	37803748	-122415151	12/15/2015	835000	\$	FALSE	1	1	1048 sqft	N/	
1500000000	1000000000	37802400	-122412405	12/4/2015	2.83	\$M	FALSE	3	2	2115 sqft	N/	
1500000000	1000000000	37801889	-122410485	11/16/2015	4.05	\$M	TRUE	N/A	N/A	4102 sqft	N/	
1500000000	1000000000	37801873	-122418873	11/16/2015	2.19	\$M	FALSE	2	3	1182 sqft	N/	
1500000000	1000000000	37803470	-122418873	11/16/2015	800000	\$	FALSE	1	1	1000 sqft	N/	
1500000000	1000000000	37802225	-122412826	11/28/2016	976000	\$	FALSE	1	1	1000 sqft	N/	
1500000000	1000000000	37801802	-122411616	11/16/2015	720000	\$	FALSE	1	1	552 sqft	N/	
1500000000	1000000000	37800260	-122406123	11/25/2015	2.25	\$M	FALSE	N/A	4	2658 sqft	N/	
1500000000	1000000000	37799474	-122414835	11/30/2015	1.29	\$M	FALSE	2	2	1165 sqft	N/	



# Activity | Subsetting with *pandas* (10 minutes)



## EXERCISE

	DataFrame	Series
Column subsetting		
by name		
by location		
Row subsetting		
by index label		
by location		
Cell subsetting/scalar lookup		
By index label/column name		
By location		

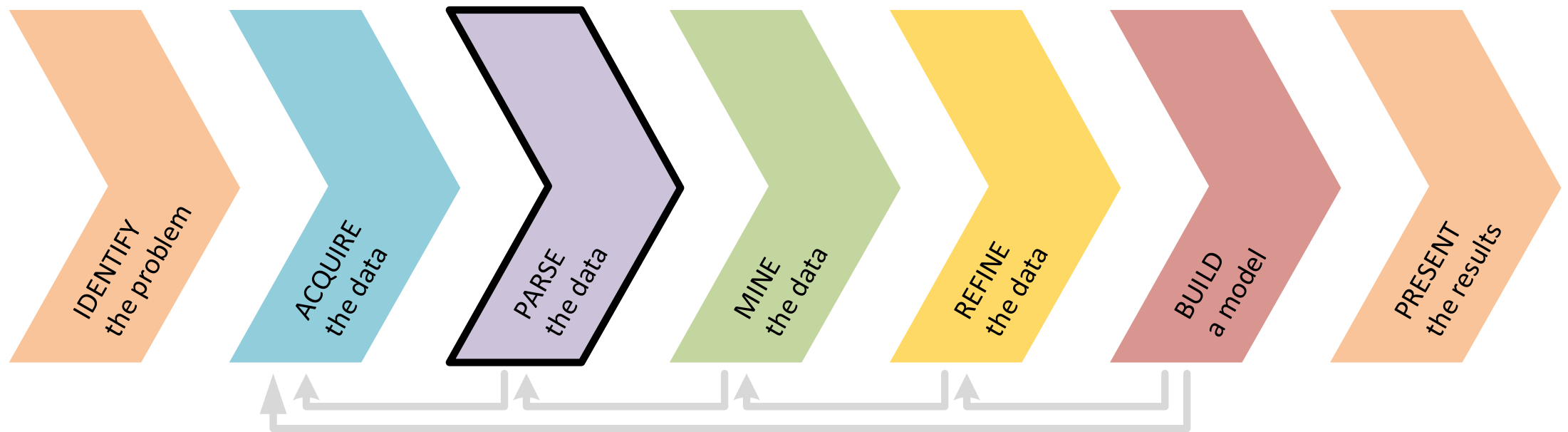
	DataFrame	Series
Column subsetting		
<b>by name</b>  (Columns names are stored in df.columns) (df.columns.get_loc('X1') returns X1's column index)	# New DataFrame with column named X1 df[ ['X1'] ]  # 2+ columns (in the order listed) df[ ['X1', 'X2', ...] ]	df['X1']  df.X1
<b>by location</b>	# New DataFrame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

A black circle containing the white text "DS".

DS

# Today

Today we'll keep our focus on **PARSE** the data



# Today, we are covering Research Design and introducing the *pandas* library

<b>Research Design and Data Analysis</b>	Research Design	Data Visualization in <i>pandas</i>	Statistics	Exploratory Data Analysis in <i>pandas</i>
<b>Foundations of Modeling</b>	Linear Regression	Classification Models	Evaluating Model Fit	Presenting Insights from Data Models
<b>Data Science in the Real World</b>	Decision Trees and Random Forests	Time Series Data	Natural Language Processing	Databases

# Here's what's happening today:

- Announcements and Exit Tickets
- Review
- **③ Parse the Data**
  - Types of Data and Types of Measurement Scales
  - Populations and Samples; Descriptive vs. Inferential Statistics
  - Measures of Central Tendency and Measures of Dispersion
- Boxplots
- Outliers
- Histograms
- Correlation
- Review
- Exit Tickets

A black circle containing the white text "DS".

DS

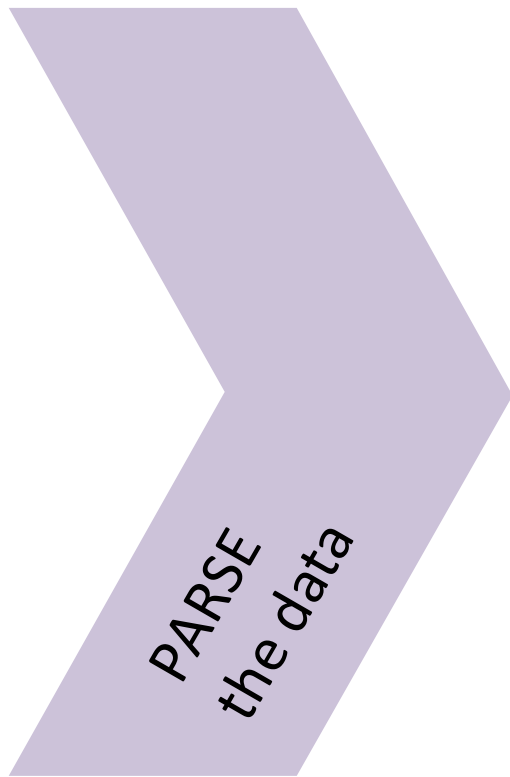
Q & A

DS

## ③ PARSE the Data

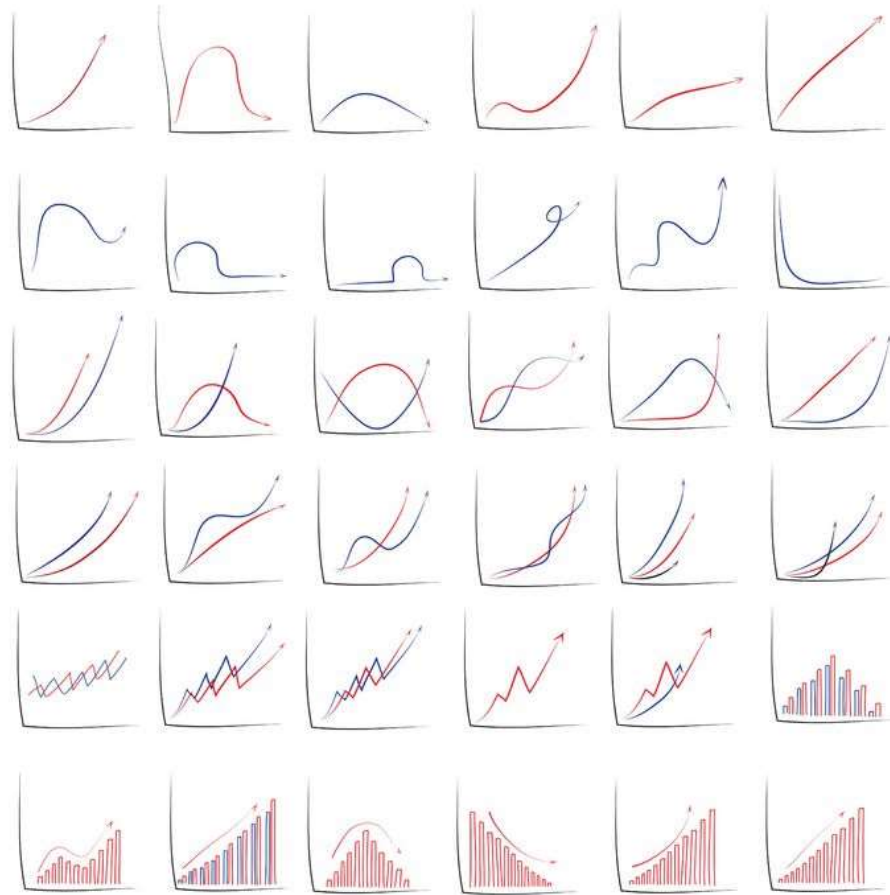


### ③ Parse the Data



- Parse the Data
  - *Read any documentation provided with the data (session 2)*
  - **Perform exploratory data analysis (session 3)**
    - *Verify the quality of the data (sessions 2/3)*

# The main theme today is to have enough statistics knowledge to perform Exploratory Data Analysis



Napat Polchoke © 123RF.com

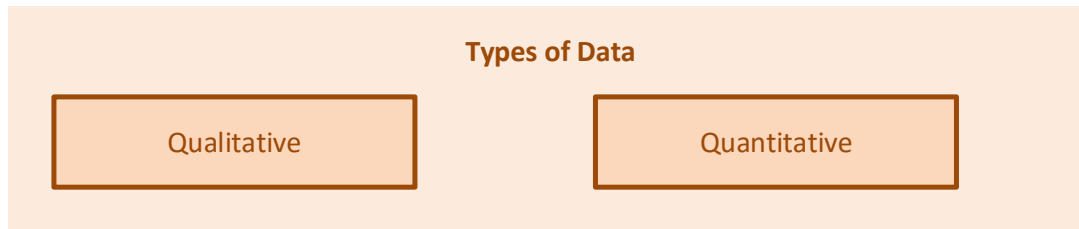
- Types of Data and Types of Measurement Scales
- Populations and Samples; Descriptive vs. Inferential Statistics
- Measures of Central Tendency and Measures of Dispersion
- Boxplots
- Outliers
- Histograms
- Correlation

DS

## ③ PARSE the Data

*Types of Data and  
Types of Measurement Scales*

# Types of Data



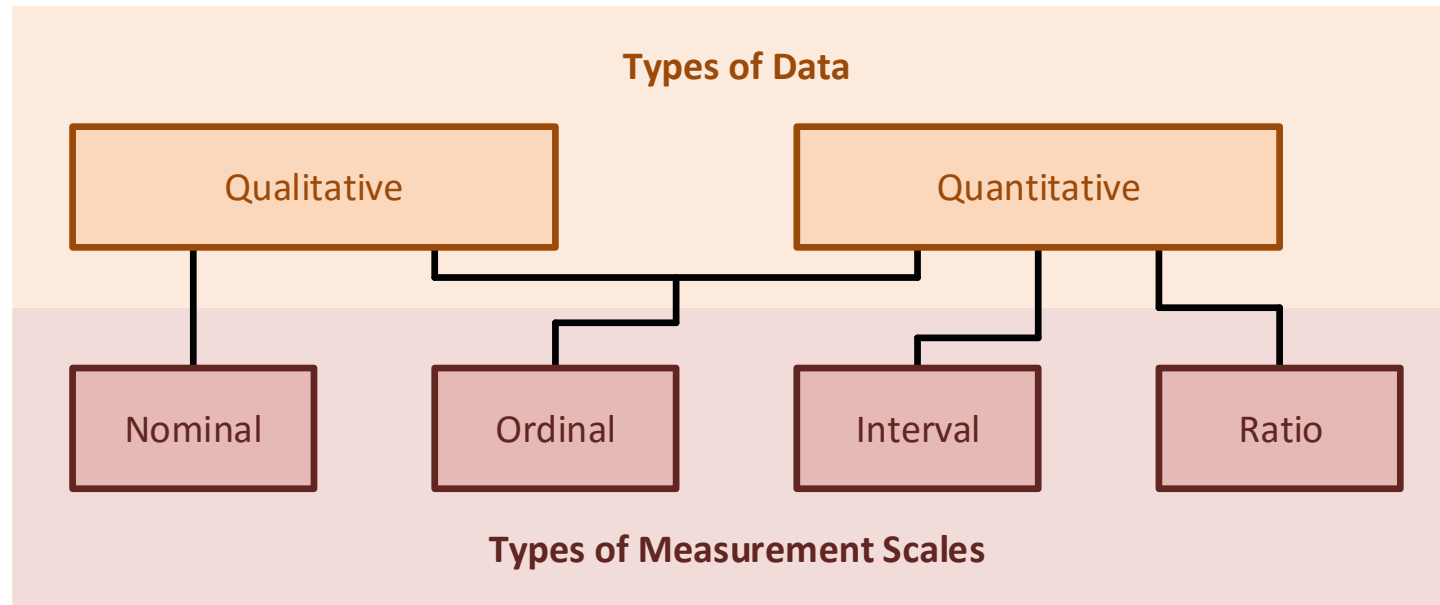
- Qualitative Data

- 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.)

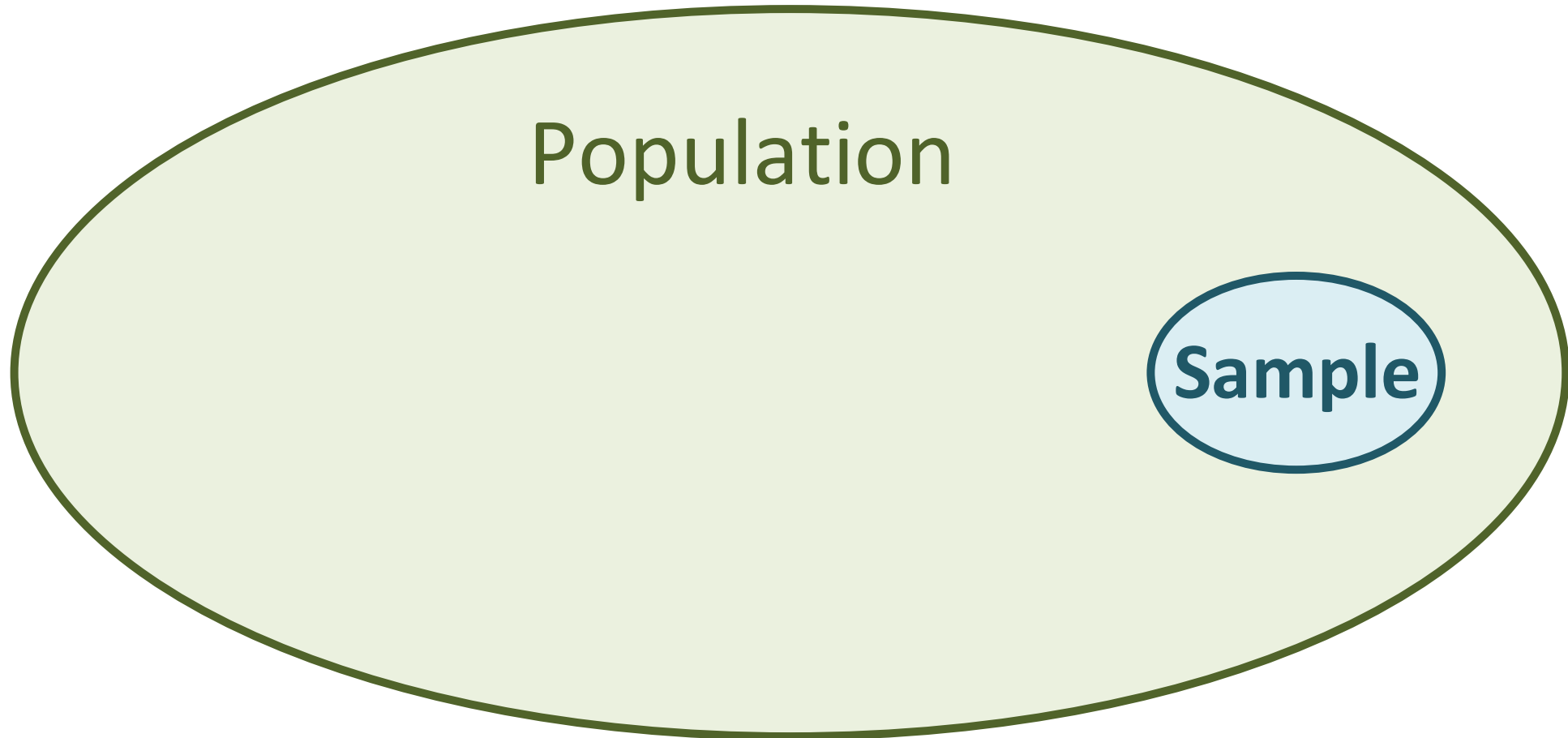
	Nominal	Ordinal	Interval	Ratio
e.g.	Gender	Movie ratings	Temperature	Salary
<b>Categorize?</b>	✓ (male, female)	✓	✓	✓
<b>Rank-order?</b>	✗	✓ (★ < 2★ < 3★ < 4★)	✓	✓
<b>Add and subtract?</b>	✗	✗ (4★ - 3★ ≠ ★)	✓ (75°C is 50°C warmer than 25°C)	✓
<b>Multiply and divide?</b>	✗	✗ (4★ not 4× better than 1★)	✗ (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

## ③ PARSE the Data

*Populations and Samples*

# Populations and Samples

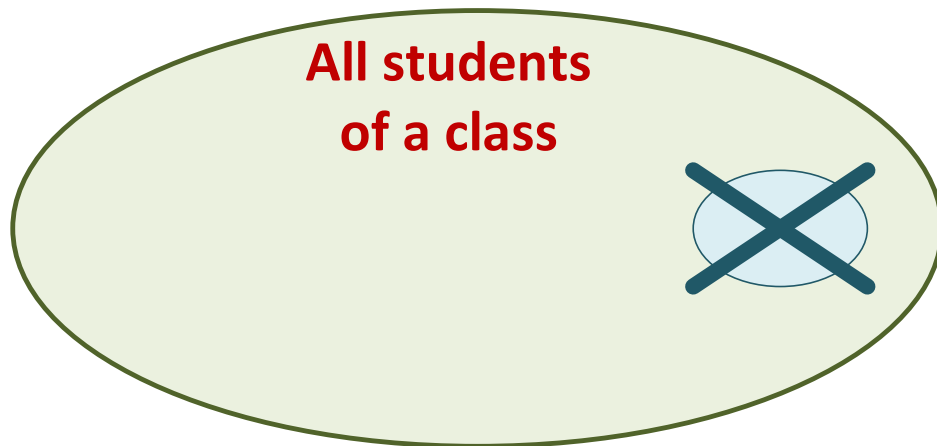




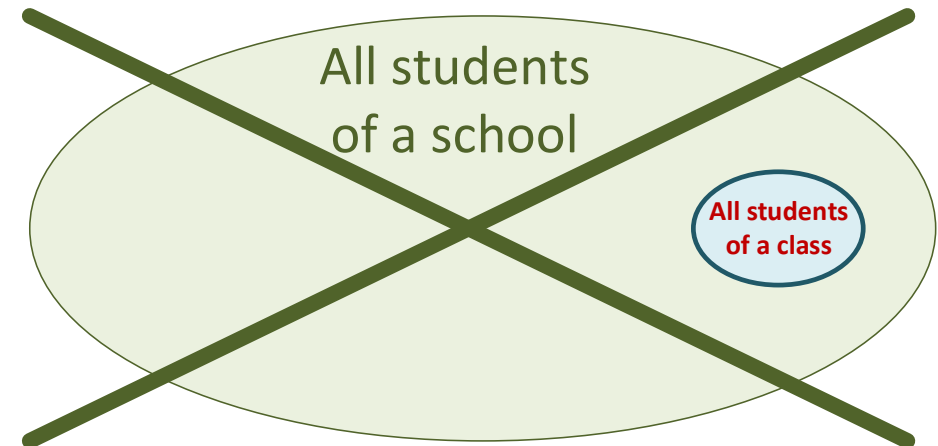
# A dataset may be considered either as a population or a sample, depending on the reason for its collection and analysis

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

## Descriptive Statistics



## Inferential Statistics



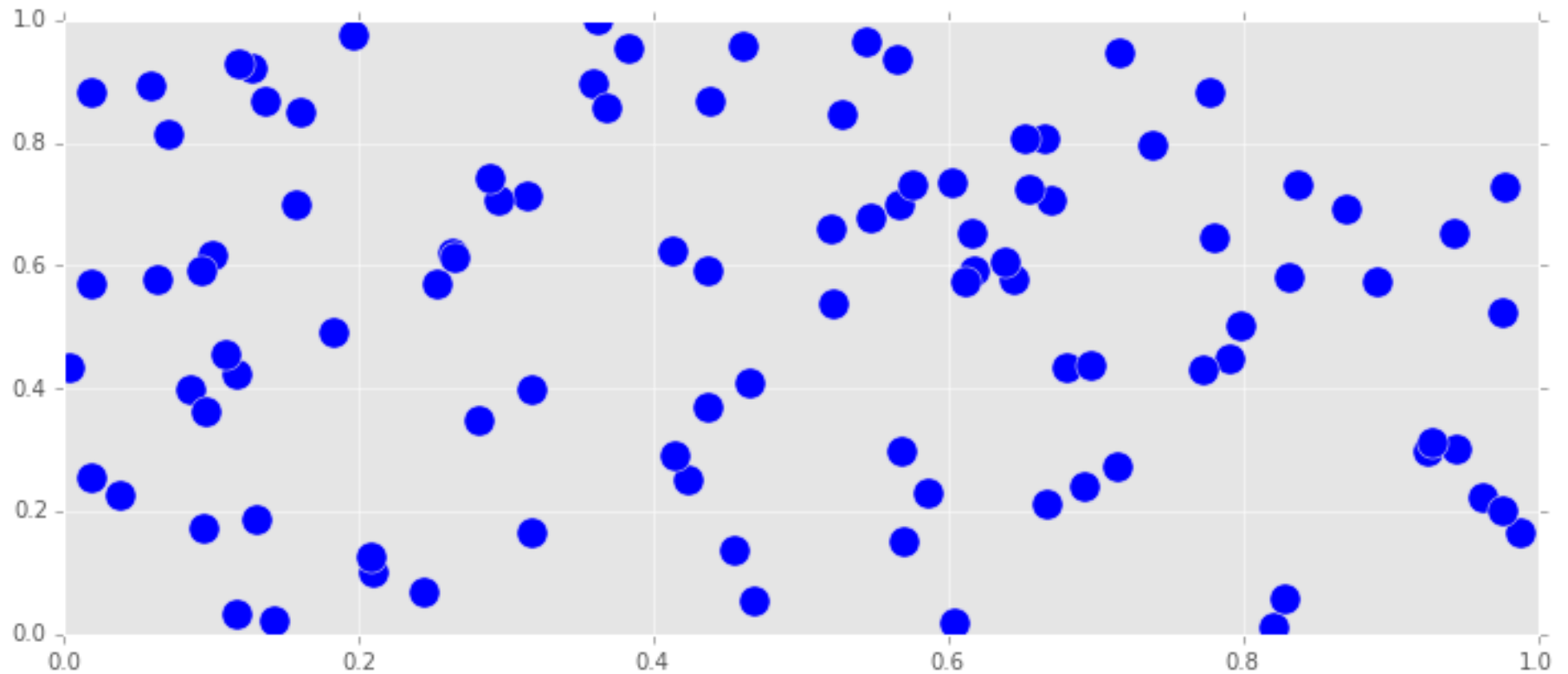
DS

## ③ PARSE the Data

*Activity / Summarizing Data*

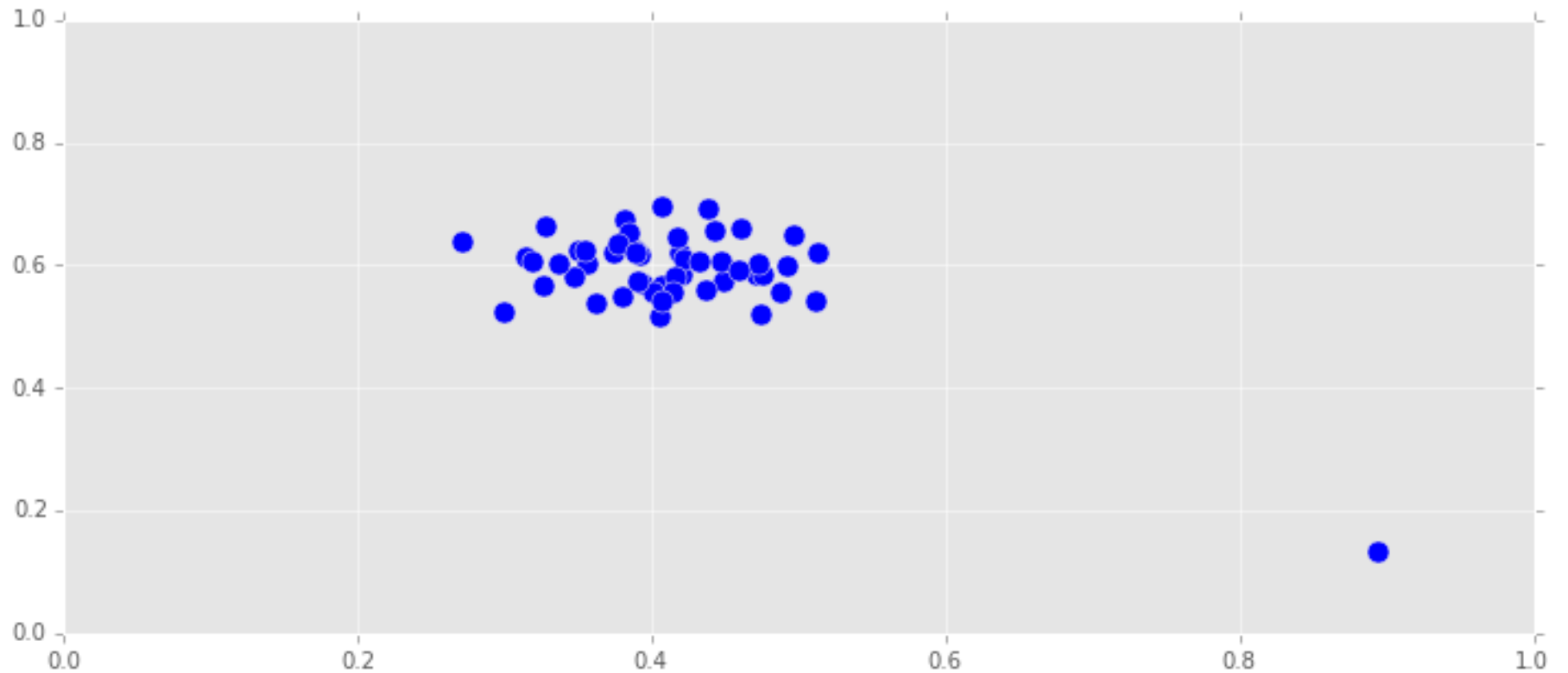
# Activity | How would you summarize this data?

EXERCISE



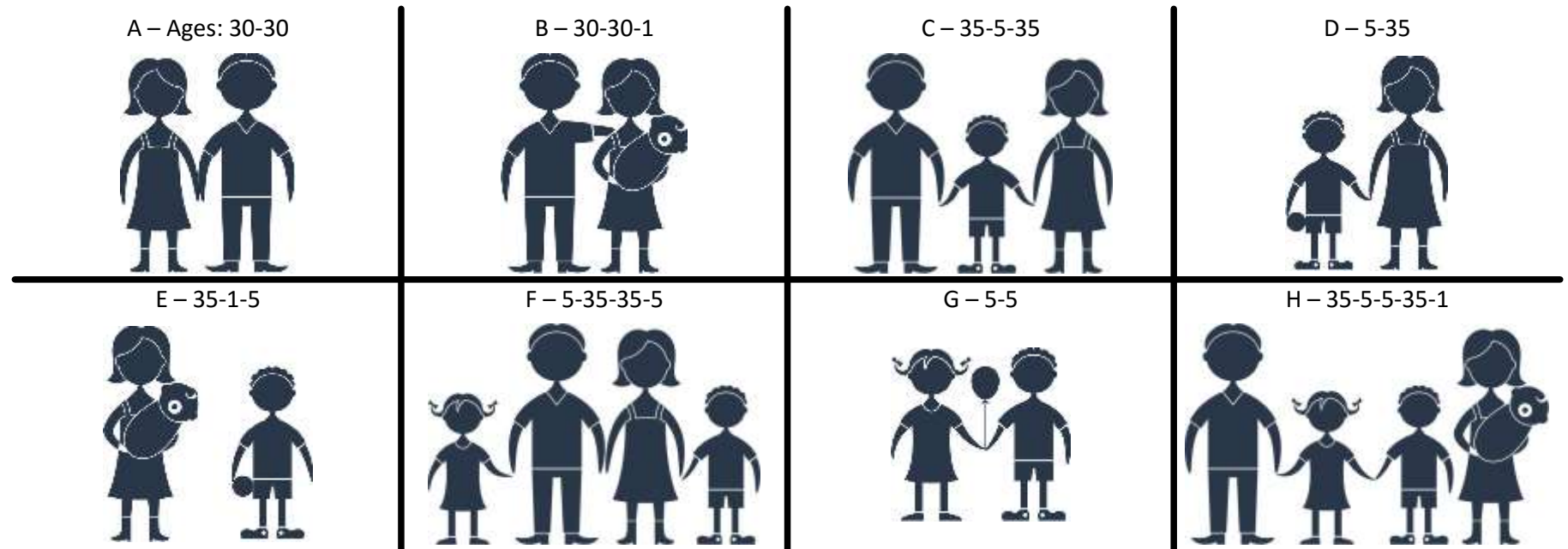
# Activity | How would you summarize this data? (cont.)

## EXERCISE



Activity | Measures of Central Tendency. What is the typical age for each of these 8 groups of people? (10 minutes)

EXERCISE



macrovector © 123RF.com

# Activity | What is the typical age for each of these 8 groups of people? (cont.)

Group	Mean	Median	Mode
A (30-30)	30 <sup>(1)</sup>	30 <sup>(1)</sup>	30 <sup>(1)</sup>
B (30-30-1)	20.3 <sup>(2)</sup> (i.e., no 20-year-olds in the group)	30 <sup>(3)</sup>	30 <sup>(3)</sup>
C (35-5-35)	25 <sup>(2)</sup>	35 <sup>(3)</sup>	35 <sup>(3)</sup>
D (5-35)	20 <sup>(2)</sup>	20 <sup>(2)</sup>	None <sup>(4)</sup>
E (35-1-5)	13.6 <sup>(2)</sup>	5 <sup>(2)</sup>	None <sup>(4)</sup>
F (5-35-35-5)	20 <sup>(2)</sup>	20 <sup>(2)</sup>	5 and 35 <sup>(5)</sup>
G (5-5)	5 <sup>(1)</sup>	5 <sup>(1)</sup>	5 <sup>(1)</sup>
H (35-5-5-35-1)	16.2 <sup>(2)</sup>	5 <sup>(6)</sup>	5 and 35 <sup>(5)</sup>

<sup>(1)</sup> All values are equal

<sup>(2)</sup> Value not representative
















<sup>(3)</sup> Follow the “majority”

<sup>(4)</sup> All values are different

<sup>(5)</sup> Follow the “majorities”

<sup>(6)</sup> Partially correct

# Mean, Median, and Mode: There is no “Winner-Take-All”

	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)

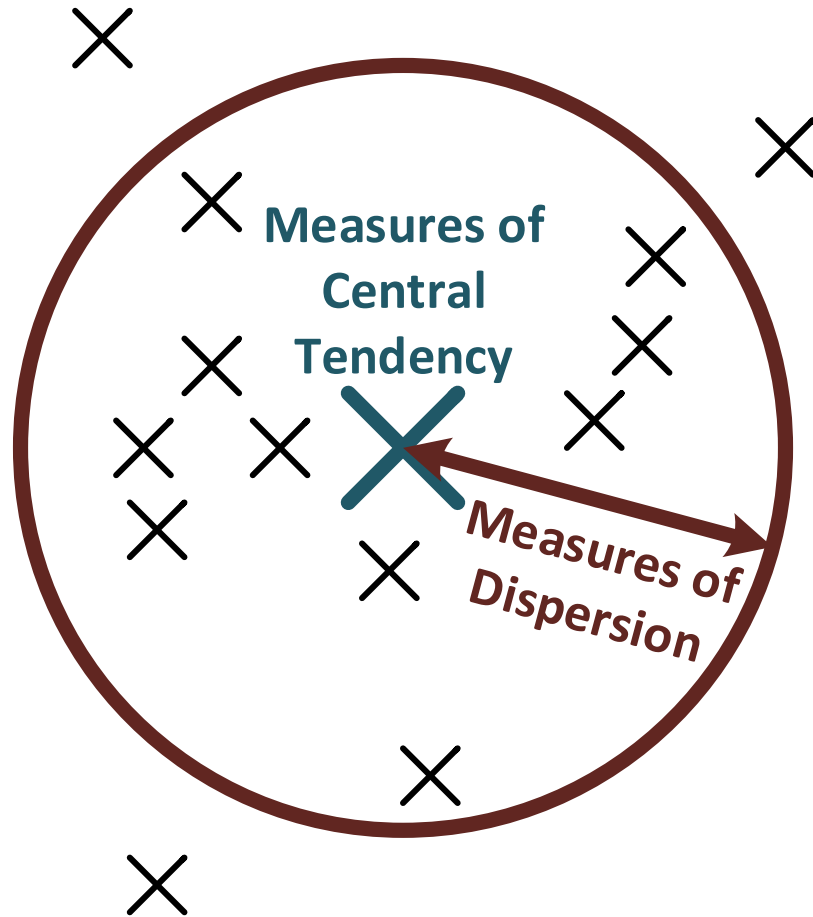
A black circle containing the white text 'DS'.

## ③ PARSE the Data

*Measures of Central Tendency and Measures of Dispersion*



# 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	
<b>(Arithmetic) Mean</b> <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)	
<b>Variance</b> <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)		$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$	
<b>Standard Deviation</b> $[X]$	$\sigma = \sqrt{\sigma^2}$ (sigma)		$s = \sqrt{s^2}$	

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

## ③ PARSE the Data

*Codealong – Part A*

*.mean()*

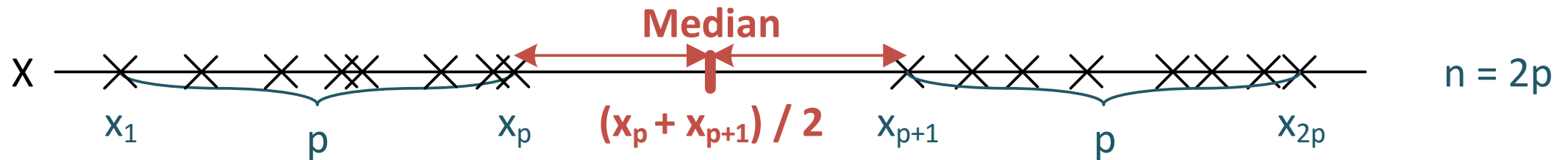
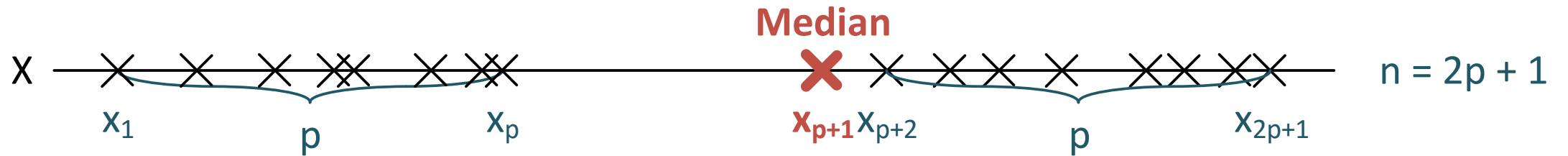
*.var(), .std()*

DS

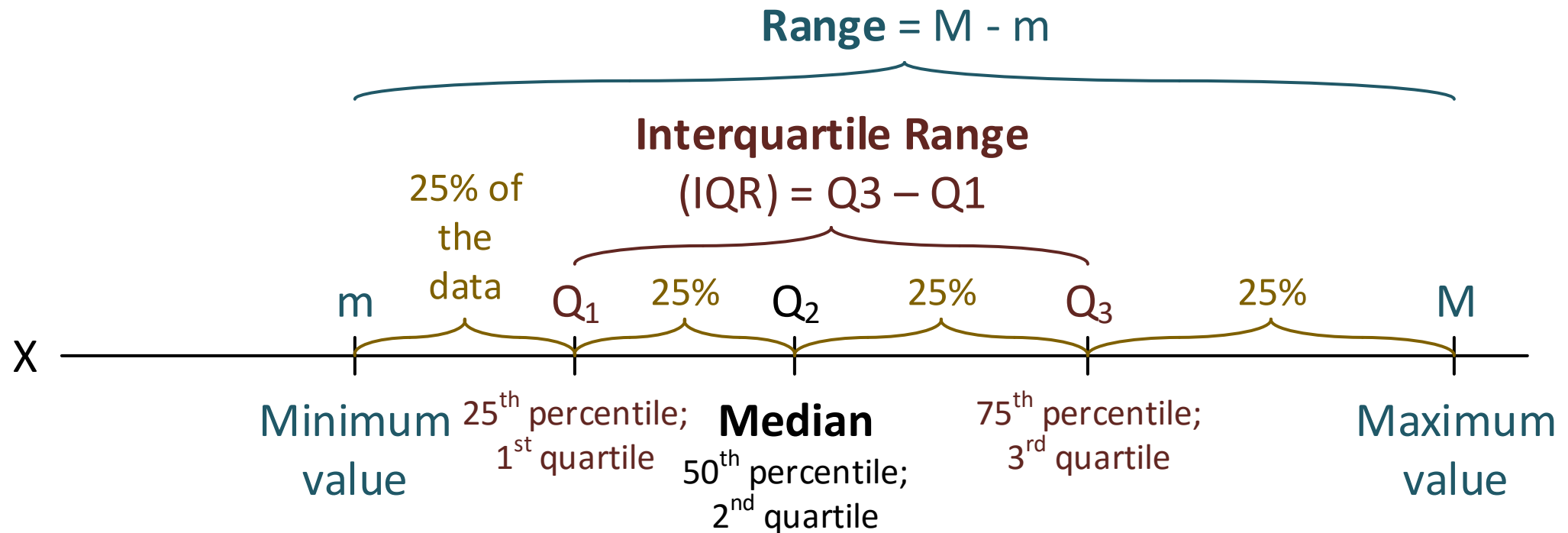
## ③ PARSE the Data

*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)

## ③ PARSE the Data

*Codealong – Part B*

```
.mean(), .median()  
.count(), .dropna(), .isnull()  
.min(), .max()  
.quantile()  
.describe()
```

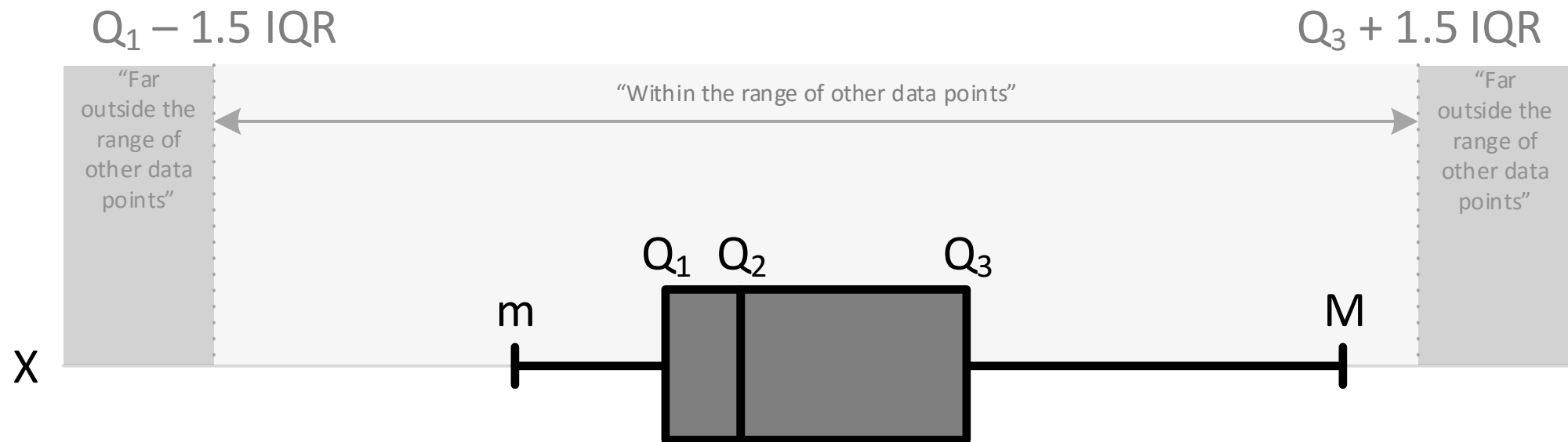


DS

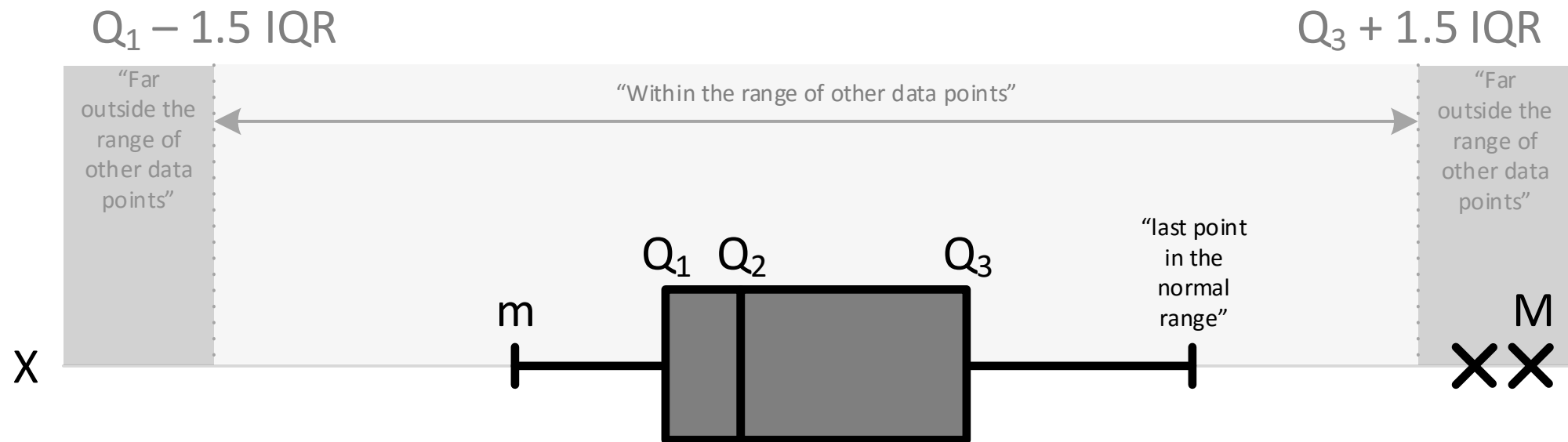
## ③ PARSE the Data

*Median, Range, Interquartile Range, and Boxplots*

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



# Boxplot #2 | Median, Range, Interquartile Range, and Outliers



DS

## ③ PARSE the Data

*Codealong – Part C*

*Boxplots*

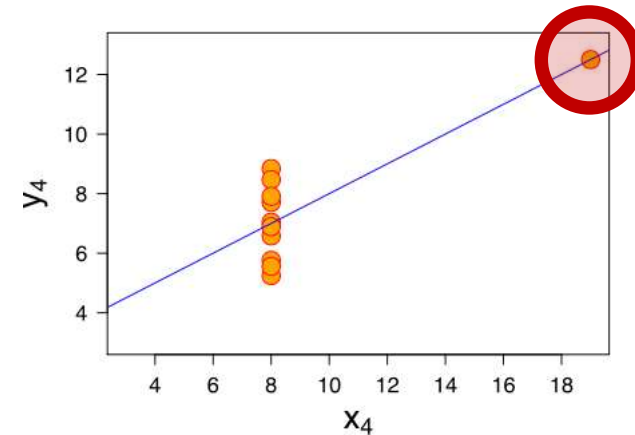
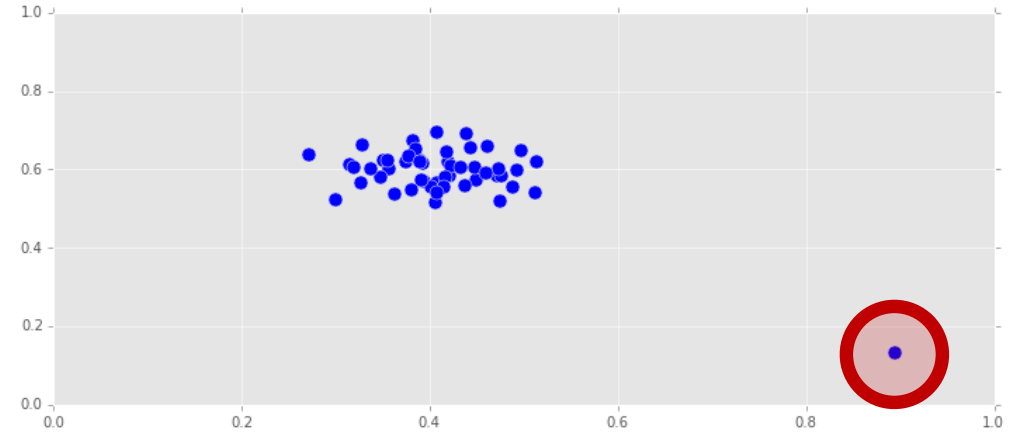
DS

## ③ PARSE the Data

*Outliers*

# Think twice before discarding outliers; they might be the most important points

- 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

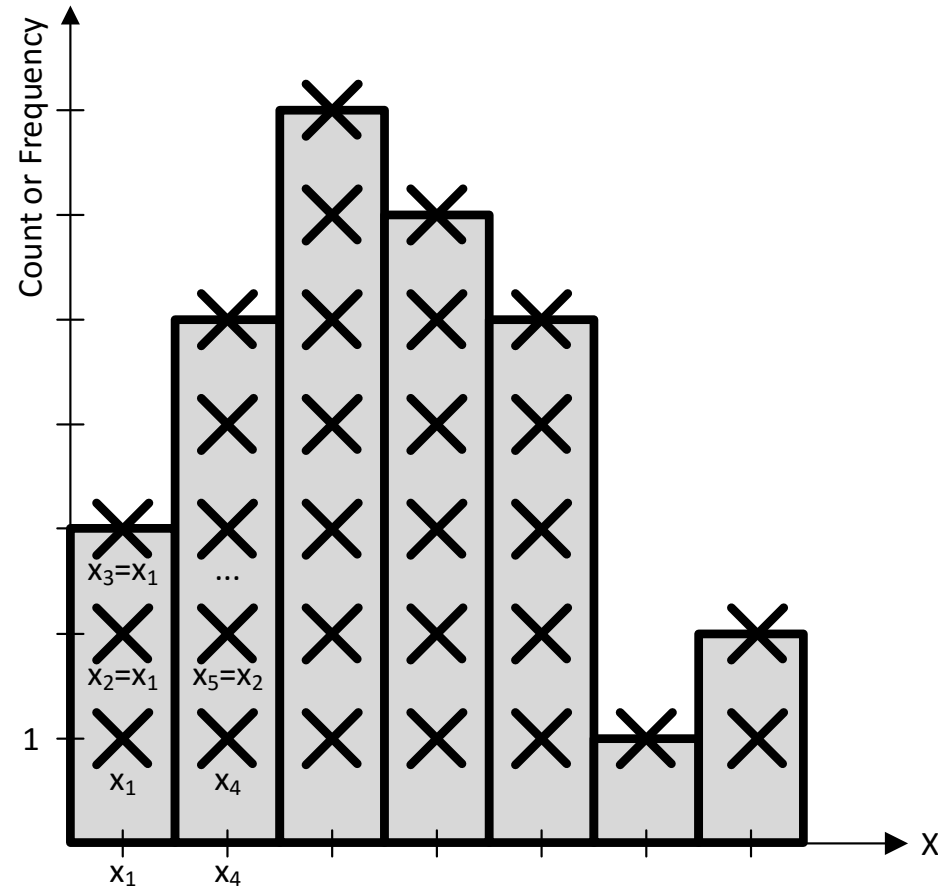


DS

## ③ PARSE the Data

*Histograms*

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





DS

## ③ PARSE the Data

*Codealong – Part D*

*Histograms*

DS

## ③ PARSE the Data

*Mode*

# Modes and Histograms

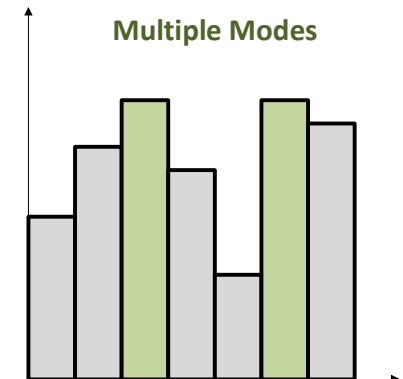
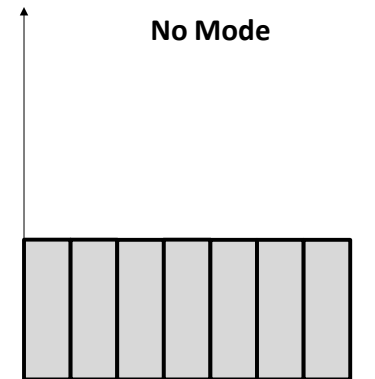
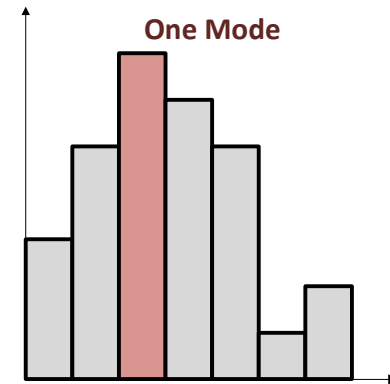
Nominal ✓

Ordinal ✓

Interval ✓

Ratio ✓

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



DS

## ③ PARSE the Data

*Codealong – Part E*

*.mode()*

A black circle containing the white text "DS".

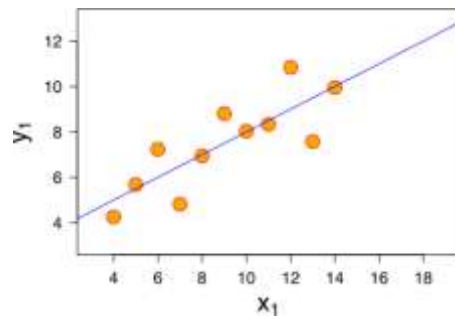
DS

## ③ PARSE the Data

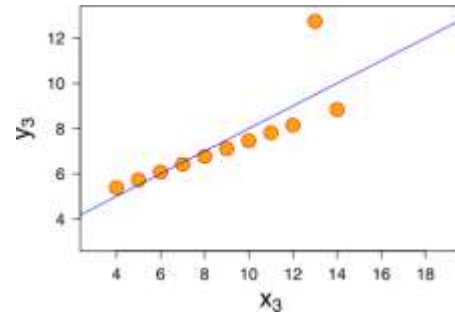
*Plot the Data!*

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

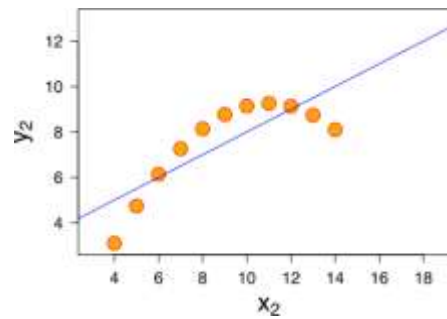
Scatter plot appears to be a simple linear relationship, corresponding to two variables correlated and following the assumption of normality.



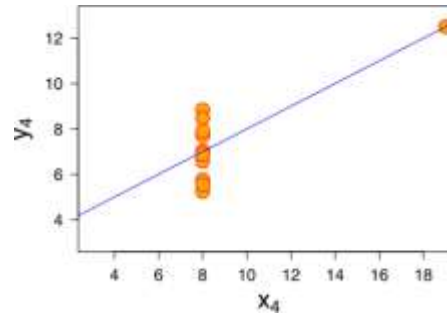
Distribution is linear, but with a different regression line, which is offset by the one outlier which exerts enough influence to alter the regression line.



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.



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

## ③ PARSE the Data

*(Linear) Correlation*

# Correlation

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

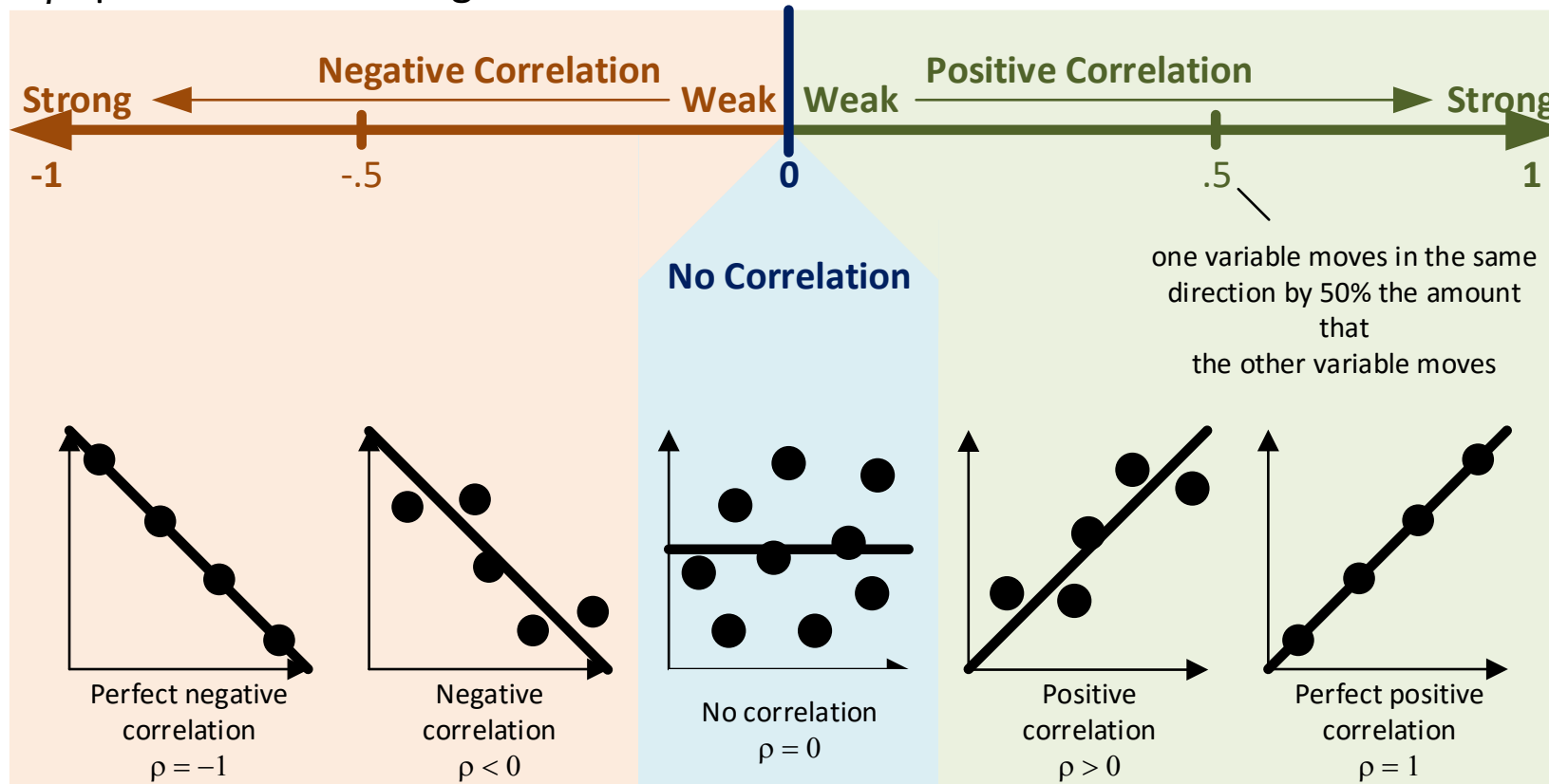
$$\rho_{X,Y} = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y}$$

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



# Correlation (cont.)

$\rho$  quantifies the strength and direction of movements of two random variables





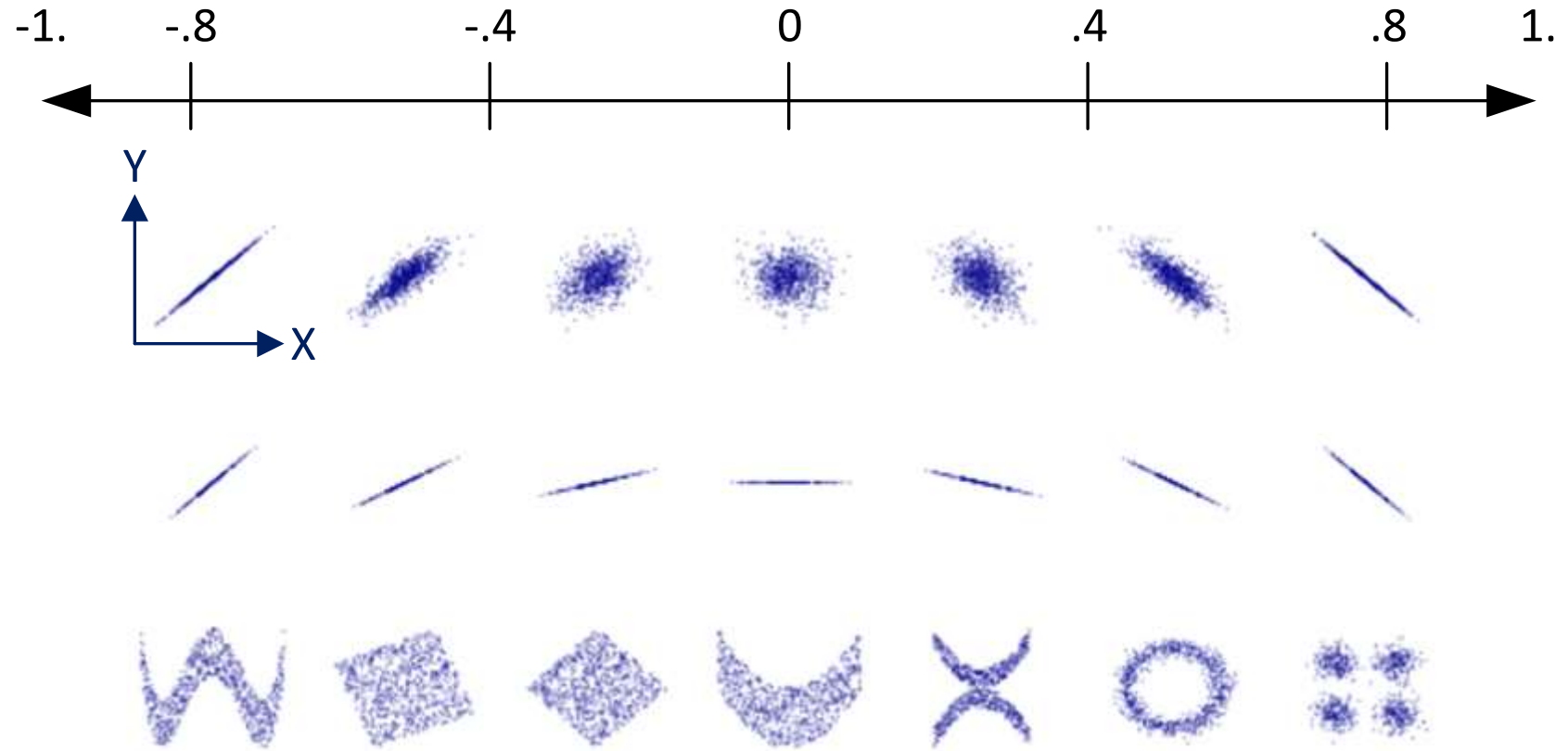
DS

## ③ PARSE the Data

*Activity / Correlations and Scatter Plots*

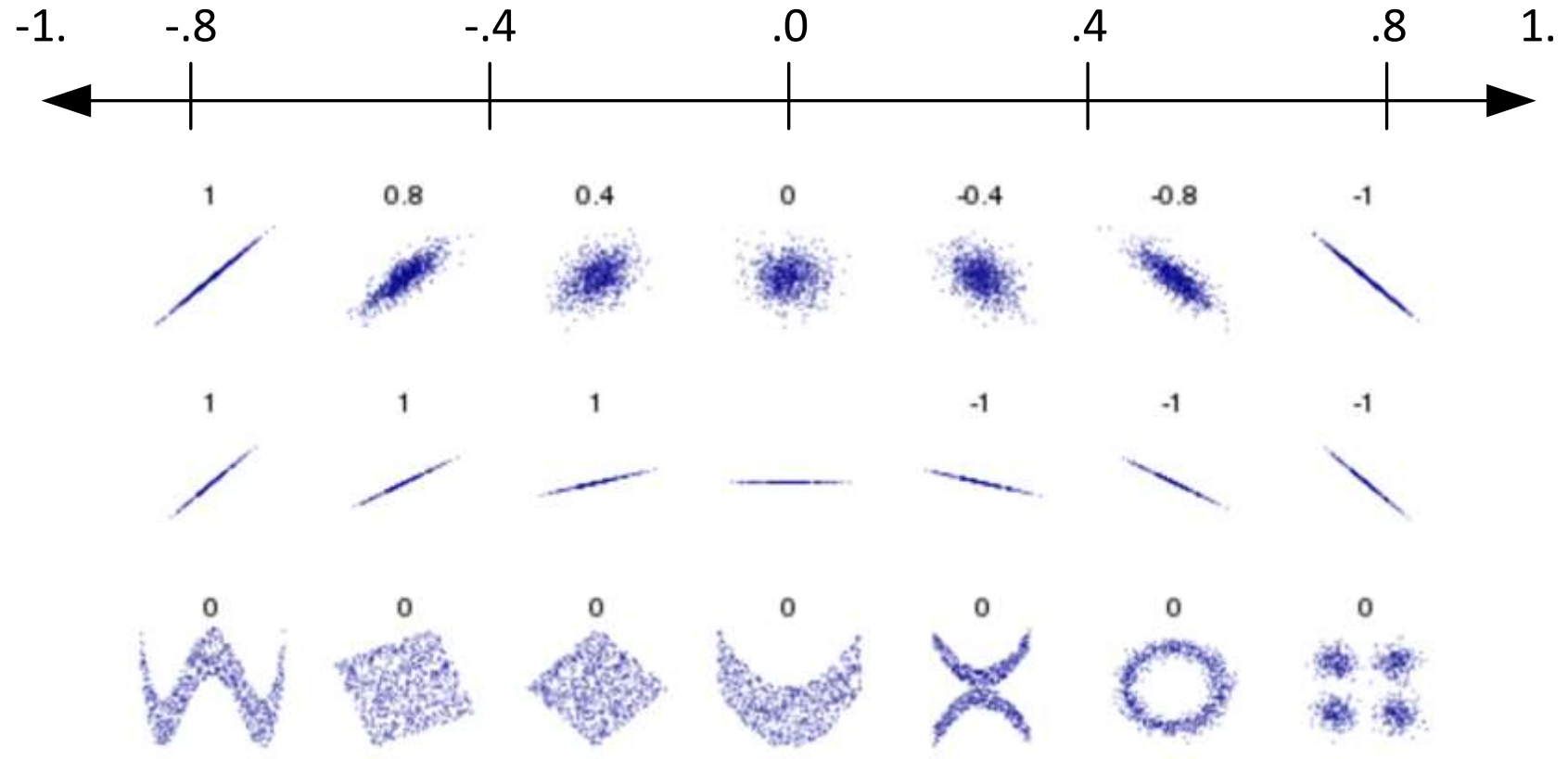
# Activity | What's the correlations for the following scatter plots (5 minutes)

## EXERCISE



# Activity | What's the correlations for the following scatter plots (cont.)

## EXERCISE



DS

## ③ PARSE the Data

*Codealong – Part F*

*.corr()*

*Heatmaps*

*Scatter plots and matrices*

## ③ PARSE the Data

*Codealong – Part G*

```
.value_counts()  
.crosstab()
```



# Lab

*Exploratory Data Analysis with pandas*

**DS**

# Review



# Review

You should now 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

A black circle containing the white text "DS".

DS

Q & A

# Next Class

*Flexible Class Session #1 | Exploratory Data Analysis*



DS

# Exit Ticket

*Don't forget to fill out your exit ticket [here](#)*

Slides © 2016 Ivan Corneillet Where Applicable  
Do Not Reproduce Without Permission