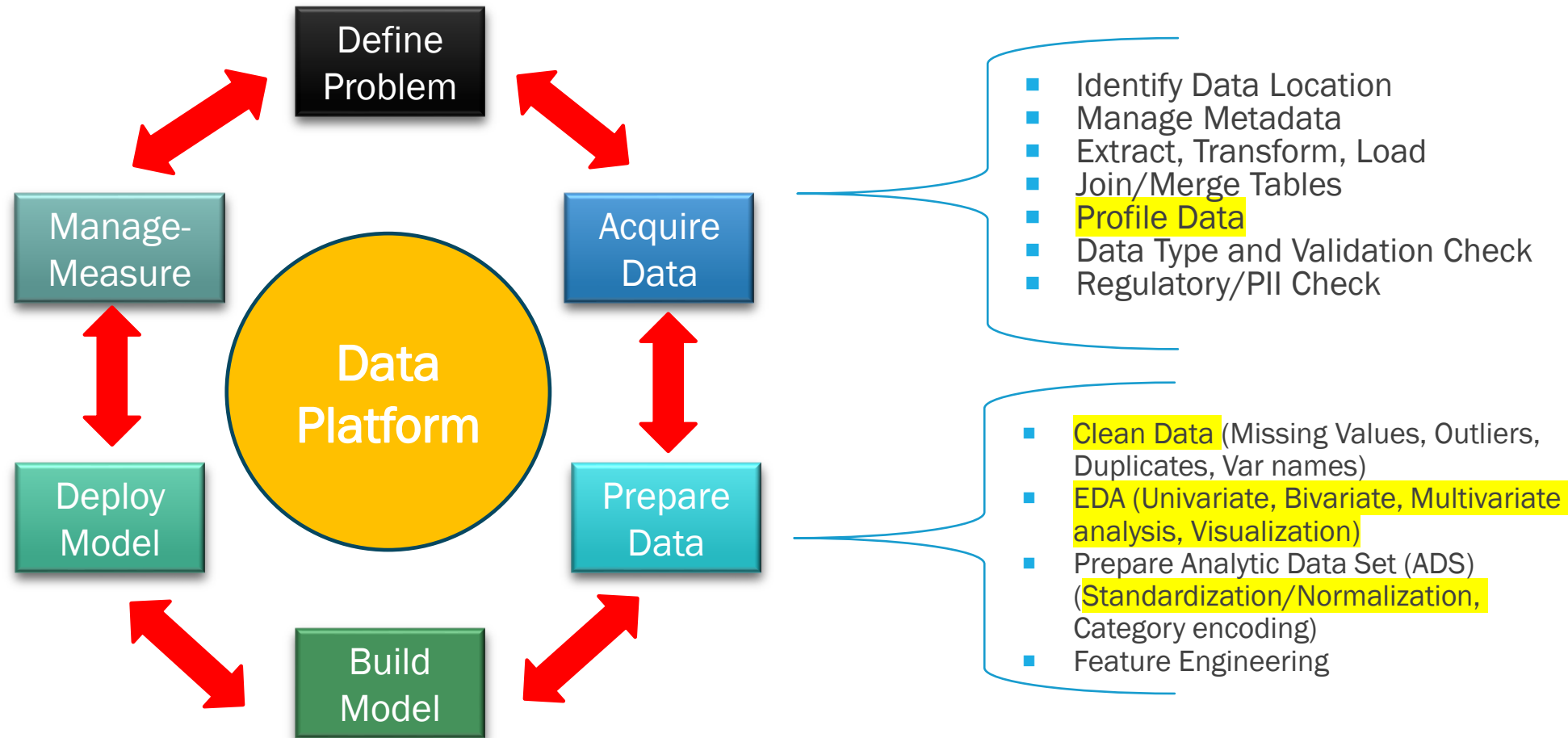


A surreal illustration depicting a data center where King penguins are the primary subjects. Several penguins stand on desks, looking at computer monitors displaying various data visualizations like bar charts and line graphs. In the background, a large window shows a snowy, mountainous landscape with a red and blue ship on the water. Numerous birds are flying in the sky. The scene is lit with a cool blue light, and a semi-transparent dark blue banner with the text "Exploratory Data Analysis" is centered over the image.

# Exploratory Data Analysis



# DATA SCIENCE LIFECYCLE



Sometimes the EDA is the objective of the project.

# SUMMARY: DATA ACQUISITION → PREPARATION → EXPLORATION

Data Analysis Processes	Tools
Acquire Data	<code>pd.read_csv()</code> , <code>pd.read_excel()</code> , <code>df.to_csv(index = False)</code> , <code>df.to_excel(index = False)</code> ,
Profile Data	<code>df.head()</code> , <code>df.tail()</code> , <code>df.info()</code> , <code>df.describe(include='all')</code> , <code>df.dtypes</code> , <code>df.shape</code>
Manipulate Data	<code>df.groupby()</code> , <code>df.pivot_table()</code> , <code>df.insert()</code>
Clean Data Missing Data Imputing Data Duplicates Outliers Data Types	<code>df.isnull().sum()</code> , <code>df.isna()</code> , <code>df.fillna(value=)</code> , <code>df.drop()</code> , <code>df.drop_duplicates()</code> , <code>df.dropna()</code> , <code>df.replace()</code> , <code>df.astype()</code> ,
Exploratory Data Analysis Univariate Multivariate Visualization	<code>df.plot()</code> , <code>plt.show()</code> , <code>df.plot.scatter()</code> , <code>df.plot.box()</code> , <code>plt.hist()</code> , <code>plt.subplots(figsize=)</code> , <code>sns.histplot()</code> , <code>sns.kdeplot()</code> , <code>sns.boxplot()</code> , <code>sns.violinplot()</code> <code>df.mean()</code> , <code>sns.scatterplot()</code> , <code>sns.swarmplot()</code> , <code>sns.stripplot()</code> , <code>plotly.express</code>
Transform Data	<code>preprocessing.scale()</code> , <code>preprocessing.MinMaxScaler().fit_transform()</code> [Part of Modeling because we need to wait until after train/validation/test split]

# EXPLORATORY DATA ANALYSIS

Process 4.5.1: Exploratory data analysis.

These four steps should be followed in conducting exploratory data analysis:

Step	Description
Step 1: Understand the data	This is done as part of the "Profile Data" step before analysis. Objective is to understanding the data types, numbers, ranges, overall cleanliness
Step 2: Detect and address Outliers and missing data	The is done in the Data Cleaning stage. Data needs to be cleaned before analysis; otherwise, analysis could be skewed by dirty data.
Step 3: Describe the shape of each feature of the data	Use descriptive univariate statistics and visualization to characterize data distributions for each feature.
Step 4: Identify and address correlations between features	Use multivariate analysis. Assess whether features with high (+/-) correlations can be dropped.

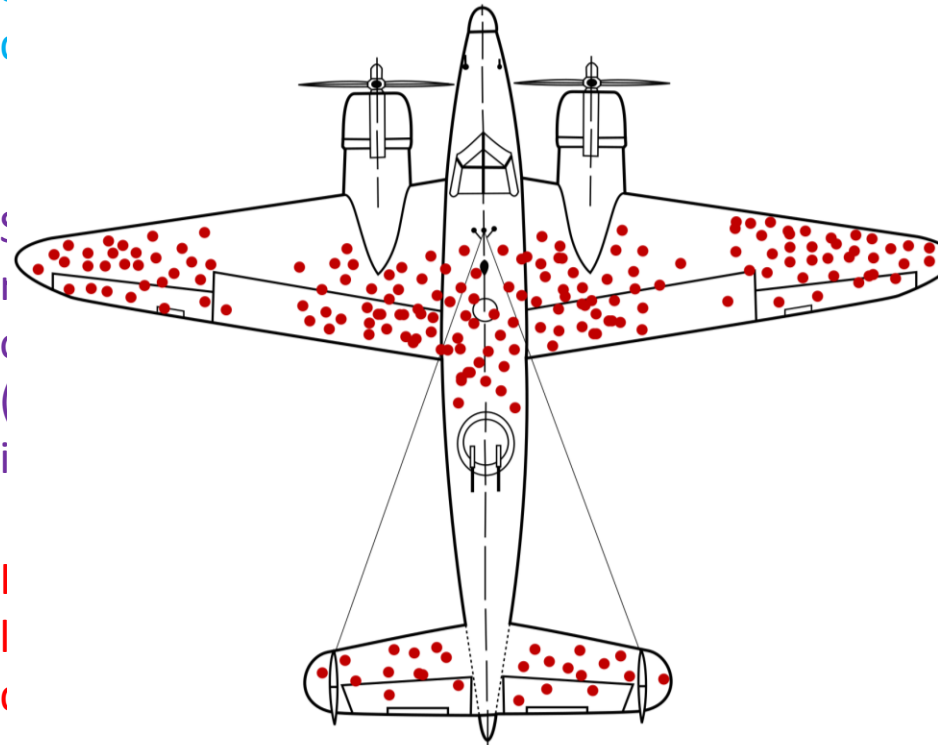
# MISSING DATA

**“Missing Completely at Random” (MCAR)** refers to where the missingness occurs purely by chance and has no relationship with any of the variables in the dataset, either the missing ones or the observed ones.

**“Missing at Random” (MAR)** refers to where the missingness can be explained by other known (observed) variables in the dataset, but not by the value of the missing variable.

**“Missing Not at Random” (MNAR)** refers to when the missingness is not random and cannot be explained by the observed data. Instead, it is systematically related to the unobserved, missing values.

Imagine a survey where respondents are asked to provide their age and income, and some survey forms are accidentally lost



me. Some  
the likelihood  
location level  
itself (which

natically less  
variables, the

# MISSING DATA QUIZ 1

## 1. Which of the following best defines “Missing Completely at Random” (MCAR)?

- a) The missingness is related to both observed and missing data.
- b) The missingness is related to observed data but not missing data.
- c) The missingness is unrelated to both observed and missing data.
- d) The missingness is related to the missing data only.

## 2. In which of the following scenarios is data considered “Missing Not at Random” (MNAR)?

- a) Data is missing due to a random system failure.
- b) Patients with higher incomes are less likely to report their income.
- c) Missingness is related to other observed variables, like age or gender.
- d) Missing data occurs randomly across the dataset.

## 3. Which type of missing data is easiest to handle without introducing bias in analyses?

- a) Missing Not at Random (MNAR)
- b) Missing Completely at Random (MCAR)
- c) Missing at Random (MAR)
- d) None of the above

## 4. Which of the following is true about “Missing at Random” (MAR)?

- a) Missingness is explained by variables that are not observed.
- b) Missingness is random and unrelated to any variables in the dataset.
- c) Missingness is related to observed variables but not the missing data itself.
- d) Missingness is dependent on the value of the missing data.

## MISSING DATA QUIZ 2

**5. Which of the following techniques can help handle data that is “Missing at Random” (MAR)?**

- a) Removing all missing data rows.
- b) Imputing missing values based on observed variables.
- c) Ignoring the missing data and proceeding with the analysis.
- d) Using the mean of the missing variable to fill in the gaps.

**6. If data is Missing Completely at Random (MCAR), removing missing data points will not introduce bias.**  
(True / False)

**7. Data that is Missing Not at Random (MNAR) can be safely ignored without any effect on the analysis.**  
(True / False)

**8. Missing at Random (MAR) means the missingness is related to the unobserved, missing data itself.**  
(True / False)

# OUTLIER DETECTION

Method	Description
Tukey's Fences	<p>Often used to determine outliers in box plots.</p> <ol style="list-style-type: none"><li>1. Calculate the interquartile range, <math>IQR = Q_3 - Q_1</math> for a feature.</li><li>2. Classify all points that fall <math>1.5IQR</math> above <math>Q_3</math> or <math>1.5IQR</math> below <math>Q_1</math> as outliers.</li></ol>
z-scores	<ol style="list-style-type: none"><li>1. Calculate the z-score <math>z = \frac{\text{value} - \text{mean}}{\text{standard deviation}}</math> for each value.</li><li>2. Classify all points that have a z-score of <math> z  &gt; 3</math> as outliers.</li></ol>



# MEAN, VARIANCE, STANDARD DEVIATION

The formula for the **population mean** is:

$$\mu = \frac{\sum_{i=1}^N X_i}{N}$$

Where:

- $\mu$  = population mean,
- $X_i$  = each individual data point in the population,
- $N$  = the total number of data points in the population.

The formula for **population variance** is:

$$\sigma^2 = \frac{\sum_{i=1}^N (X_i - \mu)^2}{N}$$

Where:

- $\sigma^2$  = population variance,
- $X_i$  = each individual data point in the population,
- $\mu$  = population mean,
- $N$  = total number of data points in the population.

The formula for **population standard deviation** is:

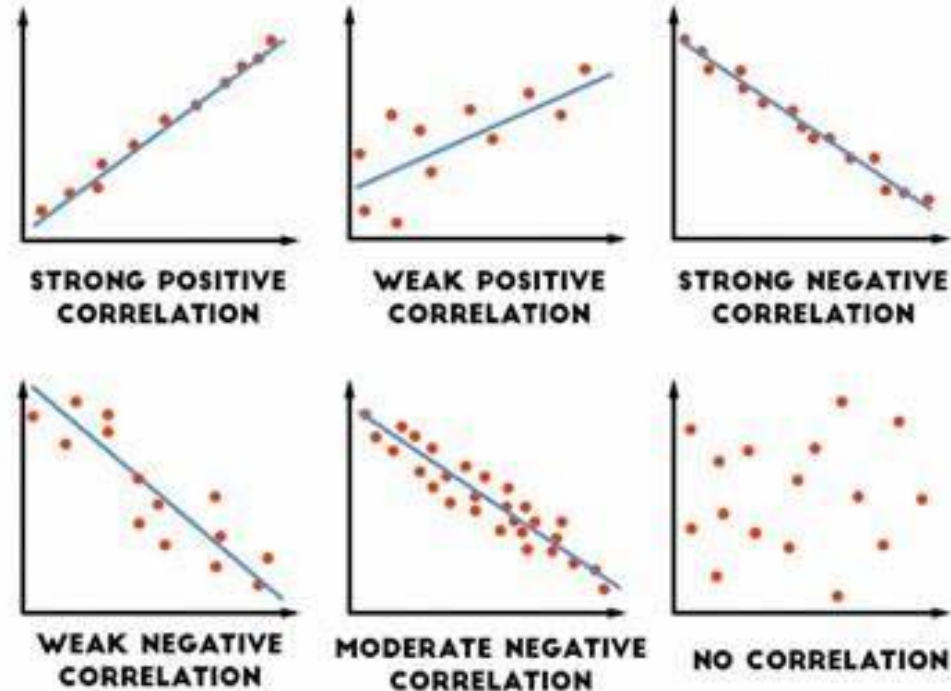
$$\sigma = \sqrt{\sigma^2} = \sqrt{\frac{\sum_{i=1}^N (X_i - \mu)^2}{N}}$$

Where:

- $\sigma$  = population standard deviation,
- $\sigma^2$  = population variance,
- $X_i$  = each individual data point in the population,
- $\mu$  = population mean,
- $N$  = total number of data points in the population.

# COVARIANCE VS CORRELATION

- Covariance:** Measures the direction of the relationship between two variables; units depend on the variables.
- Covariance:** Can range from negative to positive infinity.
- Correlation:** Measures both direction and strength; ranges from -1 to 1, and is unitless.
- Correlation:** Normalized, making it easier to compare across datasets.



# PYTHON CODE

UPDATE LINK:

<https://colab.research.google.com/github/rhodes-byu/cs180-winter25/blob/main/notebooks/05a-eda.ipynb>

EDA with Palmer Penguins:

<https://colab.research.google.com/github/rhodes-byu/cs180-winter25/blob/main/notebooks/05b-eda-penguins.ipynb>

