

# Data Prep & EDA Benchmark Assessment



## Review Results

Assessment PYTHON-1697674-BENCHMARK

Thinkific User ID: 195642843

Full Name: Onyebuchi, Augustine

Email Address: augustinesopuluonyebuchi@gmail.com

Respondent ID: 126877

Date Started: 10/24/2024 09:09:36 AM

Date Completed: 10/24/2024 09:23:21 AM

Attempt: 1

### Question Topic

Num	Question	
	Respondent's Answer	Correct Answer

#### Intro to Data Science

Earned 0 of 1 points (0%).

1. The first step of a data science project is to \_\_\_\_\_.

☒ (X) gather data  
☐ ( ) clean the data  
☐ ( ) model the data  
☐ ( ) scope the project  
☐ ( ) I don't know yet

☐ ( ) gather data  
☐ ( ) clean the data  
☐ ( ) model the data  
☒ (X) scope the project  
☐ ( ) I don't know yet

Explanation: Even though data science has data in the name, data science projects don't start with data. Instead, data scientists start by scoping a project and identifying end users, problems, etc. before diving into actual data.

Related Lecture: Data Science Workflow

#### Scoping a Project

Earned 0 of 1 points (0%).

2. Which of the following is an example of unsupervised learning?

☐ ( ) Predicting house prices over the next year  
☐ ( ) Flagging which customers are most likely to cancel their membership  
☐ ( ) Identifying the main themes mentioned in customer reviews  
☒ (X) Estimating how many customers will visit your website next week  
☐ ( ) I don't know yet

☐ ( ) Predicting house prices over the next year  
☐ ( ) Flagging which customers are most likely to cancel their membership  
☒ (X) Identifying the main themes mentioned in customer reviews  
☐ ( ) Estimating how many customers will visit your website next week  
☐ ( ) I don't know yet

Explanation: Supervised learning is all about using historical data to make future predictions. When identifying the main themes in customer reviews, no predictions are made. Finding themes is an unsupervised learning problem.

Related Lecture: Supervised vs Unsupervised Learning

#### Gathering Data

Earned 2 of 4 points (50%).

3. A Pandas DataFrame has which characteristics?

- |  |  |
|--|--|
| <input checked="" type="checkbox"/> (X) Index starting at 0, each column containing a single data type | <input checked="" type="checkbox"/> (X) Index starting at 0, each column containing a single data type |
| <input type="checkbox"/> ( ) Index starting at 1, each column containing a single data type            | <input type="checkbox"/> ( ) Index starting at 1, each column containing a single data type            |
| <input type="checkbox"/> ( ) Index starting at 0, each column containing multiple data types           | <input type="checkbox"/> ( ) Index starting at 0, each column containing multiple data types           |
| <input type="checkbox"/> ( ) Index starting at 1, each column containing multiple data types           | <input type="checkbox"/> ( ) Index starting at 1, each column containing multiple data types           |
| <input type="checkbox"/> ( ) I don't know yet  | <input type="checkbox"/> ( ) I don't know yet  |

Explanation: Counting in Python starts at 0 and each column in a DataFrame can only contain a single data type.

Related Lecture: The Pandas DataFrame

4. Which file formats can Python read in with one line of code?

- |  |  |
|--|--|
| <input type="checkbox"/> ( ) .csv files                  | <input type="checkbox"/> ( ) .csv files                  |
| <input type="checkbox"/> ( ) .xlsx files                 | <input type="checkbox"/> ( ) .xlsx files                 |
| <input type="checkbox"/> ( ) .json files                 | <input type="checkbox"/> ( ) .json files                 |
| <input checked="" type="checkbox"/> (X) All of the above | <input checked="" type="checkbox"/> (X) All of the above |
| <input type="checkbox"/> ( ) I don't know yet            | <input type="checkbox"/> ( ) I don't know yet            |

Explanation: Using `pd.read_csv()`, `pd.read_excel()` and `pd.read_json()`, you can read those file types into Python with just one line of code.

Related Lecture: Reading Flat Files

5. What will happen when the code above is executed?

- |   |  |
|---|--|
| <input type="checkbox"/> ( ) The "Instructors" tab will be read into Python                                     | <input type="checkbox"/> ( ) The "Instructors" tab will be read into Python                          |
| <input checked="" type="checkbox"/> (X) The "Instructors" tab will be read into Python and saved as a DataFrame | <input type="checkbox"/> ( ) The "Instructors" tab will be read into Python and saved as a DataFrame |
| <input type="checkbox"/> ( ) The "Courses" tab will be read into Python   | <input checked="" type="checkbox"/> (X) The "Courses" tab will be read into Python                   |
| <input type="checkbox"/> ( ) The "Courses" tab will be read into Python and saved as a DataFrame                | <input type="checkbox"/> ( ) The "Courses" tab will be read into Python and saved as a DataFrame     |
| <input type="checkbox"/> ( ) I don't know yet   | <input type="checkbox"/> ( ) I don't know yet  |

Explanation: Because Python is zero-indexed, reading the 0th tab would mean reading the first tab and reading the 1st tab would mean reading the second tab. This code only reads in the data, but does not save it as a DataFrame.

Related Lecture: Reading Excel Files

6. Which method will return the range of values within each column?

- |   |   |
|---|---|
| <input type="checkbox"/> ( ) .head()            | <input type="checkbox"/> ( ) .head()                |
| <input type="checkbox"/> ( ) .count()           | <input type="checkbox"/> ( ) .count()               |
| <input type="checkbox"/> ( ) .describe()        | <input checked="" type="checkbox"/> (X) .describe() |
| <input checked="" type="checkbox"/> (X) .info() | <input type="checkbox"/> ( ) .info()                |
| <input type="checkbox"/> ( ) I don't know yet   | <input type="checkbox"/> ( ) I don't know yet       |

Explanation: The describe method will return summary statistics including the min and max of each column.

Related Lecture: Quickly Exploring a DataFrame

## Cleaning Data

Earned 1 of 6 points (17%).

7. What will happen when the code above is executed?

- |   |  |
|---|--|
| <input checked="" type="checkbox"/> df.Income will change from an object to a numeric data type | <input type="checkbox"/> df.Income will change from an object to a numeric data type |
| <input type="checkbox"/> df.Income will be set equal to numeric                                 | <input type="checkbox"/> df.Income will be set equal to numeric                      |
| <input type="checkbox"/> Both 1 and 2   | <input type="checkbox"/> Both 1 and 2  |
| <input type="checkbox"/> Neither 1 or 2, you would get an error                                 | <input checked="" type="checkbox"/> Neither 1 or 2, you would get an error           |
| <input type="checkbox"/> I don't know yet   | <input type="checkbox"/> I don't know yet  |

Explanation: You would get an error because text with \$ values cannot be converted to numeric data types. You would need to use str.replace() to remove the \$ values before converting.

Related Lecture: Converting to Numeric

8. Which of the following is NOT a way that missing data is represented in Python?

- |   |  |
|---|--|
| <input type="checkbox"/> np.NaN           | <input type="checkbox"/> np.NaN            |
| <input type="checkbox"/> pd.NaN           | <input checked="" type="checkbox"/> pd.NaN |
| <input checked="" type="checkbox"/> pd.NA | <input type="checkbox"/> pd.NA             |
| <input type="checkbox"/> None             | <input type="checkbox"/> None              |
| <input type="checkbox"/> I don't know yet | <input type="checkbox"/> I don't know yet  |

Explanation: Various ways to represent missing values in Python are Numpy's NaN, Pandas' NA and base Python's None.

Related Lecture: Finding Missing Data

9. Which approach allows you to replace the "New York" value with "NY" in the State column?

- |  |  |
|--|--|
| <input checked="" type="checkbox"/> .loc[] | <input type="checkbox"/> .loc[]                      |
| <input type="checkbox"/> np.where()        | <input type="checkbox"/> np.where()                  |
| <input type="checkbox"/> .map()            | <input type="checkbox"/> .map()                      |
| <input type="checkbox"/> All of the above  | <input checked="" type="checkbox"/> All of the above |
| <input type="checkbox"/> I don't know yet  | <input type="checkbox"/> I don't know yet            |

Explanation: These are all ways that you can replace a value within a column, .loc[] to replace a specific value, np.where to replace a values based on a conditional and .map() to map a set of values to another set of values.

Related Lecture: Handling Inconsistent Text & Typos

10. What is considered an outlier?

- |   |  |
|---|--|
| <input type="checkbox"/> Data that is negative                                      | <input type="checkbox"/> Data that is negative   |
| <input checked="" type="checkbox"/> Data that is greater than three times the mean  | <input type="checkbox"/> Data that is greater than three times the mean                        |
| <input type="checkbox"/> Data that is more than 3 standard deviations from the mean | <input checked="" type="checkbox"/> Data that is more than 3 standard deviations from the mean |
| <input type="checkbox"/> All of the above   | <input type="checkbox"/> All of the above  |
| <input type="checkbox"/> I don't know yet   | <input type="checkbox"/> I don't know yet  |

Explanation: A rule of thumb in statistics is that outliers are data points that are more than ~3 standard deviations from the mean.

Related Lecture: Finding Outliers

11. How would you extract the day of the week from the Run Date column?

- |   |  |
|---|--|
| <input type="checkbox"/> dayofweek(run_times['Run Date'])               | <input type="checkbox"/> dayofweek(run_times['Run Date'])              |
| <input checked="" type="checkbox"/> dt.dayofweek(run_times['Run Date']) | <input type="checkbox"/> dt.dayofweek(run_times['Run Date'])           |
| <input type="checkbox"/> run_times['Run Date'].dayofweek                | <input type="checkbox"/> run_times['Run Date'].dayofweek               |
| <input type="checkbox"/> run_times['Run Date'].dt.dayofweek             | <input checked="" type="checkbox"/> run_times['Run Date'].dt.dayofweek |
| <input type="checkbox"/> I don't know yet                               | <input type="checkbox"/> I don't know yet                              |

Explanation: To extract the day of the week, you would use a datetime method, dt.dayofweek, and chain it on to the Run Date column.

Related Lecture: Creating DateTime Columns

12. If you applied `str[:6]` to the `run_notes` data above, which characters in the text would be returned?

- |   |   |
|---|---|
| <input checked="" type="radio"/> Characters 0-5 | <input checked="" type="radio"/> Characters 0-5 |
| <input type="radio"/> Characters 0-6            | <input type="radio"/> Characters 0-6            |
| <input type="radio"/> Characters 1-5            | <input type="radio"/> Characters 1-5            |
| <input type="radio"/> Characters 1-6            | <input type="radio"/> Characters 1-6            |
| <input type="radio"/> I don't know yet          | <input type="radio"/> I don't know yet          |

Explanation: Extracting a portion of a string value using `str` returns the first location (if blank, then 0) up to the last location, non-inclusive. So in this case, it would be characters 0 through 5.

Related Lecture: Creating Text Columns

## Exploratory Data Analysis

Earned 2 of 4 points (50%).

13. What does `df.groupby('col').head()` do?

- |  |   |
|--|---|
| <input checked="" type="radio"/> Groups the data by <code>col</code> and returns the first 5 rows of the results | <input type="radio"/> Groups the data by <code>col</code> and returns the first 5 rows of the results |
| <input type="radio"/> Groups the columns and returns the first 5 rows of the results                             | <input type="radio"/> Groups the columns and returns the first 5 rows of the results                  |
| <input type="radio"/> Returns the first 5 rows within each group   | <input checked="" type="radio"/> Returns the first 5 rows within each group                           |
| <input type="radio"/> Returns an error   | <input type="radio"/> Returns an error  |
| <input type="radio"/> I don't know yet   | <input type="radio"/> I don't know yet  |

Explanation: For each column value, this will return the first 5 rows for each value.

Related Lecture: Grouping

14. Why would you put parentheses around multiple chained methods?

- |   |   |
|---|---|
| <input checked="" type="radio"/> To be able to put each method on a separate line | <input checked="" type="radio"/> To be able to put each method on a separate line |
| <input type="radio"/> To allow the chained code to run                            | <input type="radio"/> To allow the chained code to run                            |
| <input type="radio"/> To create a function from the code                          | <input type="radio"/> To create a function from the code                          |
| <input type="radio"/> To comment the code   | <input type="radio"/> To comment the code   |
| <input type="radio"/> I don't know yet  | <input type="radio"/> I don't know yet  |

Explanation: When chaining multiple methods together, if you wrap the code in parentheses, you can place each method on a separate line, which makes reading the code easier.

Related Lecture: Grouping

15. What data is typically right skewed (picture above)?

- |   |   |
|---|---|
| <input type="radio"/> Student grades              | <input type="radio"/> Student grades              |
| <input type="radio"/> Household income            | <input checked="" type="radio"/> Household income |
| <input type="radio"/> Human lifespan              | <input type="radio"/> Human lifespan              |
| <input type="radio"/> Men and women's heights     | <input type="radio"/> Men and women's heights     |
| <input checked="" type="radio"/> I don't know yet | <input type="radio"/> I don't know yet            |

Explanation: Right skewed data is data where there are very few large values, and this example of household income is one of them.

Related Lecture: Distributions

16. Which one of these charts shows a correlation of 0?

- ☐ Chart 1
- ☒ Chart 2
- ☐ Chart 3
- ☐ None of them
- ☐ I don't know yet

- ☐ Chart 1
- ☒ Chart 2
- ☐ Chart 3
- ☐ None of them
- ☐ I don't know yet

Explanation: The middle chart shows a correlation of 0, meaning there is no relationship between Hours Talking to Friends and their Grade on the Test.

Related Lecture: Correlations

## Preparing for Modeling

Earned 2 of 4 points (50%).

17. What things need to be changed in the table above to put it into a model to predict house prices?

- ☐ Remove the Address column
- ☐ Change City and Color to numeric values
- ☒ Both 1 and 2
- ☐ Neither 1 or 2, the data is ready for modeling
- ☐ I don't know yet

- ☐ Remove the Address column
- ☐ Change City and Color to numeric values
- ☒ Both 1 and 2
- ☐ Neither 1 or 2, the data is ready for modeling
- ☐ I don't know yet

Explanation: All data must be non-null and numeric before inputting it into a model. Unique identifiers like address are not useful for modeling and should be removed before modeling.

Related Lecture: Data Prep for EDA vs Modeling

18. Which of the following could you use to vertically stack two DataFrames?

- ☒ .append()
- ☐ .join()
- ☐ .merge()
- ☐ .concat()
- ☐ I don't know yet

- ☐ .append()
- ☐ .join()
- ☐ .merge()
- ☒ .concat()
- ☐ I don't know yet

Explanation: `pd.concat()` allows you to both vertically stack and horizontally combine two DataFrames.

Related Lecture: Creating a Single Table

19. What can you do to a right skewed histogram to make it normally distributed?

- ☒ Apply a log transformation
- ☐ Apply a supervised learning algorithm
- ☐ Add a mirror image of the right skewed data
- ☐ Nothing, you cannot change skewed data
- ☐ I don't know yet

- ☒ Apply a log transformation
- ☐ Apply a supervised learning algorithm
- ☐ Add a mirror image of the right skewed data
- ☐ Nothing, you cannot change skewed data
- ☐ I don't know yet

Explanation: Log transformations turn skewed data into more normally-distributed data.

Related Lecture: Feature Transformations

20. Why can't you input the above DataFrame directly into a model to predict house prices?

- |   |  |
|---|--|
| <input type="radio"/> The columns are not all the same data type                    | <input type="radio"/> The columns are not all the same data type                               |
| <input type="radio"/> The columns all need to be on the same scale                  | <input type="radio"/> The columns all need to be on the same scale                             |
| <input type="radio"/> One house price seems much lower than the rest                | <input type="radio"/> One house price seems much lower than the rest                           |
| <input type="radio"/> The zip code column is saying that 60202 is better than 60201 | <input checked="" type="radio"/> The zip code column is saying that 60202 is better than 60201 |
| <input checked="" type="radio"/> I don't know yet                                   | <input type="radio"/> I don't know yet   |

Explanation: Sometimes data seems ready to be input into a model because it's numeric, but in cases like this, a higher zip code doesn't necessarily mean a higher / lower house price. In situations like this, you can either create a dummy variable or identify a proxy variable.

Related Lecture: Proxy Variables

Time Used: 00:13:45

Final Score: 35%