

V2_full_osemn

December 25, 2025

1 Activity: Full OSEMN

1.1 Introduction

In this assignment, you will work on a data analysis project. This project will let you practice the skills you have learned in this course and write real code in Python.

You will perform the following steps of the OSEMN framework:

- Section 1.2 - Section 1.3 - Section 1.5

```
[1]: # We'll import the libraries you'll likely use for this activity
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Data
df = pd.read_csv('transactions-pet_store.csv')
df_orig = df.copy()
```

1.2 Scrub

You will scrub the data. It's important that you follow the directions as stated. Doing more or less than what is asked might lead to not getting full points for the question.

If while you're working on the scrubbing phase you need to reset the DataFrame, you can restart the kernel (in the toolbar: "Kernel" > "Restart").

Question 1 Remove all rows that have are missing either the `Product_Name` or the `Product_Category`. Assign the cleaned DataFrame to the variable `df` (overwriting the original DataFrame.).

```
[2]: mask_1 = (df['Product_Name'].isnull() | df['Product_Category'].isnull())
df = df[~mask_1]
```

```
[3]: # Question 1 Grading Checks

assert df.shape[0] <= 2874, 'Did you remove all the rows with missing values,
↳for the columns Product_Name & Product_Category?'
assert df.shape[0] >= 2700, 'Did you remove too many the rows with missing
↳values?'
assert len(df.columns) == 10, 'Make sure you do not drop any columns.'
```

Question 2 Find any clearly “incorrect” values in the Price column and “clean” the DataFrame to address those values.

Ensure you make the changes to the DataFrame assigned to the variable df.

```
[4]: mask_2 = df['Price'] < 0
mask_3 = (df['Price'] < df.Price.quantile(0.999)) & (df['Price'] > df.Price.
↳quantile(0.001))
df.loc[mask_2, 'Price'] = df.loc[mask_2, 'Price'] * -1
df = df[mask_3]
```

```
[5]: # Question 2 Grading Checks

assert (df.Price < df.Price.quantile(0.0001)).sum() == 0, 'Check for very small
↳values'
assert (df.Price > df.Price.quantile(0.999)).sum() == 0, 'Check for very large
↳values'
```

Question 3 After you’ve done the cleaning above, remove any column that has more than 500 missing values.

Ensure you make the changes to the DataFrame assigned to the variable df.

```
[6]: df.isnull().sum()
df = df.drop('Size',axis=1)
```

```
[7]: # Question 3 Grading Checks

assert len(df.columns) < 10, 'You should have dropped 1 or more columns (with
↳more than 500 missing values)'
```

Question 4 Address the other missing values. You can replace the values or remove them, but whatever method you decide to clean the DataFrame, you should no longer have any missing values.

Ensure you make the changes to the DataFrame assigned to the variable df.

```
[8]: df = df.dropna()
```

```
[9]: # Question 4 Grading Checks

assert df.Customer_ID.isna().sum() == 0, 'Did you address all the missing_
      ↳values?'
```

1.3 Explore

You will explore the data. It's important that you follow the directions as stated. Doing more or less than what is asked might lead to not getting full points for the question.

You may use either exploratory statistics or exploratory visualizations to help answer these questions.

Note that the DataFrame loaded for this section (in the below cell) is different from the data you used in the Section 1.2 section.

If while you're working on the scrubbing phase you need to reset the DataFrame, you can restart the kernel (in the toolbar: "Kernel" > "Restart").

```
[10]: df = pd.read_csv('transactions-pet_store-clean.csv')
```

Question 5 Create a Subtotal column by multiplying the Price and Quantity values. This represents how much was spent for a given transaction (row).

```
[11]: df['Subtotal'] = df['Price'] * df['Quantity']
```

```
[12]: # Question 5 Grading Checks

assert 'Subtotal' in df.columns, ''
```

Question 6 Determine most common category (Product_Category) purchases (number of total items) for both Product_Line categories. Assign the (string) name of these categories to their respective variables common_category_cat & common_category_dog.

```
[13]: categories = df.groupby(['Product_Line', 'Product_Category'])['Quantity'].sum()
common_category_cat = categories['cat'].idxmax()
common_category_dog = categories['dog'].idxmax()
```

```
[14]: # Question 6 Grading Checks

assert isinstance(common_category_dog, str), 'Ensure you assign the name of the_
      ↳category (string) to the variable common_category_dog'
assert isinstance(common_category_cat, str), 'Ensure you assign the name of the_
      ↳category (string) to the variable common_category_cat'
```

Question 7 Determine which categories (Product_Category), by Product_Line have the *median* highest Price. Assign the (string) name of these categories to their respective variables priciest_category_cat & priciest_category_dog.

```
[15]: categories = df.groupby(['Product_Line', 'Product_Category'])['Price'].median()
      priciest_category_cat = categories['cat'].idxmax()
      priciest_category_dog = categories['dog'].idxmax()
```

```
[16]: # Question 7 Grading Checks

      assert isinstance(priciest_category_dog, str), 'Ensure you assign the name of_
      ↳the category (string) to the variable priciest_category_dog'
      assert isinstance(priciest_category_cat, str), 'Ensure you assign the name of_
      ↳the category (string) to the variable priciest_category_cat'
```

1.4 Modeling

This is the point of the framework where we'd work on modeling with our data. However, in this activity, we're going to move straight to interpreting.

1.5 Interpret

You will interpret the data based on what you found so far. It's important that you follow the directions as stated. Doing more or less than what is asked might lead to not getting full points for the question.

Note that the DataFrame loaded for this section (in the below cell) is the same as the data you used in the Section 1.3 section.

If while you're working on the scrubbing phase you need to reset the DataFrame, you can restart the kernel (in the toolbar: "Kernel" > "Restart").

Question 8 You want to emphasize to your stakeholders that the total number of product categories sold differ between the two Product_Line categories ('cat' & 'dog').

Create a **horizontal bar plot** that has Product_Category on the y-axis and the total number of that category sold (using the Quantity) by each Product_Line category. Also **change the axis labels** to something meaningful and add a title.

You will likely want to use Seaborn. Make sure you set the result to the variable ax like the following:

```
ax = # code to create a bar plot
```

```
[17]: ax = sns.barplot(data=df, x='Quantity',
                     y='Product_Category',
```

```

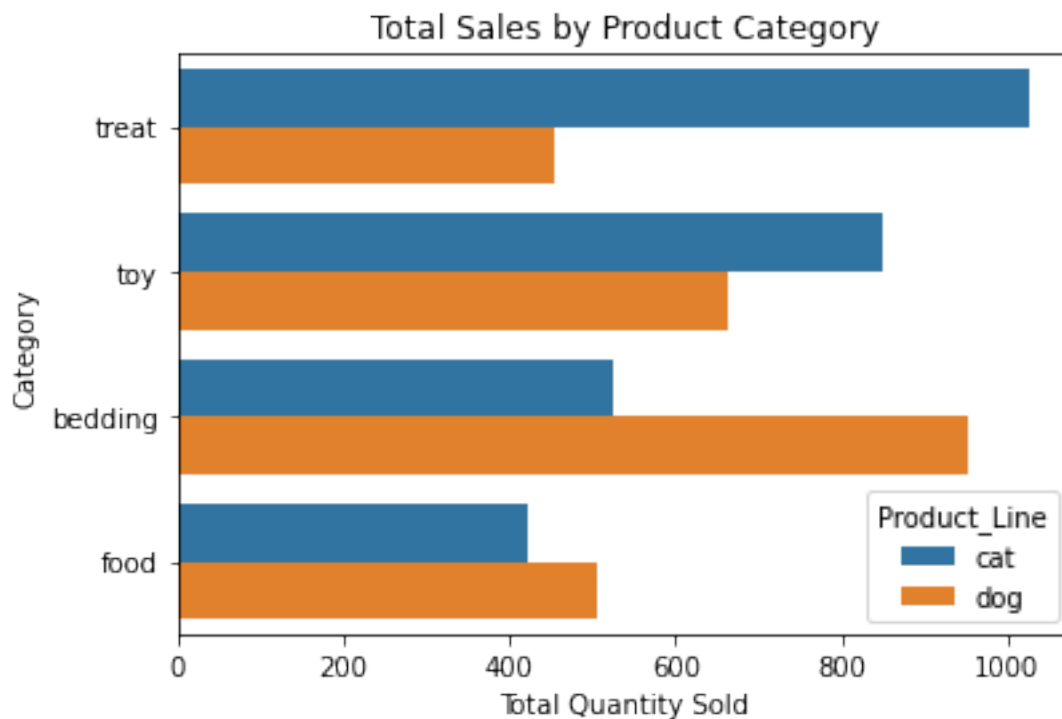
        hue='Product_Line',
        estimator='sum',
        errorbar=None)
ax.set(
    title='Total Sales by Product Category',
    xlabel='Total Quantity Sold',
    ylabel='Category'
)

```

```

[17]: [Text(0.5, 1.0, 'Total Sales by Product Category'),
      Text(0.5, 0, 'Total Quantity Sold'),
      Text(0, 0.5, 'Category')]

```



```

[18]: # Question 8 Grading Checks

assert isinstance(ax, plt.Axes), 'Did you assign the plot result to the_
↪variable ax?'

```

Question 9 Based on the plot from Section 1.5, what would you conclude for your stakeholders about what products they should sell? What would be the considerations and/or caveats you'd communicate to your stakeholders?

Write at least a couple sentences of your thoughts in a string assigned to the variable `answer_to_9`.

The cell below should look something like this:

```
answer_to_9 = '''
I think that based on the visualization that ****.
Therefore I would communicate with the stakeholders that ****
'''
```

```
[19]: # Your code here
answer_to_9 = '''
I think that based on the Total Sales by Product Category graph, there is a
    ↳ large discrepancy between dog and cat owner purchasing behaviors.
Dog owners seem to spend somewhat more on food, and overwhelmingly more on
    ↳ bedding, while cat owners spend somewhat more on toys, and overwhelmingly
    ↳ more on treats.
With this in mind, I would communicate to stakeholders that there is a
    ↳ noticeable difference in purchasing behavior, and that they should study
    ↳ (surveys, polling etc.) why there is a difference in purchasing behavior,
    ↳ and let that drive future business strategies.'''
print(len(answer_to_9))
```

572

```
[20]: # Question 9 Grading Checks

assert isinstance(answer_to_9, str), 'Make sure you create a string for your
    ↳ answer.'
```

Question 10 The plot you created for Section 1.5 is good but could be modified to emphasize which products are important for the business.

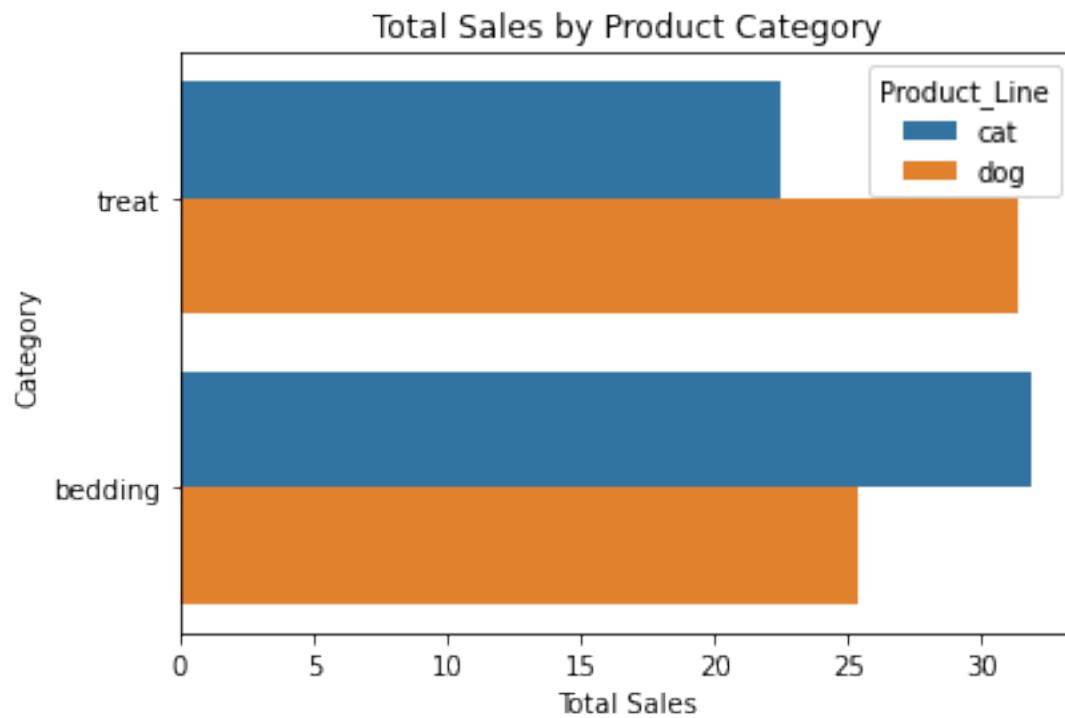
Create an explanatory visualization that emphasizes the insight you about the product category. This would be a visualization you'd share with the business stakeholders.

Make sure you set the result to the variable `ax` like the following:

```
ax = # code to create explanatory visualization
```

```
[25]: treat_bed = (df['Product_Category'] == 'bedding') | (df['Product_Category'] ==
    ↳ 'treat')
ax = sns.barplot(data=df[treat_bed], x='Price',
                y='Product_Category',
                hue='Product_Line',
                estimator='mean',
                errorbar=None)
ax.set(
    title='Total Sales by Product Category',
    xlabel='Total Sales',
    ylabel='Category'
)
```

```
[25]: [Text(0.5, 1.0, 'Total Sales by Product Category'),  
      Text(0.5, 0, 'Total Sales'),  
      Text(0, 0.5, 'Category')]
```



```
[26]: # Question 10 Grading Checks  
  
assert isinstance(ax, plt.Axes), 'Did you assign the plot result to the_  
→variable ax?'
```