

Sales Analysis (Lecture Designed by Syed Umaid Ahmed)

```
import os
```

```
import pandas as pd
```

```
#### Merge data from each month into one CSV
```

```
path = "./Sales_Data"
```

```
files = [file for file in os.listdir(path) if not file.startswith('.')]
```

```
all_months_data = pd.DataFrame()
```

```
for file in files:
```

```
    current_data = pd.read_csv(path+"/"+file)
```

```
    all_months_data = pd.concat([all_months_data, current_data])
```

```
all_months_data.to_csv("all_data_copy.csv", index=False)
```

```
##### Read in updated dataframe
```

```
all_data = pd.read_csv("all_data.csv")
```

```
all_data.head()
```

```
# ### Clean up the data!
```

```
# The first step in this is figuring out what we need to clean. I have found in practice, that you find things you need to clean as you perform operations and get errors. Based on the error, you decide how you should go about cleaning the data
```

```
##### Drop rows of NAN
```

Find NAN

```
nan_df = all_data[all_data.isna().any(axis=1)]  
display(nan_df.head())  
all_data = all_data.dropna(how='all')  
all_data.head()
```

Get rid of text in order date column

```
all_data = all_data[all_data['Order Date'].str[0:2]!='Or']
```

Make columns correct type

```
all_data['Quantity Ordered'] = pd.to_numeric(all_data['Quantity Ordered'])  
all_data['Price Each'] = pd.to_numeric(all_data['Price Each'])
```

Augment data with additional columns

Add month column

```
all_data['Month'] = all_data['Order Date'].str[0:2]  
all_data['Month'] = all_data['Month'].astype('int32')  
all_data.head()
```

Add month column (alternative method)

```
all_data['Month 2'] = pd.to_datetime(all_data['Order Date']).dt.month
all_data.head()
```

Add city column

```
def get_city(address):
    return address.split(",")[1].strip(" ")
```

```
def get_state(address):
    return address.split(",")[2].split(" ")[1]
```

```
all_data['City'] = all_data['Purchase Address'].apply(lambda x: f"{get_city(x)} {get_state(x)}")
all_data.head()
```

Data Exploration!

Question 1: What was the best month for sales? How much was earned that month?

```
all_data['Sales'] = all_data['Quantity Ordered'].astype('int') * all_data['Price Each'].astype('float')
all_data.groupby(['Month']).sum()
```

##PLOT

```
import matplotlib.pyplot as plt
months = range(1,13)
print(months)
```

```
plt.bar(months,all_data.groupby(['Month']).sum()['Sales'])
plt.xticks(months)
plt.ylabel('Sales in USD ($)')
plt.xlabel('Month number')
plt.show()
```

Question 2: What city sold the most product?

```
all_data.groupby(['City']).sum()
```

```
import matplotlib.pyplot as plt
```

```
keys = [city for city, df in all_data.groupby(['City'])]
```

```
plt.bar(keys,all_data.groupby(['City']).sum()['Sales'])
```

```
plt.ylabel('Sales in USD ($)')
```

```
plt.xlabel('Month number')
```

```
plt.xticks(keys, rotation='vertical', size=8)
```

```
plt.show()
```

Question 3: What time should we display advertisements to maximize likelihood of customer's buying product?

Add hour column

```
all_data['Hour'] = pd.to_datetime(all_data['Order Date']).dt.hour
```

```
all_data['Minute'] = pd.to_datetime(all_data['Order Date']).dt.minute
```

```
all_data['Count'] = 1
```

```
all_data.head()
```

```
keys = [pair for pair, df in all_data.groupby(['Hour'])]
plt.plot(keys, all_data.groupby(['Hour']).count()['Count'])
plt.xticks(keys)
plt.grid()
plt.show()
```

My recommendation is slightly before 11am or 7pm

Question 4: What products are most often sold together?

<https://stackoverflow.com/questions/43348194/pandas-select-rows-if-id-appear-several-time>

```
df = all_data[all_data['Order ID'].duplicated(keep=False)]
```

Referenced: <https://stackoverflow.com/questions/27298178/concatenate-strings-from-several-rows-using-pandas-groupby>

```
df['Grouped'] = df.groupby('Order ID')['Product'].transform(lambda x: ', '.join(x))
df2 = df[['Order ID', 'Grouped']].drop_duplicates()
```

Referenced: <https://stackoverflow.com/questions/52195887/counting-unique-pairs-of-numbers-into-a-python-dictionary>

```
from itertools import combinations
```

```
from collections import Counter
```

```
count = Counter()
```

```
for row in df2['Grouped']:
    row_list = row.split(',')
    count.update(Counter(combinations(row_list, 2)))
```

```
for key,value in count.most_common(10):
    print(key, value)
```

What product sold the most? Why do you think it sold the most?

```
product_group = all_data.groupby('Product')
quantity_ordered = product_group.sum()['Quantity Ordered']
```

```
keys = [pair for pair, df in product_group]
plt.bar(keys, quantity_ordered)
plt.xticks(keys, rotation='vertical', size=8)
plt.show()
```

Referenced: <https://stackoverflow.com/questions/14762181/adding-a-y-axis-label-to-secondary-y-axis-in-matplotlib>

```
prices = all_data.groupby('Product').mean()['Price Each']
fig, ax1 = plt.subplots()
ax2 = ax1.twinx()
ax1.bar(keys, quantity_ordered, color='g')
ax2.plot(keys, prices, color='b')
ax1.set_xlabel('Product Name')
ax1.set_ylabel('Quantity Ordered', color='g')
ax2.set_ylabel('Price ($)', color='b')
ax1.set_xticklabels(keys, rotation='vertical', size=8)
fig.show()
```