

Data Analysis, RIT Business Competition

```
import folium
import requests
import json
from geopy.geocoders import Nominatim
import pandas as pd
import numpy as np
import seaborn as sns
import typing
from datetime import datetime
from scipy.stats import pearsonr, spearmanr
from matplotlib import pyplot as plt
from collections import Counter

import warnings
from IPython.display import HTML
warnings.simplefilter(action='ignore', category=FutureWarning)
```

Data cleaning

Shops data

```
df_info = pd.read_csv('pharm_info.csv')

df_info.head()

storenum storename telephone address1 address2 city state zipcode zipcode4 longitude latitude website

0 28500 CVS-1142 585-671-5665 5900 EMPRE NaN WEBSTER NY 14580 19340 -77.591365 43.191229 HTTP://WWW.CVS.COM

1 28490 CVS-2067 585-321-2581 2580 EAST HENRIETTA ROAD NaN ROCHESTER NY 14623 45260 -77.606816 43.069285 HTTP://WWW.CVS.COM

2 28520 CVS-2217 585-723-3051 878 LONG ROAD NaN GREECE NY 14612 30490 -77.695905 43.238923 HTTP://WWW.CVS.COM

3 46030 CVS-2218 585-581-5101 3750 MOUNT BOULEVARD NaN ROCHESTER NY 14616 34360 -77.658331 43.227241 HTTP://WWW.CVS.COM

4 60570 CVS-2501 585-426-2091 2709 CHU AVENUE NaN ROCHESTER NY 14624 41230 -77.709106 43.122328 HTTP://WWW.CVS.COM

df_info.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 69 entries, 0 to 68
Data columns (total 13 columns):
# Column Non-Null Count Dtype
---
0 storenum 69 non-null object
1 storename 69 non-null object
2 telephone 69 non-null object
3 address1 69 non-null object
4 address2 2 non-null object
5 city 69 non-null object
6 state 69 non-null object
7 zipcode 69 non-null object
8 zipcode4 69 non-null float64
9 longitude 69 non-null float64
10 latitude 69 non-null float64
11 website 69 non-null object
12 orgname 69 non-null object
dtypes: float64(4), object(9)
memory usage: 7.4+ MB

df_info.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 69 entries, 0 to 68
Data columns (total 13 columns):
# Column Non-Null Count Dtype
---
0 storenum 69 non-null object
1 storename 69 non-null object
2 city 69 non-null object
3 zipcode 69 non-null object
4 longitude 69 non-null float64
5 latitude 69 non-null float64
6 orgname 69 non-null object
dtypes: float64(3), object(4)
memory usage: 4.3+ MB
```

Sales data

```
df_sales = pd.read_csv('pharm_sales.csv')

df_sales.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 390746 entries, 0 to 390745
Data columns (total 10 columns):
# Column Non-Null Count Dtype
---
0 invnum 390746 non-null object
1 date 390746 non-null datetime64[ns]
2 storenum 390746 non-null int64
3 vendornum 390746 non-null float64
4 itemnum 390746 non-null int64
5 itemcategory 390746 non-null object
6 productcost 390746 non-null float64
7 product_retail_price 390746 non-null float64
8 units_sold 390746 non-null int64
9 revenue_original 390746 non-null float64
dtypes: datetime64[ns](1), float64(4), int64(3), object(2)
memory usage: 75.6+ MB
```

Merging the two dataframes

```
# Merging both of them using the storenumber field as that is unique for stores.
df_merged = pd.merge(df_sales, df_info, on='storenum', how='outer')

df_merged.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 397418 entries, 0 to 390745
Data columns (total 16 columns):
# Column Non-Null Count Dtype
---
0 invnum 390746 non-null object
1 date 390746 non-null datetime64[ns]
2 storenum 390746 non-null int64
3 vendornum 390746 non-null float64
4 itemnum 390746 non-null int64
5 itemcategory 390746 non-null object
6 productcost 390746 non-null float64
7 product_retail_price 390746 non-null float64
8 units_sold 390746 non-null int64
9 revenue_original 390746 non-null float64
10 storename 390746 non-null object
11 city 390746 non-null object
12 longitude 390746 non-null float64
13 latitude 390746 non-null float64
14 orgname 390746 non-null object
dtypes: datetime64[ns](1), float64(6), int64(3), object(6)
memory usage: 128.5+ MB
```

Renaming columns

```
df_merged = df_merged.rename(columns = {'sales_num':'revenue_original'})

df_merged.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 397418 entries, 0 to 390745
Data columns (total 16 columns):
# Column Non-Null Count Dtype
---
0 invnum 390746 non-null object
1 date 390746 non-null datetime64[ns]
2 storenum 390746 non-null int64
3 vendornum 390746 non-null float64
4 itemnum 390746 non-null int64
5 itemcategory 390746 non-null object
6 productcost 390746 non-null float64
7 product_retail_price 390746 non-null float64
8 units_sold 390746 non-null int64
9 revenue_original 390746 non-null float64
10 storename 390746 non-null object
11 city 390746 non-null object
12 longitude 390746 non-null float64
13 latitude 390746 non-null float64
14 orgname 390746 non-null object
dtypes: datetime64[ns](1), float64(6), int64(3), object(6)
memory usage: 128.5+ MB
```

Adding new fields and fixing types

```
# Original sales figures are off, this is probably a safer bet, might ignore taxes
df_merged['revenue_original'] = df_merged['product_retail_price'] * df_merged['units_sold']

# Total cost per sale.
df_merged['inventory_cost'] = df_merged['productcost'] * df_merged['units_sold']

# Inventory costs are 68% of a businesses cost
#https://blog.shelvingdesignsystems.com/what-does-it-really-cost-to-start-a-pharmacy
df_merged['total_cost'] = df_merged['inventory_cost'] * 0.68

# Cost of non inventory overhead
df_merged['overhead_cost'] = df_merged['total_cost'] - df_merged['inventory_cost']

# Profit per sale
df_merged['profit'] = df_merged['revenue_original'] - df_merged['total_cost']

# Margin per unit
df_merged['margin'] = df_merged['product_retail_price'] - df_merged['productcost']

# Making store number field a string
df_merged['storenum'] = df_merged['storenum'].astype(str)

df_merged.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 390744 entries, 0 to 390743
Data columns (total 21 columns):
# Column Non-Null Count Dtype
---
0 date 390744 non-null datetime64[ns]
1 storenum 390744 non-null object
2 vendornum 390744 non-null float64
3 itemnum 390744 non-null int64
4 itemcategory 390744 non-null object
5 productcost 390744 non-null float64
6 product_retail_price 390744 non-null float64
7 units_sold 390744 non-null int64
8 revenue_original 390744 non-null float64
9 storename 390744 non-null object
10 city 390744 non-null object
11 longitude 390744 non-null float64
12 latitude 390744 non-null float64
13 orgname 390744 non-null object
14 inventory_cost 390744 non-null float64
15 total_cost 390744 non-null float64
16 overhead_cost 390744 non-null float64
17 profit 390744 non-null float64
18 margin 390744 non-null float64
dtypes: datetime64[ns](1), float64(11), int64(2), object(7)
memory usage: 159.7+ MB
```

The plot below will help me see whether I should worry too much about dropping shops or not. If there is a huge range in variance between sales, then I should.

```
df_merged.groupby(['date'])['units_sold'].sum().plot(color='green')

<AxesSubplot: xlabel='date'>
```



Sale patterns are random, IE not that much yearly variation, thus I will drop pharmacies which made no sales in 2021, which I will assume is an indication for being closed.

Dropping shops that probably closed.

Do we know the stores closed beyond a shadow of a doubt? No.

However doing any analysis would be moot if I don't have the data normalized in some way shape or form. It seems to me the best way to do that would be to normalize over a specific period of time. And to further normalize the data, I'll drop stores that did not have any sales past the year of 2021, then analyze with only those stores over a reasonable timeframe, where the loss of data would not be too high.

```
# Making a new dataframe with only sales past 2021
df_merged_open = df_merged.loc[df_merged['date'] >= '2018-01-01']
df_merged_open = df_merged_open.reset_index(drop=True)

# Finding the shops that had sales in 2021, and keeping the values in a variable to assign to another dataframe
shops_still_open = df_merged_open['storenum'].unique()
shops_still_open
```

```
array(['13695', '15248', '15503', '15557', '15338', '1126', '12673', '4796',
       '15533', '4997', '14641', '12561', '14140', '15051', '12527', '4594',
       '15537', '1263', '12627', '14361', '12626', '1495', '15131', '12527',
       '15507', '15127', '15128', '14622', '12696', '14601', '4802', '45967',
       '15375', '4599', '14603', '14377', '14801', '14603', '15539', '12852',
       '12850', '13821', '1404', '1479', '14800', '1564', '4799', '2635',
       '15970', '5891', '16557', '16064'], dtype=object)
```

```
# Number of stores still open.
df_merged_open['storenum'].nunique()

52
```

```
52/69

0.7536231884057971
```

I lost 25% of the stores by only keeping the stores that were open in 2021, I need to see how much of the total dataset I will lose.

```
df_merged_final = df_merged[(df_merged['storenum'].isin(shops_still_open)) & (df_merged['date'] >= '2018-01-01')]
df_merged_final = df_merged_final.reset_index(drop=True)

df_info_final = df_info.loc[df_info['storenum'].isin(shops_still_open)]
df_info_final = df_info_final.reset_index(drop=True)
df_info_final.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 52 entries, 0 to 51
Data columns (total 7 columns):
# Column Non-Null Count Dtype
---
0 storenum 52 non-null float64
1 storename 52 non-null object
2 city 52 non-null object
3 zipcode 52 non-null object
4 longitude 52 non-null float64
5 latitude 52 non-null float64
6 orgname 52 non-null object
dtypes: float64(3), object(4)
memory usage: 3.0+ KB
```

```
df_merged_final.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 371147 entries, 0 to 371146
Data columns (total 21 columns):
# Column Non-Null Count Dtype
---
0 invnum 371147 non-null object
1 date 371147 non-null datetime64[ns]
2 storenum 371147 non-null object
3 vendornum 371147 non-null float64
4 itemnum 371147 non-null int64
5 itemcategory 371147 non-null object
6 productcost 371147 non-null float64
7 product_retail_price 371147 non-null float64
8 units_sold 371147 non-null int64
9 revenue_original 371147 non-null float64
10 storename 371147 non-null object
11 city 371147 non-null object
12 longitude 371147 non-null float64
13 latitude 371147 non-null float64
14 orgname 371147 non-null object
15 inventory_cost 371147 non-null float64
16 total_cost 371147 non-null float64
17 overhead_cost 371147 non-null float64
18 profit 371147 non-null float64
19 margin 371147 non-null float64
dtypes: datetime64[ns](1), float64(11), int64(2), object(7)
memory usage: 59.5+ MB
```

```
# Saving as csv
df_merged_final.to_csv('pharm_sales_final.csv')
df_info_final.to_csv('pharm_info_final.csv')
```

EDA

Total sales over time

```
df_merged_final.groupby(['date'])['units_sold'].sum().plot(color='green', title='Total units sold, time')

<AxesSubplot: title='center': 'Total units sold, time', xlabel='date'>
```



Zipcode

```
total_profit_zipcode_groups = df_merged_final.groupby(['zipcode']).sum().sort_values(['profit'], ascending=True)
total_profit_zipcode_groups.sort_values().head(10)
```

```
zipcode
14615 1307.171765
14667 1613.135294
14609 2322.936471
14625 2439.005888
14650 3041.305882
14611 3220.129412
14616 3723.339359
14607 6343.448824
14559 7186.052353
14526 11422.208224
Name: profit, dtype: float64
```

```
total_profit_zipcode_groups.plot(kind='bar',
                                ylabel='Total ($)',
                                title='Profit per zipcode')

<AxesSubplot: title='center': 'Profit per zipcode', xlabel='zipcode', ylabel='Total ($)'>
```



```
number_stores_zipcode = df_info_final.groupby(['zipcode']).count()['storenum']
number_stores_zipcode.head(10)
```

```
zipcode
14650 1
14667 1
14514 2
14526 2
14559 1
14580 5
14596 3
14607 1
14609 2
14611 1
14612 1
Name: storenum, dtype: int64
```

```
profit_per_zipcode = total_profit_zipcode_groups/number_stores_zipcode
profit_per_zipcode.sort_values().head()
```

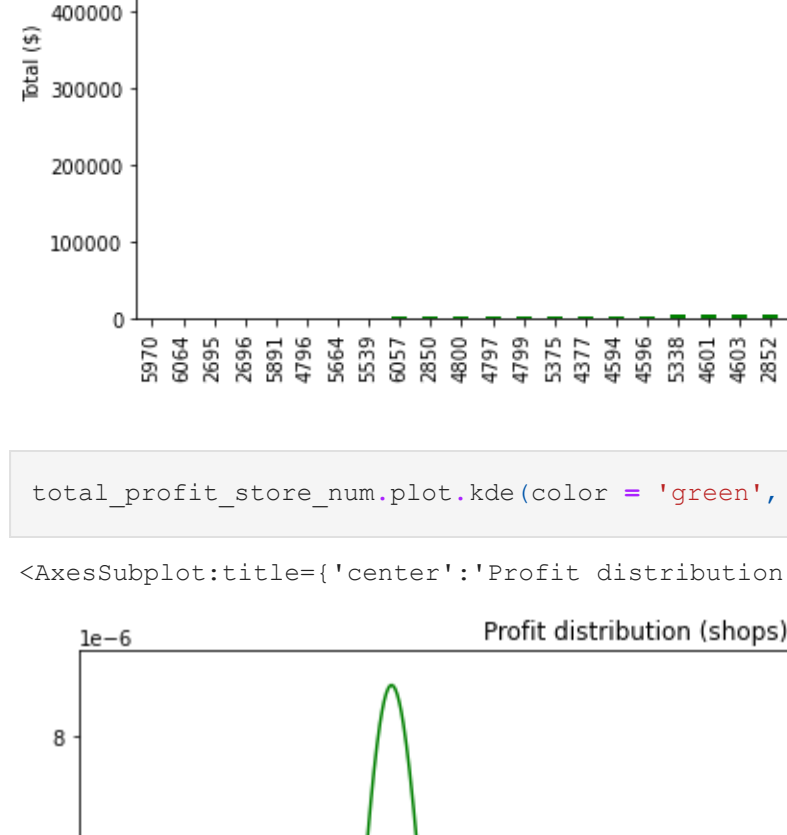
```
zipcode
14609 161.468235
14615 1307.171765
14667 1613.135294
14625 2439.005888
14612 2862.668529
dtype: float64
```

```
profit_per_zipcode.max()
profit_per_zipcode.mean()

12.17106385903507
```

```
profit_per_zipcode.sort_values().plot(kind='bar',
                                       color='green',
                                       ylabel='Total ($)',
                                       title='Average Profit/ Zipcode')

<AxesSubplot: title='center': 'Average Profit/ Zipcode', xlabel='zipcode', ylabel='Total ($)'>
```



Choropleth map

```
# Link below has all the zipcode boundaries in NY state, I'll only keep the ones I need
url="https://github.com/OpenDataDB/State-zip-code-GeoJSON/raw/master/ny_new_york_zip_codes_geo_min.json"
g = requests.get(url).json()

# Making new geojson file
zip_list = df_merged_final['zipcode'].unique() #list of zips in your dataframe

inlist = []
for i in zip_list:
    if i in g['features']:
        inlist.append(i)

new_zip_json = {}
new_zip_json['type'] = 'FeatureCollection'
new_zip_json['features'] = inlist
```

```
#Running the series above into a dataframe
df_zip = pd.DataFrame({'zipcode': [i for i in profit_per_zipcode.index],
                      'average_profit': [x for x in profit_per_zipcode.values]})

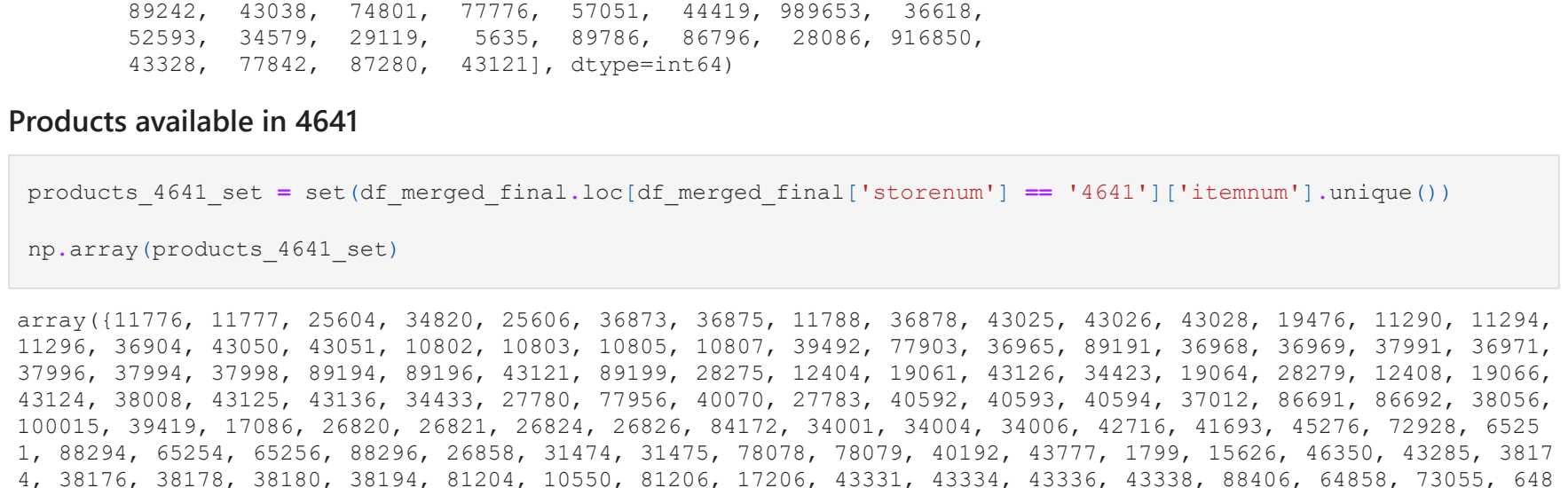
# Casting type as string so that parser can parse connect it to geojson file
df_zip['ECTA5CE10'] = df_zip['zipcode']

m = folium.Map(location=[43.191229,-77.5], zoom_start=11, tile='stamenwatercolor')

style_function = {"font-size": 15px; font-weight: bold"

geo_data = new_zip_json
name = "choropleth",
columns = ['ECTA5CE10', 'average_profit'],
key_on = "feature.properties.ECTA5CE10",
fill_color = "color",
fill_opacity = 2,
line_opacity = 2,
legend_name = "Average profit",
labels = True
).add_to(m)

# Add labels indicating zipcode
style_function = {"font-size": 15px; font-weight: bold"
folium.features.GeoJsonTooltip(['ECTA5CE10'], style=style_function, labels=False)
folium.LayerControl().add_to(m)
m
```



Heatmap

```
df_zip = df_zip.rename(columns = {'ECTA5CE10': 'zipcode', 'average_profit': 'AverageProfit'})
df_zip.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 22 entries, 0 to 21
Data columns (total 2 columns):
# Column Non-Null Count Dtype
---
0 zipcode 22 non-null object
1 AverageProfit 22 non-null float64
dtypes: float64(1), object(1)
memory usage: 480.0 bytes
```

```
df_demographics = pd.read_csv('zip_demographics.csv')
df_demographics['zipcode'] = df_demographics['zipcode'].astype(str)
df_demographics['populationdensity'] = df_demographics['CurrentPopulation']/df_demographics['Area']
df_demographics.head()
```

```
zipcode CurrentPopulation RacialMajority UnemploymentRate MedianHouseholdIncome AverageAdjustedGrossIncome SchoolTestPerformance
0 14580 50587 White 5.3 75618.6 815100 Above Average
1 14623 27173 White 5.5 54283.0 491300 Above Average
2 14612 28514 White 7.3 62148.0 615100 Average
3 14616 34535 White 10.8 54538.0 462700 Above Average
4 14624 36296 White 6.5 69312.0 623600 Average
```

```
#Label encoding
encode_dict = {'SchoolTestPerformance': {'Poor': 1, 'Below Average': 2, 'Average': 3, 'Above Average': 4, 'Excellent': 5}, 'RacialMajority': {'White': 0, 'Black': 1, 'Hispanic': 2, 'Asian': 3, 'Other': 4}}
df_demographics = df_demographics.replace(encode_dict)
```

```
df_demographics.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 23 entries, 0 to 24
Data columns (total 12 columns):
# Column Non-Null Count Dtype
---
0 zipcode 23 non-null object
1 CurrentPopulation 23 non-null int64
2 RacialMajority 23 non-null object
3 UnemploymentRate 23 non-null float64
4 MedianHouseholdIncome 23 non-null float64
5 AverageAdjustedGrossIncome 23 non-null float64
6 SchoolTestPerformance 23 non-null object
7 Area 23 non-null float64
8 SchoolRacialMajorityPercentage 23 non-null object
9 SchoolHouseholdIncomePercentage 23 non-null object
10 RacialMajorityPercentage 23 non-null float64
11 PopulationDensity 23 non-null float64
dtypes: float64(7), int64(3), object(2)
memory usage: 2.5+ KB
```

```
df_zip_merged = pd.merge(df_demographics, df_zip, on='zipcode', how='outer')
df_zip_num = df_zip_merged.select_dtypes(include=['float64', 'int64'])

plt.subplots(figsize=(15, 10))
ax = sns.heatmap(df_zip_num.corr(), annot=True, cmap='Greens')
```


Store, most profitable

```
total_profit_store_num = df_merged_final.groupby(['storenum']).sum().sort_values(['profit'], ascending=True)
total_profit_store_num.sort_values().head()
```

```
#Map 4641 in mind
storenum
5970 257.497467
6064 328.004706
6096 486.239412
6096 1270.866471
5891 1107.171765
Name: profit, dtype: float64
```

```
total_profit_store_num.describe()

count      22.000000
mean      28282.520687
std       90299.770394
min        257.497467
25%       3175.423529
50%       6484.453824
75%       18233.935483
max       646563.885882
Name: profit, dtype: float64
```

```
total_profit_store_num.sort_values().plot(kind='bar',
                                           color='green',
                                           ylabel='Total ($)',
                                           title='Profit / Store (number)',
                                           figsize=(15,5))

<AxesSubplot: title='center': 'Profit / Store (number)', xlabel='storenum', ylabel='Total ($)'>
```



```
total_profit_store_num.plot.kde(color='green', title='Profit distribution (shops)')

<AxesSubplot: title='center': 'Profit distribution (shops)', ylabel='Density'>
```


Most profitable products & product sets

```
profit_products = df_merged_final.groupby(['itemnum'])['profit'].sum().sort_values().tail(500)
profit_products.set = set(np.array(profit_products.index))

#numpy arrays for well formatted output
np.array(profit_products.index)[0:100]
```

```
array([13695, 64513, 34783, 42703, 34819, 81208, 81124, 48107,
       86018, 52314, 36447, 11588, 5696, 43136, 918010, 34003,
       77174, 80456, 33663, 15582, 57125, 59154, 5446, 43053,
       27780, 5347, 19063, 7632, 82147, 43387, 68340, 10548,
       34014, 43095, 34116, 82636, 87589, 37346, 65200, 37886,
       35354, 35416, 65258, 36900, 76478, 89866, 34881, 89025,
       100413, 89445, 17766, 86691, 12408, 27605, 27544, 38179,
       21598, 5008, 42444, 2788, 85200, 22228, 36989, 3430,
       33658, 101187, 27410, 89121, 52318, 77852, 46658, 68022,
       37655, 89796, 64986, 82846, 77708, 58875, 87511, 32238,
       89242, 43038, 74803, 17776, 57053, 44419, 98953, 36618,
       35693, 34579, 29113, 5635, 89788, 86796, 28086, 618650,
       43528, 77842, 87280, 43121], dtype=int64)
```

```
Products available in 4641
products_4641_set = set(df_merged_final.loc[df_merged_final['storenum'] == '4641']['itemnum'].unique())
np.array(products_4641_set)
```

```
array([111776, 11777, 25604, 34820, 25606, 36873, 36875, 11788, 36878, 43025, 43026, 43028, 19476, 11290, 11294,
       37996, 37994, 37998, 89194, 89196, 43121, 89199, 28275, 12404, 19061, 43126, 34423, 19064, 28279, 12408, 19066,
       43124, 38028, 43125, 43136, 34515, 27780, 77558, 40070, 27783, 40502, 40535, 40594, 37012, 86691, 86692, 38036,
       1, 88294, 65254, 65256, 86296, 26858, 31474, 31475, 78078, 78079, 40132, 43777, 1799, 15626, 46350, 43285, 3817
       4, 28175, 38178, 38180, 30194, 81204, 81205, 81206, 81706, 43331, 43334, 43356, 43358, 88406, 44856, 73055, 648
       64, 64865, 64866, 52579, 66884, 66885, 86886, 64870, 64865, 87403, 87408, 52594, 77714, 27544, 61009, 77
       70, 75183, 28081, 86466, 75208, 36301, 36304, 36306, 36308, 68052, 58838, 66519, 26585, 37338, 53214, 27102, 4
       8101, 32232, 32232, 48105, 32234, 32235, 32236, 39017, 39018, 27628, 87026, 58668, 65013, 27125, 58672, 11771,
       11774], dtype=object)
```

```
len(products_4641_set)

159
```

```
recommended_products_set = profit_products.set - products_4641_set

total_set_4641 = products_4641_set | recommended_products_set
np.array(recommended_products_set)
```



```
array([36819, 34828, 69637, 36874, 78842, 77843, 36886, 36887, 69657, 40306, 77852, 40308, 49185, 49186, 36900,
49189, 36903, 36908, 88167, 34871, 100113, 34881, 100423, 26710, 43095, 88162, 36967, 36970, 34876, 34876,
0, 41077, 41078, 41079, 41080, 41081, 41082, 41083, 41084, 41085, 41086, 41087, 41088, 41089, 41090, 41091, 41092, 41093,
27, 45277, 45278, 45279, 45280, 45281, 45282, 45283, 45284, 45285, 45286, 45287, 45288, 45289, 45290, 45291, 45292, 45293, 45294,
45295, 45296, 45297, 45298, 45299, 45300, 45301, 45302, 45303, 45304, 45305, 45306, 45307, 45308, 45309, 45310, 45311,
646, 1885, 49333, 49334, 49335, 49336, 49337, 49338, 49339, 49340, 49341, 49342, 49343, 49344, 49345, 49346, 49347, 49348,
646, 1885, 3337, 3338, 88536, 33374, 33374, 33374, 88536, 33374, 33374, 33374, 33374, 33374, 33374, 33374, 33374, 33374,
35354, 10784, 10789, 10792, 10784, 10804, 10808, 10825, 10809, 40406, 35416, 35418, 29287, 19063, 19067, 48886,
19069, 41604, 41614, 41616, 41617, 41618, 41619, 41620, 41621, 41622, 41623, 41624, 41625, 41626, 41627, 41628, 41629,
41630, 41631, 41632, 41633, 41634, 41635, 41636, 41637, 41638, 41639, 41640, 41641, 41642, 41643, 41644, 41645, 41646,
41647, 41648, 41649, 41650, 41651, 41652, 41653, 41654, 41655, 41656, 41657, 41658, 41659, 41660, 41661, 41662, 41663,
41664, 41665, 41666, 41667, 41668, 41669, 41670, 41671, 41672, 41673, 41674, 41675, 41676, 41677, 41678, 41679, 41680,
41681, 41682, 41683, 41684, 35638, 35663, 506, 15248, 82846, 82847, 64418, 31658, 506, 307, 82867, 26055, 48,
979, 48102, 31719, 48106, 48107, 48108, 48102, 31708, 64513, 25067, 25068, 19467, 19466, 25068, 11297, 89121, 1,
1299, 12998, 44884, 74801, 37939, 45673, 81954, 81954, 81954, 81954, 81954, 81954, 81954, 81954, 81954, 81954, 81954,
21596, 21597, 21598, 21599, 21600, 21601, 21602, 21603, 21604, 21605, 21606, 21607, 21608, 21609, 21610, 21611, 21612, 21613,
21614, 21615, 21616, 21617, 21618, 21619, 21620, 21621, 21622, 21623, 21624, 21625, 21626, 21627, 21628, 21629, 21630,
46, 5347, 8124, 30929, 34030, 81240, 81247, 50316, 34036, 10836, 89339, 40131, 64776, 15627, 64776, 15627, 64776,
54, 3042, 30279, 30428, 64816, 15667, 50316, 81248, 34036, 10836, 89339, 40131, 64776, 15627, 64776, 15627, 64776,
8, 8945, 17567, 89447, 5486, 52599, 91650, 34262, 87409, 69635, 44419, 40327, 64904, 26244, 64914, 75155, 65,
04, 17811, 20806, 28067, 87485, 67608, 71210, 36305, 36307, 87510, 87511, 45516, 64969, 32233, 89579, 32238, 58,
979, 48102, 31719, 48106, 48107, 48108, 48102, 31708, 64513, 25067, 25068, 19467, 19466, 25068, 11297, 89121, 1,
1299, 12998, 44884, 74801, 37939, 45673, 81954, 81954, 81954, 81954, 81954, 81954, 81954, 81954, 81954, 81954, 81954,
21596, 21597, 21598, 21599, 21600, 21601, 21602, 21603, 21604, 21605, 21606, 21607, 21608, 21609, 21610, 21611, 21612, 21613,
21614, 21615, 21616, 21617, 21618, 21619, 21620, 21621, 21622, 21623, 21624, 21625, 21626, 21627, 21628, 21629, 21630,
46, 5347, 8124, 30929, 34030, 81240, 81247, 50316, 34036, 10836, 89339, 40131, 64776, 15627, 64776, 15627, 64776,
54, 3042, 30279, 30428, 64816, 15667, 50316, 81248, 34036, 10836, 89339, 40131, 64776, 15627, 64776, 15627, 64776,
8, 8945, 17567, 89447, 5486, 52599, 91650, 34262, 87409, 69635, 44419, 40327, 64904, 26244, 64914, 75155, 65,
04, 17811, 20806, 28067, 87485, 67608, 71210, 36305, 36307, 87510, 87511, 45516, 64969, 32233, 89579, 32238, 58,
979, 48102, 31719, 48106, 48107, 48108, 48102, 31708, 64513, 25067, 25068, 19467, 19466, 25068, 11297, 89121, 1,
1299, 12998, 44884, 74801, 37939, 45673, 81954, 81954, 81954, 81954, 81954, 81954, 81954, 81954, 81954, 81954, 81954,
21596, 21597, 21598, 21599, 21600, 21601, 21602, 21603, 21604, 21605, 21606, 21607, 21608, 21609, 21610, 21611, 21612, 21613,
21614, 21615, 21616, 21617, 21618, 21619, 21620, 21621, 21622, 21623, 21624, 21625, 21626, 21627, 21628, 21629, 21630,
46, 5347, 8124, 30929, 34030, 81240, 81247, 50316, 34036, 10836, 89339, 40131, 64776, 15627, 64776, 15627, 64776,
54, 3042, 30279, 30428, 64816, 15667, 50316, 81248, 34036, 10836, 89339, 40131, 64776, 15627, 64776, 15627, 64776,
8, 894
```

[illegible]

```
Out[53]: array(['3696', '2521',
                '5533', '499',
                '2633', '2633'])
```

```
'5507', '5127',
'5375', '4599',
'2850', '3825'
```

```

In [54]: # Getting time delta of all stores
stores_delta_list = np.array([store,sale_timedelta(store,df_merged_final)] for store in storenums])

In [55]: stores_delta_list[0:5]

Out[55]: array([[ '3696', '1182'],
 ['2528', '1180'],
 ['5501', '1176'],
 ['5057', '1178'],
 ['5338', '1173']], dtype=<O4>)

In [56]: #Turning array to series, for index matching
idx, values = zip(*stores_delta_list)

stores_timedelta_series = pd.Series([int(k) for k in values], idx)

In [57]: # Finding profit per day
store_profit_perday= total_profit_store_num/stores_timedelta_series

store_profit_perday.sort_values().tail()

Out[57]: 2626    54.743972
2527    60.304628
2528    65.698499
2561    90.813035
2633    546.545973
dtype: float64

In [58]: for index in store_profit_perday.tail().index:
print(store_profit_perday[index])

2.189442190669386
3.458126361655788
0.8903952185426444
10.059784313725558
3.904817927170887

Comparison plots

In [176]: # class else write 100 lines of code
class Opti:
    def __init__(self, name: str, cost: float, daily_profit: float, fixed_cost: float, n: int, days: int, investment: float):
        self.name = name
        self.cost = cost
        self.daily_profit = daily_profit
        self.fixed_cost = fixed_cost
        self.n = n
        self.days = days
        self.investment = investment

    def cashplot(self):
        if self.n != 1 or self.n == 0:
            self.cost = (self.cost*self.n)-(self.fixed_cost*(self.n-1))
            self.daily_profit = self.daily_profit*self.n

        cashsum_pharmacy = 0

```

ph
ph
el

```

summed_list = []
for v in pharmacy_cash:
    cashsum_pharmacy += v
    summed_list.append(cashsum_pharmacy)

return summed_list

def breakeven(self):
    if self.investment:
        return 0
    return self.cost//self.daily_profit

def profit(self):
    if self.investment:
        return self.cashplot()[-1] - self.cost
    return self.cashplot()[-1]

def plot(self):
    interval = range(1,self.days+1)
    plt.rcParams['figure.figsize'] = (15, 10)
    plt.plot(interval, self.cashplot(), label=self.name)
    plt.ylabel('$')
    plt.xlabel('days')
    plt.title('Net positive cashflow over 10 years')
    plt.legend()
    plt.grid(True)

```

```

In [302]: consise_list = []

consise_list.append(Option('Buy 4641',0,store_profit_perday('4641'),0.1,3650,False))
consise_list.append(Option('Building pharmacy',0,store_profit_perday.mean(),0.1,3650,False))

```

```

In [303]: store_profit_perday.mean()

```

```

Out[303]: 24.3177582684731

```

```

In [304]: for v in consise_list:
          v.plot()

```

Vendors

```

In [324]: df_merged_final.groupby(['vendornum'])\
          .sum().sort_values(['profit'], ascending = True) ['profit'].tail(15)\
          .plot(kind = 'bar',color = 'green')

```

```

Out[324]: <AxesSubplot: xlabel='vendornum'>

```

```

In [346]: df_4641['vendornum'].unique()

Out[346]: 26

In [325]: 161*0.2

Out[325]: 32.2

In [327]: vendors = df_merged_final.groupby(['vendornum'])\
            .sum().sort_values(['profit'], ascending = True)['profit'].tail(64)\
            np.array(vendors.index)

Out[327]: array([479., 357., 342., 216., 497., 285., 308., 482., 564., 239., 557.,
                287., 267., 451., 391., 978., 116., 31., 402., 384., 306., 214.,
                409., 346., 521., 130., 566., 388., 492., 195., 154., 266.,
                297., 322., 277., 125., 619., 461., 330., 305., 626., 229., 389.,
                267., 235., 192., 300., 410., 380., 205., 259., 55., 115., 420.,
                305., 35., 85., 434., 301., 65., 370., 421., 260.])

In [341]: vendors_4641 = set(df_4641['vendornum'].unique())

In [344]: np.array(set(vendors.index) | vendors_4641)

Out[344]: array([259.0, 260.0, 521.0, 246.0, 267.0, 277.0, 285.0, 287.0, 35.0, 297.0, 300.0, 557.0, 301.0, 305.0, 306.0,
                30.0, 564.0, 566.0, 95.0, 65.0, 322.0, 330.0, 85.0, 342.0, 346.0, 91.0, 357.0, 419.0, 348.0, 626.0, 115.0, 11.
                6.0, 370.0, 380.0, 125.0, 384.0, 130.0, 389.0, 391.0, 395.0, 402.0, 154.0, 410.0, 420.0, 421.0, 434.0, 184.0, 1
                92.0, 451.0, 195.0, 461.0, 205.0, 978.0, 469.0, 214.0, 216.0, 479.0, 482.0, 229.0, 492.0, 239.0, 240.0, 497.0,
                255.0])
dtype=object)

Random calculations

In [278]: stores_perzipcode = df_info_final.groupby(['zipcode'])['storenum'].count()

          area_zipcode = df_demographics.groupby(['zipcode']).first()['Area']

          store_perarea = stores_perzipcode/area_zipcode

          store_perarea.store(perfit_per_zipcode, method='pearson')

Out[278]: -0.05322828236009978

In [279]: df_4641 = df_merged_final.loc[(df_merged_final['storenum'] == '4641') & (df_merged_final['date'] <= '2021-01-01')
df_4641.head()

Out[279]:
   invnum    date  storenum  vendornum  itemnum  itemcategory  productcost  product_retail_price  units_sold  revenue_original  ...
75633    INV- 2020- 02-07    4641      2590      11776      OTC/VMS             5.23                7.85              12              94.20  ...

```

75634 INV
2507310000

	INV- 2020- 25073100010	4641	2590	17788	Personal hygiene	1045	15.68	6	94.08 ...
75636	INV- 2020- 25073100003	4641	421.0	86884	Skin care items	4.78	7.17	10	71.70 ...
75637	INV- 2020- 25073100008	4641	192.0	65254	Skin care items	6.52	9.78	10	97.80 ...

5 rows × 21 columns

```
In [280]: df_4641['month'] = df_4641['date'].dt.month

<ipython-input-280-27bf6e2656d2>:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
df_4641['month'] = df_4641['date'].dt.month

In [281]: df_4641['month']

Out[281]:
75633    2
75634    2
75635    2
75636    2
75637    2
...
79276   12
79277   12
79278   12
79279   12
79280   12
Name: month, Length: 3648, dtype: int64

In [282]: df_4641['hour'] = df_4641['date'].dt.hour
df_4641['minute'] = df_4641['date'].dt.minute
#df_4641['count'] = 1

<ipython-input-282-99625344d591>:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
df_4641['hour'] = df_4641['date'].dt.hour
<ipython-input-282-99625344d591>:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

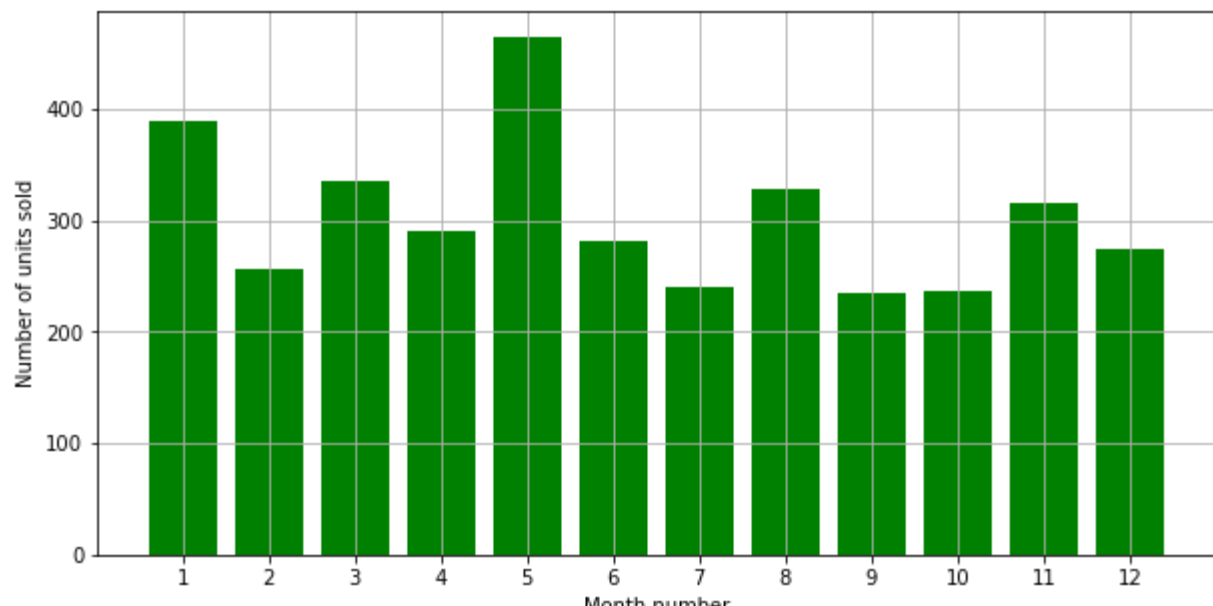
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
df_4641['minute'] = df_4641['date'].dt.minute

In [286]: months = range(1,13)

plt.rcParams['figure.figsize'] = (10, 5)
plt.bar(months, df_4641.groupby('month').count()['itemnum'], color = 'green')
plt.xticks(months)
```

```
plt.ylabel('Nu')
plt.xlabel('Mo')
plt.grid()
```

```
plt.show()
```



Month number	Number of units sold
1	380
2	260
3	330
4	290
5	450
6	280
7	240
8	320
9	240
10	240
11	310
12	280

```
In [287]: store_profit_perday.describe()
```

```
Out[287]: count    52.000000
          mean    24.317776
          std     76.238075
          min      0.857565
          25%     3.417513
          50%     5.538680
          75%    15.509820
          max    546.545973
          dtype: float64
```

```
In [288]: percentilescore(store_profit_perday, store_profit_perday['6641'])
```

```
Out[288]: 48.07692307692308
```

```
In [ ]:
```