

```
In [21]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from datetime import datetime, relativedelta

from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import LabelEncoder
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score
from sklearn.preprocessing import LabelEncoder, OneHotEncoder
from sklearn.compose import ColumnTransformer
from sklearn.model_selection import train_test_split

import plotly.express as px
from mpl_toolkits.mplot3d import Axes3D

from statsmodels import cycle
from statsmodels import product

from googletrans import Translator

import plotlylib inline

In [22]: #Read in .csv files into pandas dataframes
train = pd.read_csv('sales_train.csv')
test = pd.read_csv('test.csv').set_index('ID')
submission = pd.read_csv('sample_submission.csv')
items = pd.read_csv('items.csv')
item_cats = pd.read_csv('item_categories.csv')
shops = pd.read_csv('shops.csv')
shops_t = pd.read_csv('shops_translated.csv')
items_t = pd.read_csv('items_translated_test.csv')

In [23]: #Down casts the data entries from int64 to int32 and float64 to float32
#This reduces the size of the records by almost half. (from 134mb to 61mb)
def downcast_dtypes(df):
    float_cols = [c for c in df if df[c].dtype == "float64"]
    int_cols = [c for c in df if df[c].dtype in ["int64", "int32"]]
    df[float_cols] = df[float_cols].astype(np.float32)
    df[int_cols] = df[int_cols].astype(np.int16)
    return df

In [24]: #Calls the downcasting function
train = downcast_dtypes(train)
test = downcast_dtypes(test)
submission = downcast_dtypes(submission)
items = downcast_dtypes(items)
item_cats = downcast_dtypes(item_cats)
shops = downcast_dtypes(shops)
shops_t = downcast_dtypes(shops_t)
items_t = downcast_dtypes(items_t)

In [25]: #Adda item id to main train dataframe and drops the string name column
train = train.merge(items, on='item_id')
train = train.drop(columns = 'item_name')

In [26]: train

Out[26]:
```

|         | date       | date_block_num | shop_id | item_id | item_price | item_cnt_day | item_category_id |
|---------|------------|----------------|---------|---------|------------|--------------|------------------|
| 0       | 02.01.2013 | 0              | 59      | 22154   | 999.0      | 1.0          | 37               |
| 1       | 23.01.2013 | 1              | 24      | 22154   | 999.0      | 1.0          | 37               |
| 2       | 20.01.2013 | 0              | 27      | 22154   | 999.0      | 1.0          | 37               |
| 3       | 02.01.2013 | 0              | 25      | 22154   | 999.0      | 1.0          | 37               |
| 4       | 03.01.2013 | 0              | 25      | 22154   | 999.0      | 1.0          | 37               |
| ...     | ...        | ...            | ...     | ...     | ...        | ...          | ...              |
| 2935844 | 17.10.2015 | 33             | 25      | 8428    | 249.0      | 1.0          | 40               |
| 2935845 | 01.10.2015 | 33             | 25      | 7903    | 12198.0    | 1.0          | 15               |
| 2935846 | 29.10.2015 | 33             | 25      | 7610    | 2890.0     | 1.0          | 64               |
| 2935847 | 22.10.2015 | 33             | 25      | 7635    | 2100.0     | 1.0          | 64               |
| 2935848 | 01.10.2015 | 33             | 25      | 7840    | 4040.0     | 1.0          | 64               |

2935849 rows x 7 columns

```
In [8]: ## Exploratory Merging and graphing cell for preliminary data analysis ##

# group data by month and shop_id, return sum of items sold per shop per month
month_group = pd.DataFrame(train.groupby(['date_block_num', 'shop_id'])['item_cnt_day'].sum().reset_index())

# added the item category into sales_train
train = pd.merge(train, items[['item_id', 'item_category_id']], on = 'item_id')

# group data by month and item category_id, return sum of items sold per category per month
category_group = pd.DataFrame(merged.groupby(['date_block_num', 'item_category_id'])['item_cnt_day'].sum().reset_index())

In [9]: # plotting data for visual reference

# Plot for items sold per shop per month
fig1, axes = plt.subplots(30, 2, figsize = (18,80), sharex = True, sharey = True)
shop1 = 0
axis1 = []

for i in range(int(np.max(month_group['shop_id'])/2+1)):
    axis.append([[],[],[],[],[]])

for row in range(30):
    for col in range(42):
        for date in range(len(month_group['shop_id'])):
            if month_group['shop_id'][date] == shop1:
                axis1[row][col][0].append(month_group['date_block_num'][date])
                axis1[row][col][1].append(month_group['item_cnt_day'][date])
                shop1+=1

shop = 0
cyclo = cycle('bgrcmk')
for row in range(42):
    for col in range(2):
        axes1[row, col].plot(axis1[row][col][0], axis1[row][col][1], color = next(cyclo))
        axes1[row, col].set_title('shop ' + str(shop))
        shop += 1

In [10]: # plotting data for visual reference

# Plot for items sold per category per month
fig2, axes2 = plt.subplots(42, 2, figsize = (18,100), sharex = True, sharey = True)
category = 0
axis2 = []

for i in range(int(np.max(category_group['item_category_id'])/2+1)):
    axis2.append([[],[],[],[],[]])

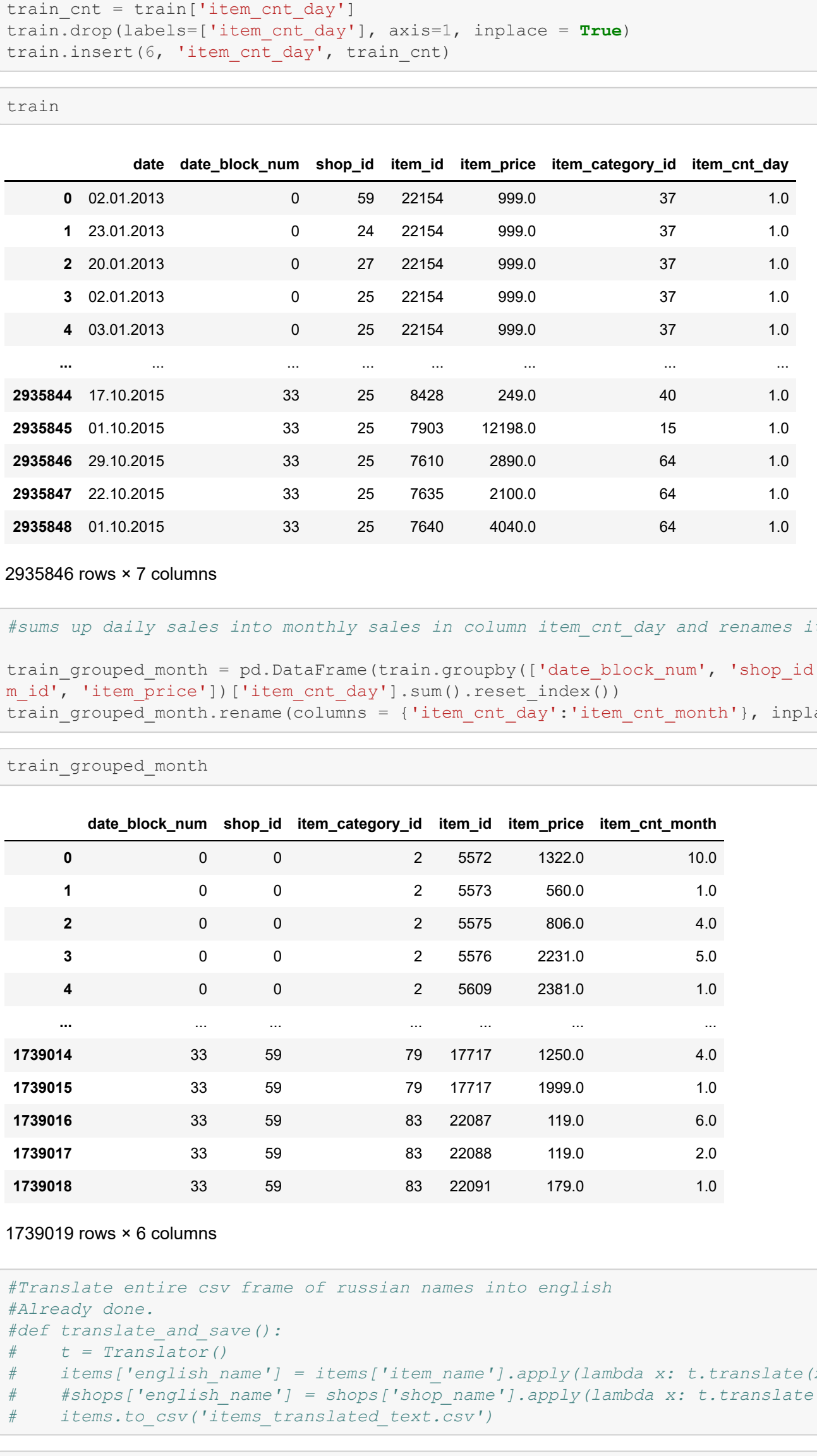
for row in range(42):
    for col in range(2):
        for date in range(len(category_group['item_category_id'])):
            if category_group['item_category_id'][date] == category:
                axis2[row][col][0].append(category_group['date_block_num'][date])
                axis2[row][col][1].append(category_group['item_cnt_day'][date])
                category+=1

category = 0
cyclo = cycle('bgrcmk')
for row in range(42):
    for col in range(2):
        axes2[row, col].plot(axis2[row][col][0], axis2[row][col][1], color = next(cyclo))
        axes2[row, col].set_title('category ID ' + str(category))
        category += 1

In [14]: # Boxplot to show outliers in item_cnt_day and item_price
# This shows the outlier in price and the one instance of an unusually high sales count
plt.figure(figsize=(12,4))
plt.xlim(-100, train.item_cnt_day.max()*1.2)
sns.boxplot(x=train.item_cnt_day)

plt.figure(figsize=(12,4))
plt.xlim(train.item_price.min()*1.1, train.item_price.max()*1.2)
sns.boxplot(x=train.item_price)

Out[14]:
```



```
In [30]: #Removes outliers from train
train = train[train.item_price < 90000]
train = train[train.item_cnt_day < 999]

#Replaces the negative price item with the median item price of all items with the id of 2973 and in sh
op id 32
train.loc[train.shop_id==32](train.item_id==2973)(train.date_block_num=4)(train.item_price=0
),item_price.median()

In [31]: # moves item_cnt_day to the last column
train_cnt = train['item_cnt_day']
train.drop(labels=['item_cnt_day'], axis=1, inplace = True)
train.insert(6, 'item_cnt_day', train_cnt)

In [32]: train

Out[32]:
```

|         | date       | date_block_num | shop_id | item_id | item_price | item_category_id | item_cnt_day |
|---------|------------|----------------|---------|---------|------------|------------------|--------------|
| 0       | 02.01.2013 | 0              | 59      | 22154   | 999.0      | 37               | 1.0          |
| 1       | 23.01.2013 | 0              | 24      | 22154   | 999.0      | 37               | 1.0          |
| 2       | 20.01.2013 | 0              | 27      | 22154   | 999.0      | 37               | 1.0          |
| 3       | 02.01.2013 | 0              | 25      | 22154   | 999.0      | 37               | 1.0          |
| 4       | 03.01.2013 | 0              | 25      | 22154   | 999.0      | 37               | 1.0          |
| ...     | ...        | ...            | ...     | ...     | ...        | ...              | ...          |
| 2935844 | 17.10.2015 | 33             | 25      | 8428    | 249.0      | 40               | 1.0          |
| 2935845 | 01.10.2015 | 33             | 25      | 7903    | 12198.0    | 15               | 1.0          |
| 2935846 | 29.10.2015 | 33             | 25      | 7610    | 2890.0     | 64               | 1.0          |
| 2935847 | 22.10.2015 | 33             | 25      | 7635    | 2100.0     | 64               | 1.0          |
| 2935848 | 01.10.2015 | 33             | 25      | 7840    | 4040.0     | 64               | 1.0          |

2935849 rows x 7 columns

```
In [36]: #sums up daily sales into monthly sales in column item_cnt_month and renames it item_cnt_month
train_grouped_month = pd.DataFrame(train.groupby(['date_block_num', 'shop_id', 'item_category_id', 'item_id', 'item_price'])['item_cnt_day'].sum().reset_index())
train_grouped_month.rename(columns = ['item_cnt_day':'item_cnt_month'], inplace = True)

In [37]: train_grouped_month

Out[37]:
```

|         | date_block_num | shop_id | item_category_id | item_id | item_price | item_cnt_month |
|---------|----------------|---------|------------------|---------|------------|----------------|
| 0       | 0              | 0       | 0                | 2       | 5572       | 1322.0         |
| 1       | 1              | 0       | 0                | 2       | 5573       | 560.0          |
| 2       | 2              | 0       | 0                | 2       | 5575       | 806.0          |
| 3       | 3              | 0       | 0                | 2       | 5576       | 2231.0         |
| 4       | 4              | 0       | 0                | 2       | 5609       | 2381.0         |
| ...     | ...            | ...     | ...              | ...     | ...        | ...            |
| 1739014 | 33             | 59      | 79               | 17717   | 1250.0     | 4.0            |
| 1739015 | 33             | 59      | 79               | 17717   | 1999.0     | 1.0            |
| 1739016 | 33             | 59      | 83               | 22087   | 119.0      | 6.0            |
| 1739017 | 33             | 59      | 83               | 22088   | 119.0      | 2.0            |
| 1739018 | 33             | 59      | 83               | 22091   | 179.0      | 1.0            |

1739019 rows x 6 columns

```
In [39]: #Translate entire .csv frame of russian names into english
#Already done.
def translate_and_save():
    t = Translator()
    # items['english_name'] = items['item_name'].apply(lambda x: t.translate(x).text)
    # #shops['english_name'] = shops['shop_name'].apply(lambda x: t.translate(x).text)
    # items.to_csv('items_translated_test.csv')

In [40]: # Changes numerical, categorical features into strings to properly be represented as categorical in one
hotencoding
# nominal integers can not be converted to binary encoding, convert to string
train_grouped_month['date_block_num'] = [[('month ' + str(i)) for i in train_grouped_month['date_block_num']]
train_grouped_month['shop_id'] = [[('shop ' + str(i)) for i in train_grouped_month['shop_id']]
train_grouped_month['item_category_id'] = [[('item_category ' + str(i)) for i in train_grouped_month['item_category_id']]
train_grouped_month['item_id'] = [[('item ' + str(i)) for i in train_grouped_month['item_id']]

In [41]: train_grouped_month

Out[41]:
```

|         | date_block_num | shop_id | item_category_id | item_id    | item_price | item_cnt_month |
|---------|----------------|---------|------------------|------------|------------|----------------|
| 0       | month 0        | shop 0  | item_category 2  | item 5572  | 1322.0     | 10.0           |
| 1       | month 0        | shop 0  | item_category 2  | item 5573  | 560.0      | 1.0            |
| 2       | month 0        | shop 0  | item_category 2  | item 5575  | 806.0      | 4.0            |
| 3       | month 0        | shop 0  | item_category 2  | item 5576  | 2231.0     | 5.0            |
| 4       | month 0        | shop 0  | item_category 2  | item 5609  | 2381.0     | 1.0            |
| ...     | ...            | ...     | ...              | ...        | ...        | ...            |
| 1739014 | month 33       | shop 59 | item_category 79 | item 17717 | 1250.0     | 4.0            |
| 1739015 | month 33       | shop 59 | item_category 79 | item 17717 | 1999.0     | 1.0            |
| 1739016 | month 33       | shop 59 | item_category 83 | item 22087 | 119.0      | 6.0            |
| 1739017 | month 33       | shop 59 | item_category 83 | item 22088 | 119.0      | 2.0            |
| 1739018 | month 33       | shop 59 | item_category 83 | item 22091 | 179.0      | 1.0            |

1739019 rows x 6 columns

```
In [44]: # stores all feature columns in the x variable, and stores the target variable as the last column in y
variable
x = train_grouped_month.iloc[:, :-1].values
y = train_grouped_month.iloc[:, -1].values

In [45]: # Encoding categorical data
# provides a value to data that can be then used in regression equations, eg. Friday = 1 etc.
ct = ColumnTransformer([('encoder', OneHotEncoder([0, 1, 2, 3])), remainder = 'passthrough'])
x = ct.fit_transform(x)

In [46]: # split data set into training set and testing set
X_train, X_test, y_train, y_test = train_test_split(x, y, test_size = 0.2, random_state = 0)

In [55]: # fitting multiple linear regression to the training set
regressor = LinearRegression()
regressor.fit(X_train, y_train)

# prediction of test split from original data, train_grouped_month
y_prediction = regressor.predict(X_test)

# the output of this is the r squared
regressor.score(X_test, y_test)

Out[55]: 0.24532067463255955

In [ ]:

In [48]: # This is for playing around with different items and shops
# sample format of test prediction = ['month 0', 'shop 0', 'item_category 2', 'item 5572', 1322]

# A simple program to see how the predictor would work
# from scipy import stats
def auto_predictor():
    global ct
    month = input("Please enter the month you want to predict. Enter the month, eg. Jan = 0")
    shop = input("Please enter the shop ID, eg. 2 : ")
    item = input("Please enter the item ID you wish to predict, eg. for item 55, enter '55' : ")
    item_cat = (items.loc[items['item_id'] == int(item)]['item_category_id'].values[0][0])
    prices = train.loc[train['item_id'] == int(item), ['item_price']].values
    price = (stats.mode(prices))[0][0][0]

    z = ['month ' + str(month), 'shop ' + str(shop), 'item_category ' + str(item_cat), 'item ' + str(item), price]
    z = np.array(z, dtype = object).reshape(1, -1)
    z = ct.transform(z)
    z_pred = regressor.predict(z)

    # Things to check!!!
    # ***there are multiple prices available for the item in that specific shop, provide an option to choose
```



