

Pythonによるドーナツの売上予測

山形東高校H8班 2年 木野紗花

研究動機 目的

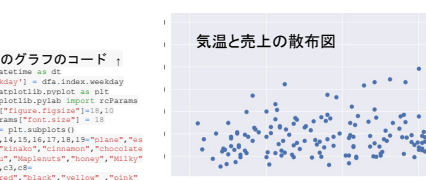
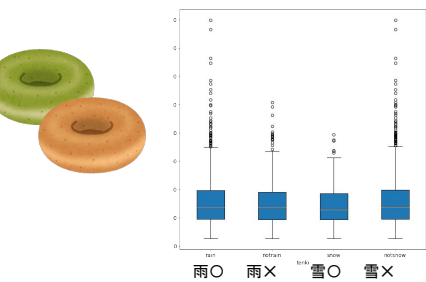
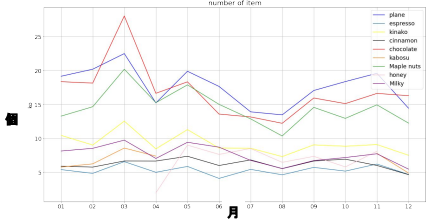
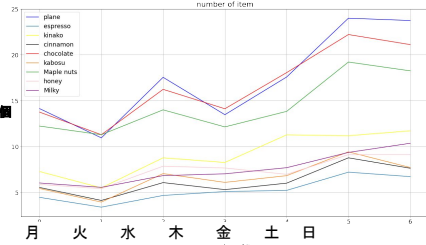
学校近くにドーナツ店があるが、私が行くころには売り切れている事が多い。需要と供給のバランスを導き出せれば私と同じような利用者にも、なおかつ御店主にも貢献できるのではないかと思い、探究活動で学習中のPythonを用い、売上予測を行っている。

研究(予測)方法

店舗は七日町商店街近くの道路沿いにあり、専用の駐車場はなく、徒歩でのお客さんが多い。客層は主婦が多いようだ。売上は過去3年分、商品数は過去1年分のデータを加工し、また外部からもデータを取得し、Pythonのライブラリのモデルで説明変数(週や降雪量、コロナ...)と目的変数(金額)の対応関係を学習させ、一月の予測をする

データの可視化

Pythonのライブラリmatplotlib.pyplotやseabornを用いて売上と週、月、天気等との関係を調べる。



```
一番上のグラフのコード
import datetime as dt
def [weekly] = df.index.weekday
import matplotlib.pyplot as plt
from matplotlib.pyplot import rcParams
rcParams['figure.figsize'] = 10, 10
plt.rcParams['font.size'] = 10
fig, ax = plt.subplots(1, 1)
ax = plt.subplot(1, 1)
ax.plot(week, index, caseweek, 'plane', color='red', label='plane')
ax.set_xlabel('week') # 日付をみる
ax.set_ylabel('no') # 箱ひげ図
ax.grid()
ax.plot(week, index, caseweek, 'espresso', color='blue', label='espresso')
ax.plot(week, index, caseweek, 'kinako', color='green', label='kinako')
ax.plot(week, index, caseweek, 'chocolate', color='purple', label='chocolate')
ax.plot(week, index, caseweek, 'honey', color='orange', label='honey')
ax.plot(week, index, caseweek, 'maple_nuts', color='pink', label='maple_nuts')
ax.plot(week, index, caseweek, 'honey_lemon', color='yellow', label='honey_lemon')
ax.plot(week, index, caseweek, 'milk_ee', color='cyan', label='milk_ee')
ax.legend(loc='0')
fig.tight_layout() # レイアウトの設定
plt.show()
```

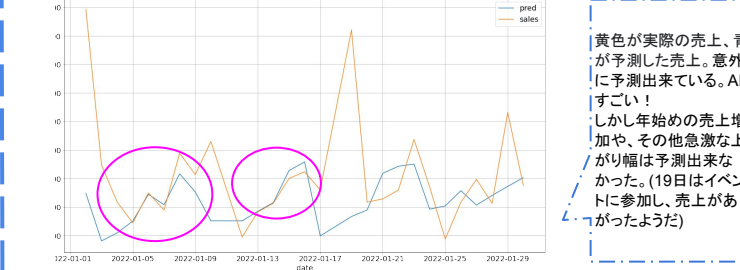
参考文献
kaggleで勝つデータ分析の技術
Pythonで始めるkaggleスタートブック
SIGNATE Qiita

謝辞
売上データ提供していただいたnicoドーナツ山形店、
助言をしてくださった先生方に感謝の意を表します。

予測

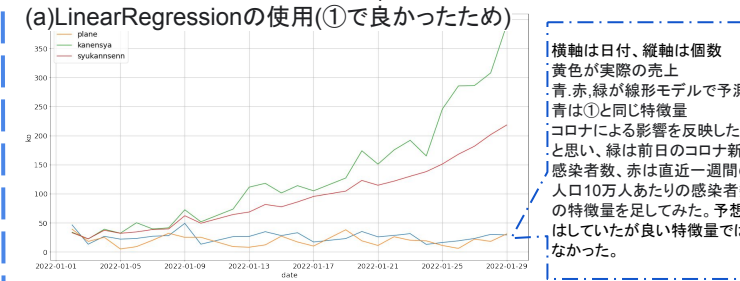
①合計の売上

3年分の過去データを基に、特徴量を日、月、週、降水量、降雪量にしてLinearRegressionモデル(使いやすかったため)に当てはめ、売上予測

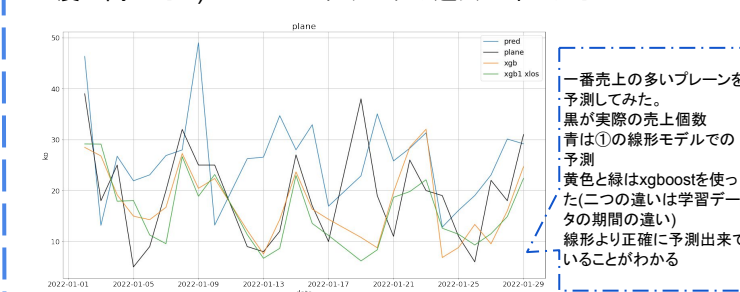


②商品(プレーン)の売上個数

目的である個数の売上予測をするため1年分の過去データをトレーニングデータとした。(商品ごとのデータは一日ごとで3年分のデータをダウンロードする時間がなかったため)



(b)XGboostの使用(不要な特徴量を追加しても影響が小さい上、精度が高いため) トレーニングデータは過去一年とした



③ほかの商品の売上個数の予測



考察 今後の展望
Pythonで売上をある程度予測できた。
しかし、これをもとにドーナツを作るとすると危険があるので、より良い特徴量とハイパラメータを見つける必要がある。
また、天気などは今回は過去データを使ったが、実際は明日の天気予報を入力しなければならない。すべての特徴量を入力するのは大変なので、影響力の少ない特徴量を消す必要がある。一日後であったら、ARIMAモデルなども役に立ちそう。祝日や感染警戒レベル等も特徴量に加えてみたい。

同じように、xgboostを使った
黒が実際の売上個数
青が予測した売上個数
scikit-learnのmetricsモジュール
mean_squared_log_errorでRMSLEを求めた。1日、19日は例外として、全体的に惜しい。

データ元
売上データ:nicoドーナツ山形店
気象データ:気象庁
コロナ関連:NHK