

## Data Modeling 數據建模 (sns regplot&scatterplot, statmodels.formula.api)

The model before transform 模型轉換前



將 x 軸用 np.sqrt 平方根後，數值修正，得到正確的分佈

The model after transform 模型轉換後

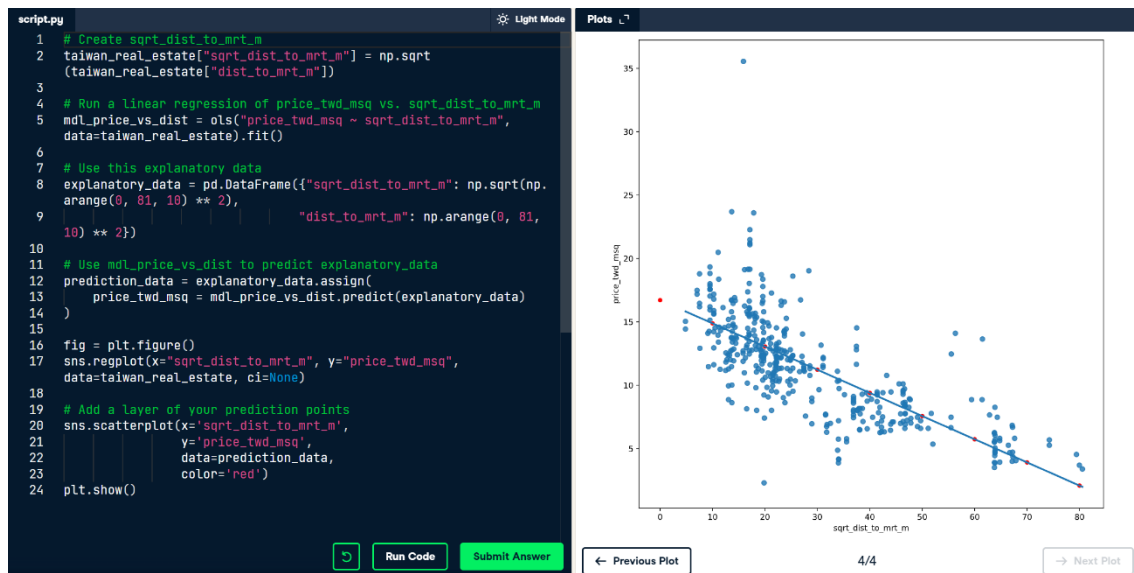


## The Comparison of Original and Predict model

利用上述模型，與預測模型(線上紅點)做比較，觀察結果跟預測數值差異

紅點皆在點上，表示預測值與實際結果相符

當周圍的便利商店數量多，房地坪數價格會相對高



RSE of original model and transformed model  
RSE 在模型轉換後，值變小，代表模型準確率提升

```
script.py
1 # Calculate mse_orig for mdl_click_vs_impression_orig
2 mse_orig = mdl_click_vs_impression_orig.mse_resid
3
4 # Calculate rse_orig for mdl_click_vs_impression_orig and print it
5 rse_orig = np.sqrt(mse_orig)
6 print("RSE of original model: ", rse_orig)
7
8 # Calculate mse_trans for mdl_click_vs_impression_trans
9 mse_trans = mdl_click_vs_impression_trans.mse_resid
10
11 # Calculate rse_trans for mdl_click_vs_impression_trans and print it
12 rse_trans = np.sqrt(mse_trans)
13 print("RSE of transformed model: ", rse_trans)
```

Run Code Submit Answer

```
IPython Shell Slides Notes
# calculate and print the specificity
specificity = TN / (TN + FP)
print('specificity: ', specificity)

# Calculate mse_orig for mdl_click_vs_impression_orig
mse_orig = mdl_click_vs_impression_orig.mse_resid

# Calculate rse_orig for mdl_click_vs_impression_orig and print it
rse_orig = np.sqrt(mse_orig)
print("RSE of original model: ", rse_orig)

# Calculate mse_trans for mdl_click_vs_impression_trans
mse_trans = mdl_click_vs_impression_trans.mse_resid

# Calculate rse_trans for mdl_click_vs_impression_trans and print it
rse_trans = np.sqrt(mse_trans)
print("RSE of transformed model: ", rse_trans)

RSE of original model: 19.905838862478138
RSE of transformed model: 0.19690640896875722
```

## The Comparison of Original and Logistic Predict Model

預測模型(藍色)在 logistic 後變成曲線，尾端拉升  
代表客戶距離第一次購買產品(time\_since\_first\_purchase)，時間隔得越久  
客戶越容易流失 (has\_churned = 0~1 未流失~流失)



將預測模型轉為 scatterplot



將模型改為另一種簡單預測(紅點)

### Most Likely Outcome

(has\_churned, 即 $<0.5$  客戶不會流失,  $>0.5$  客戶更有可能會流失)

可得出距離首次購買時間間隔越久, 流失機率越高

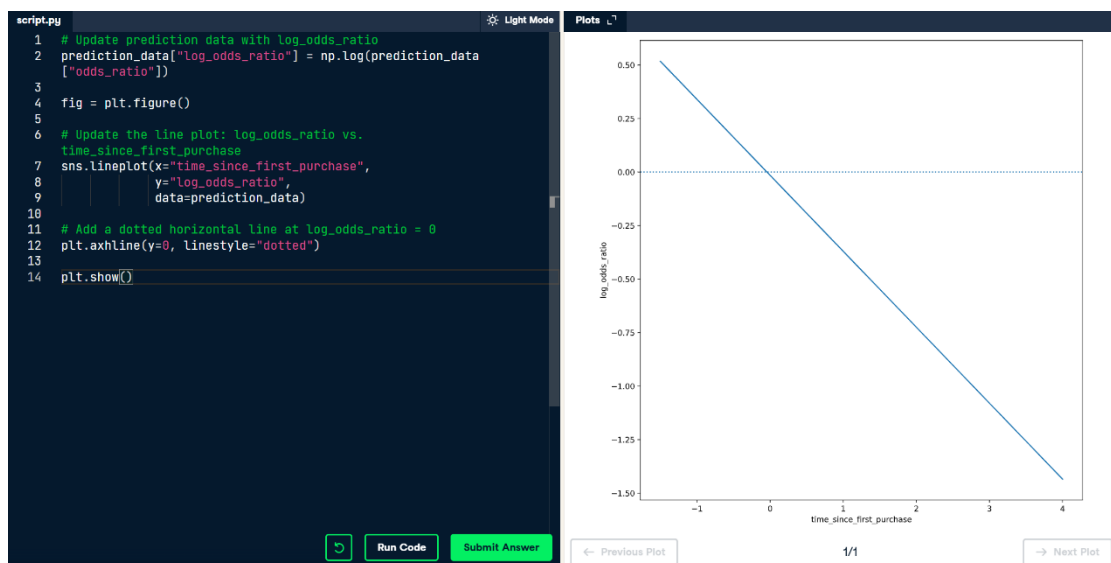


## The Comparison of Odds Ratio and the Time First Purchase

利用 Odds Ratio 賠率比與首次購買時間做比較  
距離首購時間越久，賠率比越接近 1，流失率越高



將賠率比 log，使模型穩定  
沒有 log 的賠率比(即曲線)，模型波動較大，不易穩定



## Linear Regression 線性回歸分類比較

```
script.py
1 # Update the model formula to remove the intercept
2 mdl_price_vs_age0 = ols("price_twd_msq ~ house_age_years + 0", data=taiwan_real_estate).fit()
3
4 # Print the parameters of the fitted model
5 print(mdl_price_vs_age0.params)
```

Run Code Submit Answer

```
IPython Shell Slides Notes
house_age_years[1.30 to 45] -1.244
dtype: float64

<script.py> output:
Intercept 12.637
house_age_years[1.15 to 30] -2.761
house_age_years[1.30 to 45] -1.244
dtype: float64

# Update the model formula to remove the intercept
mdl_price_vs_age0 = ols("price_twd_msq ~ house_age_years + 0", data=taiwan_real_estate).fit()

# Print the parameters of the fitted model
print(mdl_price_vs_age0.params)

house_age_years[0 to 15] 12.637
house_age_years[15 to 30] 9.877
house_age_years[30 to 45] 11.393
dtype: float64

In [1]:
```

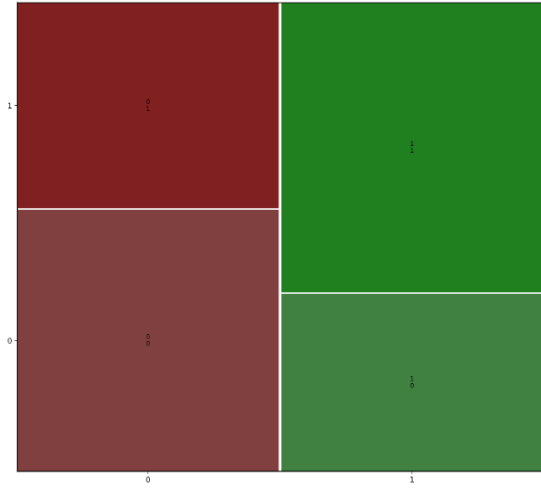
## Confusion Matrix 模糊矩陣

```
script.py  Light Mode
1 # Get the actual responses
2 actual_response = churn['has_churned']
3
4 # Get the predicted responses
5 predicted_response = np.round(mdL_churn_vs_relationship.predict())
6
7 # Create outcomes as a DataFrame of both Series
8 outcomes = pd.DataFrame({"actual_response": actual_response,
9                           "predicted_response": predicted_response})
10
11 # Print the outcomes
12 print(outcomes.value_counts(sort = False))
```

Run Code Submit Answer

```
script.py solution.py  Light Mode
1 # Import mosaic from statsmodels.graphics.mosaicplot
2 from statsmodels.graphics.mosaicplot import mosaic
3
4 # Calculate the confusion matrix conf_matrix
5 conf_matrix = mdL_churn_vs_relationship.pred_table()
6
7 # Print it
8 print(conf_matrix)
9
10 # Draw a mosaic plot of conf_matrix
11 mosaic(conf_matrix)
12 plt.show()
```

Plots 1/1



Previous Plot 1/1 Next Plot