

Analyze_ab_test_results_notebook

July 31, 2020

0.1 Analyze A/B Test Results

You may either submit your notebook through the workspace here, or you may work from your local machine and submit through the next page. Either way assure that your code passes the project [RUBRIC](#). **Please save regularly.**

This project will assure you have mastered the subjects covered in the statistics lessons. The hope is to have this project be as comprehensive of these topics as possible. Good luck!

0.2 Table of Contents

- Section ??
- Section ??
- Section ??
- Section ??

Introduction

A/B tests are very commonly performed by data analysts and data scientists. It is important that you get some practice working with the difficulties of these

For this project, you will be working to understand the results of an A/B test run by an e-commerce website. Your goal is to work through this notebook to help the company understand if they should implement the new page, keep the old page, or perhaps run the experiment longer to make their decision.

As you work through this notebook, follow along in the classroom and answer the corresponding quiz questions associated with each question. The labels for each classroom concept are provided for each question. This will assure you are on the right track as you work through the project, and you can feel more confident in your final submission meeting the criteria. As a final check, assure you meet all the criteria on the [RUBRIC](#).

Part I - Probability

To get started, let's import our libraries.

```
In [1]: import pandas as pd
import numpy as np
import random
import matplotlib.pyplot as plt
%matplotlib inline
#We are setting the seed to assure you get the same answers on quizzes as we set up
random.seed(42)
```

1. Now, read in the `ab_data.csv` data. Store it in `df`. Use your dataframe to answer the questions in Quiz 1 of the classroom.

a. Read in the dataset and take a look at the top few rows here:

```
In [2]: df = pd.read_csv("ab_data.csv")
        df.head()
```

```
Out[2]:
```

	user_id	timestamp	group	landing_page	converted
0	851104	2017-01-21 22:11:48.556739	control	old_page	0
1	804228	2017-01-12 08:01:45.159739	control	old_page	0
2	661590	2017-01-11 16:55:06.154213	treatment	new_page	0
3	853541	2017-01-08 18:28:03.143765	treatment	new_page	0
4	864975	2017-01-21 01:52:26.210827	control	old_page	1

b. Use the cell below to find the number of rows in the dataset.

```
In [3]: df.describe()
```

```
Out[3]:
```

	user_id	converted
count	294478.000000	294478.000000
mean	787974.124733	0.119659
std	91210.823776	0.324563
min	630000.000000	0.000000
25%	709032.250000	0.000000
50%	787933.500000	0.000000
75%	866911.750000	0.000000
max	945999.000000	1.000000

c. The number of unique users in the dataset.

```
In [4]: df.nunique()
```

```
Out[4]:
```

user_id	290584
timestamp	294478
group	2
landing_page	2
converted	2
dtype:	int64

d. The proportion of users converted.

```
In [5]: df_converted = df[df['converted'] == 1]
        df_converted.describe()
```

```
Out[5]:
```

	user_id	converted
count	35237.000000	35237.0
mean	788394.376962	1.0
std	91398.565565	0.0
min	630001.000000	1.0

25%	709555.000000	1.0
50%	787633.000000	1.0
75%	867831.000000	1.0
max	945991.000000	1.0

```
In [6]: proportion = (35237.0/294478.000000)*100
        proportion
```

```
Out[6]: 11.96591935560551
```

e. The number of times the new_page and treatment don't match.

```
In [7]: df_treat = df[df['group'] == 'treatment']
        df_treatnotnew = df_treat[df_treat['landing_page'] != 'new_page']
        df_treatnotnew.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 1965 entries, 308 to 294252
Data columns (total 5 columns):
user_id      1965 non-null int64
timestamp    1965 non-null object
group        1965 non-null object
landing_page  1965 non-null object
converted     1965 non-null int64
dtypes: int64(2), object(3)
memory usage: 92.1+ KB
```

```
In [8]: df_new = df[df['landing_page'] == 'new_page' ]
        df_newnottreat = df_new[df_new['group'] != 'treatment']
        df_newnottreat.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 1928 entries, 22 to 294331
Data columns (total 5 columns):
user_id      1928 non-null int64
timestamp    1928 non-null object
group        1928 non-null object
landing_page  1928 non-null object
converted     1928 non-null int64
dtypes: int64(2), object(3)
memory usage: 90.4+ KB
```

```
In [9]: df_newnottreat.set_index('user_id', inplace=True)
        df_newnottreat.head()
```

```
Out[9]:
```

	timestamp	group	landing_page	converted
user_id				

767017	2017-01-12 22:58:14.991443	control	new_page	0
733976	2017-01-11 15:11:16.407599	control	new_page	0
808613	2017-01-10 21:44:01.292755	control	new_page	0
637639	2017-01-11 23:09:52.682329	control	new_page	1
793580	2017-01-08 03:25:33.723712	control	new_page	1

```
In [10]: together = 1928 + 1965
         together
```

```
Out[10]: 3893
```

f. Do any of the rows have missing values?

```
In [11]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 294478 entries, 0 to 294477
Data columns (total 5 columns):
user_id      294478 non-null int64
timestamp    294478 non-null object
group        294478 non-null object
landing_page  294478 non-null object
converted     294478 non-null int64
dtypes: int64(2), object(3)
memory usage: 11.2+ MB
```

2. For the rows where **treatment** does not match with **new_page** or **control** does not match with **old_page**, we cannot be sure if this row truly received the new or old page. Use **Quiz 2** in the classroom to figure out how we should handle these rows.

a. Now use the answer to the quiz to create a new dataset that meets the specifications from the quiz. Store your new dataframe in **df2**.

```
In [12]: df2 = df.loc[(df['group']=='treatment') & (df['landing_page']=='new_page') | \
                    (df['group']=='control') & (df['landing_page']=='old_page'),:]
```

```
In [ ]:
```

```
In [13]: # Double Check all of the correct rows were removed - this should be 0
         df2[((df2['group'] == 'treatment') == (df2['landing_page'] == 'new_page')) == False].sh
```

```
Out[13]: 0
```

3. Use **df2** and the cells below to answer questions for **Quiz3** in the classroom.

a. How many unique **user_ids** are in **df2**?

```
In [14]: df2.loc[:, 'user_id'].count()
```

```
Out[14]: 290585
```

b. There is one **user_id** repeated in **df2**. What is it?

```
In [15]: df2.loc[df2['user_id'].duplicated(),['user_id']]
```

```
Out[15]:      user_id
        2893    773192
```

c. What is the row information for the repeat **user_id**?

```
In [16]: df2.loc[df2['user_id'].duplicated(keep=False)]
```

```
Out[16]:      user_id      timestamp      group landing_page  converted
        1899    773192  2017-01-09 05:37:58.781806  treatment    new_page         0
        2893    773192  2017-01-14 02:55:59.590927  treatment    new_page         0
```

d. Remove **one** of the rows with a duplicate **user_id**, but keep your dataframe as **df2**.

```
In [17]: df2 = df2.drop_duplicates(subset=['user_id'], keep='first')
```

```
In [18]: df2.shape
```

```
Out[18]: (290584, 5)
```

4. Use **df2** in the cells below to answer the quiz questions related to **Quiz 4** in the classroom.

a. What is the probability of an individual converting regardless of the page they receive?

```
In [19]: df2.loc[:, 'converted'].mean()
```

```
Out[19]: 0.11959708724499628
```

b. Given that an individual was in the control group, what is the probability they converted?

```
In [20]: df2.loc[df2['group']=='control', 'converted'].mean()
```

```
Out[20]: 0.1203863045004612
```

c. Given that an individual was in the treatment group, what is the probability they converted?

```
In [21]: df2.loc[df2['group']=='treatment', 'converted'].mean()
```

```
Out[21]: 0.11880806551510564
```

d. What is the probability that an individual received the new page?

```
In [22]: df2.loc[df2['landing_page']=='new_page', :].count()[0]/df2.shape[0]
```

```
Out[22]: 0.50006194422266881
```

- e. Consider your results from parts (a) through (d) above, and explain below whether you think there is sufficient evidence to conclude that the new treatment page leads to more conversions.

According to the previous lines, there's 50% percent of chance of any individual receives a new page, and a probability of conversion of control and treatment group of 12.0% and 11.9%, respectively, therefore there's no evidence that the new treatment page leads to more conversions.

Part II - A/B Test

Notice that because of the time stamp associated with each event, you could technically run a hypothesis test continuously as each observation was observed.

However, then the hard question is do you stop as soon as one page is considered significantly better than another or does it need to happen consistently for a certain amount of time? How long do you run to render a decision that neither page is better than another?

These questions are the difficult parts associated with A/B tests in general.

1. For now, consider you need to make the decision just based on all the data provided. If you want to assume that the old page is better unless the new page proves to be definitely better at a Type I error rate of 5%, what should your null and alternative hypotheses be? You can state your hypothesis in terms of words or in terms of p_{old} and p_{new} , which are the converted rates for the old and new pages.

H1: $p_{new} - p_{old} > 0$.

2. Assume under the null hypothesis, p_{new} and p_{old} both have "true" success rates equal to the **converted** success rate regardless of page - that is p_{new} and p_{old} are equal. Furthermore, assume they are equal to the **converted** rate in **ab_data.csv** regardless of the page.

Use a sample size for each page equal to the ones in **ab_data.csv**.

Perform the sampling distribution for the difference in **converted** between the two pages over 10,000 iterations of calculating an estimate from the null.

Use the cells below to provide the necessary parts of this simulation. If this doesn't make complete sense right now, don't worry - you are going to work through the problems below to complete this problem. You can use **Quiz 5** in the classroom to make sure you are on the right track.

- a. What is the **conversion rate** for p_{new} under the null?

```
In [23]: p_new = df2.loc[:, 'converted'].mean()
         p_new
```

```
Out[23]: 0.11959708724499628
```

- b. What is the **conversion rate** for p_{old} under the null?

```
In [24]: p_old = df2.loc[:, 'converted'].mean()
         p_old
```

```
Out[24]: 0.11959708724499628
```

- c. What is n_{new} , the number of individuals in the treatment group?

```
In [25]: n_new = df2.loc[df2['group']=='treatment'].shape[0]
         n_new
```

```
Out[25]: 145310
```

d. What is n_{old} , the number of individuals in the control group?

```
In [26]: n_old = df2.loc[df2['group']=='control'].shape[0]
         n_old
```

```
Out[26]: 145274
```

e. Simulate n_{new} transactions with a conversion rate of p_{new} under the null. Store these n_{new} 1's and 0's in **new_page_converted**.

```
In [27]: new_page_converted = np.random.binomial(1, p_new, n_new)
         new_page_converted.sum()
```

```
Out[27]: 17416
```

f. Simulate n_{old} transactions with a conversion rate of p_{old} under the null. Store these n_{old} 1's and 0's in **old_page_converted**.

```
In [28]: old_page_converted = np.random.binomial(1, p_old, n_old)
         old_page_converted.sum()
```

```
Out[28]: 17609
```

g. Find $p_{new} - p_{old}$ for your simulated values from part (e) and (f).

```
In [29]: new_page_converted.shape, old_page_converted.shape
```

```
Out[29]: ((145310,), (145274,))
```

```
In [30]: new_page_converted = new_page_converted[:145274]
         new_page_converted.mean() - old_page_converted.mean()
```

```
Out[30]: -0.0013560582072497523
```

h. Create 10,000 $p_{new} - p_{old}$ values using the same simulation process you used in parts (a) through (g) above. Store all 10,000 values in a NumPy array called **p_diffs**.

```
In [31]: p_diffs = []
```

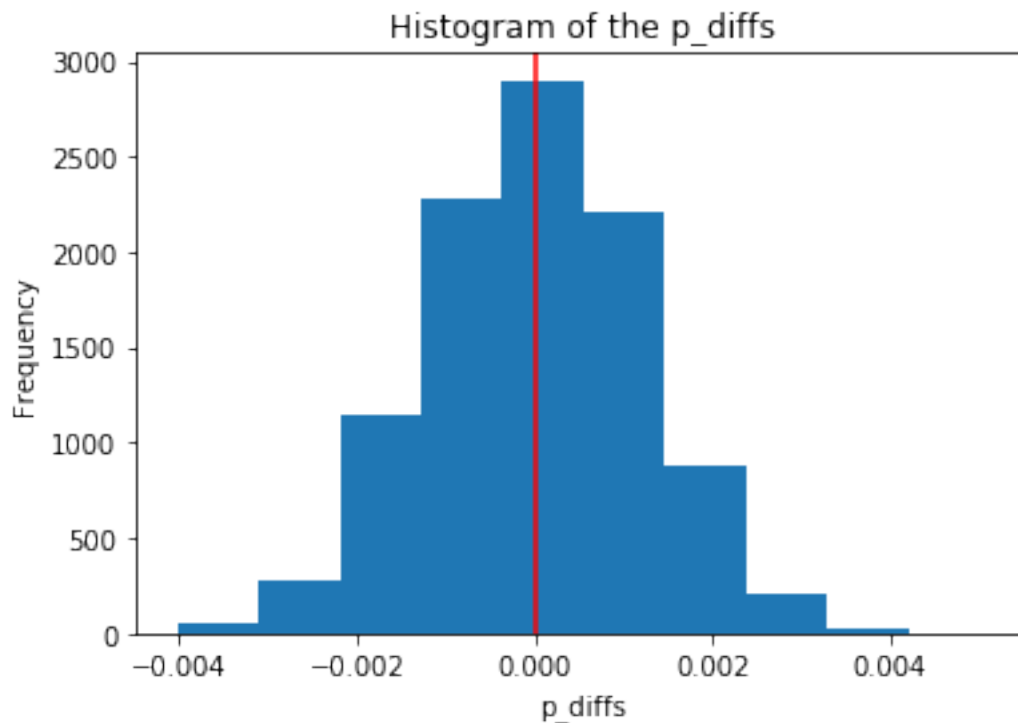
```
for _ in range(10000):
    new_page_converted = np.random.binomial(1, p_new, n_new)
    old_page_converted = np.random.binomial(1, p_old, n_old)
    diffs = new_page_converted.mean() - old_page_converted.mean()
    p_diffs.append(diffs)
```

```
In [32]: p_diffs = np.array(p_diffs)
         p_diffs.mean()
```

```
Out[32]: 5.3491899280244675e-07
```

- i. Plot a histogram of the **p_diffs**. Does this plot look like what you expected? Use the matching problem in the classroom to assure you fully understand what was computed here.

```
In [33]: plt.hist(p_diffs);  
         plt.axvline(np.array(p_diffs).mean(), color='red');  
         plt.xlabel('p_diffs');  
         plt.ylabel('Frequency');  
         plt.title('Histogram of the p_diffs');
```



- j. What proportion of the **p_diffs** are greater than the actual difference observed in **ab_data.csv**?

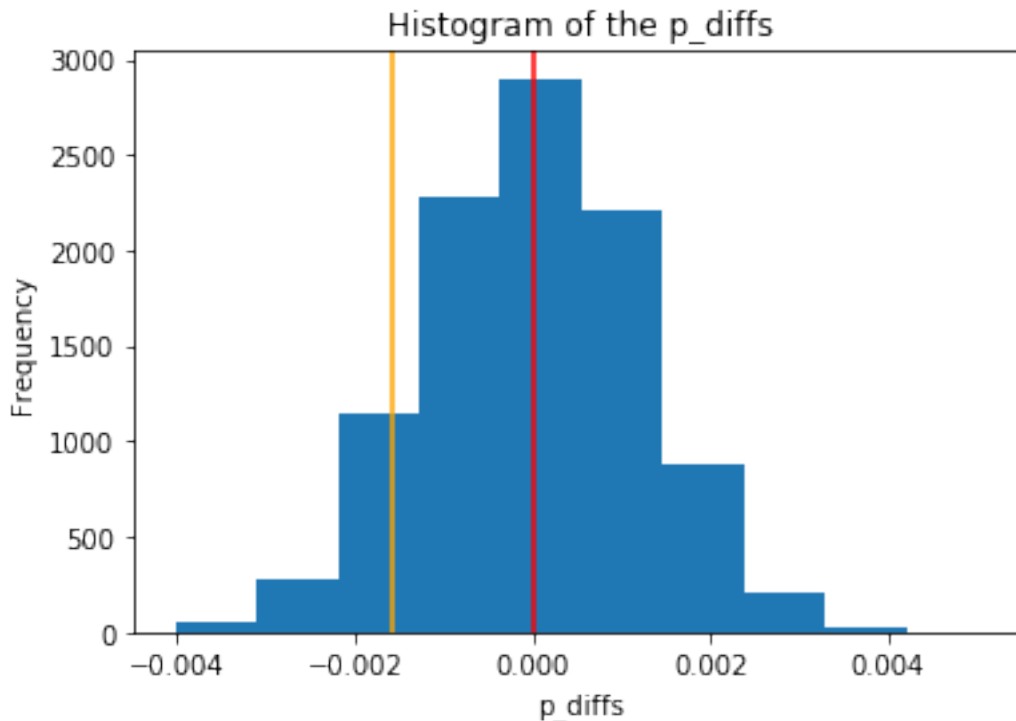
```
In [34]: actual_diff = df2.loc[df2['group']=='treatment', 'converted'].mean()-df2.loc[df2['group']=='control', 'converted'].mean()  
         actual_diff
```

```
Out[34]: -0.0015782389853555567
```

```
In [35]: (actual_diff < p_diffs).mean()
```

```
Out[35]: 0.90090000000000003
```

```
In [36]: plt.hist(p_diffs);  
         plt.axvline(p_diffs.mean(), color='red');  
         plt.axvline(np.array(actual_diff).mean(), color='orange');  
         plt.xlabel('p_diffs');  
         plt.ylabel('Frequency');  
         plt.title('Histogram of the p_diffs');
```

In []:

- k. Please explain using the vocabulary you've learned in this course what you just computed in part j. What is this value called in scientific studies? What does this value mean in terms of whether or not there is a difference between the new and old pages?

In scientific studies, this value is called as P-value, which is the probability of finding the observed, or more extreme results when the null hypothesis (H_0) of a study question is true. In our study, we find a very large P-value, which indicates weak evidence against the null hypothesis so we can conclude that there's no conversion advantage utilizing the new page, in other words, the null hypothesis is true..

- l. We could also use a built-in to achieve similar results. Though using the built-in might be easier to code, the above portions are a walkthrough of the ideas that are critical to correctly thinking about statistical significance. Fill in the below to calculate the number of conversions for each page, as well as the number of individuals who received each page. Let `n_old` and `n_new` refer the the number of rows associated with the old page and new pages, respectively.

In [37]: `import statsmodels.api as sm`

```
df2.loc[(df2['landing_page']=='new_page')&(df2['converted']==1)].shape[0]
```

/opt/conda/lib/python3.6/site-packages/statsmodels/compat/pandas.py:56: FutureWarning: The pandas from pandas.core import datetools

Out[37]: 17264

```
In [38]: convert_old = df2.loc[(df2['landing_page']=='old_page')&(df2['converted']==1)].shape[0]
        convert_new = df2.loc[(df2['landing_page']=='new_page')&(df2['converted']==1)].shape[0]
        n_old = df2.loc[df2['group']=='control'].shape[0]
        n_new = df2.loc[df2['group']=='treatment'].shape[0]
```

- m. Now use `stats.proportions_ztest` to compute your test statistic and p-value. [Here](#) is a helpful link on using the built in.

```
In [39]: z_score, p_value = sm.stats.proportions_ztest([convert_old, convert_new], [n_old, n_new])
        print('z_score: {}'.format(z_score))
        print('p_value: {}'.format(p_value))
```

```
z_score: 1.3109241984234394
p_value: 0.9050583127590245
```

- n. What do the z-score and p-value you computed in the previous question mean for the conversion rates of the old and new pages? Do they agree with the findings in parts j. and k.?

```
In [49]: import scipy.stats
        # How significant our z-score is?
        cdf_z_score = norm.cdf(z_score)
        print('z_score: {}'.format(cdf_z_score))

        # Define the critical value at the confidence interval of single-sided test
        conf_int = 0.95
        z_crit_value = norm.ppf(conf_int)
        print('z_crit_value: {}'.format(z_crit_value))
```

```
z_score: 0.9050583127590245
z_crit_value: 1.6448536269514722
```

The z-critical value is greater than z-score, it indicates that we fail to reject the null hypothesis, which suggest the new page conversion rate is higher than the old page. These values agree with the findings in parts j. and k..

Part III - A regression approach

1. In this final part, you will see that the result you achieved in the A/B test in Part II above can also be achieved by performing regression.

- a. Since each row is either a conversion or no conversion, what type of regression should you be performing in this case?

Logistic Regression.

- b. The goal is to use **statsmodels** to fit the regression model you specified in part a. to see if there is a significant difference in conversion based on which page a customer receives. However, you first need to create in `df2` a column for the intercept, and create a dummy variable column for which page each user received. Add an **intercept** column, as well as an **ab_page** column, which is 1 when an individual receives the **treatment** and 0 if **control**.

```
In [50]: df2.loc[:, 'intercept'] = 1
         df2[['control', 'treatment']] = pd.get_dummies(df2.loc[:, ['group']])
         df2['ab_page'] = df2['treatment']
         df2 = df2.drop(columns=['treatment', 'control'])
```

```
In [51]: df2.head()
```

```
Out[51]:
```

	user_id	timestamp	group	landing_page	converted	\
0	851104	2017-01-21 22:11:48.556739	control	old_page	0	
1	804228	2017-01-12 08:01:45.159739	control	old_page	0	
2	661590	2017-01-11 16:55:06.154213	treatment	new_page	0	
3	853541	2017-01-08 18:28:03.143765	treatment	new_page	0	
4	864975	2017-01-21 01:52:26.210827	control	old_page	1	

	intercept	ab_page
0	1	0
1	1	0
2	1	1
3	1	1
4	1	0

- c. Use **statsmodels** to instantiate your regression model on the two columns you created in part b., then fit the model using the two columns you created in part b. to predict whether or not an individual converts.

```
In [52]: logistic = sm.Logit(df2['converted'], df2[['intercept', 'ab_page']])
         results = logistic.fit()
```

```
Optimization terminated successfully.
Current function value: 0.366118
Iterations 6
```

- d. Provide the summary of your model below, and use it as necessary to answer the following questions.

```
In [55]: print(results.summary2())
```

```
Results: Logit
=====
Model:                Logit                No. Iterations:    6.0000
Dependent Variable:    converted            Pseudo R-squared:    0.000
Date:                 2020-07-31 20:11 AIC:                212780.3502
```

```

No. Observations:  290584          BIC:          212801.5095
Df Model:          1              Log-Likelihood: -1.0639e+05
Df Residuals:      290582          LL-Null:       -1.0639e+05
Converged:         1.0000          Scale:         1.0000
-----
                Coef.   Std.Err.    z      P>|z|    [0.025   0.975]
-----
intercept    -1.9888    0.0081  -246.6690  0.0000   -2.0046   -1.9730
ab_page      -0.0150    0.0114   -1.3109  0.1899   -0.0374    0.0074
=====

```

- e. What is the p-value associated with **ab_page**? Why does it differ from the value you found in **Part II**? **Hint:** What are the null and alternative hypotheses associated with your regression model, and how do they compare to the null and alternative hypotheses in **Part II**?

The P-value associated with **ab_page** is 0.190. Its value differ because of the methods utilized, on the first we calculated P-value with one-sided test and this logistic regression approach consider the two-sided test. In this case, we don't test anymore for not greater than or equal, instead we test for not equal in our hypothesis, as following: H1: pnew unknown character. pold!=0

- f. Now, you are considering other things that might influence whether or not an individual converts. Discuss why it is a good idea to consider other factors to add into your regression model. Are there any disadvantages to adding additional terms into your regression model?

By adding other factors to our regression model, we may be leading to unreliable and unstable estimates of regression coefficients (multicollinearity) in our model, which may affect our predictions. Every time we include a new predictor variable we need to take the necessary steps in order to assure the reliability of the model/prediction.

- g. Now along with testing if the conversion rate changes for different pages, also add an effect based on which country a user lives in. You will need to read in the **countries.csv** dataset and merge together your datasets on the appropriate rows. [Here](#) are the docs for joining tables.

Does it appear that country had an impact on conversion? Don't forget to create dummy variables for these country columns - **Hint: You will need two columns for the three dummy variables.** Provide the statistical output as well as a written response to answer this question.

```

In [58]: countries_df = pd.read_csv('countries.csv')
         df_new = countries_df.set_index('user_id').join(df2.set_index('user_id'), how='inner')
         df_new.head()

```

```

Out[58]:
   country  timestamp  group landing_page \
user_id
834778    UK  2017-01-14 23:08:43.304998  control    old_page
928468    US  2017-01-23 14:44:16.387854  treatment    new_page

```

822059	UK	2017-01-16 14:04:14.719771	treatment	new_page
711597	UK	2017-01-22 03:14:24.763511	control	old_page
710616	UK	2017-01-16 13:14:44.000513	treatment	new_page

	converted	intercept	ab_page
user_id			
834778	0	1	0
928468	0	1	1
822059	1	1	1
711597	0	1	0
710616	0	1	1

```
In [59]: df_new['country'].value_counts()
```

```
Out[59]: US    203619
UK       72466
CA       14499
Name: country, dtype: int64
```

```
In [60]: df_new[['CA', 'UK', 'US']] = pd.get_dummies(df_new['country'])
df_new = df_new.drop(columns=['US'])
df_new.head()
```

```
Out[60]:
```

	country	timestamp	group	landing_page	\
user_id					
834778	UK	2017-01-14 23:08:43.304998	control	old_page	
928468	US	2017-01-23 14:44:16.387854	treatment	new_page	
822059	UK	2017-01-16 14:04:14.719771	treatment	new_page	
711597	UK	2017-01-22 03:14:24.763511	control	old_page	
710616	UK	2017-01-16 13:14:44.000513	treatment	new_page	

	converted	intercept	ab_page	CA	UK
user_id					
834778	0	1	0	0	1
928468	0	1	1	0	0
822059	1	1	1	0	1
711597	0	1	0	0	1
710616	0	1	1	0	1

- h. Though you have now looked at the individual factors of country and page on conversion, we would now like to look at an interaction between page and country to see if there significant effects on conversion. Create the necessary additional columns, and fit the new model.

Provide the summary results, and your conclusions based on the results.

```
In [66]: logistic_2 = sm.Logit(df_new['converted'], df_new[['intercept', 'ab_page', 'CA', 'UK']])
results_2 = logistic_2.fit()
print(results_2.summary2())
```

Optimization terminated successfully.
 Current function value: 0.366113
 Iterations 6

Results: Logit

```
=====
Model:                Logit                No. Iterations:    6.0000
Dependent Variable: converted                Pseudo R-squared: 0.000
Date:                2020-07-31 20:33 AIC:                212781.1253
No. Observations:    290584                BIC:                212823.4439
Df Model:            3                    Log-Likelihood:    -1.0639e+05
Df Residuals:        290580                LL-Null:            -1.0639e+05
Converged:            1.0000                Scale:            1.0000
-----
```

	Coef.	Std.Err.	z	P> z	[0.025	0.975]
intercept	-1.9893	0.0089	-223.7628	0.0000	-2.0067	-1.9718
ab_page	-0.0149	0.0114	-1.3069	0.1912	-0.0374	0.0075
CA	-0.0408	0.0269	-1.5161	0.1295	-0.0934	0.0119
UK	0.0099	0.0133	0.7433	0.4573	-0.0162	0.0359

```
=====
```

0.2.1 The Conclusions

the performance of the old page was better. then we can accept the null hypothesis and reject the alternate hypothesis.

Finishing Up

Congratulations! You have reached the end of the A/B Test Results project! You should be very proud of all you have accomplished!

Tip: Once you are satisfied with your work here, check over your report to make sure that it satisfies all the areas of the rubric (found on the project submission page at the end of the lesson). You should also probably remove all of the "Tips" like this one so that the presentation is as polished as possible.

0.3 Directions to Submit

Before you submit your project, you need to create a .html or .pdf version of this notebook in the workspace here. To do that, run the code cell below. If it worked correctly, you should get a return code of 0, and you should see the generated .html file in the workspace directory (click on the orange Jupyter icon in the upper left).

Alternatively, you can download this report as .html via the **File > Download as** sub-menu, and then manually upload it into the workspace directory by clicking on the orange Jupyter icon in the upper left, then using the Upload button.

Once you've done this, you can submit your project by clicking on the "Submit Project" button in the lower right here. This will create and submit a zip file with this .ipynb doc and the .html or .pdf version you created. Congratulations!

```
In [ ]: from subprocess import call
        call(['python', '-m', 'nbconvert', 'Analyze_ab_test_results_notebook.ipynb'])
```