

Analyze_ab_test_results_notebook

December 28, 2019

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.

```
[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:

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

```
[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

```
[3]: df.drop_duplicates(inplace = True)
```

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

```
[4]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
Int64Index: 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: 13.5+ MB
```

294478 = Number of Rows

c. The number of unique users in the dataset.

```
[5]: df.nunique()
```

```
[5]: user_id          290584  
timestamp          294478  
group              2  
landing_page       2  
converted          2  
dtype: int64
```

Number of Unique Users = 290584

d. The proportion of users converted.

```
[6]: df.query('converted == 1').count()
```

```
[6]: user_id      35237
      timestamp   35237
      group       35237
      landing_page 35237
      converted    35237
      dtype: int64
```

```
[7]: 35237/290584
```

```
[7]: 0.12126269856564711
```

0.12126(approx) = Proportion of Users Converted

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

```
[8]: df.query('landing_page == "new_page" and group != "treatment" ').count() + df.
      ↪query('landing_page != "new_page" and group == "treatment" ').count()
```

```
[8]: user_id      3893
      timestamp   3893
      group       3893
      landing_page 3893
      converted    3893
      dtype: int64
```

3893

f. Do any of the rows have missing values?

```
[9]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 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: 13.5+ MB
```

There are no missing values.

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**.

```
[10]: op = df['landing_page'] == 'old_page'

np = df['landing_page'] == 'new_page'

ct = df['group'] == 'control'

t = df['group'] == 'treatment'

df_ = (df[(t) & (op)] + df[(ct) & (np)]).index

df2 = df.drop(df_)
```

```
[11]: df2.head()
```

```
[11]:
```

	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

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

```
[12]: 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**?

```
[13]: df2['user_id'].nunique()
```

```
[13]: 290584
```

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

```
[14]: df2.loc[df2.duplicated('user_id')]
```

```
[14]:
```

	user_id	timestamp	group	landing_page	converted
2893	773192	2017-01-14 02:55:59.590927	treatment	new_page	0

773192

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

```
[15]: df2[df2.duplicated(['user_id'], keep=False)]
```

```
[15]:      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**.

```
[16]: df2.drop_duplicates(subset='user_id', keep="first", inplace=True)
```

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?

```
[17]: (df2['converted'] == 1).mean()
```

```
[17]: 0.11959708724499628
```

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

```
[18]: df2.query('group == "control" and converted == 1').count()/df2.query('group ==_
↪"control").count()
```

```
[18]: user_id      0.120386
      timestamp    0.120386
      group        0.120386
      landing_page  0.120386
      converted     0.120386
      dtype: float64
```

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

```
[19]: df2.query('group == "treatment" and converted == 1').count()/df2.query('group_
↪== "treatment").count()
```

```
[19]: user_id      0.118808
      timestamp    0.118808
      group        0.118808
      landing_page  0.118808
      converted     0.118808
      dtype: float64
```

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

```
[20]: len(df2.query("landing_page == 'new_page')) / df2.shape[0]
```

```
[20]: 0.5000619442226688
```

```
0.5001
```

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.

Since, neither the overall conversion rate nor the conversion rate of the people in treatment group was high, also 50% people already had new page, we can say that **there is not much evidence that treatment page lead to significant conversions.**

```
[21]: df.head(10)
```

```
[21]:
```

	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
5	936923	2017-01-10 15:20:49.083499	control	old_page	0
6	679687	2017-01-19 03:26:46.940749	treatment	new_page	1
7	719014	2017-01-17 01:48:29.539573	control	old_page	0
8	817355	2017-01-04 17:58:08.979471	treatment	new_page	1
9	839785	2017-01-15 18:11:06.610965	treatment	new_page	1

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.

$H_0 - P_{old} \geq P_{new}$. . $H_1 - P_{old} < P_{new}$

Put your answer here.

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?

```
[22]: p_new = df2['converted'].mean()  
p_new
```

```
[22]: 0.11959708724499628
```

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

```
[23]: p_old = df2['converted'].mean()  
p_old
```

```
[23]: 0.11959708724499628
```

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

```
[24]: n_new = df2[df2['group'] == 'treatment'].shape[0]  
n_new
```

```
[24]: 145310
```

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

```
[25]: n_old = df2[df2['group'] == 'control'].shape[0]  
n_old
```

```
[25]: 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**.

```
[26]: import numpy as np  
new_page_converted = np.random.binomial(n_new, p_new)
```

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**.

```
[27]: old_page_converted = np.random.binomial(n_old, p_old)
```

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

```
[28]: new_page_converted/n_new - old_page_converted/n_old
```

```
[28]: -0.0003804852409692444
```

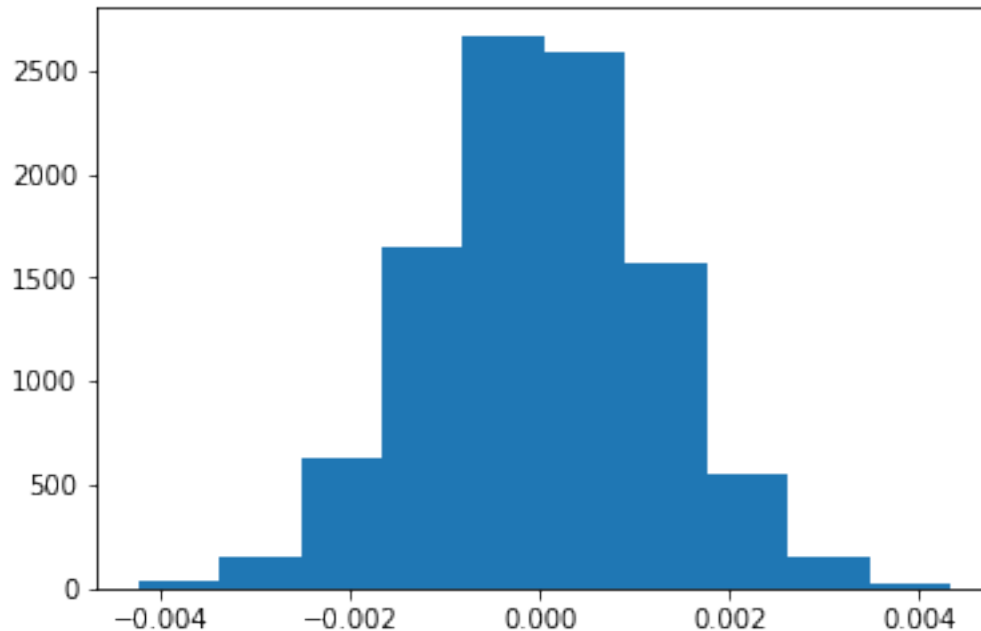
h. Create 10,000 $p_{new} - p_{old}$ values using the same simulation process you used in parts (a) through (g) above. Store all 10,000 values in a NumPy array called **p_diffs**.

```
[29]: p_diffs = []  
for _ in range(10000):  
    new_page_converted = np.random.binomial(n_new, p_new)
```

```
old_page_converted = np.random.binomial(n_old, p_old)
diff = new_page_converted/n_new - old_page_converted/n_old
p_diffs.append(diff)
```

- 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.

```
[30]: plt.hist(p_diffs);
```



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

```
[31]: actual_diff = df2[df2['group'] == 'treatment']['converted'].mean() -
↳ df2[df2['group'] == 'control']['converted'].mean()
p_diffs = np.array(p_diffs)
(actual_diff < p_diffs).mean()
```

```
[31]: 0.90569999999999995
```

- 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?

The Value obtained in part j is **P-value** . P - value value is the probability of obtaining the observed statistic or one more extreme in favour of alternative hypothesis (H1) , given that the null hypothesis (Ho) is true. Since P-value is large, we will Stay with the Ho that old page still is recieving more traffic than new page.

- 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 to the number of rows associated with the old page and new pages, respectively.

```
[32]: import statsmodels.api as sm

convert_old = df2.query(" landing_page == 'old_page' and converted == 1").
    ↪shape[0]
convert_new = df2.query(" landing_page == 'new_page' and converted == 1").
    ↪shape[0]
n_old = df2[df2['group'] == 'control'].shape[0]
n_new = df2[df2['group'] == 'treatment'].shape[0]
convert_old
```

```
/opt/conda/lib/python3.6/site-packages/statsmodels/compat/pandas.py:56:
FutureWarning: The pandas.core.datetools module is deprecated and will be
removed in a future version. Please use the pandas.tseries module instead.
    from pandas.core import datetools
```

[32]: 17489

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

```
[33]: z_score, p_value = sm.stats.proportions_ztest([convert_old, convert_new],
    ↪[n_old, n_new], alternative='smaller')
print(z_score, p_value)
```

1.31092419842 0.905058312759

z-score = 1.311, p-value = 0.905

- 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.**?

```
[34]: from scipy.stats import norm
# significance of z-score
print(norm.cdf(z_score))

# for single-sides test, assumed at 95% confidence level, we calculate:
print(norm.ppf(1-(0.05)))
```

0.905058312759

1.64485362695

Ans. Since z-score of 1.311 is less than the critical value of 1.644, we **fail to reject the null hypothesis(H_0)**. This states that conversion rate of old page is more than new page. **Yes** it

agrees with part 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

ANS. We will be performing logistics regression in this case because there are two outcomes for the response variable.

- 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**.

```
[35]: df2['intercept']=1
df2[['control', 'ab_page']]=pd.get_dummies(df2['group'])
df2.drop(labels=['control'], axis=1, inplace=True)
df2.head()
```

```
[35]:
```

	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.

```
[36]: # Now we will run our logistic regression model.
import statsmodels.api as sm

logit_mod = sm.Logit(df2['converted'], df2[['intercept', 'ab_page']])
results = logit_mod.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.

```
[37]: results.summary()
```

```

[37]: <class 'statsmodels.iolib.summary.Summary'>
      """
                                Logit Regression Results
      =====
Dep. Variable:                converted    No. Observations:                290584
Model:                        Logit       Df Residuals:                290582
Method:                       MLE        Df Model:                  1
Date:                         Sat, 23 Nov 2019    Pseudo R-squ.:                8.077e-06
Time:                         20:12:04    Log-Likelihood:               -1.0639e+05
converged:                    True         LL-Null:                   -1.0639e+05
                                      LLR p-value:                0.1899
      =====
                                coef      std err          z      P>|z|      [0.025      0.975]
      -----
intercept          -1.9888         0.008   -246.669      0.000      -2.005      -1.973
ab_page            -0.0150         0.011    -1.311      0.190      -0.037       0.007
      =====
      """

```

- 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**?

ANS. P-Value associated with ab page is **0.190**.

It differs from that in part II because **here the test is a two sided test with $H_0 : p_{\text{new}} = p_{\text{old}}$** (where p_{new} and p_{old} are the conversion rates of new and old pages respectively . In part II the test was a one sided test with $H_0 : p_{\text{new}} \leq p_{\text{old}}$.

- 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?

ANS. Yes it is a good idea to consider other aspects that might have influenced the decision to convert from old to new page. The factors related to the phenomenons on **Change Aversion** and **Novelty Effect** can have a great impact on the results. Some of the considerable factors are as follows :- 1. If there are more older people in the dataset , we know that old people are change averse and they would be biased towards chosing the old page while younger people would be baised to chose the new page as they welcome change contributing to the novelty effect. 2. Also factors like people belonging to the more mordernised areas or rural and slum areas may also effect the decision as modern people like to have greater exploration and visual enhancement whi;le people

from backward areas prefer simplicity.

Disadvantages with adding additional factors to the regression model :- 1. Adding additional factors to the model might lead to a **problem of multicollinearity**. It happens when there is high correlation among the explanatory variables leading to unreliable regression coefficients. 2. Also accurate data regarding the factors that we want to consider might not be available which would make our regression study unreliable.

- 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.

```
[38]: countries_df = pd.read_csv('./countries.csv')
df_new = countries_df.set_index('user_id').join(df2.set_index('user_id'),
→how='inner')
```

```
[39]: #Create dummy variables
df_new['intercept'] = 1
df_new[['CA', 'US']] = pd.get_dummies(df_new['country'])[['CA', 'US']]
```

- 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.

```
[40]: mod = sm.Logit(df_new['converted'], df_new[['intercept', 'US', 'CA']])

results = mod.fit()

results.summary()
```

Optimization terminated successfully.

Current function value: 0.366116

Iterations 6

```
[40]: <class 'statsmodels.iolib.summary.Summary'>
      """
```

```

                                Logit Regression Results
=====
Dep. Variable:                  converted    No. Observations:                  290584
Model:                            Logit      Df Residuals:                  290581
Method:                           MLE        Df Model:                        2
Date:                Sat, 23 Nov 2019    Pseudo R-squ.:                  1.521e-05
Time:                20:12:05      Log-Likelihood:                  -1.0639e+05
```

```

converged:                True    LL-Null:                -1.0639e+05
                               LLR p-value:                0.1984
=====
               coef      std err          z      P>|z|      [0.025      0.975]
-----
intercept    -1.9868      0.011    -174.174      0.000     -2.009     -1.964
US           -0.0099      0.013     -0.746      0.456     -0.036      0.016
CA           -0.0507      0.028     -1.786      0.074     -0.106      0.005
=====
"""

```

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** submenu, 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!

```

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

```

[41]: 0