

Analyze_ab_test_results_notebook

October 30, 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.

```
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.shape[0]
```

```
Out[3]: 294478
```

c. The number of unique users in the dataset.

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

```
Out[4]: 290584
```

d. The proportion of users converted.

```
In [5]: df.converted.mean()
```

```
Out[5]: 0.11965919355605512
```

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

```
In [6]: treat_old = df[(df.group == 'treatment') & (df.landing_page == 'old_page')]
        treat_old.shape[0]
```

```
Out[6]: 1965
```

```
In [7]: ctl_new = df[(df.group == 'control') & (df.landing_page == 'new_page')]
        ctl_new.shape[0]
```

```
Out[7]: 1928
```

```
In [8]: treat_old.shape[0] + ctl_new.shape[0]
```

```
Out[8]: 3893
```

f. Do any of the rows have missing values?

```
In [9]: df.isnull().sum()
```

```
Out[9]: user_id      0
        timestamp    0
        group        0
        landing_page  0
        converted     0
        dtype: int64
```

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 [10]: remove = ctl_new.append(treat_old).index
        remove
```

```
Out[10]: Int64Index([    22,    240,    490,    846,    850,    988,   1198,   1354,
                   1474,   1877,
                   ...,
                   293240, 293302, 293391, 293443, 293530, 293773, 293817, 293917,
                   294014, 294252],
                  dtype='int64', length=3893)
```

```
In [11]: df2 = df.drop(remove)
        df2.head()
```

```
Out[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
```

```
In [12]: # 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[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**?

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

```
Out[13]: 290584
```

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

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

```
Out[14]:
```

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

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

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

```
Out[15]:
```

	user_id	timestamp	group	landing_page	converted
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 [16]: df2[df2.user_id == 773192]
```

```
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

```
In [17]: df2.drop(2893, inplace=True)
```

```
In [18]: df2[df2.user_id == 773192]
```

```
Out[18]:
```

	user_id	timestamp	group	landing_page	converted
1899	773192	2017-01-09 05:37:58.781806	treatment	new_page	0

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.converted.mean()
```

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

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

```
In [20]: control_convert=df2.query('group == "control").converted.mean()
control_convert
```

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

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

```
In [21]: treat_convert =df2.query('group == "treatment").converted.mean()
treat_convert
```

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

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

```
In [22]: p_new_page= df2.query('landing_page == "new_page").count()[0] / df2.shape[0]
p_new_page
```

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.

The results in the previous two portions are very close and there is no evidence in my opinion that one page will lead 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.

Put your answer here.

$$H_0 : p_{new} - p_{old} \leq 0$$

$$H_1 : p_{new} - p_{old} > 0$$

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]: convert_mean = df2.converted.mean()
        convert_mean
```

Out [23]: 0.11959708724499628

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

```
In [24]: n_new = df2.query('landing_page == "new_page"').shape[0]
        n_new
```

Out [24]: 145310

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

```
In [25]: n_old = df2.query('landing_page == "old_page").shape[0]
        n_old
```

```
Out[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**.

```
In [26]: new_page_converted = np.random.choice([0, 1], size=n_new, p=[(1 - convert_mean), convert_mean])
```

- 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 [27]: old_page_converted = np.random.choice([0, 1], size=n_old, p=[(1 - convert_mean), convert_mean])
```

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

```
In [28]: new_page_converted.mean() - old_page_converted.mean()
```

```
Out[28]: -0.00038057050956437355
```

- 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 [29]: p_diffs = []

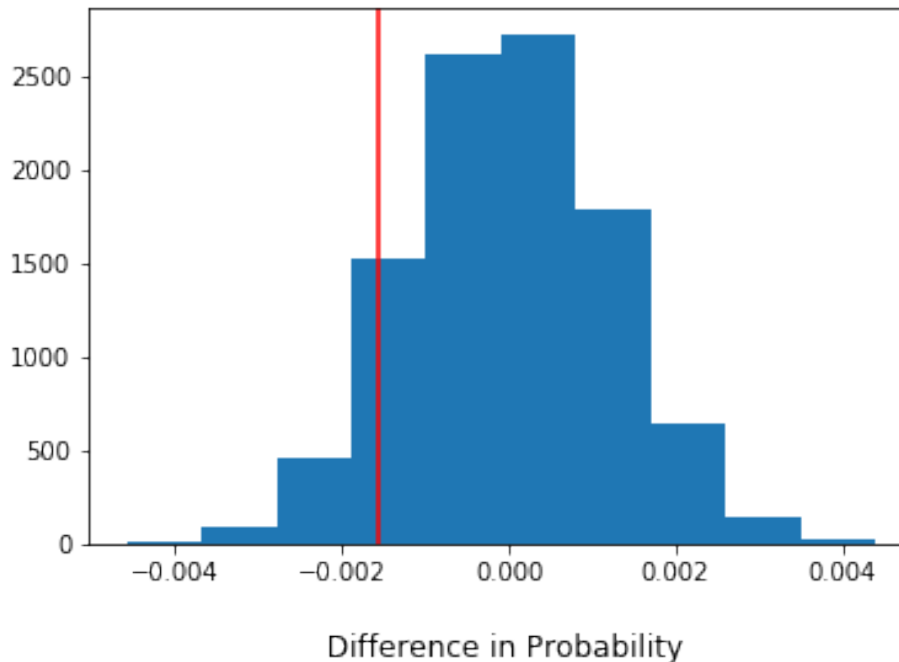
        for i in range(10000):
            new_page_converted = np.random.choice([0, 1], size=n_new, p=[(1 - convert_mean), convert_mean])
            old_page_converted = np.random.choice([0, 1], size=n_old, p=[(1 - convert_mean), convert_mean])
            p_diffs.append(new_page_converted.mean() - old_page_converted.mean())
```

```
In [30]: p_diffs = np.asarray(p_diffs)
```

- 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 [31]: plt.hist(p_diffs)
        plt.title("Simulated Differences in Conversion Rates for Null Hypothesis \n", fontsize=12)
        plt.xlabel("\n Difference in Probability", fontsize=12)
        plt.axvline(treat_convert - control_convert, color='r');
```

Simulated Differences in Conversion Rates for Null Hypothesis



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

```
In [32]: obs_diff = treat_convert - control_convert  
  
         (p_diffs > obs_diff).mean()
```

```
Out [32]: 0.9052999999999999
```

- 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 p-value calculated is 0.9065. This is far greater than the typical α level of 0.05 in business studies. As such, we would fail to reject the null and conclude that there is not sufficient evidence to say that there is a difference between the two values.

- 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.

```
In [33]: import statsmodels.api as sm
```

```
convert_old = df2.query('landing_page == "old_page" & converted == "1").count()[0]
convert_new = df2.query('landing_page == "new_page" & converted == "1").count()[0]
```

```
n_old = df2.query('landing_page == "old_page").count()[0]
n_new = df2.query('landing_page == "new_page").count()[0]
```

```
convert_old, convert_new, n_old, n_new
```

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

```
Out[33]: (17489, 17264, 145274, 145310)
```

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 [34]: z_score, p_value= sm.stats.proportions_ztest([convert_new, convert_old], [n_new, n_old])
z_score, p_value
```

```
Out[34]: (-1.3109241984234394, 0.90505831275902449)
```

```
In [35]: from scipy.stats import norm
```

```
print(norm.cdf(z_score))
```

```
print(norm.ppf(1-(0.05)))
```

```
0.094941687241
```

```
1.64485362695
```

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.?

Since the z-score of 1.31 less than the critical value of 1.64485362695, we fail to reject the null hypothesis which suggest the new page conversion rate is higher than the old rate. Since they are different, And Yes I Agree with 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 [36]: df2[['ab_page', 'old_page']] = pd.get_dummies(df2['landing_page'])
df2['intercept'] = 1
df2.head()
```

```
Out[36]:
```

	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

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

- 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 [37]: log_mod = sm.Logit(df2['converted'], df2[['intercept', 'ab_page']])
result = log_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.

```
In [38]: from scipy import stats
stats.chisqprob = lambda chisq, df: stats.chi2.sf(chisq, df)

result.summary()
```

```
Out[38]: <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:                 Wed, 30 Oct 2019    Pseudo R-squ.:          8.077e-06
Time:                 11:22:49    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**?

The p-value (0.190) here remains above an α level of 0.05 but is different because this is a two tailed test. We will still reject the null in this situation.

- 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?

it is important to take into considerations the factors that might affect the conversion rate for any given case. Having additional terms to our model is great as long as they are relevant to the case. More terms can provide more insights and increase/decrease our confidence when either rejecting the null hypothesis or the failure of rejecting the null hypothesis

- 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 [39]: 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[39]:
   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	ab_page	old_page	intercept
user_id					
834778		0	0	1	1
928468		0	1	0	1
822059		1	1	0	1
711597		0	0	1	1
710616		0	1	0	1

```
In [40]: df_new['country'].unique()
```

```
Out[40]: array(['UK', 'US', 'CA'], dtype=object)
```

```
In [41]: df_new[['CA', 'UK', 'US']] = pd.get_dummies(df_new['country'])
df_new.head()
```

```
Out[41]:
```

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

```
In [42]: log_mod = sm.Logit(df_new['converted'], df_new[['intercept', 'CA', 'UK']])
result = log_mod.fit()
result.summary()
```

```
Optimization terminated successfully.
Current function value: 0.366116
Iterations 6
```

```
Out[42]: <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:                Wed, 30 Oct 2019    Pseudo R-squ.:          1.521e-05
Time:                11:22:51           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    -1.9967      0.007    -292.314      0.000     -2.010     -1.983
CA           -0.0408      0.027     -1.518      0.129     -0.093      0.012
UK            0.0099      0.013      0.746      0.456     -0.016      0.036
=====
"""

```

- 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.

Pages column is already included as per exercise in part b); hence, model may be made similar to previous part while including pages column.

```
In [43]: log_mod = sm.Logit(df_new['converted'], df_new[['intercept', 'UK', 'US', 'ab_page']])
```

```
In [44]: results = log_mod.fit()
         results.summary()
```

```

Optimization terminated successfully.
Current function value: 0.366113
Iterations 6

```

```

Out[44]: <class 'statsmodels.iolib.summary.Summary'>
"""
                Logit Regression Results
=====
Dep. Variable:          converted    No. Observations:          290584
Model:                Logit        Df Residuals:              290580
Method:                MLE         Df Model:                  3
Date:                Wed, 30 Oct 2019    Pseudo R-squ.:          2.323e-05
Time:                11:22:52           Log-Likelihood:         -1.0639e+05
converged:            True              LL-Null:                -1.0639e+05
                                      LLR p-value:                0.1760
=====
              coef      std err          z      P>|z|      [0.025      0.975]
-----
intercept    -2.0300      0.027    -76.249      0.000     -2.082     -1.978
UK            0.0506      0.028      1.784      0.074     -0.005      0.106
US            0.0408      0.027      1.516      0.130     -0.012      0.093
ab_page      -0.0149      0.011     -1.307      0.191     -0.037      0.007

```

```
=====
"""
```

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

```
Out[45]: 0
```