

Analyze_ab_test_results_notebook

July 24, 2022

0.1 Analyze A/B Test Results

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

- Introduction
- Part I - Probability
- Part II - A/B Test
- Part III - Regression

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(3)
```

```
[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
```

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

```
[3]: len(df)
```

```
[3]: 294478
```

c. The number of unique users in the dataset.

```
[4]: unique_users = df['user_id'].drop_duplicates()
len(unique_users)
```

```
[4]: 290584
```

d. The proportion of users converted.

```
[5]: converted = df.where(df['converted'] == 1).dropna()
len(converted)/len(df['converted'])*100
```

```
[5]: 11.96591935560551
```

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

```
[6]: controldf = df.drop(axis=1, columns=['user_id', 'timestamp', 'converted'])

wrong1 = len(controldf.where((controldf['group'] == 'control') &
    ↪ (controldf['landing_page'] == 'new_page')).dropna())
wrong2 = len(controldf.where((controldf['group'] == 'treatment') &
    ↪ (controldf['landing_page'] != 'new_page')).dropna())
wrong1 + wrong2
```

```
[6]: 3893
```

f. Do any of the rows have missing values?

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

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

2. For the rows where **treatment** is not aligned with **new_page** or **control** is not aligned with **old_page**, we cannot be sure if this row truly received the new or old page. Use **Quiz 2** in the classroom to provide 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**.

```
[8]: df01 = df.where((df['group'] == 'control') & (df['landing_page'] == 'old_page')).dropna()
      df02 = df.where((df['group'] == 'treatment') & (df['landing_page'] == 'new_page')).dropna()
      df2 = df01.merge(df02,how='outer',sort=False)
```

```
[9]: # 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]
```

```
[9]: 0
```

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

```
[10]:   user_id          timestamp   group landing_page  converted
0  851104.0  2017-01-21 22:11:48.556739  control      old_page        0.0
1  804228.0  2017-01-12 08:01:45.159739  control      old_page        0.0
2  864975.0  2017-01-21 01:52:26.210827  control      old_page        1.0
3  936923.0  2017-01-10 15:20:49.083499  control      old_page        0.0
4  719014.0  2017-01-17 01:48:29.539573  control      old_page        0.0
```

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

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

```
[11]: uniq_users = len(df2['user_id'].drop_duplicates())
      uniq_users
```

```
[11]: 290584
```

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

```
[12]: df2.where(df2['user_id'].duplicated() == True).dropna()
```

```
[12]:      user_id      timestamp      group landing_page \
146678  773192.0  2017-01-14 02:55:59.590927  treatment    new_page

      converted
146678      0.0
```

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

```
[13]: df2.where(df2['user_id'].duplicated() == True).dropna()
```

```
[13]:      user_id      timestamp      group landing_page \
146678  773192.0  2017-01-14 02:55:59.590927  treatment    new_page

      converted
146678      0.0
```

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

```
[14]: df2.drop_duplicates(subset=['user_id'],inplace=True)
```

4. Use `df2` in the below cells 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?

```
[15]: converted = df2.where(df2['converted'] == 1).dropna()
Pg = len(converted)/len(df2['converted'])
Pg
```

```
[15]: 0.11959708724499628
```

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

```
[16]: controlConv = df2.where((df2['converted'] == 1) & (df2['group'] == 'control')).
      ↪dropna()
controlall = df2.where(df2['group'] == 'control').dropna()
Pcontrol = len(controlConv)/len(controlall)
Pcontrol
```

```
[16]: 0.1203863045004612
```

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

```
[17]: controlConv = df2.where((df2['converted'] == 1) & (df2['group'] ==
      ↪'treatment')).dropna()
controlall = df2.where(df2['group'] == 'treatment').dropna()
Ptreatment = len(controlConv)/len(controlall)
```

```
Ptreatment
```

```
[17]: 0.11880806551510564
```

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

```
[18]: newpage = df2.where(df2['landing_page'] == 'new_page').dropna()  
P_of_new_page = len(newpage)/len(df2['landing_page'])  
P_of_new_page
```

```
[18]: 0.5000619442226688
```

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

The control group was more likely to make a conversion, but it was not a huge difference. In control group we had 12.04% conversion rate while in treatment group 11.88%, so that is a no!

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.

$H_0 : P_{old} \geq P_{new}$,

$H_1 : P_{old} < P_{new}$,

$\alpha = 0.05$

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 **convert rate** for p_{new} under the null?

```
[19]: Pnew = len(converted)/len(df2['converted'])
Pnew
```

```
[19]: 0.11959708724499628
```

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

```
[20]: Pold = len(converted)/len(df2['converted'])
Pold
```

```
[20]: 0.11959708724499628
```

c. What is n_{new} ?

```
[21]: Nnew = len(df2.where(df2['group'] == 'treatment').dropna())
Nnew
```

```
[21]: 145310
```

d. What is n_{old} ?

```
[22]: Nold = len(df2.where(df2['group'] == 'control').dropna())
Nold
```

```
[22]: 145274
```

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

```
[23]: new_page_converted = np.random.choice(a=[0,1],p=[1-Pnew,Pnew],size=1000)
```

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

```
[24]: old_page_converted = np.random.choice(a=[0,1],p=[1-Pold,Pold],size=1000)
```

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

```
[25]: Pnew_Pold = new_page_converted.mean() - old_page_converted.mean()
Pnew_Pold
```

```
[25]: -0.00400000000000000036
```

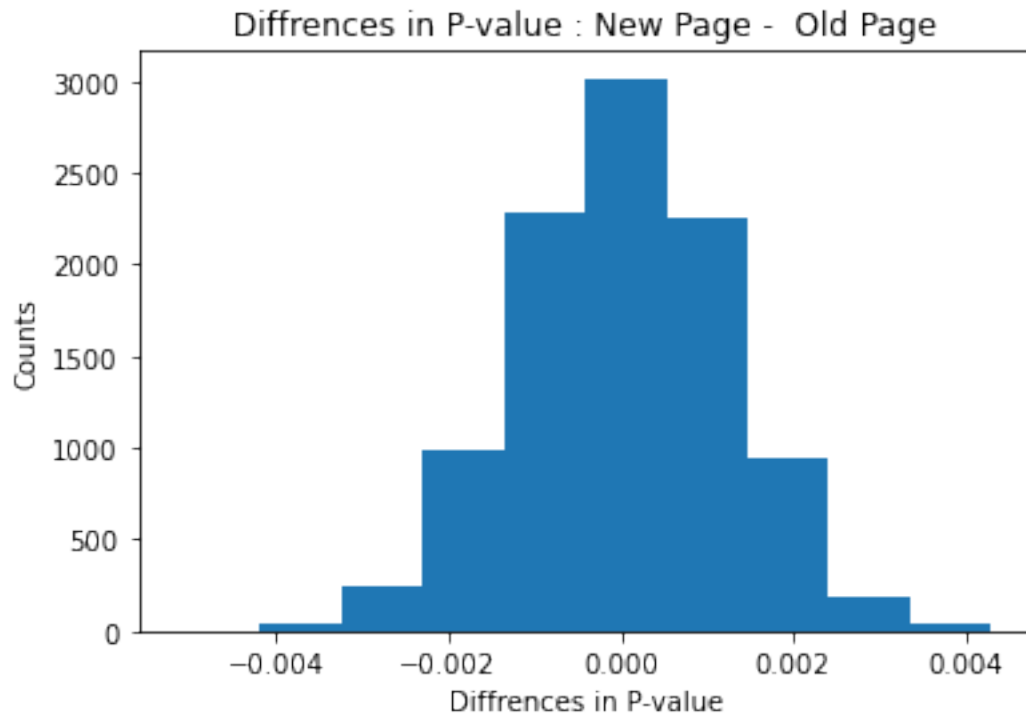
h. Simulate 10,000 $p_{new} - p_{old}$ values using this same process similarly to the one you calculated in parts **a. through g.** above. Store all 10,000 values in a numpy array called **p_diffs**.

```
[46]: p_diffs = []
new_converted_simulation = np.random.binomial(Nnew, Pnew, 10000)/Nnew
old_converted_simulation = np.random.binomial(Nold, Pold, 10000)/Nold
```

```
p_diffs = new_converted_simulation - old_converted_simulation
```

- 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.

```
[47]: plt.hist(p_diffs);  
plt.xlabel('Differences in P-value')  
plt.ylabel('Counts')  
plt.title('Differences in P-value : New Page - Old Page');
```



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

```
[48]: obs_diffs = 0.11880724790277405 - 0.1203863045004612
```

```
[49]: obs_diffs
```

```
[49]: -0.0015790565976871451
```

```
[50]: box = []  
for val in p_diffs:  
    if(val) > obs_diffs:  
        box.append(val)
```

```
[51]: len(box)/10000
```

```
[51]: 0.9076
```

- k. In words, explain 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?

We computed the P value since its around 0.9 there is reason for converting since it had to be 0.05 or less to be relevant.

- 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.

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

convert_old = sum(df2.query("group == 'control'")['converted'])
convert_new = sum(df2.query("group == 'treatment'")['converted'])
```

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

```
[53]: sm.stats.proportions_ztest([convert_new,convert_old], [Nnew, Nold], alternative='larger')
```

```
[53]: (-1.3109241984234394, 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**.?

Z-score is standard deviations above the mean and p-value is probability of getting 'random' values while simulating from a sample. No they do not agree since second one is two-tailed and value in parts j,k. is one-tailed.

Part III - A regression approach

1. In this final part, you will see that the result you acheived in the previous A/B test can also be acheived 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 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**.

```
[34]: df2['intercept'] = 1
      df2['ab_page'] = 0
```

```
[35]: groupa = df2.where(df2['group'] == 'treatment').dropna()
      groupa['ab_page'] = 1
      groupb = df2.where(df2['group'] == 'control').dropna()
      df2 = groupb.merge(groupa,how='outer')
```

```
[36]: df2.head(3)
```

```
[36]:   user_id      timestamp  group landing_page  converted \
0  851104.0  2017-01-21 22:11:48.556739  control    old_page      0.0
1  804228.0  2017-01-12 08:01:45.159739  control    old_page      0.0
2  864975.0  2017-01-21 01:52:26.210827  control    old_page      1.0

   intercept  ab_page
0          1.0      0.0
1          1.0      0.0
2          1.0      0.0
```

c. Use **statsmodels** to import your regression model. Instantiate the model, and fit the model using the two columns you created in part **b.** to predict whether or not an individual converts.

```
[37]: stats = sm.Logit(df2['converted'],df2[['intercept','ab_page']])
      results = stats.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.

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

```
[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:                   Sun, 24 Jul 2022    Pseudo R-squ.:          8.077e-06
Time:                   18:06:18          Log-Likelihood:         -1.0639e+05
converged:               True            LL-Null:                -1.0639e+05
Covariance Type:         nonrobust        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 the **Part II**?

The $P = 0.19$, $H_0 : B_0 = 0$, $H_1: B_0 \neq 0$,

So again this is two-tailed since $H_1: B_0 \neq 0$ and in part II we had, one-tailed one.

- 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?

Yes We might consider other factors. It will be good idea to have overall good understanding why peoples are converting, but yes there will be few disadvantages like: A linear relationship may not exist, Correlated errors, Non-constant variance, Outliers, Multicollinearity

- g. Now along with testing if the conversion rate changes for different pages, also add an effect based on which country a user lives. 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.

```
[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(3)
```

```
[39]:      country      timestamp      group landing_page \
user_id
834778.0      UK  2017-01-14 23:08:43.304998    control    old_page
928468.0      US  2017-01-23 14:44:16.387854  treatment    new_page
822059.0      UK  2017-01-16 14:04:14.719771  treatment    new_page

      converted  intercept  ab_page
user_id
834778.0        0.0         1.0     0.0
928468.0        0.0         1.0     1.0
822059.0        1.0         1.0     1.0
```

```
[40]: # US - 0
      # UK - 1
      # CA - 2

      df_new['is_US'] = 0
      df_new['is_UK'] = 0
      df_new['is_CA'] = 0
      df_new.reset_index(inplace=True)

[41]: gro1 = df_new.where(df_new['country'] == 'UK').dropna()
      gro1['is_UK'] = 1

      gro2 = df_new.where(df_new['country'] == 'CA').dropna()
      gro2['is_CA'] = 1

      gro0 = df_new.where(df_new['country'] == 'US').dropna()
      gro0['is_US'] = 1

      df_new = gro0.merge(gro1,how='outer')
      df_new = df_new.merge(gro2,how='outer')

[42]: stats2 = sm.
      ↪Logit(df_new['converted'],df_new[['intercept','ab_page','is_UK','is_CA']])
      results2 = stats2.fit()
      results2.summary()
```

Optimization terminated successfully.
 Current function value: 0.366113
 Iterations 6

```
[42]: <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:	Sun, 24 Jul 2022	Pseudo R-squ.:	2.323e-05			
Time:	18:06:26	Log-Likelihood:	-1.0639e+05			
converged:	True	LL-Null:	-1.0639e+05			
Covariance Type:	nonrobust	LLR p-value:	0.1760			
=====						
	coef	std err	z	P> z	[0.025	0.975]

intercept	-1.9893	0.009	-223.763	0.000	-2.007	-1.972
ab_page	-0.0149	0.011	-1.307	0.191	-0.037	0.007
is_UK	0.0099	0.013	0.743	0.457	-0.016	0.036

```
is_CA          -0.0408      0.027      -1.516      0.130      -0.093      0.012
=====
"""
```

- 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.

```
[43]: df_new['UK_ab_page'] = df_new['is_UK'] * df_new['ab_page']
df_new['CA_ab_page'] = df_new['is_CA'] * df_new['ab_page']
stats2 = sm.
    ↪Logit(df_new['converted'],df_new[['intercept','ab_page','is_UK','is_CA','UK_ab_page','CA_ab
results2 = stats2.fit()
results2.summary()
```

Optimization terminated successfully.

Current function value: 0.366109

Iterations 6

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

```

                                Logit Regression Results
=====
Dep. Variable:                converted    No. Observations:                290584
Model:                        Logit        Df Residuals:                290578
Method:                        MLE         Df Model:                    5
Date:                Sun, 24 Jul 2022    Pseudo R-squ.:                3.482e-05
Time:                18:06:29            Log-Likelihood:               -1.0639e+05
converged:                True            LL-Null:                    -1.0639e+05
Covariance Type:            nonrobust      LLR p-value:                0.1920
=====
               coef      std err          z      P>|z|      [0.025      0.975]
-----
intercept    -1.9865      0.010   -206.344      0.000      -2.005      -1.968
ab_page      -0.0206      0.014    -1.505      0.132      -0.047      0.006
is_UK        -0.0057      0.019    -0.306      0.760      -0.043      0.031
is_CA        -0.0175      0.038    -0.465      0.642      -0.091      0.056
UK_ab_page     0.0314      0.027     1.181      0.238      -0.021      0.084
CA_ab_page    -0.0469      0.054    -0.872      0.383      -0.152      0.059
=====
"""
```

Conclusions:

We found out there is almost the same probability for user to get either new or old page.

There is not enough evidences to support the hypothesis that new_page will lead to more conver-

sions.

After adding countries in the equation, We found that user country of origin is also irrelevant in terms of conversion.

Conclusions

Congratulations on completing the project!

0.2.1 Gather Submission Materials

Once you are satisfied with the status of your Notebook, you should save it in a format that will make it easy for others to read. You can use the **File -> Download as -> HTML (.html)** menu to save your notebook as an .html file. If you are working locally and get an error about “No module name”, then open a terminal and try installing the missing module using `pip install <module_name>` (don’t include the “<” or “>” or any words following a period in the module name).

You will submit both your original Notebook and an HTML or PDF copy of the Notebook for review. There is no need for you to include any data files with your submission. If you made reference to other websites, books, and other resources to help you in solving tasks in the project, make sure that you document them. It is recommended that you either add a “Resources” section in a Markdown cell at the end of the Notebook report, or you can include a `readme.txt` file documenting your sources.

0.2.2 Submit the Project

When you’re ready, click on the “Submit Project” button to go to the project submission page. You can submit your files as a .zip archive or you can link to a GitHub repository containing your project files. If you go with GitHub, note that your submission will be a snapshot of the linked repository at time of submission. It is recommended that you keep each project in a separate repository to avoid any potential confusion: if a reviewer gets multiple folders representing multiple projects, there might be confusion regarding what project is to be evaluated.

It can take us up to a week to grade the project, but in most cases it is much faster. You will get an email once your submission has been reviewed. If you are having any problems submitting your project or wish to check on the status of your submission, please email us at dataanalyst-project@udacity.com. In the meantime, you should feel free to continue on with your learning journey by beginning the next module in the program.

[]: