

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!

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 (<https://review.udacity.com/#!/projects/37e27304-ad47-4eb0-a1ab-8c12f60e43d0/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 s
et 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 below cell to find the number of rows in the dataset.

```
In [3]: 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
```

c. The number of unique users in the dataset.

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

Out[4]: 290584

d. The proportion of users converted.

```
In [5]: n_con = df[df['converted']==1]['user_id'].count()
p=n_con/n
p
```

Out[5]: 0.12126269856564711

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

```
In [6]: df[((df['group'] == 'treatment') == (df['landing_page'] != 'new_page'))==True]
        ['user_id'].count()
```

Out[6]: 3893

f. Do any of the rows have missing values?

```
In [7]: # check if any columns have null values; it should print false
        df.isnull().sum().any()
```

Out[7]: False

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**.

```
In [8]: remove=df[((df['group'] == 'treatment') == (df['landing_page'] != 'new_page'))
        ==True]
        remove.head()
```

Out[8]:

	user_id	timestamp	group	landing_page	converted
22	767017	2017-01-12 22:58:14.991443	control	new_page	0
240	733976	2017-01-11 15:11:16.407599	control	new_page	0
308	857184	2017-01-20 07:34:59.832626	treatment	old_page	0
327	686623	2017-01-09 14:26:40.734775	treatment	old_page	0
357	856078	2017-01-12 12:29:30.354835	treatment	old_page	0

```
In [9]: df2=df.drop(remove.index,axis=0)
        df2.head()
```

Out[9]:

	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 [10]: df2['user_id'].count()
```

```
Out[10]: 290585
```

```
In [11]: # 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]
```

```
Out[11]: 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 [12]: df2['user_id'].nunique()
```

```
Out[12]: 290584
```

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

```
In [13]: df2[df2.duplicated(['user_id'], keep=False)]
```

```
Out[13]:
```

	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

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

```
In [14]: # The landing_page for the non-unique id is "new_page". The group for the non-
unique id is "treatment".
# The value of converted column for the non-unique id is 0
```

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

```
In [15]: df2.drop_duplicates(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?

```
In [16]: npc=df2[df2['converted']==1]['user_id'].count()  
prc = npc/n  
prc
```

Out[16]: 0.11959708724499628

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

```
In [17]: cong=df2[df2['group']=='control']  
ng=cong.user_id.count()  
cong_conv=cong[cong['converted']==1]['user_id'].count()  
pc=cong_conv/ng  
pc
```

Out[17]: 0.1203863045004612

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

```
In [18]: cont=df2[df2['group']=='treatment']  
nt=cont.user_id.count()  
cont_conv=cont[cont['converted']==1]['user_id'].count()  
pt=cont_conv/nt  
pt
```

Out[18]: 0.11880724790277405

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

```
In [19]: new_page=df2[df2['landing_page']=='new_page']['user_id'].count()  
new_page/n
```

Out[19]: 0.50006538556837266

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.

From the metrics above, although half of the population received the new page, there is no significant evidence to conclude that the new treatment page led to more conversions. There were as many conversions with those who landed on the old page as were with the new page.

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.

Null Hypothesis: $p_{new} - p_{old} \leq 0$

Alternative Hypothesis: $p_{new} - p_{old} \geq 0$

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

```
In [20]: npc=df2[df2['converted']==1]['user_id'].count()  
p_new = npc/n  
p_new
```

```
Out[20]: 0.11959708724499628
```

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

```
In [21]: npc=df2[df2['converted']==1]['user_id'].count()  
p_old = npc/n  
p_old
```

```
Out[21]: 0.11959708724499628
```

c. What is n_{new} ?

```
In [22]: new_df=df2.query('group=="treatment"')  
nnew=new_df.query('landing_page=="new_page"').user_id.nunique()  
nnew
```

```
Out[22]: 145310
```

d. What is n_{old} ?

```
In [23]: old_df=df2.query('group=="control"')  
nold=old_df.query('landing_page=="old_page"').user_id.nunique()  
nold
```

```
Out[23]: 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**.

```
In [32]: new_page_converted = np.random.binomial(1, p_new,nnew)
```

```
In [25]: n_new=len(new_page_converted)  
n_new
```

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

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**.

```
In [31]: old_page_converted = np.random.binomial(1, p_old,nold)
old_page_converted
```

```
Out[31]: array([1, 0, 0, ..., 0, 0, 0])
```

```
In [27]: n_old=len(old_page_converted)
n_old
```

```
Out[27]: 145274
```

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

```
In [33]: n_old_converted=(old_page_converted==1).sum()
n_new_converted=(new_page_converted==1).sum()
obs_diff=n_new_converted/n_new-n_old_converted/n_old
obs_diff
```

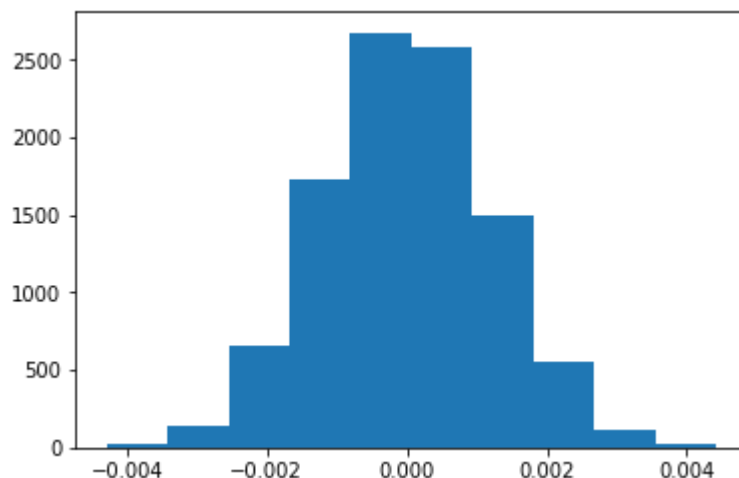
```
Out[33]: 8.0715631123454146e-05
```

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**.

```
In [40]: p_diffs = []
for _ in range(10000):
    new_pages_converted = np.random.binomial(1, p_new,nnew)
    n_news=len(new_pages_converted)
    old_pages_converted = np.random.binomial(1, p_old,nold)
    n_old=len(old_pages_converted)
    n_old_converted=(old_pages_converted==1).sum()
    n_news_converted=(new_pages_converted==1).sum()
    obs_diffs=n_news_converted/n_news-n_old_converted/n_old
    p_diffs.append(obs_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 [41]: plt.hist(p_diffs);
```



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

```
In [42]: ndiff=len(p_diffs)
```

```
In [43]: obsdiff=pt-pc
grt = [i for i in p_diffs if i>float(obsdiff)]
n_grt=len(grt)
p_grt = n_grt/ndiff
p_grt
```

```
Out[43]: 0.9031
```

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?

The index calculated in "j" above is the p-value, which is the probability of observing the statistic in favor of the alternative hypothesis if the null hypothesis is true. The null hypothesis suggests that there is no difference in conversion based on the landing page. With a p-value of 0.9, suggests the null hypothesis should not be rejected.

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 [ ]: 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 = n_old
n_new = n_new
```

m. Now use `stats.proportions_ztest` to compute your test statistic and p-value. [Here](http://knowledgegetack.com/python/statsmodels/proportions_ztest/) (http://knowledgegetack.com/python/statsmodels/proportions_ztest/) is a helpful link on using the built in.

```
In [ ]: z_score, p_value = sm.stats.proportions_ztest([convert_new, convert_old], [nnew, nold], alternative='larger')
```

```
In [ ]: p_value
```

```
In [ ]: z_score
```

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.**?

The p-value, as in indicated in the response to j above, is the probability of observing our statistic provided the null hypothesis is true. Similar to parts j and k, the p-value obtained in part l suggests that the landing page has no impact on the conversion rates. The z-score is a measure of how many standard deviations are below or above the population mean.

Part III - A regression approach

1. In this final part, you will see that the result you achieved in the previous A/B test 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 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 [45]: import statsmodels.api as sm
        ##df2['group'].replace('treatment', 'ab_page', inplace=True)
```

```
C:\Users\User\Anaconda3\lib\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
```

```
In [46]: ## Getting duration column
        import datetime as dt
        from datetime import datetime
        df2['timestamp'] = df2['timestamp'].apply(lambda x : pd.to_datetime(str(x)))
```

```
In [47]: df2['hours'] = df2['timestamp'].dt.hour + (df2['timestamp'].dt.minute)/60 + (df2['timestamp'].dt.second)/3600
```

```
In [48]: df2[['control', 'treatment']] = pd.get_dummies(df2['group'])
        df2['ab_column'] = df2['treatment']
        df2 = df2.drop(['control', 'treatment'], axis=1)
        df2.head()
```

Out[48]:

	user_id	timestamp	group	landing_page	converted	hours	ab_column
0	851104	2017-01-21 22:11:48.556739	control	old_page	0	22.196667	0
1	804228	2017-01-12 08:01:45.159739	control	old_page	0	8.029167	0
2	661590	2017-01-11 16:55:06.154213	treatment	new_page	0	16.918333	1
3	853541	2017-01-08 18:28:03.143765	treatment	new_page	0	18.467500	1
4	864975	2017-01-21 01:52:26.210827	control	old_page	1	1.873889	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.

```
In [49]: df2['intercept']=1
logit_mod=sm.Logit(df2['converted'],df2[['intercept','ab_column']])
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.

```
In [50]: results.summary()
```

Out[50]: Logit Regression Results

Dep. Variable:	converted	No. Observations:	290585
Model:	Logit	Df Residuals:	290583
Method:	MLE	Df Model:	1
Date:	Wed, 24 Jan 2018	Pseudo R-squ.:	8.085e-06
Time:	11:05:35	Log-Likelihood:	-1.0639e+05
converged:	True	LL-Null:	-1.0639e+05
		LLR p-value:	0.1897

	coef	std err	z	P> z	[0.025	0.975]
intercept	-1.9888	0.008	-246.669	0.000	-2.005	-1.973
ab_column	-0.0150	0.011	-1.312	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-value associated with **ab_column** is $0.19 > \text{new threshold rate}$. Thus suggesting that the null hypothesis should not be rejected. The null hypothesis suggests that there is no significant relationship between conversion and page, while the alternative hypothesis indicates that there is significant relationship between the new page and conversion rates; the two hypotheses are similar to Part II. And the results of this analysis similar to what was observed in Part II above. The simulation and the z-test performed above were implemented as one-tailed tests while the regression is a two-tailed test; this resulted in the difference in the value of p-value obtained in this regression part.

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?

Other possible factors what considering include the country of residents of the individual converts, profession, age, and gender. With more factors to consider, there could be possible problems such as multicollinearity, nonlinear relationship existence among explanatory variables, outliers, and many others.

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 \(https://pandas.pydata.org/pandas-docs/stable/generated/pandas.DataFrame.join.html\)](https://pandas.pydata.org/pandas-docs/stable/generated/pandas.DataFrame.join.html) 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 [51]: 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[51]:

	country	timestamp	group	landing_page	converted	hours	ab_col
user_id							
630000	US	2017-01-19 06:26:06.548941	treatment	new_page	0	6.435000	1
630001	US	2017-01-16 03:16:42.560309	treatment	new_page	1	3.278333	1
630002	US	2017-01-19 19:20:56.438330	control	old_page	0	19.348889	0
630003	US	2017-01-12 10:09:31.510471	treatment	new_page	0	10.158611	1
630004	US	2017-01-18 20:23:58.824994	treatment	new_page	0	20.399444	1

```
In [52]: ### Create the necessary dummy variables
df_new[['US','UK','CA']]=pd.get_dummies(df_new['country'])
df_new=df_new.drop('UK',axis=1)
df_new.head()
```

Out[52]:

	country	timestamp	group	landing_page	converted	hours	ab_col
user_id							
630000	US	2017-01-19 06:26:06.548941	treatment	new_page	0	6.435000	1
630001	US	2017-01-16 03:16:42.560309	treatment	new_page	1	3.278333	1
630002	US	2017-01-19 19:20:56.438330	control	old_page	0	19.348889	0
630003	US	2017-01-12 10:09:31.510471	treatment	new_page	0	10.158611	1
630004	US	2017-01-18 20:23:58.824994	treatment	new_page	0	20.399444	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 [53]: ### Fit Your Linear Model And Obtain the Results
df_new['intercept']=1
logit_mod2=sm.Logit(df_new['converted'],df_new[['intercept','US','CA']])
results2=logit_mod2.fit()
```

```
Optimization terminated successfully.
Current function value: 0.366115
Iterations 6
```

In [54]: `results2.summary()`

Out[54]: Logit Regression Results

Dep. Variable:	converted	No. Observations:	290585
Model:	Logit	Df Residuals:	290582
Method:	MLE	Df Model:	2
Date:	Wed, 24 Jan 2018	Pseudo R-squ.:	1.521e-05
Time:	11:05:48	Log-Likelihood:	-1.0639e+05
converged:	True	LL-Null:	-1.0639e+05
		LLR p-value:	0.1983

	coef	std err	z	P> z 	[0.025	0.975]
intercept	-1.9868	0.011	-174.174	0.000	-2.009	-1.964
US	-0.0507	0.028	-1.786	0.074	-0.106	0.005
CA	-0.0099	0.013	-0.746	0.456	-0.036	0.016

In [55]: `np.exp(-0.0507)`

Out[55]: 0.95056379690040338

In [56]: `np.exp(-0.0099)`

Out[56]: 0.99014884368295719

From the results above, the conversion rates from each country are equally possible. Thus suggesting that conversion rate is not tilted in favor of a particular country.

Conclusions

Congratulations on completing the project!

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.

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.