

Malak Alshedokhi

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 \(https://review.udacity.com/#!/projects/37e27304-ad47-4eb0-a1ab-8c12f60e43d0/rubric\)](https://review.udacity.com/#!/projects/37e27304-ad47-4eb0-a1ab-8c12f60e43d0/rubric).

Part I - Probability

To get started, let's import our libraries.

In [2]:

```
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 [3]:

```
df = pd.read_csv('ab_data.csv')
df.head(10)
```

Out[3]:

	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

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

In [4]:

```
df.shape
```

Out[4]:

(294478, 5)

In []:

```
# number of rows = 294478
```

c. The number of unique users in the dataset.

In [5]:

```
len(pd.unique(df['user_id']))
```

Out[5]:

290584

In []:

```
#The number of unique users in the dataset = 290584
```

d. The proportion of users converted.

In [6]:

```
df['converted'].mean()
```

Out[6]:

0.11965919355605512

In []:

```
# The proportion of users converted = 0.11965919355605512
```

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

In [7]:

```
df.query("(group == 'treatment' and landing_page == 'old_page' or group == 'control' and landing_page == 'new_page')").shape
```

Out[7]:

(3893, 5)

In []:

```
# The number of times the new_page and treatment don't line u  
p = 3893
```

f. Do any of the rows have missing values?

In [8]:

```
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
```

In [9]:

```
df.isnull().sum()
```

Out[9]:

```
user_id          0
timestamp        0
group            0
landing_page     0
converted        0
dtype: int64
```

In []:

```
# No missing values in any of the rows
```

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 [10]:

```
#Dropping the rows where group is treatment and landing page is old page
df.drop(df.query("group == 'treatment' and landing_page == 'old_page'").index, inplace=True)
#Dropping the rows where group is control and landing page is new page
df.drop(df.query("group == 'control' and landing_page == 'new_page'").index, inplace=True)
```

In [11]:

```
#Checking the updated data frame
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 290585 entries, 0 to 294477
Data columns (total 5 columns):
user_id          290585 non-null int64
timestamp        290585 non-null object
group            290585 non-null object
landing_page     290585 non-null object
converted        290585 non-null int64
dtypes: int64(2), object(3)
memory usage: 13.3+ MB
```

In [12]:

```
#Saving the new data in a modified CSV file
df.to_csv('ab_modified.csv', index=False)
```

In [13]:

```
#Creating df2 where it reads from the modified file
df2 = pd.read_csv('ab_modified.csv')
```

In [14]:

```
#Checking
df2[((df2['group'] == 'treatment') == (df2['landing_page'] == 'new_page')) == False].shape[0]
```

Out[14]:

0

In [15]:

```
df2.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 290585 entries, 0 to 290584
Data columns (total 5 columns):
user_id          290585 non-null int64
timestamp        290585 non-null object
group            290585 non-null object
landing_page     290585 non-null object
converted        290585 non-null int64
dtypes: int64(2), object(3)
memory usage: 11.1+ MB
```

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

In [16]:

```
df2.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 290585 entries, 0 to 290584
Data columns (total 5 columns):
user_id          290585 non-null int64
timestamp        290585 non-null object
group            290585 non-null object
landing_page     290585 non-null object
converted        290585 non-null int64
dtypes: int64(2), object(3)
memory usage: 11.1+ MB
```

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

In [17]:

```
len(pd.unique(df2['user_id']))
```

Out[17]:

290584

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

In [18]:

```
sum(df2['user_id'].duplicated())
```

Out[18]:

1

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

In [19]:

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

Out[19]:

	user_id	timestamp	group	landing_page	converted
1876	773192	2017-01-09 05:37:58.781806	treatment	new_page	0
2862	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 [20]:

```
timestamp_duplicate = "2017-01-09 05:37:58.781806"  
df2 = df2[df2.timestamp != timestamp_duplicate]
```

In [21]:

```
#Check if it was removed  
sum(df2['user_id'].duplicated())
```

Out[21]:

0

In [22]:

```
df2.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 290584 entries, 0 to 290584
Data columns (total 5 columns):
user_id      290584 non-null int64
timestamp    290584 non-null object
group        290584 non-null object
landing_page  290584 non-null object
converted     290584 non-null int64
dtypes: int64(2), object(3)
memory usage: 13.3+ MB
```

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 [23]:

```
df2['converted'].mean()
```

Out[23]:

```
0.11959708724499628
```

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

In [24]:

```
df_group = df2.groupby('group')
df_group.describe()
```

Out[24]:

user_id					
	count	mean	std	min	25%
group					
control	145274.0	788164.072594	91287.914601	630002.0	709279.5
treatment	145310.0	787845.719290	91161.564429	630000.0	708745.7

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

In [25]:

```
#0.118808
```

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

In [26]:

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

Out[26]:

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.

I think the old page and new page conversion rate are almost similar

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.

$$H_0 : p_{(new)} \leq p_{(old)}$$

$$H_1 : p_{(new)} > p_{(old)}$$

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 [42]:

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

```
0.11959708724499628
```

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

In [44]:

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

```
0.11959708724499628
```

c. What is n_{new} ?

In [45]:

```
n_new = df2[df2['group']=='treatment'].shape[0]  
print(n_new)
```

```
145310
```

d. What is n_{old} ?

In [46]:

```
n_old = df2[df2['group']=='control'].shape[0]  
print(n_old)
```

```
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 [35]:

```
new_page_converted = np.random.binomial(n_new,p_new)
```

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 [36]:

```
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).

In [37]:

```
new_page_converted/n_new - old_page_converted/n_old
```

Out[37]:

```
-5.688533261327677e-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**.

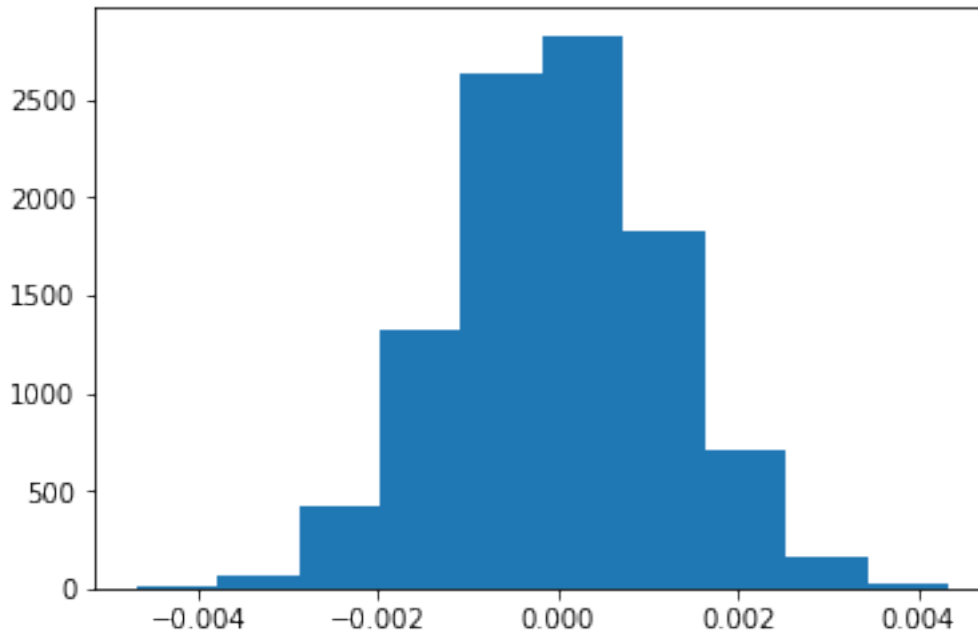
In [74]:

```
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)
    #new_page_converted/n_new - old_page_converted/n_old
    #add more simulatio for old page converted
    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.

In [75]:

```
plt.hist(p_diffs);
```



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

In [76]:

```
act_diff = df[df['group'] == 'treatment']['converted'].mean()  
- df[df['group'] == 'control']['converted'].mean()  
(act_diff < p_diffs).mean()
```

Out[76]:

0.9024

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 p-value is 0.9024, so that value indicate a weak evidence toward the null hypothesis.

Reference :

https://www.researchgate.net/post/how_to
(https://www.researchgate.net/post/how_to

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 [77]:

```
import statsmodels.api as sm

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

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 [73]:

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

```
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.?

P-value is = 0.9, z-score is = 1.31. Since the value of z-score is larger than P-value, we fail to reject the null hypothesis.

Reference :

[http://www.cs.uni.edu/~campbell/stat/inf3.](http://www.cs.uni.edu/~campbell/stat/inf3)
(<http://www.cs.uni.edu/~campbell/stat/inf3>)

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 [53]:

```
df['intercept'] = 1
df[['control', 'treatment']] = pd.get_dummies(df['group'])
```


In [54]:

```
df.head()
```

Out[54]:

	user_id	timestamp	group	landing_page	converted	inter
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	

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 [55]:

```
df.columns
```

Out[55]:

```
Index(['user_id', 'timestamp', 'group', 'landing_p  
age', 'converted',  
      'intercept', 'control', 'treatment'],  
      dtype='object')
```

In [56]:

```
import statsmodels.api as sm
```

In [57]:

```
df.columns
```

Out[57]:

```
Index(['user_id', 'timestamp', 'group', 'landing_p  
age', 'converted',  
      'intercept', 'control', 'treatment'],  
      dtype='object')
```

In [58]:

```
logit = sm.Logit(df['converted'], df[['intercept', 'treatment'  
]])  
  
results = logit.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 [59]:

```
results.summary()
```

Out[59]:

Logit Regression Results

Dep. Variable:	converted	No. Observations:	290585
Model:	Logit	Df Residuals:	290583
Method:	MLE	Df Model:	1
Date:	Wed, 19 Jun 2019	Pseudo R-squ.:	8.085e-06
Time:	11:48:26	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
treatment	-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**?

p-value is 0.216

AB testing is one tail while regression is 2 tails.

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 would be powerful to add more factors, however that might affect p-value.

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 [60]:

```
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[60]:

	country	timestamp	group	landing_page	converted
user_id					
834778	UK	2017-01-14 23:08:43.304998	control	old_page	0
928468	US	2017-01-23 14:44:16.387854	treatment	new_page	0
822059	UK	2017-01-16 14:04:14.719771	treatment	new_page	1
711597	UK	2017-01-22 03:14:24.763511	control	old_page	0
710616	UK	2017-01-16 13:14:44.000513	treatment	new_page	0

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 [65]:

```
df_new['country'].value_counts()  
df_new[['CA', 'US']] = pd.get_dummies(df_new['country'])[['CA',  
'US']]  
df_new.head()
```

Out[65]:

	country	timestamp	group	landing_page	converted
user_id					
834778	UK	2017-01-14 23:08:43.304998	control	old_page	0
928468	US	2017-01-23 14:44:16.387854	treatment	new_page	0
822059	UK	2017-01-16 14:04:14.719771	treatment	new_page	1
711597	UK	2017-01-22 03:14:24.763511	control	old_page	0
710616	UK	2017-01-16 13:14:44.000513	treatment	new_page	0

In [71]:

```
logit_model = sm.Logit(df_new['converted'], df_new[['CA', 'US']])  
logit_model  
results = logit_model.fit()  
results.summary()
```

Optimization terminated successfully.
Current function value: 0.447174
Iterations 6

Out[71]:

Logit Regression Results

Dep. Variable:	converted	No. Observations:	290584
Model:	Logit	Df Residuals:	290582
Method:	MLE	Df Model:	1
Date:	Wed, 19 Jun 2019	Pseudo R-squ.:	-0.2214
Time:	12:01:12	Log-Likelihood:	-1.2994e+05
converged:	True	LL-Null:	-1.0639e+05
		LLR p-value:	1.000

	coef	std err	z	P> z	[0.025	0.975]
CA	-2.0375	0.026	-78.364	0.000	-2.088	-1.987
US	-1.9967	0.007	-292.314	0.000	-2.010	-1.983

Conclusions

AB test fail to reject the null hypothesis Regression model showed that the old page is better than the new page.

References:

Udacity's video "Practical Stats Project Walkthrough Final"

In []: