

## 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 not page. Either way assume 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!

### Table of Contents

- [Introduction](#)
- [Part I - Probability](#)
- [Part II - A/B Test](#)
- [Part III - Regression](#)

### Introduction

A/B testing 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 [2]: import pandas as pd
import numpy as np
import random
import matplotlib.pyplot as plt
matplotlib inline
# We will 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)
```

	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 01:52:26.218027	control	old_page	1
5	936923	2017-01-10 15:20:49.083499	control	old_page	0
6	676697	2017-01-19 03:26:46.040749	treatment	new_page	1
7	719014	2017-01-17 01:48:29.639573	control	old_page	0
8	817355	2017-01-04 17:50:05.879471	treatment	new_page	1
9	839785	2017-01-15 18:11:06.610965	treatment	new_page	1

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

```
In [4]: df.shape[0]
Out[4]: 294478
```

c. The number of unique users in the dataset.

```
In [5]: df.user_id.nunique()
Out[5]: 290584
```

d. The proportion of users converted.

```
In [6]: uniq_users=df.user_id.nunique()
df.query('converted==1').user_id.nunique()/uniq_users
Out[6]: 0.1219425424669237
```

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

```
In [7]: len(df[(df['group']!='treatment') & (df['landing_page']!='new_page') | (df['group']=='treatm
ent') & (df['landing_page']!='new_page')])
Out[7]: 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 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
```

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 [9]: df_remove=df[(df['group']=='treatment') & (df['landing_page']!='new_page') | (df['group']=='treatm
ent') & (df['landing_page']!='new_page')]
In [10]: df2=df.drop(df_remove.index, inplace=False)
df2.shape[0]
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 **Quiz 3** in the classroom.

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

```
In [12]: df2.shape[0]
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')].user_id
Out[13]: 2893
Name: user_id, dtype: int64
```

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

```
In [14]: duplicate_row=df2[df2['user_id']=='773192']
duplicate_row
Out[14]:
```

	user_id	timestamp	group	landing_page	converted
1899	773192	2017-01-09 05:37:58.781806	treatment	new_page	0
2693	773192	2017-01-14 02:55:59.590927	treatment	new_page	0

d. Remove one of the rows with a duplicate `user_id`. Keep your dataframe as `df2`.

```
In [15]: # Since the timestamp for both the entries are different, we need to filter
# the row using user_id column
df2.drop_duplicates(subset='user_id', inplace=True)
In [16]: df2.drop_duplicates(subset='user_id', inplace=True)
In [17]: df2[df2['user_id']=='773192']
Out[17]:
```

	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 [20]: # This also works fine but better use the minimized ones : len(df2[df2['converted']==1])/len
(df2['landing_page'])
# rates equal to the converted success rate regardless of page - that is p_new
# and p_old are equal.
p_converted = df2.converted.mean()
p_converted
Out[20]: 0.11959788724499628
```

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

```
In [21]: len(df2[(df2['group']=='control') & (df2['converted']==1)])/len(df2[(df2['group']=='control'
)])
Out[21]: 0.1263853845064612
```

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

```
In [22]: len(df2[(df2['group']=='treatment') & (df2['converted']==1)])/len(df2[(df2['group']=='treatm
ent')])
Out[22]: 0.1188896055151564
```

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

```
In [23]: len(df2[df2['landing_page']=='new_page']/len(df2['landing_page']))
Out[23]: 0.500601942226688
```

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

Your answer goes here. **Discourage not!** The results generated above rather suggest us that the control page conversion is higher than the treatment page. However, the results generated upon the above calculations are not enough to support that the control page conversions are better because, often these predictions are delusional.

### 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 to 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 you 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.

$H_0: p_{new} = p_{old}$

$H_1: p_{new} \neq p_{old}$

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** success 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 project. 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 [21]: df2.head(5)
Out[21]:
```

	user_id	timestamp	group	landing_page	converted
0	851104	2017-01-21 22:11:48.556739	control	old_page	0
1	804228	2017-01-12 08:01:45.159739	control	old_page	0
2	661590	2017-01-11 16:55:06.154213	treatment	new_page	0
3	853541	2017-01-08 18:28:03.143765	treatment	new_page	0
4	864975	2017-01-21 01:52:26.218027	control	old_page	1

```
In [24]: # 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.
p_new=df2['converted'].mean()
p_old
Out[24]: 0.11959788724499628
```

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

```
In [25]: p_old=df2['converted'].mean()
p_old
Out[25]: 0.11959788724499628
```

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

```
In [26]: n_new=df2[df2['group']=='treatment'].count()
n_new
Out[26]:
```

	user_id	timestamp	group	landing_page	converted	dtype
0	145310	145310	145310	145310	145310	int64

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

```
In [27]: df2[df2['group']=='control'].count()
Out[27]:
```

	user_id	timestamp	group	landing_page	converted	dtype
0	145274	145274	145274	145274	145274	int64

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 [28]: new_page_converted=np.random.choice([0,1],size=145310,p=[1-p_new,p_new])
p_new_converted=new_page_converted.mean()
p_new_converted
Out[28]: 0.1186497832267693
```

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 [29]: old_page_converted=np.random.choice([0,1],size=145274,p=[1-p_old,p_old])
p_old_converted=old_page_converted.mean()
p_old_converted
Out[29]: 0.12612472982089917
```

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

```
In [30]: p_new_converted-p_old_converted
Out[30]: -0.0614749465088132399
```

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 [35]: p_diffs=[]
for i in range(10000):
    new_page_converted=np.random.choice([0,1],size=145310,p=[1-p_new,p_new])
    p_new_converted=new_page_converted.mean()
    old_page_converted=np.random.choice([0,1],size=145274,p=[1-p_old,p_old])
    p_old_converted=old_page_converted.mean()
    p_diffs.append(p_new_converted-p_old_converted)
p_diffs=array(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 [36]: plt.hist(p_diffs);
```

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

```
In [37]: obs_diff=df2.query('group=="treatment"')['converted'].mean()-df2.query('group=="control"')['
converted'].mean()
obs_diff
Out[37]: -0.061578238965355567
```

In [38]: df2.info()

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 290584 entries, 0 to 294477
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 int64
 converted     290584 non-null int64
 dtypes: int64(2), object(3)
memory usage: 13.3+ MB
```

```
In [39]: (p_diffs-obs_diff).mean()
Out[39]: 0.98229999999999999
```

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?

Put your answer here. The value in part j is the **p-value** that we obtain after comparing the mean of observed difference mean and the mean obtained by simulated difference between new and old converted pages). **diffs**.

Observed difference `mean(obs_diff)` is the what we get by observing the actual values from the dataset whereas `p_diffs` we get by simulating the transactions for 10000 times and then calculating the difference in means of `new_page_converted` and `old_page_converted`.

Now, By `obs_diff` being negative, it suggests us that the null hypothesis can't be rejected implying that old page is better to keep than implementing a new ones.

Secondly, the **p-value** being much greater than 0.05 suggests the same that we should stick to Null hypothesis which states that the old page is better than the new ones.

Moreover, we did not build our histogram here to simulate distribution under null hypothesis unlike before due to the fact that, **p\_new** and **p\_old** are calculated assuming to be under null hypothesis.

1. 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 walk-through of the idea so that you can be 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 number of rows associated with the old page and new pages, respectively.

```
In [40]: import statsmodels.api as sm
convert_old = len(df2[(df2['landing_page']=='old_page') & (df2['converted']==1)])
convert_new = len(df2[(df2['landing_page']=='new_page') & (df2['converted']==1)])
n_old = len(df2.query('landing_page=="old_page"'))
n_new = len(df2.query('landing_page=="new_page"'))

/opt/conda/lib/python3.6/site-packages/statsmodels/compat/pandas.py:56: FutureWarning: The pa
ndas.core.iat[i] attribute is deprecated and will be removed in a future version. Please use
the pandas.iat[i] attribute instead.
from pandas.core import iatools
```

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 [41]: from scipy.stats import norm
z_score, p_value = sm.stats.proportions_ztest([convert_old, convert_new], [n_old, n_new], al
ternative="smaller")
print(z_score)
print(p_value)
```

```
1.31892419842
0.965959312759
```

n. What do the z-score and p-value in the previous question mean for the conversion rates of the old and new pages? Do they do the z-score and p-values in parts j, and k?

Put your answer here.

Yes, they agree to the findings in j and k. We know that when a p-value is greater than the value of alpha that is 0.05, we fail to reject the Null hypothesis.

Also, considering the z-score, if the z-score has between -1.96 and +1.96 (1.31 then), the p-value is likely to be larger than 0.05(0.905 then) implies that you can reject the null hypothesis.

### Part III - A regression approach

1. In this first 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?

Put your answer here.

**Undoubtedly, 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`, which is 1 when an individual receives the `treatment` and 0 if `control`.

```
In [42]: df2.head(4)
Out[42]:
```

	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 01:52:26.218027	control	old_page	1

```
In [43]: df2.head(5)
Out[43]:
```

	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 01:52:26.218027	control	old_page	1

```
In [44]: df2['intercept'] = 1
df2['ab_page'] = pd.get_dummies(df2['group'])['treatment']
Out[44]:
```

	user_id	timestamp	group	landing_page	converted	intercept	ab_page
0	851104	2017-01-21 22:11:48.556739	control	old_page	0	1	0
1	804228	2017-01-12 08:01:45.159739	control	old_page	0	1	0
2	661590	2017-01-11 16:55:06.154213	treatment	new_page	0	1	1
3	853541	2017-01-08 18:28:03.143765	treatment	new_page	0	1	1
4	864975	2017-01-21 01:52:26.218027	control	old_page	1	1	0

```
In [47]: df2.head(5)
Out[47]:
```

	user_id	timestamp	group	landing_page	converted	intercept	ab_page
0	851104	2017-01-21 22:11:48.556739	control	old_page	0	1	0
1	804228	2017-01-12 08:01:45.159739	control	old_page	0	1	0
2	661590	2017-01-11 16:55:06.154213	treatment	new_page	0	1	1
3	853541	2017-01-08 18:28:03.143765	treatment	new_page	0	1	1
4	864975	2017-01-21 01:52:26.218027	control	old_page	1	1	0

```
In [48]: df2.head(6)
Out[48]:
```

	user_id	timestamp	group	landing_page	converted	intercept	ab_page
0	851104	2017-01-21 22:11:48.556739	control	old_page	0	1	0
1	804228	2017-01-12 08:01:45.159739	control	old_page	0	1	0
2	661590	2017-01-11 16:55:06.154213	treatment	new_page	0	1	1
3	853541	2017-01-08 18:28:03.143765	treatment	new_page	0	1	1
4	864975	2017-01-21 01:52:26.218027	control	old_page	1	1	0
5	936923	2017-01-10 15:20:49.083499	control	old_page	0	1	0

c. Use `statsmodels` to instantiate your regression model on the two columns you created in part b., then fit the model using the two columns you created in part b. to predict whether or not an individual converts.

```
In [49]: import statsmodels.api as sm
log_model=sm.Log
```