

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](https://review.udacity.com/#!/projects/37e27304-ad47-4eb0-a1ab-8c12f60e43d0/rubric) (<https://review.udacity.com/#!/projects/37e27304-ad47-4eb0-a1ab-8c12f60e43d0/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 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 [1]: import pandas as pd
import numpy as np
import random
import matplotlib.pyplot as plt
import statsmodels.api as sm

%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(5)
```

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

```
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
```

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

```
In [4]: df.shape
```

Out[4]: (294478, 5)

c. The number of unique users in the dataset.

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

Out[5]: 290584

```
In [6]: df[df.duplicated(subset = ["user_id"])].shape
```

Out[6]: (3894, 5)

```
In [7]: df[df.duplicated(subset = ["user_id"])].info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 3894 entries, 2656 to 294355
Data columns (total 5 columns):
user_id      3894 non-null int64
timestamp    3894 non-null object
group        3894 non-null object
landing_page  3894 non-null object
converted     3894 non-null int64
dtypes: int64(2), object(3)
memory usage: 182.5+ KB
```

d. The proportion of users converted.

```
In [8]: df.converted.mean()
```

```
Out[8]: 0.11965919355605512
```

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

```
In [9]: df.groupby(["group", "landing_page"]).count()
```

```
Out[9]:
```

		user_id	timestamp	converted
group	landing_page			
control	new_page	1928	1928	1928
	old_page	145274	145274	145274
treatment	new_page	145311	145311	145311
	old_page	1965	1965	1965

f. Do any of the rows have missing values?

```
In [10]: df.isnull().values.any()
```

```
Out[10]: False
```

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 [11]: #df2 = df.drop(df[((df['group'] == 'treatment') == (df['landing_page'] == 'new_page')) == False].index)
df2=df
```

```
In [12]: df2 = df[((df.group=='treatment') & (df.landing_page=='new_page')) | ((df.group=='control') & (df.landing_page=='old_page'))]
```

```
In [13]: # 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[13]: 0

```
In [14]: #df2 = df[((df['group'] == 'control') == (df['landing_page'] == 'new_page')) == False].count()
#df2
```

```
In [15]: #df = df.drop(df[df.duplicated(subset = ["user_id"])].index)
#df[df.duplicated(subset = ["user_id"], keep = False)]
```

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 [16]: df2.shape
```

Out[16]: (290585, 5)

```
In [17]: df2.user_id.nunique()
```

Out[17]: 290584

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

```
In [18]: df2.user_id[df2.user_id.duplicated()]
```

Out[18]: 2893 773192
Name: user_id, dtype: int64

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

```
In [19]: df2[((df2['user_id'] == 773192))]
```

```
Out[19]:
```

	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

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

```
In [20]: df2 = df2.drop_duplicates()

df2.drop_duplicates(subset = "user_id",
                    keep = False, inplace = True)
df2.user_id[df2.user_id.duplicated()]
```

```
Out[20]: Series([], Name: user_id, dtype: int64)
```

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 [21]: converted = df2.converted.mean()
converted
```

```
Out[21]: 0.11959749882133504
```

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

```
In [22]: controlledConverted = df2.query("group == 'control'").converted.mean()
controlledConverted
```

```
Out[22]: 0.1203863045004612
```

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

```
In [23]: treatedConverted = df2.query("group == 'treatment'").converted.mean()
treatedConverted
```

```
Out[23]: 0.11880888313869065
```

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

```
In [24]: convertedProbability = len(df2.query("landing_page == 'new_page'))/len(df2)
convertedProbability
```

```
Out[24]: 0
```

```
In [25]: num_conv_treat = df2.query("group == 'treatment' and converted == 1").count()[0]
num_conv_treat
```

Out[25]: 17264

```
In [26]: df2.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 290583 entries, 0 to 294477
Data columns (total 5 columns):
user_id      290583 non-null int64
timestamp    290583 non-null object
group        290583 non-null object
landing_page  290583 non-null object
converted     290583 non-null int64
dtypes: int64(2), object(3)
memory usage: 13.3+ MB
```

```
In [27]: treatedConvertedProbability = len(df2.query("group == 'treatment' and converted == 1"))/len(df2)
treatedConvertedProbability
```

Out[27]: 0

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.

****Probability that a treated individual is converted is 0.059. Clearly if new page is used chances are more that they will click.**

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.

$$\begin{aligned} H_0 : p_{new} - p_{old} &\leq 0 \\ H_1 : p_{new} - p_{old} &> 0 \end{aligned}$$

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 [28]: p_new = len(df2.query( 'converted==1'))/len(df2.index)
p_new
```

```
Out[28]: 0
```

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

```
In [29]: p_old = len(df2.query('converted==1'))/len(df2.index)
p_old
```

Out[29]: 0

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

```
In [30]: p_diff=p_new-p_old
```

```
In [31]: n_new = len(df2.query('landing_page=="new_page"'))
n_new
```

Out[31]: 145309

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

```
In [32]: n_old = len(df2.query('landing_page=="old_page"'))
n_old
```

Out[32]: 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 [33]: new_page_converted = np.random.choice([0, 1], n_new, p = [p_new, 1-p_new])
new_page_converted
```

Out[33]: array([1, 1, 1, ..., 1, 1, 1])

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 [34]: old_page_converted = np.random.choice([0, 1], n_old, p = [p_old, 1-p_old])
old_page_converted
```

Out[34]: array([1, 1, 1, ..., 1, 1, 1])

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


```
In [35]: obs_diff= new_page_converted.mean() - old_page_converted.mean()# differences c
omputed in from p_new and p_old
obs_diff
```

Out[35]: 0.0

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**.

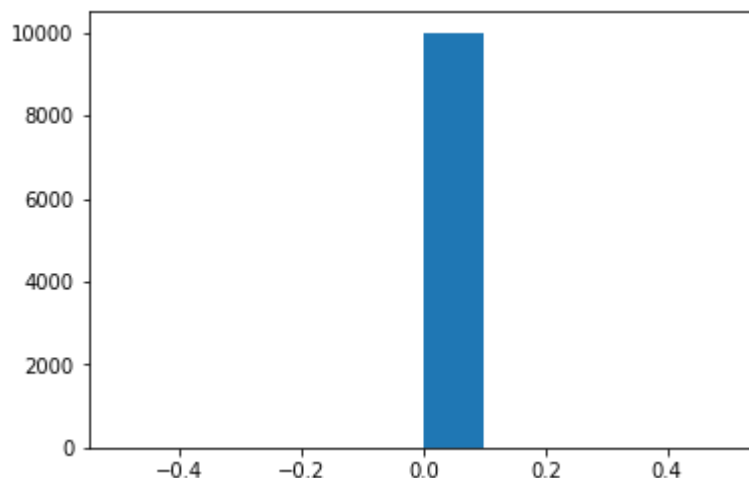
```
In [36]: p_diffs = []
for i in range(10000):
    p_new1 = np.random.choice([1, 0],n_new,replace = True,p = [p_new, 1-p_new
])
    p_old1 = np.random.choice([1, 0],n_old,replace = True,p = [p_old, 1-p_old
])
    p_new2 = p_new1.mean()
    p_old2 = p_old1.mean()
    p_diffs.append(p_new2-p_old2)
p_diffs=np.array(p_diffs)
p_diffs
```

Out[36]: array([0., 0., 0., ..., 0., 0., 0.])

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.

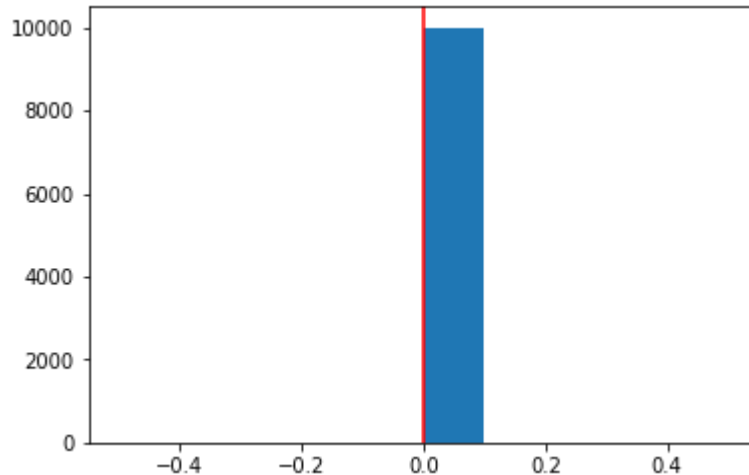
```
In [37]: plt.hist(p_diffs)
```

```
Out[37]: (array([ 0.,  0.,  0.,  0.,  0., 10000.,  0.,  0.,
 0.,  0.]),
array([-0.5, -0.4, -0.3, -0.2, -0.1,  0. ,  0.1,  0.2,  0.3,  0.4,  0.5]),
<a list of 10 Patch objects>)
```



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

```
In [38]: plt.hist(p_diffs)
plt.axvline(x= obs_diff, color='red');
```



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?

Null hypothesis is true and the new page is not performing better than old page.

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 [39]: convert_old = len(df2.query('converted==1 and landing_page=="old_page"'))
convert_new = len(df2.query('converted==1 and landing_page=="new_page"'))
n_old = len(df2.query('landing_page=="old_page"'))
n_new = len(df2.query('landing_page=="new_page"'))
n_new
```

Out[39]: 145309

m. Now use `stats.proportions_ztest` to compute your test statistic and p-value. [Here \(https://docs.w3cub.com/statsmodels/generated/statsmodels.stats.proportion.proportions_ztest/\)](https://docs.w3cub.com/statsmodels/generated/statsmodels.stats.proportion.proportions_ztest/) is a helpful link on using the built in.

```
In [40]: z_score, p_value = sm.stats.proportions_ztest([convert_old,convert_new], [n_old, n_new],alternative='smaller')
p_value
```

```
Out[40]: 0.9049428161159749
```

```
In [41]: z_score
```

```
Out[41]: 1.3102408579271012
```

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.**?

Again Z score is less than 95% confidence interval. Hence we reject the hypothesis.

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 [42]: df2['intercept']=1
ab_page = ['treatment', 'control']
df2['ab_page'] = pd.get_dummies(df2.group)['treatment']
```

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 [43]: import statsmodels.api as sm
model=sm.Logit(df2['converted'],df2[['intercept','ab_page']])
results=model.fit()
```

Optimization terminated successfully.
Current function value: 0.366119
Iterations 6

d. Provide the summary of your model below, and use it as necessary to answer the following questions.

```
In [44]: from scipy import stats
stats.chisqprob = lambda chisq, df: stats.chi2.sf(chisq, df)
results.summary()
```

Out[44]: Logit Regression Results

Dep. Variable:	converted	No. Observations:	290583
Model:	Logit	Df Residuals:	290581
Method:	MLE	Df Model:	1
Date:	Wed, 18 Mar 2020	Pseudo R-squ.:	8.068e-06
Time:	20:59:42	Log-Likelihood:	-1.0639e+05
converged:	True	LL-Null:	-1.0639e+05
Covariance Type:	nonrobust	LLR p-value:	0.1901

	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	-1.310	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**?

p-value = 0.19 which is less than the p-value calculated using the z-score function.

Lets define the hypothesis

In Logistic regression

$$H_0 : p_{new} - p_{old} = 0$$

$$H_1 : p_{new} - p_{old} \neq 0$$

Part 2

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

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

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?

**The conversion rate may be dependednt on users age group, sex, social status etc. With the current dataset we such parameters are not avaiable.

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 \(https://pandas.pydata.org/pandas-docs/stable/generated/pandas.DataFrame.join.html\)](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 [45]: countries_df = pd.read_csv('./countries.csv')
df_new = countries_df.set_index('user_id').join(df2.set_index('user_id'))
df_new.head()
```

Out[45]:

	country	timestamp	group	landing_page	converted	intercept	ab_page
user_id							
834778	UK	2017-01-14 23:08:43.304998	control	old_page	0.0	1.0	0.0
928468	US	2017-01-23 14:44:16.387854	treatment	new_page	0.0	1.0	1.0
822059	UK	2017-01-16 14:04:14.719771	treatment	new_page	1.0	1.0	1.0
711597	UK	2017-01-22 03:14:24.763511	control	old_page	0.0	1.0	0.0
710616	UK	2017-01-16 13:14:44.000513	treatment	new_page	0.0	1.0	1.0

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.

```
df_new.country.unique()
```

```
In [46]: df_new[['CA', 'US']] = pd.get_dummies(df_new['country'])[['CA', 'US']]
df_new.head()
```

Out[46]:

	country	timestamp	group	landing_page	converted	intercept	ab_page	CA	U
user_id									
834778	UK	2017-01-14 23:08:43.304998	control	old_page	0.0	1.0	0.0	0	
928468	US	2017-01-23 14:44:16.387854	treatment	new_page	0.0	1.0	1.0	0	
822059	UK	2017-01-16 14:04:14.719771	treatment	new_page	1.0	1.0	1.0	0	
711597	UK	2017-01-22 03:14:24.763511	control	old_page	0.0	1.0	0.0	0	
710616	UK	2017-01-16 13:14:44.000513	treatment	new_page	0.0	1.0	1.0	0	

```
In [49]: df_new['intercept'] = 1
log_mod = sm.Logit(df_new['converted'], df_new[['CA', 'US', 'intercept', 'ab_p
age']], missing='drop')
results = log_mod.fit()
results.summary()
```

Optimization terminated successfully.
Current function value: 0.366114
Iterations 6

Out[49]:

Logit Regression Results

Dep. Variable:	converted	No. Observations:	290583
Model:	Logit	Df Residuals:	290579
Method:	MLE	Df Model:	3
Date:	Wed, 18 Mar 2020	Pseudo R-squ.:	2.322e-05
Time:	21:07:55	Log-Likelihood:	-1.0639e+05
converged:	True	LL-Null:	-1.0639e+05
Covariance Type:	nonrobust	LLR p-value:	0.1761

	coef	std err	z	P> z	[0.025	0.975]
CA	-0.0506	0.028	-1.784	0.074	-0.106	0.005
US	-0.0099	0.013	-0.743	0.458	-0.036	0.016
intercept	-1.9794	0.013	-155.415	0.000	-2.004	-1.954
ab_page	-0.0149	0.011	-1.306	0.191	-0.037	0.007

**The p-value is higher in US than in Canada, which means that users in the US are more likely to convert, but as the p-value is less than .5 we can not reject the hypothesis.

```
In [52]: from subprocess import call  
         call(['python', '-m', 'nbconvert', 'Analyze_ab_test_results_notebook.ipynb'])
```

```
Out[52]: -1
```