

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](#).

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 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 [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]:

```
#Checking the number of rows of dataset
df.shape[0]
```

Out[3]:

294478

c. The number of unique users in the dataset.

In [4]:

```
#Checking the unique users
df['user_id'].nunique()
```

Out[4]:

290584

d. The proportion of users converted.

In [5]:

```
# The converted Rate
df['converted'].mean()
```

Out[5]:

0.11965919355605512

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

In [6]:

```
(df.query("group == 'treatment' & landing_page == 'old_page'").shape[0]) + (df.query("group == 'control' & landing_page == 'new_page'").shape[0])
```

Out[6]:

3893

f. Do any of the rows have missing values?

In [7]:

```
#Checking null Values
df.isnull().sum()
```

Out[7]:

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

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]:

```
#Gathering the data with the correct data.
```

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

In [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]
```

Out[9]:

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 [10]:

```
#checking the change on dataset
df2['user_id'].nunique()
```

Out[10]:

290584

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

In [11]:

```
#checking repeated id from users.
df2.user_id.astype(str).value_counts()
```

Out[11]:

773192	2
775180	1
740710	1
848318	1
924914	1
712246	1
735110	1
926205	1
660220	1
870432	1
872120	1
760981	1
675903	1
884337	1
638410	1
926711	1
835964	1
747944	1
731600	1
631119	1
832489	1
922540	1
651439	1
747623	1
935032	1
812129	1
867930	1
776442	1
895518	1
676836	1
..	
667752	1
786776	1
712119	1
704107	1
911532	1
723414	1
845914	1

```

843917 1
918521 1
641042 1
789476 1
886737 1
939244 1
743959 1
677622 1
668124 1
799734 1
824050 1
796090 1
923785 1
909700 1
635496 1
648323 1
815832 1
712707 1
680345 1
832649 1
726592 1
817987 1
860740 1
897371 1
Name: user_id, Length: 290584, dtype: int64

```

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

In [12]:

```

#Checking the repeated user ID
df2.query("user_id == '773192'")

```

Out[12]:

	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 [13]:

```

#removing the repeated id and checking the change
c = df2.query("timestamp == '2017-01-14 02:55:59.590927'").index
df2 = df2.drop(c)

df2.shape

```

Out[13]:

(290584, 5)

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 [14]:

```

# The Conversion rate regardless of the page they receive
df2.converted.mean()

```

Out[14]:

0.11959708724499628

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

In [15]:

```
#Conversion rate of control group
p_old = df2.query("group == 'control'").converted.mean()
p_old
```

Out[15]:

0.1203863045004612

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

In [16]:

```
#Conversion rate of treatment group
p_new = df2.query("group == 'treatment'").converted.mean()
p_new
```

Out[16]:

0.11880806551510564

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

In [17]:

```
#The probability that an individual received the new page
df2.query("landing_page == 'new_page'").user_id.count()/df2.user_id.count()
```

Out[17]:

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.

There isn't any evidence that treatment page leads to more conversions, since the conversion rate of group control is slightly higher

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.

$\$H_0: P_{new} - P_{old} \leq 0$ $\$H_1: P_{new} - P_{old} > 0$

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 [18]:

```
# that is pnewpnew and poldpold are equal. Furthermore, assume they are equal to the converted
rate in ab_data.csv regardless of the page.
p_new = df2.converted.mean()
p_new
```

Out[18]:

0.11959708724499628

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

In [19]:

```
# that is pnewpnew and poldpold are equal. Furthermore, assume they are equal to the converted
rate in ab_data.csv regardless of the page.
p_old = df2.converted.mean()
p_old
```

Out[19]:

0.11959708724499628

c. What is n_{new} ?

In [20]:

```
# Nnew is equal to ab_data csv
n_new = df2.query("landing_page == 'new_page'").shape[0]
```

d. What is n_{old} ?

In [21]:

```
# Nold is equal to ab_data csv
n_old = df2.query("landing_page == 'old_page'").shape[0]
```

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 [22]:

```
#Since it's a binomial choice we will use np.random.choice
new_page_converted = np.random.choice([0,1],n_new, p=(1 - p_new, p_new))
new_page_converted.mean()
```

Out[22]:

0.1200330328263712

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 [23]:

```
#Since it's a binomial choice we will use np.random.choice
old_page_converted = np.random.choice([0,1],n_old, p=( 1- p_old, p_old))
old_page_converted.mean()
```

Out[23]:

0.1192987045169817

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

In [24]:

```
obs_diff = new_page_converted.mean() - old_page_converted.mean()
```

h. Simulate 10,000 $p_{\text{new}} - p_{\text{old}}$ 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 [25]:

```
p_diffs = []
size = df.shape[0]

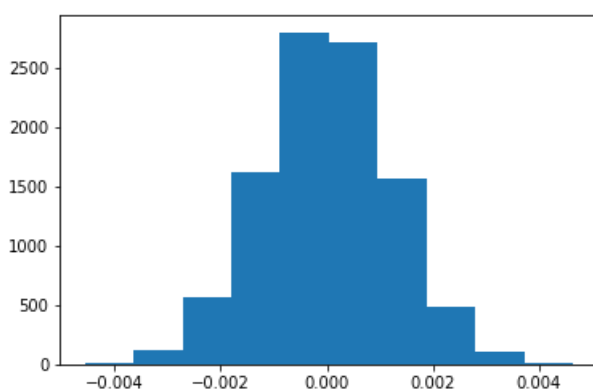
for _ in range(10000):

    new_page_converted = np.random.choice([0,1],n_new, p=(p_new, 1-p_new))
    old_page_converted = np.random.choice([0,1],n_old, p=(p_old, 1-p_old))
    p_diffs.append(new_page_converted.mean() - old_page_converted.mean())
```

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 [26]:

```
p_diffs = np.array(p_diffs)
plt.hist(p_diffs);
```



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

In [27]:

```
(p_diffs > obs_diff).mean()
```

Out[27]:

0.273

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?

This Value is called P-value, since the value is 0.6228, it suggests we shouldn't move away from the null hypothesis, in

This value is called a **p-value**, since the value is 0.0220, it suggests we shouldn't move away from the null hypothesis, in other words, we fail to reject the null

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 [28]:

```
import statsmodels.api as sm

convert_old = df2.query("landing_page == 'old_page'").converted.sum()
convert_new = df2.query("landing_page == 'new_page'").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](#) is a helpful link on using the built in.

In [29]:

```
count = convert_old
nobs = n_old
value = convert_old/n_old

stats, pval = sm.stats.proportions_ztest(count,nobs,value)
print('{0:0.3f}'.format(pval))
```

1.000

In [30]:

```
count = convert_new
nobs = n_new
value = convert_new/n_new

stats, pval = sm.stats.proportions_ztest(count,nobs,value)
print('{0:0.3f}'.format(pval))
```

1.000

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.?

Since the p-value is 1.0, it means, this data agree with the "k" question, we fail to reject the null hypothesis

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?

It must be a logistic regression, since we have 0 from not converted and 1 of converted.

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 [31]:

```
#checking the dataset
```



```
df2.head()
```

Out[31]:

	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 [32]:

```
#Creating Dummies Variable
df2[['new_page','old_page']] = pd.get_dummies(df2['landing_page'])
df2.head()
df2 = df2.drop('old_page',axis=1)
```

In [33]:

```
#Rename to "ab_page"
df2.rename(columns = {'new_page':'ab_page'}, inplace = True)
df2.head()
```

Out[33]:

	user_id	timestamp	group	landing_page	converted	ab_page
0	851104	2017-01-21 22:11:48.556739	control	old_page	0	0
1	804228	2017-01-12 08:01:45.159739	control	old_page	0	0
2	661590	2017-01-11 16:55:06.154213	treatment	new_page	0	1
3	853541	2017-01-08 18:28:03.143765	treatment	new_page	0	1
4	864975	2017-01-21 01:52:26.210827	control	old_page	1	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 [34]:

```
df2['intercept'] = 1
log_mod = sm.Logit(df2['converted'], df2[['intercept', 'ab_page']])
results = log_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 [35]:

```
results.summary()
```

Out[35]:

Logit Regression Results

Dep. Variable:	converted	No. Observations:	290584
Model:	Logit	Df Residuals:	290582

Method:	MLE	Df Model:	1
Date:	Wed, 09 Oct 2019	Pseudo R-squ.:	8.077e-06
Time:	23:31:36	Log-Likelihood:	-1.0639e+05
converged:	True	LL-Null:	-1.0639e+05
		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

In [36]:

```
#converted is 1.01 times as likely on old page than new page holding all else constant
np.exp(-0.0150) , 1/np.exp(-0.0150)
```

Out[36]:

```
(0.9851119396030626, 1.015113064615719)
```

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 coefficient **ab_page** is -0.0150, or 0.98 on exponential, this means converted is 0.98 times as likely on new_page than old page holding all else constant, it means the old page slightly better than new page, like as we found on part two. This means we fail to reject the null hypothesis. About the p-value, the value is 0.1899, it's higher to accept the alternative hypothesis.

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's good to add some new features, because we can find other relevant variables to our model. The problems with adding new features can be multicollinearity, correlated errors, and even overfitting.

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.

In [37]:

```
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[37]:

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

In [38]:

```
df_new['country'].value_counts()
```

Out[38]:

```
US    203619
UK     72466
CA     14499
Name: country, dtype: int64
```

In [39]:

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

Out[39]:

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

In [40]:

```
### Fit Your Linear Model And Obtain the Results
df_new['intercept'] = 1
log_mod = sm.Logit(df_new['converted'], df_new[['intercept', 'ab_page', 'US', 'UK']])
results = log_mod.fit()
results.summary()
```

```
Optimization terminated successfully.
      Current function value: 0.366113
      Iterations 6
```

Out[40]:

Logit Regression Results

Dep. Variable:	converted	No. Observations:	290584
Model:	Logit	Df Residuals:	290580
Method:	MLE	Df Model:	3
Date:	Wed, 09 Oct 2019	Pseudo R-squ.:	2.323e-05
Time:	23:31:40	Log-Likelihood:	-1.0639e+05
converged:	True	LL-Null:	-1.0639e+05
		LLR p-value:	0.1760

	coef	std err	z	P> z	[0.025	0.975]
intercept	0.0000	0.0000	76.040	0.000	0.000	1.070

Intercept	-2.0300	0.027	-76.249	0.000	-2.082	-1.978
ab_page	-0.0149	0.011	-1.307	0.191	-0.037	0.007
US	0.0408	0.027	1.516	0.130	-0.012	0.093
UK	0.0506	0.028	1.784	0.074	-0.005	0.106

In [41]:

```
np.exp(-0.0155) , np.exp(0.0357) , np.exp(0.0449)
```

Out[41]:

```
(0.9846195067517329, 1.036344896381825, 1.045923262352691)
```

The countries didn't make a great change on p-value and, the coefficients and ab_data change very slightly. We still fail to reject the null. The alternative hypothesis- the new page- isn't better than the old page.