

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

import seaborn as sns
sns.set_style('darkgrid')

#We are setting the seed to assure you get the same answers on quizzes as we set up
random.seed(1)
```

1. Now, read in the `ab_data.csv` data. Store it in `df`. Use your **Quiz 1** of 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.shape()
```

```
Out[3]:
```

c. The number of unique users in the dataset.

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

```
Out[4]:
```

d. The proportion of users converted.

```
In [5]: converted = df.query('converted == 1')
converted.shape()[0]/df.shape()[0]
```

```
Out[5]:
```

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

```
In [6]: df.query('group == "treatment" and landing_page != "new_page").shape()[0] + \
df.query('group == "treatment" and landing_page == "new_page").shape()[0]
```

```
Out[6]:
```

f. Do any of the rows have missing values?

```
In [7]: df.isna().sum()
```

```
Out[7]:
```

```
class: pandas.core.frame.DataFrame
RangeIndex: 294478 entries, 0 to 294477
Data columns (total 5 columns):
#   Column                Non-Null Count  Dtype
---  --
0   user_id                294478 non-null  int64
1   timestamp              294478 non-null  object
2   group                  294478 non-null  object
3   landing_page           294478 non-null  object
4   converted              294478 non-null  int64
dtypes: int64(2), object(3)
memory usage: 11.2+ MB
```

No, there isn't missing values.

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 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]: df2 = df
df2.drop(df.query('group == "treatment" and landing_page == "old_page") or (group == "control" and landing_page == "new_page")').index, inplace=True)
```

```
In [9]: #A Double Check all of the correct rows were removed - this should be 0
df2[(df2['group'] == 'treatment') == (df2['landing_page'] == 'new_page')] ==
df2.shape()[0]
```

```
Out[9]:
```

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

```
Out[10]:
```

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

```
In [11]: df2[df2.duplicated(['user_id'], keep=False)][['user_id']].unique()
```

```
Out[11]:
```

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

```
In [12]: df2[df2.duplicated(['user_id'], keep=False)]
```

```
Out[12]:
```

	user_id	timestamp	group	landing_page	converted
1899	773192	2017-01-09 05:57:58.781806	treatment	new_page	0
2893	773192	2017-01-14 02:55:59.595927	treatment	new_page	0

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

```
In [13]: df2.drop(df2.loc[[1899]].index, 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 [14]: prob_convert = df2.query('converted == 1').shape()[0] / df2.shape()[0]
prob_convert
```

```
Out[14]:
```

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

```
In [15]: control_population = df2.query('group == "control"')
control_converted = df2.query('converted == 1 and group == "control"')
control_converted.shape()[0] / control_population.shape()[0]
```

```
Out[15]:
```

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

```
In [16]: treatment_population = df2.query('group == "treatment"')
treatment_converted = control_converted = df2.query('converted == 1 and group == "treatment"')
treatment_converted.shape()[0] / treatment_population.shape()[0]
```

```
Out[16]:
```

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

```
In [17]: #For further analysis, we are querying new individual who received old page
received_new_page = df2.query('landing_page == "new_page"')
received_old_page = df2.query('landing_page == "old_page"')

received_new_page.shape()[0] / df2.shape()[0]
```

```
Out[17]:
```

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.

Once control group conversion rate is higher than treatment group, we fail to reject the null hypothesis.

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.

$$H_0 : p_{old} \geq p_{new}$$
$$H_1 : p_{old} < p_{new}$$

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

```
Out[18]:
```

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

```
In [19]: p_old = df2['converted'].mean()
p_old
```

```
Out[19]:
```

c. What is n_{new} ?

```
In [20]: n_new = received_new_page.shape()[0]
n_new
```

```
Out[20]:
```

d. What is n_{old} ?

```
In [21]: n_old = received_old_page.shape()[0]
n_old
```

```
Out[21]:
```

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]: np.random.binomial(n_new, p_new)
```

```
Out[22]:
```

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]: np.random.binomial(n_old, p_old)
```

```
Out[23]:
```

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

```
In [24]: p_new - p_old
```

```
Out[24]:
```

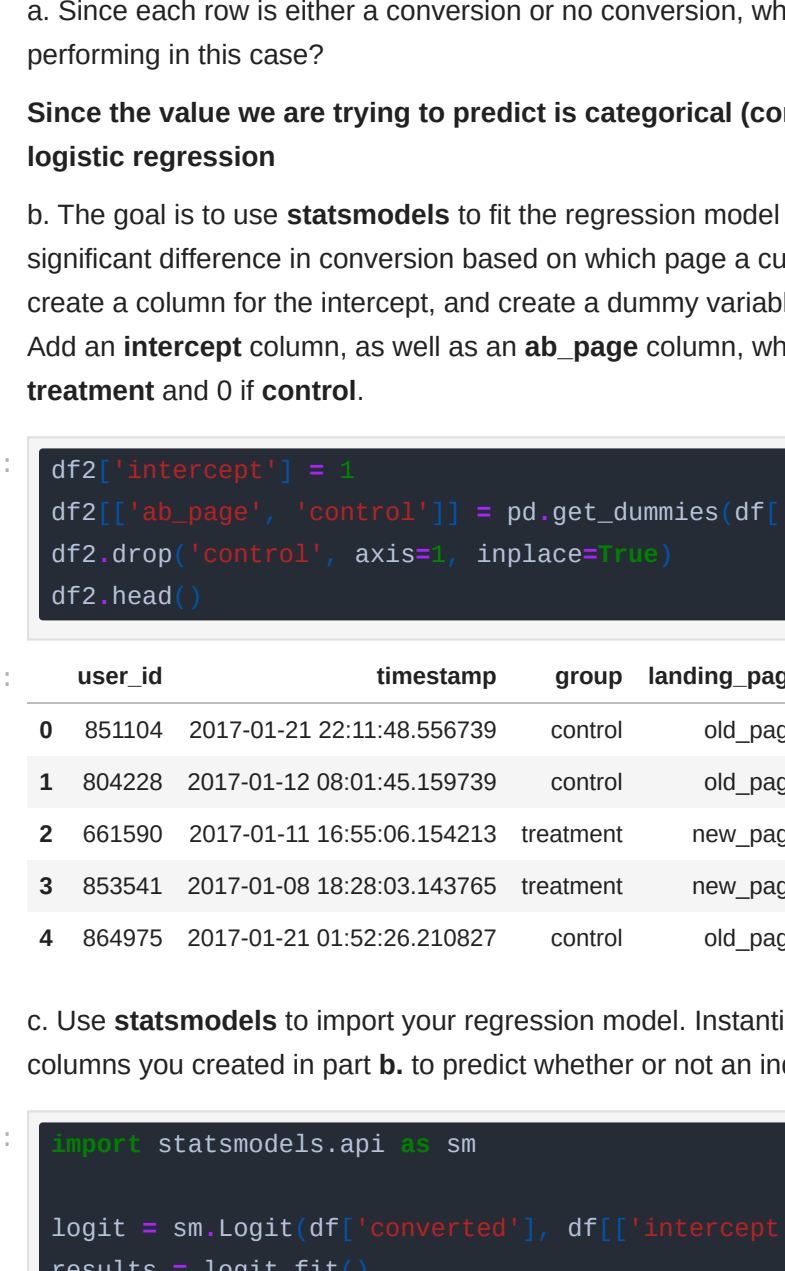
h. Simulate 10,000 $p_{new} - p_{old}$ values using this same process similarly to the one you calculated in parts a. through g. above. Store all 10,000 values in a numpy array called `p_diffs`.

```
In [25]: p_diffs = []

for _ in range(int(10000)):
    new_page_converted = np.random.binomial(n_new, p_new)
    old_page_converted = np.random.binomial(n_old, p_old)
    p_diff = new_page_converted/n_new - old_page_converted/n_old
    p_diffs.append(p_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 [26]: plt.title('Simulated p-values')
plt.xlabel('Frequency')
plt.ylabel('p-values')
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_diff_orig = (treatment_population['converted'].mean() -
control_population['converted'].mean())
p_diff_orig
```

```
Out[27]:
```

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?

A p-value is the probability of observing your statistic if the null hypothesis is true.

The null hypothesis was that the difference in means would be equal or less than 0, and the alternative was the difference would be greater than 0. However, the difference is less than zero, and the p-value is very large. We do not have evidence 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 [29]: import statsmodels.api as sm

convert_old = df2.query('landing_page == "old_page" & converted == True').shape[0]
convert_new = df2.query('landing_page == "new_page" & converted == True').shape[0]

#These value was already calculated above
n_old = n_old
n_new = n_new
```

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

```
stats: 1.310241094243494
p_value: 0.905058517596945
```

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?

```
In [31]: from scipy.stats import norm

#calculate the critical z term
z_critical = norm.ppf(1-(0.05))
print('Zscore: ', stats)
print('Critical Zscore: ', z_critical)
```

```
stats: 1.310241094243494
critical Zscore: -1.6448536269514722
```

Since both p-values are equal and larger than 5% (0.05), and our Z-score are smaller than critical z-score, then we fail to reject the null, that is: old page is more efficient to convert than new page.

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?

Since the value we are trying to predict is categorical (convert or non converted), we should use a 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 [32]: df2['intercept'] = 1
df2[['ab_page', 'control']] = pd.get_dummies(df['group'])
df2.drop('control', axis=1, inplace=True)
df2.head()
```

```
Out[32]:
```

	user_id	timestamp	group	landing_page	converted	intercept	ab_page
0	851104	2017-01-21 22:11:48.556739	control	old_page	0	1	1
1	804228	2017-01-12 08:01:45.159739	control	old_page	0	1	1
2	661590	2017-01-11 16:55:06.154213	treatment	new_page	0	1	0
3	853541	2017-01-08 18:28:03.143765	treatment	new_page	0	1	0
4	864975	2017-01-21 01:52:26.210827	control	old_page	1	1	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 [33]: import statsmodels.api as sm

logit = sm.Logit(df['converted'], df[['intercept', 'ab_page']])
results = logit.fit()
```

```
Optimization terminated successfully.
Current function value: 8.366118
Iterations: 6
```

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

```
In [34]: results.summary()
```

```
Out[34]:
```

Logit Regression Results						
Dep. Variable:	converted	No. Observations:	290584			
Model:	Logit	DF Residuals:	290582			
Method:	MLE	DF Model:	1			
Date:	Mon, 28 Dec 2020	Pseudo R-squ.:	8.077e-06			
Time:	00:10:47	Log-Likelihood:	-1.0639e+05			
converged:	True	LL-Null:	-1.0639e+05			
Covariance Type:	nonrobust	LLR p-value:	0.1899			
coef	std err	z	P> z	[0.025	0.975]	
intercept	-2.0038	0.008	-247.146	0.000	-2.020	-1.988
ab_page	0.0150	0.011	1.311	0.190	-0.007	0.037

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**?

Answer

The p-value now is 0.19. The values is different from part II (0.90) because the hypotheses were different too. In part II, the null hypothesis was new page are better to convert than old pages.

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?

Answer.

Adding more terms will increase model complexity, because can affect others variables, but can increase R-squared. One of the feature that worth to try add is duration that individual spent on the page, perhaps time of day can be useful.

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 [35]: countries_df = pd.read_csv('./countries.csv')
df_new = countries_df.set_index('user_id').join(df2.set_index('user_id'), how='inner')
```

```
In [36]: df_new['country']
```

```
Out[36]:
```

user_id	country
834778	UK
928468	US
822059	UK
711597	UK
710616	UK
653118	US
878226	UK
799368	UK
655535	CA
934996	UK

```
In [37]: df_new[['CA', 'UK', 'US']] = pd.get_dummies(df_new['country'])
df_new[['ab_page', 'control']] = pd.get_dummies(df_new['group'])

#Let's use US as baseline country, so we don't need to drop, as well as control.
df_new.drop(['US', 'control'], axis=1, inplace=True)
df_new
```

```
Out[37]:
```

	user_id	country	timestamp	group	landing_page	converted	ab_page	CA	UK		
	834778	UK	2017-01-14 23:08:43.304988	control	old_page	0	1	1	0	1	
	928468	US	2017-01-23 14:44:16.387854	treatment	new_page	0	1	0	0	0	
	822059	UK	2017-01-16 14:04:14.719771	treatment	new_page	1	1	0	0	1	
	711597	UK	2017-01-22 03:14:24.763511	control	old_page	0	1	1	0	1	
	710616	UK	2017-01-16 13:14:44.000513	treatment	new_page	0	1	0	0	1	
...	
	653118	US	2017-01-09 03:12:31.034796	control	old_page	0	1	1	1	0	0
	878226	UK	2017-01-05 15:02:50.334962	control	old_page	0	1	1	1	0	1
	799368	UK	2017-01-09 18:07:34.259353	control	old_page	0	1	1	1	0	1
	655535	CA	2017-01-09 13:30:47.524512	treatment	new_page	0	1	0	1	0	1
	934996	UK	2017-01-09 00:30:08.377677	control	old_page	0	1	1	0	1	1

290584 rows x 9 columns

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 [38]: #Creating additional column with country and page column
df_new['CA_page'] = df_new['CA'] * df_new['ab_page']
df_new['UK_page'] = df_new['UK'] * df_new['ab_page']
df_new.head()
```

```
Out[38]:
```

	user_id	country	timestamp	group	landing_page	converted	intercept	ab_page	CA	UK	CA_page	
	834778	UK	2017-01-14 23:08:43.304988	control	old_page	0	1	1	1	0	1	
	928468	US	2017-01-23 14:44:16.387854	treatment	new_page	0	1	0	0	0	0	
	822059	UK	2017-01-16 14:04:14.719771	treatment	new_page	1	1	0	0	1	0	
	711597	UK	2017-01-22 03:14:24.763511	control	old_page	0	1	1	1	0	1	
	710616	UK	2017-01-16 13:14:44.000513	treatment	new_page	0	1	0	0	1	0	
...	
	653118	US	2017-01-09 03:12:31.034796	control	old_page	0	1	1	1	1	0	0
	878226	UK	2017-01-05 15:02:50.334962	control	old_page	0	1	1	1	1	0	1
	799368	UK	2017-01-09 18:07:34.259353	control	old_page	0	1	1	1	1	0	1
	655535	CA	2017-01-09 13:30:47.524512	treatment	new_page	0	1	0	1	0	1	0
	934996	UK	2017-01-09 00:30:08.377677	control	old_page	0	1	1	0	1	0	1

290584 rows x 12 columns

i. 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 [39]: #Fitting Your Linear Model and Obtain the Results
logit = sm.Logit(df_new['converted'], df_new[['intercept', 'CA', 'CA_page', 'UK',
'UK_page', 'ab_page']])
results = logit.fit()
results.summary()
```

```
Out[39]:
```

Logit Regression Results						
Dep. Variable:	converted	No. Observations:	290584			
Model:	Logit	DF Residuals:	290578			
Method:	MLE	DF Model:	5			
Date:	Mon, 28 Dec 2020	Pseudo R-squ.:	3.482e-05			
Time:	00:10:51	Log-Likelihood:	-1.0639e+05			
converged:	True	LL-Null:	-1.0639e+05			
Covariance Type:	nonrobust	LLR p-value:	0.1920			
coef	std err	z	P> z	[0.025	0.975]	
intercept	-2.0070	0.010	-207.045	0.000	-2.026	-1.988
CA	-0.0644	0.038	-1.679	0.093	-0.140	0.011
CA_page	0.0469	0.054	0.872	0.383	-0.059	0.152
UK	0.0257	0.019	1.363	0.173	-0.011	0.063
UK_page	-0.0314	0.027	-1.181	0.238	-0.084	0.021
ab_page	0.0206	0.014	1.505	0.132	-0.006	0.047

We still fail to reject null hypothesis since conversion rate is similar. Also these interactions between countries and pages has larger p-values and could be decreasing significance of countries variables.

As reported before, add duration of user interaction would be interesting, but since in this dataset we just have timestamp, which means only the moment that user clicked on the page, we can't get that information.

Let's take a look in duration of test in dataset to see if we can discover some potential inside:

```
In [40]: #Converting timestamp column to datetime format
df_new['timestamp'] = pd.to_datetime(df_new['timestamp'])
df_new['timestamp'] = df_new['timestamp'].dt.date
print('Duration of test: ', df_new['timestamp'].max() - df_new['timestamp'].min())
```

```
Out[40]:
```

Duration of test: 22 days 05:05:00

```
In [67]: df_new['timestamp'].value_counts().sort_index()
```

```
Out[67]:
```

timestamp	count
2017-01-02	6712
2017-01-03	13268
2017-01-04	13119
2017-01-05	12932
2017-01-06	13553
2017-01-07	13213
2017-01-08	13387
2017-01-09</	