# Analyze A/B Test Results

This project will assure you have mastered the subjects covered in the statistics lessons. We have organized the current notebook into the following sections:

Specific programming tasks are marked with a **ToDo** tag.

## Introduction

A/B tests are very commonly performed by data analysts and data scientists. 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 webpage,
- Keep the old webpage, or
- Perhaps run the experiment longer to make their decision.

Each **ToDo** task below has an associated quiz present in the classroom. Though the classroom quizzes are **not necessary** to complete the project, they help ensure you are on the right track as you work through the project, and you can feel more confident in your final submission meeting the rubric specification.

## Part I - Probability

To get started, let's import our libraries.

```
In [3]:  import pandas as pd
         import numpy as np
         import random
         import matplotlib.pyplot as plt
         #We are setting the seed to assure you get the same answers on quizzes as we set up
         random.seed(42)
```

### ToDo 1.1

Now, read in the `ab_data.csv` data. Store it in `df`. Below is the description of the data, there are a total of 5 columns:

| Data columns | Purpose | Valid values |
|---|---|---|
| user_id | Unique ID | Int64 values |
| timestamp | Time stamp when the user visited the webpage | - |
| group | In the current A/B experiment, the users are categorized into two broad groups. | ['control', |

The `control` group users are expected to be served with `old_page` ; and `treatment` group users are matched with the `new_page` . However, **some inaccurate rows** are present in the initial data, such as a `control` group user is matched with a `new_page` .

`'treatment']`

| landing_page | It denotes whether the user visited the old or new webpage. | `['old_page',`<br>`'new_page']` |
| --- | --- | --- |
| converted | It denotes whether the user decided to pay for the company's product. Here, `1` means yes, the user bought the product. | `[0, 1]` |

</center> Use your dataframe to answer the questions in Quiz 1 of the classroom.

**a.** Read in the dataset from the `ab_data.csv` file and take a look at the top few rows here:

```
In [4]:  df=pd.read_csv('ab_data.csv')
         df.head()
```

Out[4]:

| | 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 cell below to find the number of rows in the dataset.

```
In [5]:  df.shape
```

Out[5]:
```
(294478, 5)
```

**c.** The number of unique users in the dataset.

```
In [6]:  df.user_id.nunique()
```

Out[6]:
```
290584
```

**d.** The proportion of users converted.

```
In [7]:  df.converted.mean()
```

Out[7]:
```
0.11965919355605512
```

**e.** The number of times when the "group" is `treatment` but "landing_page" is not a `new_page` .

```
In [8]:  Treatment_OldP=  df.query("group == 'treatment' and landing_page != 'new_page'").shape[0
         Treatment_NewP= df.query('group == "control" and landing_page == "new_page"')['landing_p
```

```
In [9]:  #The number of times when the "group" is treatment but "landing_page" is not a new_page
         Treatment_OldP
```

Out[9]:
```
1965
```

```
In [10]:  #The number of times when the "control" is treatment and "landing_page" is a new_page
          Treatment_NewP
```

`Out[10]:` 1928

**f.** Do any of the rows have missing values?

`In [11]:` `df.info()` `#No null values`

```
<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
```

## ToDo 1.2

In a particular row, the **group** and **landing_page** columns should have either of the following acceptable values:

| user_id | timestamp | group | landing_page | converted |
|---------|-----------|-------|--------------|-----------|
| XXXX | XXXX | control | old_page | X |
| XXXX | XXXX | treatment | new_page | X |

It means, the `control` group users should match with `old_page` ; and `treatment` group users should matched with the `new_page` .

However, for the rows where `treatment` does not match with `new_page` or `control` does not match with `old_page` , we cannot be sure if such rows truly received the new or old wepage.

Use **Quiz 2** in the classroom to figure out how should we handle the rows where the group and landing_page columns don't match?

**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 [12]:`
```
df2 = df #store df in df2 dataframe
df2.drop(df2.query("group == 'treatment' and landing_page == 'old_page'").index,inplace=
df2.drop(df2.query("group == 'control' and landing_page == 'new_page'").index,inplace=Tr
```

`In [13]:`
```
# Double Check all of the incorrect rows were removed from df2 -
# Output of the statement below should be 0
df2[((df2['group'] == 'treatment') == (df2['landing_page'] == 'new_page')) == False].sha
df2[((df2['group'] == 'control') == (df2['landing_page'] == 'old_page')) == False].shape
```

`Out[13]:` 0

## ToDo 1.3

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

**a.** How many unique **user_id**s are in **df2**?

```
In [14]:   df2.user_id.nunique()
```

```
Out[14]:   290584
```

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

```
In [15]:   df2[df2['user_id'].duplicated()]
```

Out[15]:

| | user_id | timestamp | group | landing_page | converted |
|---|---|---|---|---|---|
| **2893** | 773192 | 2017-01-14 02:55:59.590927 | treatment | new_page | 0 |

**c.** Display the rows for the duplicate **user_id**?

```
In [16]:   df2[df2.duplicated(['user_id'], keep=False)]
```

Out[16]:

| | 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**, from the **df2** dataframe.

```
In [17]:   df2 = df2.drop(df2[df2.duplicated(['user_id'])].index)
```

## ToDo 1.4

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 [18]:   Ppopulation = df2['converted'].mean()
           Ppopulation   #0.11959667567149027
```

```
Out[18]:   0.11959708724499628
```

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

```
In [19]:   Control_df = df2.query('group == "control"').converted.mean()
           Control_df
```

```
Out[19]:   0.1203863045004612
```

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

```
In [20]:   Treatment_df = df2.query('group == "treatment"').converted.mean()
           Treatment_df
```

```
Out[20]:   0.11880806551510564
```

```
In [21]:   # Calculate the actual difference (obs_diff) between the conversion rates for the two gr
           obs_diff =  Treatment_df - Control_df
```

```
obs_diff
```

-0.0015782389853555567

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

```python
df2['landing_page'].value_counts()[0]/len(df2)
```

0.5000619442226688

**e.** Consider your results from parts (a) through (d) above, and explain below whether the new `treatment` group users lead to more conversions.

> **The data indicates that the probability of converting individual in control group is slightly higher than the probability of converting individual in treatment group, there is not sufficient evidence to say that the new treatment page leads to more conversions as we can see in obs_diff the difference is only 0.00157, therfore there is no sufficient evidence to say that the new treatment page leads to more conversions.**

# Part II - A/B Test

Since a timestamp is associated with each event, you could run a hypothesis test continuously as long as you observe the events.

However, then the hard questions would be:

- 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.

## ToDo 2.1

For now, consider you need to make the decision just based on all the data provided.

> Recall that you just calculated that the "converted" probability (or rate) for the old page is *slightly* higher than that of the new page (ToDo 1.4.c).

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 be your null and alternative hypotheses ($H_0$ and $H_1$)?

You can state your hypothesis in terms of words or in terms of $p_{old}$ and $p_{new}$, which are the "converted" probability (or rate) for the old and new pages respectively.

> $H_0$: $p_{old}$ >= $p_{new}$
> $H_1$: $p_{old}$ < $p_{new}$

## ToDo 2.2 - Null Hypothesis $H_0$ Testing

Under the null hypothesis $H_0$, assume that $p_{new}$ and $p_{old}$ are equal. Furthermore, assume that $p_{new}$ and $p_{old}$ both are equal to the **converted** success rate in the `df2` data regardless of the page. So, our assumption is:

$$p_{new} = p_{old} = p_{population}$$

In this section, you will:

- Simulate (bootstrap) sample data set for both groups, and compute the "converted" probability $p$ for those samples.

- Use a sample size for each group equal to the ones in the `df2` data.

- Compute the difference in the "converted" probability for the two samples above.

- Perform the sampling distribution for the "difference in the converted probability" between the two simulated-samples over 10,000 iterations; and calculate an estimate.

Use the cells below to provide the necessary parts of this simulation. 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 hypothesis?

```
In [23]: p_new = df2.converted.mean()
         p_new #0.11959667567149027
```

```
Out[23]: 0.11959708724499628
```

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

```
In [24]: p_old = df2.converted.mean()
         p_old #0.11959667567149027
```

```
Out[24]: 0.11959708724499628
```

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

*Hint*: The treatment group users are shown the new page.

```
In [25]: n_new =df2.query('landing_page != "old_page"')['user_id'].nunique()
         n_new
```

```
Out[25]: 145310
```

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

```
In [30]: n_old = df2.query('landing_page != "new_page"')['user_id'].nunique()
         n_old
```

```
Out[30]: 145274
```

**e. Simulate Sample for the `treatment` Group**

Simulate $n_{new}$ transactions with a conversion rate of $p_{new}$ under the null hypothesis.

*Hint*: Use `numpy.random.choice()` method to randomly generate $n_{new}$ number of values.
Store these $n_{new}$ 1's and 0's in the `new_page_converted` numpy array.

```
In [31]:  # Simulate a Sample for H0
          new_page_converted = np.random.choice([1,0],n_new, p=(p_new,1-p_new))
          new_page_converted.mean()

Out[31]:  0.12020507879705457
```

**f. Simulate Sample for the `control` Group**

Simulate $n_{old}$ transactions with a conversion rate of $p_{old}$ under the null hypothesis.
Store these $n_{old}$ 1's and 0's in the `old_page_converted` numpy array.

```
In [32]:  # Simulate a Sample for H1
          old_page_converted = np.random.choice([1,0],n_old, p=(p_old,1-p_old))
          old_page_converted.mean()

Out[32]:  0.11826617288709611
```

**g.** Find the difference in the "converted" probability $(p'_{new} - p'_{old})$ for your simulated samples from the parts (e) and (f) above.

```
In [33]:  new_page_converted.mean() - old_page_converted.mean()

Out[33]:  0.0019389059099584605
```

**h. Sampling distribution**

Re-create `new_page_converted` and `old_page_converted` and find the $(p'_{new} - p'_{old})$ value 10,000 times using the same simulation process you used in parts (a) through (g) above.

Store all $(p'_{new} - p'_{old})$ values in a NumPy array called `p_diffs`.

```
In [34]:  p_diffs = []

          for i in range(10000):
              new_page_converted = np.random.choice([1, 0], size=n_new, p=[p_new, (1-p_new)])
              old_page_converted = np.random.choice([1, 0], size=n_old, p=[p_old, (1-p_old)])
              p_diff= new_page_converted.mean() - old_page_converted.mean()

              p_diffs.append(p_diff)
```

```
In [35]:  p_diffs = np.array(p_diffs)  #Convert to array
```
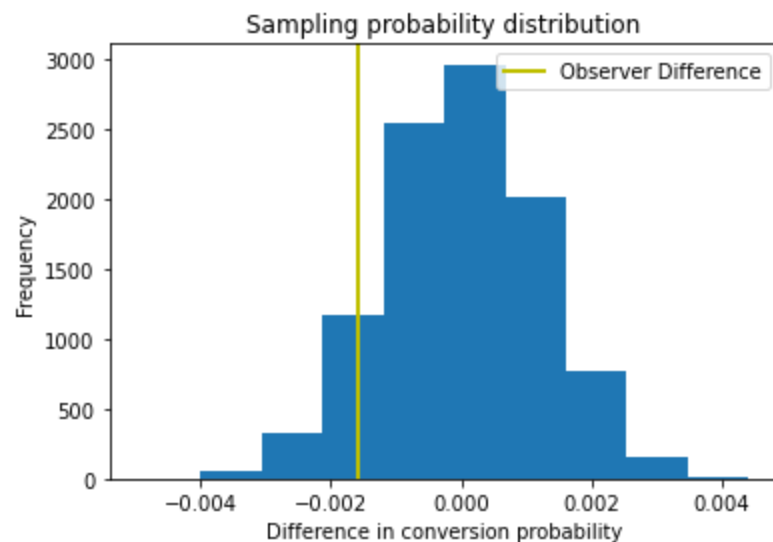
**i. Histogram**

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.

Also, use `plt.axvline()` method to mark the actual difference observed in the `df2` data (recall `obs_diff`), in the chart.

```
In [36]:  # Plot histogram
          plt.hist(p_diffs);
          plt.axvline(obs_diff, c='y', lw=2 ,label='Observer Difference')
          plt.xlabel('Difference in conversion probability')
```

```
plt.ylabel('Frequency')
plt.title('Sampling probability distribution')
plt.legend();
```



**j.** What proportion of the **p_diffs** are greater than the actual difference observed in the `df2` data?

```
In [43]:  (p_diffs>obs_diff).mean()
```
```
Out[43]:  0.9086
```

**k.** Please explain in words what you have just computed in part **j** above.

- What is this value called in scientific studies?
- What does this value signify in terms of whether or not there is a difference between the new and old pages? *Hint*: Compare the value above with the "Type I error rate (0.05)".

> **This value called as p-value,if the p-value were under 0.05 we can reject the null hypothesis,however in our case the p-value is greater than the alpha(0.05), p_value= (0.9047 > 0.05) so I have a statistical evidence to not reject the null hypothesis.**

**l. Using Built-in Methods for Hypothesis Testing**
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 statements below to calculate the:

- `convert_old` : number of conversions with the old_page
- `convert_new` : number of conversions with the new_page
- `n_old` : number of individuals who were shown the old_page
- `n_new` : number of individuals who were shown the new_page

```
In [44]:  import statsmodels.api as sm

          # number of conversions with the old page
          convert_old = convert_old = df2.query('landing_page != "new_page"').converted.sum()

          # number of conversions with the new page
          convert_new =df2.query('landing_page != "old_page"').converted.sum()
```

```
# number of individuals who were shown the old_page
n_old = df2.query('landing_page != "new_page" ').shape[0]

# number of individuals who received new_page
n_new = df2.query('landing_page != "old_page" ').shape[0]
```

In [45]:
```
print('convert_old = ', convert_old)
print('convert_new =',convert_new)
print('n_old',n_old)
print('n_new',n_new)
```

```
convert_old =  17489
convert_new = 17264
n_old 145274
n_new 145310
```

**m.** Now use `sm.stats.proportions_ztest()` to compute your test statistic and p-value. Here is a helpful link on using the built in.

The syntax is:

` proportions_ztest(count_array, nobs_array, alternative='larger')`

where,

- `count_array` = represents the number of "converted" for each group
- `nobs_array` = represents the total number of observations (rows) in each group
- `alternative` = choose one of the values from `['two-sided', 'smaller', 'larger']` depending upon two-tailed, left-tailed, or right-tailed respectively.

> **Hint**:
> It's a two-tailed if you defined $H_1$ as $(p_{new} = p_{old})$.
> It's a left-tailed if you defined $H_1$ as $(p_{new} < p_{old})$.
> It's a right-tailed if you defined $H_1$ as $(p_{new} > p_{old})$.

The built-in function above will return the z_score, p_value.

---

## About the two-sample z-test

Recall that you have plotted a distribution `p_diffs` representing the difference in the "converted" probability $(p'_{new} - p'_{old})$ for your two simulated samples 10,000 times.

Another way for comparing the mean of two independent and normal distribution is a **two-sample z-test**. You can perform the Z-test to calculate the Z_score, as shown in the equation below:

$$Z_{score} = \frac{(p'_{new} - p'_{old}) - (p_{new} - p_{old})}{\sqrt{\frac{\sigma^2_{new}}{n_{new}} + \frac{\sigma^2_{old}}{n_{old}}}}$$

where,

- $p'$ is the "converted" success rate in the sample
- $p_{new}$ and $p_{old}$ are the "converted" success rate for the two groups in the population.
- $\sigma_{new}$ and $\sigma_{new}$ are the standard deviation for the two groups in the population.

- $n_{new}$ and $n_{old}$ represent the size of the two groups or samples (it's same in our case)

> Z-test is performed when the sample size is large, and the population variance is known. The z-score represents the distance between the two "converted" success rates in terms of the standard error.

Next step is to make a decision to reject or fail to reject the null hypothesis based on comparing these two values:

- $Z_{score}$
- $Z_\alpha$ or $Z_{0.05}$, also known as critical value at 95% confidence interval. $Z_{0.05}$ is 1.645 for one-tailed tests, and 1.960 for two-tailed test. You can determine the $Z_\alpha$ from the z-table manually.

Decide if your hypothesis is either a two-tailed, left-tailed, or right-tailed test. Accordingly, reject OR fail to reject the null based on the comparison between $Z_{score}$ and $Z_\alpha$. We determine whether or not the $Z_{score}$ lies in the "rejection region" in the distribution. In other words, a "rejection region" is an interval where the null hypothesis is rejected iff the $Z_{score}$ lies in that region.

> Hint:
> For a right-tailed test, reject null if $Z_{score} > Z_\alpha$.
> For a left-tailed test, reject null if $Z_{score} < Z_\alpha$.

Reference:

- Example 9.1.2 on this page/09%3A_Two-Sample_Problems/9.01%3A_Comparison_of_Two_Population_Means-_Large_Independent_Samples), courtesy www.stats.libretexts.org

In [46]:
```python
import statsmodels.api as sm
# ToDo: Complete the sm.stats.proportions_ztest() method arguments
z_score, p_value = sm.stats.proportions_ztest([convert_old, convert_new], [n_old, n_new]
print(z_score, p_value)
```

1.3109241984234394 0.18988337448195103

**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.**?

> **The difference between the test statistics and the null hypotheis (Z-score) is 1.3109241984234394 while the p-value is equal to 0.18988337448195103 are greater than the alpha(0.05),both evidence suggest keep using the old page, we fail to reject the null hypothesis.**

# Part III - A regression approach

## ToDo 3.1

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 in the `df2` data 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** library to fit the regression model you specified in part **a.** above to see if there is a significant difference in conversion based on the page-type a customer receives. However, you first need to create the following two columns in the `df2` dataframe:

1. `intercept` - It should be `1` in the entire column.
2. `ab_page` - It's a dummy variable column, having a value `1` when an individual receives the **treatment**, otherwise `0`.

```
In [47]: df2['intercept'] = 1
         df2['ab_page'] = pd.get_dummies(df2['group'])['treatment']
         df2.head(10)
```

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 01:52:26.210827 | control | old_page | 1 | 1 | 0 |
| **5** | 936923 | 2017-01-10 15:20:49.083499 | control | old_page | 0 | 1 | 0 |
| **6** | 679687 | 2017-01-19 03:26:46.940749 | treatment | new_page | 1 | 1 | 1 |
| **7** | 719014 | 2017-01-17 01:48:29.539573 | control | old_page | 0 | 1 | 0 |
| **8** | 817355 | 2017-01-04 17:58:08.979471 | treatment | new_page | 1 | 1 | 1 |
| **9** | 839785 | 2017-01-15 18:11:06.610965 | treatment | new_page | 1 | 1 | 1 |

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

```
In [52]: X= df2.converted
         Lm = sm.Logit(X, df2[['intercept', 'ab_page']])
         Result = Lm.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 [53]: Result.summary()
```

Out[53]:

Logit Regression Results

| Dep. Variable: | converted | No. Observations: | 290584 |
|---|---|---|---|
| Model: | Logit | Df Residuals: | 290582 |
| Method: | MLE | Df Model: | 1 |

| | coef | std err | z | P>|z| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| **Date:** | Tue, 15 Nov 2022 | **Pseudo R-squ.:** | | 8.077e-06 | | |
| **Time:** | 16:06:18 | **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** | -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 |

**e.** What is the p-value associated with **ab_page**? Why does it differ from the value you found in **Part II**?

**Hints**:

- 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**?
- You may comment on if these hypothesis (Part II vs. Part III) are one-sided or two-sided.
- You may also compare the current p-value with the Type I error rate (0.05).

> **In II part the p-value was 0.18988337448195103 and the hypothesis was one-sided**
> $H_0: p_{old} >= p_{new}$
> $H_1: p_{old} < p_{new}$
> **Hence,a one-tailed test is applied.**
>
> **The p-value associated with ab_page is 0.190 which is higher than the p-value in the in part II, and the hypothesis is two-sided.**
> $H_0: p_{old} = p_{new}$
> $H_1: p_{old} != p_{new}$
> **Hence, a two-tailed test is applied. P-value > alpha (0.05), We should keep using the old page**

**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?

> **Due to treatment/control page not having much impact on the converts, I believe that adding factor will improve the model. I'll pay attention to not add many factort beacuse i don't want the model to suffer from over-fiting issue.**

**g. Adding countries**
Now along with testing if the conversion rate changes for different pages, also add an effect based on which country a user lives in.

1. You will need to read in the **countries.csv** dataset and merge together your `df2` datasets on the appropriate rows. You call the resulting dataframe `df_merged` . Here are the docs for joining tables.

2. Does it appear that country had an impact on conversion? To answer this question, consider the three unique values, `['UK', 'US', 'CA']`, in the `country` column. Create dummy variables for these country columns.

> **Hint:** Use `pandas.get_dummies()` to create dummy variables. **You will utilize two columns for the three dummy variables.**

Provide the statistical output as well as a written response to answer this question.

```
In [54]:   # Read the countries.csv
           countries_df = pd.read_csv('countries.csv')
           countries_df.head()
```

Out[54]:

| | user_id | country |
|---|---|---|
| **0** | 834778 | UK |
| **1** | 928468 | US |
| **2** | 822059 | UK |
| **3** | 711597 | UK |
| **4** | 710616 | UK |

```
In [57]:   # Join with the df2 dataframe
           df_join= df2.merge(countries_df, on ="user_id")
           df_join.head()
```

Out[57]:

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

```
In [58]:   # Create the necessary dummy variables
           df_dummies = pd.get_dummies(df_join['country'])
           df_join=df_join.join(df_dummies)
           df_join.head()
```

Out[58]:

| | user_id | timestamp | group | landing_page | converted | intercept | ab_page | country | CA | UK | US |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 851104 | 2017-01-21 22:11:48.556739 | control | old_page | 0 | 1 | 0 | US | 0 | 0 | 1 |
| **1** | 804228 | 2017-01-12 08:01:45.159739 | control | old_page | 0 | 1 | 0 | US | 0 | 0 | 1 |
| **2** | 661590 | 2017-01-11 16:55:06.154213 | treatment | new_page | 0 | 1 | 1 | US | 0 | 0 | 1 |
| **3** | 853541 | 2017-01-08 18:28:03.143765 | treatment | new_page | 0 | 1 | 1 | US | 0 | 0 | 1 |
| **4** | 864975 | 2017-01-21 01:52:26.210827 | control | old_page | 1 | 1 | 0 | US | 0 | 0 | 1 |

## h. Fit your model and obtain the results

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 are there significant effects on conversion. **Create the necessary additional columns, and fit the new model.**

Provide the summary results (statistical output), and your conclusions (written response) based on the results.

> **Hints**:
>
> - Look at all of p-values in the summary, and compare against the Type I error rate (0.05).
> - Can you reject/fail to reject the null hypotheses (regression model)?
> - Comment on the effect of page and country to predict the conversion.

```
In [59]:  # Fit your model, and summarize the results
          Lm2 = sm.Logit(df_join['converted'], df_join[['intercept', 'ab_page', 'CA', 'US']])
          Results = Lm2.fit()

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

```
In [60]:  Results.summary()
```

Out[60]:

Logit Regression Results

| | | | |
|---|---|---|---|
| **Dep. Variable:** | converted | **No. Observations:** | 290584 |
| **Model:** | Logit | **Df Residuals:** | 290580 |
| **Method:** | MLE | **Df Model:** | 3 |
| **Date:** | Tue, 15 Nov 2022 | **Pseudo R-squ.:** | 2.323e-05 |
| **Time:** | 16:06:42 | **Log-Likelihood:** | -1.0639e+05 |
| **converged:** | True | **LL-Null:** | -1.0639e+05 |
| **Covariance Type:** | nonrobust | **LLR p-value:** | 0.1760 |

| | coef | std err | z | P>\|z\| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| **intercept** | -1.9794 | 0.013 | -155.415 | 0.000 | -2.004 | -1.954 |
| **ab_page** | -0.0149 | 0.011 | -1.307 | 0.191 | -0.037 | 0.007 |
| **CA** | -0.0506 | 0.028 | -1.784 | 0.074 | -0.106 | 0.005 |
| **US** | -0.0099 | 0.013 | -0.743 | 0.457 | -0.036 | 0.016 |

> **Based on the p-values above,country does not appear to have significant impact on the conversion rate, we can't reject the null hypothesis.**

```
In [61]:  #interaction between page and country
          df_join['N_CA'] = df_join['CA']*df_join['ab_page']
          df_join['N_UK'] = df_join['UK']*df_join['ab_page']
          df_join['N_US'] = df_join['US']*df_join['ab_page']
```

```
In [62]:  log_model= sm.Logit(df_join['converted'], df_join[['intercept', 'ab_page', 'CA', 'US','N

          Results2 = log_model.fit()
```

```
In [63]:   Results2.summary()
```

```
Optimization terminated successfully.
         Current function value: 0.366109
         Iterations 6
```

Out[63]:

### Logit Regression Results

| Dep. Variable: | converted | No. Observations: | 290584 |
|---|---|---|---|
| Model: | Logit | Df Residuals: | 290578 |
| Method: | MLE | Df Model: | 5 |
| Date: | Tue, 15 Nov 2022 | Pseudo R-squ.: | 3.482e-05 |
| Time: | 16:06:46 | 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 | -1.9922 | 0.016 | -123.457 | 0.000 | -2.024 | -1.961 |
| ab_page | 0.0108 | 0.023 | 0.475 | 0.635 | -0.034 | 0.056 |
| CA | -0.0118 | 0.040 | -0.296 | 0.767 | -0.090 | 0.066 |
| US | 0.0057 | 0.019 | 0.306 | 0.760 | -0.031 | 0.043 |
| N_CA | -0.0783 | 0.057 | -1.378 | 0.168 | -0.190 | 0.033 |
| N_US | -0.0314 | 0.027 | -1.181 | 0.238 | -0.084 | 0.021 |

The p-vlaue of N_CA is 0.168 which is higher than the alpah (0.05), we conclude there is no interaction between page and country in CA.The same thing goes for N_US which its p-value equals to 0.238

## Conclusion

Ultimately, none of the A/B testing provides sufficient data to rule out the null hypothesis because all p-values are higher than the alpha level of 0.05.The old page still functions just as well, so there is no need to switch to it.

In [ ]: