

# Analyze\_ab\_test\_results\_notebook

March 23, 2021

## 0.1 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!

## 0.2 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.

```
[2]: 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:

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

```
[3]:   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.

```
[4]: df.shape[0]
```

```
[4]: 294478
```

c. The number of unique users in the dataset.

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

```
[5]: 290584
```

d. The proportion of users converted.

```
[6]: round(df['converted'].mean()*100)
```

```
[6]: 12
```

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

```
[7]: A=df.query('group == "treatment" & landing_page != "new_page").count()[0]
B=df.query('group == "control" & landing_page != "old_page").count()[0]
ct=A+B
ct
```

```
[7]: 3893
```

f. Do any of the rows have missing values?

```
[8]: df.isnull().sum()
```

```
[8]: 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**.

```
[9]: df2 = df.query("(group == 'treatment' and landing_page != 'old_page') or (group_
    => == 'control' and landing_page != 'new_page'))")
```

```
[10]: # 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]
```

```
[10]: 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**?

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

```
[11]: 290584
```

- b. There is one **user\_id** repeated in **df2**. What is it?

```
[12]: df2[df2.duplicated(['user_id'])]
```

```
[12]:      user_id      timestamp      group landing_page  converted
2893   773192  2017-01-14 02:55:59.590927  treatment      new_page          0
```

- c. What is the row information for the repeat **user\_id**?

```
[13]: df2[df2.user_id == 773192]
```

```
[13]:      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**.

```
[14]: df2.drop([1899],axis=0, inplace=True)
```

```
C:\Users\khaled\anaconda3\lib\site-packages\pandas\core\frame.py:4163:
SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame
```

```
See the caveats in the documentation: https://pandas.pydata.org/pandas-
docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
return super().drop(
```

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?

```
[15]: df2['converted'].mean()
```

```
[15]: 0.11959708724499628
```

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

```
[16]: df2[df2['group'] == "control"]['converted'].mean()
```

```
[16]: 0.1203863045004612
```

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

```
[17]: df2[df2['group'] == "treatment"]['converted'].mean()
```

```
[17]: 0.11880806551510564
```

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

```
[18]: (df2['landing_page'] == "new_page").mean()
```

```
[18]: 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.

**The results shows a very similar percentage for the covertion rate (0.1188) and (0.1203) so from the result its obvious that No there is no sufficient evidence to say that the new treatment page leads to more conversions.**

### 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.

**$H_0: P(\text{new}) - P(\text{old}) \leq 0$**

**$H_1: P(\text{new}) - P(\text{old}) > 0$**

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?

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

```
[19]: 0.11959708724499628
```

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

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

```
[20]: 0.11959708724499628
```

c. What is  $n_{new}$ ?

```
[21]: n_new = df2.query('group == "treatment"')['user_id'].count()  
n_new
```

```
[21]: 145310
```

d. What is  $n_{old}$ ?

```
[22]: n_old = df2.query('group == "control"')['user_id'].count()  
n_old
```

```
[22]: 145274
```

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**.

```
[23]: new_page_converted = np.random.binomial(n_new,p_new)
```

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**.

```
[24]: old_page_converted = np.random.binomial(n_old,p_old)
```

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

```
[25]: val=new_page_converted/n_new - old_page_converted/n_old
      val
```

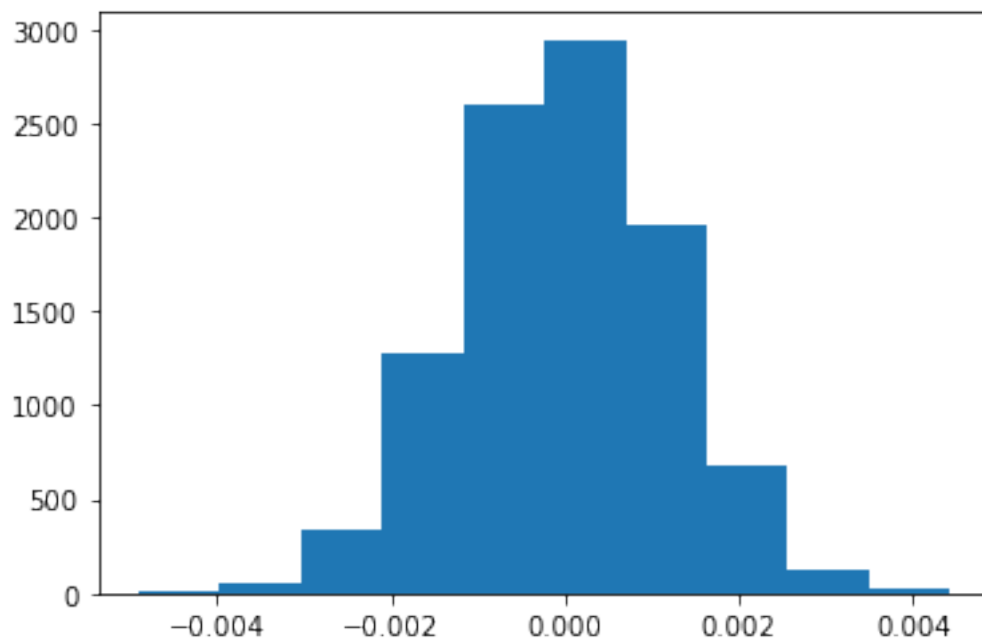
```
[25]: 0.00014940611650410274
```

- 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**.

```
[26]: p_diffs = []
      for _ in range(10000):
          new_page_converted = np.random.binomial(n_new,p_new)
          old_page_converted = np.random.binomial(n_old, p_old)
          val = new_page_converted/n_new - old_page_converted/n_old
          p_diffs.append(val)
```

- 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.

```
[27]: plt.hist(p_diffs);
```



- j. What proportion of the **p\_diffs** are greater than the actual difference observed in **ab\_data.csv**?

```
[28]: diff1=df2[df2['group'] == "treatment"]['converted'].mean()
      diff2=df2[df2['group'] == "control"]['converted'].mean()
      diff=diff1-diff2
      diff
```

[28]: -0.0015782389853555567

```
[29]: p_diffs = np.array(p_diffs)
      (diff < p_diffs).mean()
```

[29]: 0.8999

- 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?

**The values above is the Pvalues which represents the exact probability that outcomes from the statistical test and it shows that both old page and the new page has similar values and the old page has slightly better value**

- 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 the the number of rows associated with the old page and new pages, respectively.

```
[31]: import statsmodels.api as sm

convert_old = df2.query(" landing_page == 'old_page' and converted == 1").
    ↪count()[0]
convert_new = df2.query(" landing_page == 'new_page' and converted == 1").
    ↪count()[0]
n_old = df2.query('group == "control"')['user_id'].count()
n_new = df2.query('group == "treatment"')['user_id'].count()
n_old
```

[31]: 145274

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

```
[32]: z_score, p_value = sm.stats.proportions_ztest([convert_old, convert_new],
    ↪[n_old, n_new], alternative='smaller', prop_var=False)
z_score, p_value
```

[32]: (1.3109241984234394, 0.9050583127590245)

- 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 values shown in the previous question shows that the old\_page showed a better values than the new page**

**The values in the previous question do agree with the value in part j and k**

### 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?

### 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**.

```
[33]: df2['intercept'] = 1
df2[['ab_page', 'old_page']] = pd.get_dummies(df2['landing_page'])
df2.head()
```

```
<ipython-input-33-3d232eb1f5cb>:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
df2['intercept'] = 1
C:\Users\khaled\anaconda3\lib\site-packages\pandas\core\frame.py:3065:
SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
self[k1] = value[k2]
```

```
[33]:
```

	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	

	intercept	ab_page	old_page
0	1	0	1
1	1	0	1
2	1	1	0
3	1	1	0
4	1	0	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.



```
[34]: A = sm.Logit(df2['converted'], df2[['ab_page', 'intercept']])
      model = A.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.

```
[35]: model.summary()
```

```
[35]: <class 'statsmodels.iolib.summary.Summary'>
      """
                                Logit Regression Results
      =====
Dep. Variable:                converted    No. Observations:                290584
Model:                        Logit       Df Residuals:                290582
Method:                       MLE        Df Model:                  1
Date:                         Tue, 23 Mar 2021    Pseudo R-squ.:                8.077e-06
Time:                         01:25:14    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]
      -----
ab_page          -0.0150      0.011      -1.311      0.190      -0.037      0.007
intercept       -1.9888      0.008     -246.669      0.000      -2.005     -1.973
      =====
      """
```

- 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**?

**p-value=0.1899**

**different results occurred because in this model we used the two sided test unlike in part 2 when we used only one side and it was randomly chosen**

- 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?

**adding more factors will make the results more accurate but if this factors has no influence to the test the results might be useless**

- 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 vari-

ables 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.

```
[36]: 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()
```

```
[36]:
```

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

	converted	intercept	ab_page	old_page
user_id				
834778	0	1	0	1
928468	0	1	1	0
822059	1	1	1	0
711597	0	1	0	1
710616	0	1	1	0

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

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

```
[38]: df_new[['UK', 'US']] = pd.get_dummies(df_new['country'])[['UK', 'US']]
df_new.head()
```

```
[38]:
```

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

	converted	intercept	ab_page	old_page	UK	US
user_id						
834778	0	1	0	1	1	0
928468	0	1	1	0	0	1
822059	1	1	1	0	1	0
711597	0	1	0	1	1	0

```
710616          0          1          1          0          1          0
```

```
[39]: A = sm.Logit(df_new['converted'], df_new[['intercept', 'UK', 'US']])
      model = A.fit()
```

```
Optimization terminated successfully.
      Current function value: 0.366116
      Iterations 6
```

```
[40]: model.summary()
```

```
[40]: <class 'statsmodels.iolib.summary.Summary'>
      """
                                Logit Regression Results
      =====
Dep. Variable:                converted    No. Observations:                290584
Model:                        Logit       Df Residuals:                  290581
Method:                       MLE        Df Model:                      2
Date:                         Tue, 23 Mar 2021    Pseudo R-squ.:                1.521e-05
Time:                         01:25:20    Log-Likelihood:               -1.0639e+05
converged:                    True        LL-Null:                      -1.0639e+05
Covariance Type:              nonrobust    LLR p-value:                  0.1984
      =====
                                coef      std err          z      P>|z|      [0.025      0.975]
      -----
intercept      -2.0375      0.026     -78.364      0.000      -2.088      -1.987
UK              0.0507      0.028       1.786      0.074      -0.005       0.106
US              0.0408      0.027       1.518      0.129      -0.012       0.093
      =====
      """
```

As it shows above countries doesnt have a significant effect in the conversion rates

- 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.

```
[41]: A = sm.Logit(df_new['converted'],df_new[['intercept', 'ab_page', 'UK', 'US']])
      model = A.fit()
```

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

```
[42]: model.summary()
```

```
[42]: <class 'statsmodels.iolib.summary.Summary'>
      """
```

### Logit Regression Results

```
=====
Dep. Variable:          converted    No. Observations:          290584
Model:                  Logit       Df Residuals:              290580
Method:                 MLE         Df Model:                  3
Date:                  Tue, 23 Mar 2021    Pseudo R-squ.:          2.323e-05
Time:                  01:25:25    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    -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
UK           0.0506      0.028      1.784      0.074     -0.005      0.106
US           0.0408      0.027      1.516      0.130     -0.012      0.093
=====
"""
```

The results shows that the interaction between both the page and the country have no effect on the conversion rates

```
[44]: df_new['UK_new'] = df_new['UK'] * df_new['ab_page']
      df_new['US_new'] = df_new['US'] * df_new['ab_page']
      df_new.head()
```

```
[44]:
```

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

	converted	intercept	ab_page	old_page	UK	US	UK_new	US_new
user_id								
834778	0	1	0	1	1	0	0	0
928468	0	1	1	0	0	1	0	1
822059	1	1	1	0	1	0	1	0
711597	0	1	0	1	1	0	0	0
710616	0	1	1	0	1	0	1	0

```
[45]: A = sm.Logit(df_new['converted'], df_new[['intercept', 'ab_page', 'UK', 'US',
      ↪ 'UK_new', 'US_new']])
      model = A.fit()
```

Optimization terminated successfully.  
Current function value: 0.366109

Iterations 6

```
[46]: model.summary()
```

```
[46]: <class 'statsmodels.iolib.summary.Summary'>
      """
```

```

                        Logit Regression Results
=====
Dep. Variable:          converted    No. Observations:          290584
Model:                  Logit      Df Residuals:              290578
Method:                  MLE       Df Model:                  5
Date:                   Tue, 23 Mar 2021    Pseudo R-squ.:          3.482e-05
Time:                   01:27:03    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      -2.0040      0.036     -55.008      0.000      -2.075      -1.933
ab_page        -0.0674      0.052     -1.297      0.195      -0.169      0.034
UK              0.0118      0.040      0.296      0.767      -0.066      0.090
US              0.0175      0.038      0.465      0.642      -0.056      0.091
UK_new         0.0783      0.057      1.378      0.168      -0.033      0.190
US_new         0.0469      0.054      0.872      0.383      -0.059      0.152
=====
      """
```

## Conclusions > In conclusion this study showed that they should stay on the old page as the study showed that there is no evidence that the new page is any better the the old page

```
[ ]:
```