



Data Science for Marketing: A/B Testing and Statistical Analysis

Silverblue Group

April 8, 2025



Course Information and Instructor

- Course: 2025 Spring MAP 2192-001-15022 Mathematics of Data Science
- Instructor: Dr. Zhijun Yin

Group Member List

Name	Contribution
CJ Bingham	Everything



Data-Driven decision making is crucial for companies to maximize their return on investment in marketing and ad campaigns. As data analysts, we must understand how to find statistically significant relationships between different variables, and create visualizations of those relationships. We can create hypothesis tests to look at different scenarios, and Z-tests can help to show statistical significance. Logistic regression analysis is also a tool used to analyze relationships between different variables.



We take a sample dataset containing 5000 rows of marketing related data. Columns include data such as marketing data campaign id (1-20, ordinal), positive or negative response from the user (0 or 1, binary), and the percent grade for each customer's click through rate (0-100, ratio). We then perform analysis on the dataset to gather insights to better focus our marketing strategies and better understand our user.



- Definitions and principles of hypothesis testing and the Z test
- Definitions and principles of logistic regression
- Python, Jupyter, and LaTeX for implementation
- Python packages for data manipulation, analysis, and visualization



Our Hypothesis Statement and A/B Test

An A/B test takes two different campaigns shown to a user and changes a certain variable in one of the campaigns. This could be a different ad or call to action statement. In this case we look at two different sets of campaigns, one from the first half of campaigns, and one from the second half of campaigns. We are conducting a hypothesis test to compare the response rates of these two sets of marketing campaigns.

- Null Hypothesis (H_0): The response rate of the control group is equal to the response rate of the treatment group.

$$H_0 : p_{control} = p_{treatment}$$

- Alternative Hypothesis (H_1): The response rate of the treatment group is higher than the control group.

$$H_1 : p_{control} < p_{treatment}$$

Decision rule: Reject H_0 if $p\text{-value} < \alpha$ (typically 0.05)



Z-Test for Comparing Response Rates

To determine whether there is a statistically significant difference in positive response rates between the control and treatment groups, we conduct a two-proportion Z-test. The test statistic is calculated as:

$$z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}(1 - \hat{p}) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}} \quad (1)$$

- $\hat{p}_1 = \frac{x_1}{n_1}$ is the sample proportion of positive responses in the control group.
- $\hat{p}_2 = \frac{x_2}{n_2}$ is the sample proportion of positive responses in the treatment group.
- n_1, n_2 are the sample sizes of the control and treatment groups, respectively.
- x_1, x_2 are the number of successful responses in each group.
- \hat{p} is the pooled sample proportion, given by:

$$\hat{p} = \frac{x_1 + x_2}{n_1 + n_2} \quad (2)$$



Python Code

```
1 # Import necessary libraries
2 import pandas as pd
3 import numpy as np
4 from statsmodels.stats.proportion import proportions_ztest
5 from statsmodels.stats.proportion import proportion_confint
6 import matplotlib.pyplot as plt
7 # Create dataframe from dataset and create control_group and treatment_group
8 df = pd.read_csv("marketing_dataset.csv")
9 # Control_group contains campaigns 1-9
10 control_group = df[(df['campaign_id'].isin(tuple(f'CAMPAIGN_{i}' for i in range(1, 10))))]
11 # Treatment_group contains campaigns 11-19
12 treatment_group = df[(df['campaign_id'].isin(tuple(f'CAMPAIGN_{i}' for i in range(11, 20))))]
13 # Find total users in each group
14 control_total = len(control_group)
15 treatment_total = len(treatment_group)
16 print(f"Total number of users in first half of marketing campaigns: {control_total}")
17 print(f"Total number of users in second half of marketing campaigns: {treatment_total}")
18 # Compute conversions
19 control_conversions = (control_group['customer_response'] == 1).sum()
20 treatment_conversions = (treatment_group['customer_response'] == 1).sum()
21 print(f"Total number of positive responses in first half of marketing campaigns: {
    control_conversions}")
22 print(f"Total number of positive responses in second half of marketing campaigns: {
    treatment_conversions}")
23 def ab_test_analysis(control_conversions, control_total,
24                      treatment_conversions, treatment_total):
25     # Calculate proportions
26     p_control = control_conversions / control_total
27     p_treat = treatment_conversions / treatment_total
28
29     # Perform z-test
30     count = np.array([control_conversions, treatment_conversions])
31     nobs = np.array([control_total, treatment_total])
32     z_stat, p_val = proportions_ztest(count, nobs)
33
34     # Calculate confidence intervals
35     ci_control = proportion_confint(control_conversions, control_total, method='normal')
36     ci_treat = proportion_confint(treatment_conversions, treatment_total, method='normal')
```




Python Code

```
1 fig, ax = plt.subplots(figsize=(6, 5))
2     groups = ['Control Group', 'Treatment Group']
3     proportions = [p_control, p_treat]
4     ci_lower = [ci_control[0], ci_treat[0]]
5     ci_upper = [ci_control[1], ci_treat[1]]
6
7     # Create bar plot with error bars
8     ax.bar(groups, proportions, yerr=[
9         [p - l for p, l in zip(proportions, ci_lower)],
10        [u - p for u, p in zip(ci_upper, proportions)]
11    ], capsize=10, color=['blue', 'red'], alpha=0.6)
12
13    # Labels and title
14    ax.set_ylabel("Response Rate")
15    ax.set_title("Comparison of Response Rates with 95% CIs")
16    ax.axhline(y=p_control, color='blue', linestyle='dashed', alpha=0.7, label="Control Mean")
17    ax.axhline(y=p_treat, color='red', linestyle='dashed', alpha=0.7, label="Treatment Mean")
18    ax.legend()
19    plt.savefig("CI.png")
20    # Show plot
21    plt.show()
22
23    return {
24        'z_test_statistic': f"{z_stat:.2f}",
25        'p_value_z_test': f"{p_val:.2f}",
26        'control_response_rate': f"{p_control:.2f}",
27        'treatment_response_rate': f"{p_treat:.2f}",
28        'control_response_rate_confidence_interval': ci_control,
29        'treatment_response_rate_confidence_interval': ci_treat
30    }
31 ab_test_analysis(control_conversions, control_total, treatment_conversions, treatment_total)
```



Statistical Test and Results

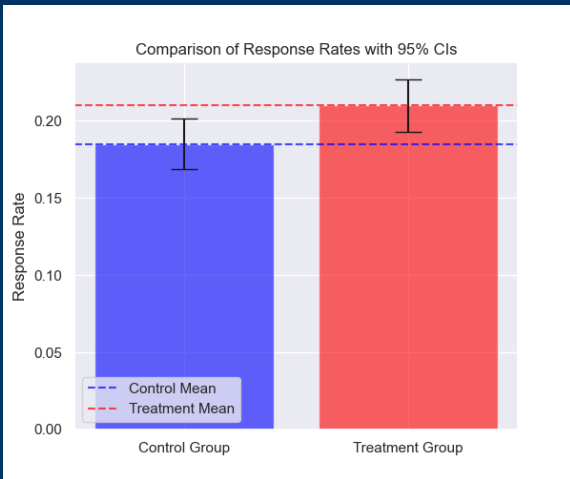
We perform a two-proportion Z-test to determine if there is a statistically significant difference between the response rates of the two groups. Below are the test results:

Metric	Control Group	Treatment Group
Sample Size (n)	2213	2266
Successful Responses (x)	409	475
Response Rate (p)	0.18	0.21
95% Confidence Interval	(0.169, 0.201)	(0.193, 0.226)
Z-test Statistic		-2.09
P-value		0.04

Table 1: Summary of A/B Test Results



Confidence Interval Overlap



Overlap of 95 percent confidence intervals



With a p-value of 0.04 (below the significance level of 0.05), we reject the null hypothesis and accept the alternative hypothesis stating that "The response rate of the treatment group is higher than the control group." However, we have yet to prove what variables influenced the improvement in the second half of marketing campaigns the most.



Introduction to Logistic Regression Analysis

Logistic regression analysis is a process used to estimate the probability of an event occurring, such as a conversion or clicking on an ad. With our marketing dataset, we use the analysis to compare two variables, click through rate and customer response. We conduct this test to determine how an increase in click through rate affects the probability of getting a positive response.

For predicting binary outcomes (positive or negative response):

$$\log \left(\frac{p}{1-p} \right) = \beta_0 + \beta_1 X_1 + \cdots + \beta_k X_k \quad (3)$$

Where:

p = Probability of positive response

β_i = Coefficients

X_i = Predictor variable (click through rate)



Python Code

The target column in the dataset is customer response, with a value of 1 for a positive response, and a 0 for a negative response. We use the following Python code to perform the logistic regression analysis:

```
1 # Import necessary libraries
2 import pandas as pd
3 import numpy as np
4 import statsmodels.api as sm
5 import seaborn as sns
6 import matplotlib.pyplot as plt
7
8 # Load data
9 df = pd.read_csv("marketing_dataset.csv")
10
11 # Define predictor (X) and target (y)
12 X = df["click_through_rate"]
13 y = df["customer_response"]
14
15 # Fit logistic regression model
16 X_const = sm.add_constant(X) # Add intercept
17 logit_model = sm.Logit(y, X_const).fit()
18 print(logit_model.summary())
19
20 # Generate prediction data
21 X_range = np.linspace(X.min(), X.max(), 300)
22 X_range_const = sm.add_constant(X_range) # Add intercept for prediction
23 y_pred = logit_model.predict(X_range_const)
```



Python Code and Logistic Regression Results

```
1 # Plot scatterplot of actual data
2 plt.figure(figsize=(8, 6))
3 sns.scatterplot(x=X, y=y, alpha=0.3, color="blue", label="Data")
4
5 # Plot logistic regression curve
6 plt.plot(X_range, y_pred, color="red", linewidth=2, label="Logistic Regression Curve")
7
8 # Labels and title
9 plt.xlabel("Click Through Rate Rating for Customer (%)")
10 plt.ylabel("Customer Response (0 = No, 1 = Yes)")
11 plt.title("Logistic Regression Test: Click Through Rate vs Customer Response")
12 plt.legend()
13 plt.savefig("log_regression.png")
14 plt.show()
```

Logit Regression Results

=====						
Dep. Variable:	customer_response	No. Observations:	5000			
Model:	Logit	Df Residuals:	4998			
Method:	MLE	Df Model:	1			
Date:	Tue, 08 Apr 2025	Pseudo R-squ.:	0.2215			
Time:	18:20:49	Log-Likelihood:	-1942.3			
converged:	True	LL-Null:	-2495.1			
Covariance Type:	nonrobust	LLR p-value:	2.054e-242			
=====						
	coef	std err	z	P> z	[0.025	0.975]

const	-4.4954	0.137	-32.890	0.000	-4.763	-4.227
click_through_rate	0.0508	0.002	27.127	0.000	0.047	0.054
=====						



Logistic Regression Test

Logistic regression used on response rate and click through rate for customer:





The analysis suggests that there is a statistically significant positive relationship between the click-through rate and the likelihood of a positive response to the marketing campaign. Specifically, an increase in the click-through ratio (0-100) per customer is associated with a higher probability of a positive response.



Conclusion

This project examined the use of hypothesis testing, Z tests, and A/B tests to take marketing campaign related data out of the dataset and examine if the second half of campaigns (11-19) performed better than the first half of campaigns (1-9) with a one month buffer in between. This test showed that there was statistical significance in an increase in response rates between the two campaigns and showcase the ability of A/B tests to inform data-driven decision making in marketing departments.

Additionally, we used a logistic regression test to examine relationships between specific variables for further analysis. This found that as a user's click through rating rose in percentage, so did the probability of a positive response to the marketing campaign. Further analysis of click through rate differences existing in both sets of marketing campaigns could give additional insight.



In future research, we could explore the following directions:

- Use other types of data in the logistic regression model
- Other descriptive, inferential, or predictive statistical analysis techniques to draw further conclusions
- Other datasets that contain real data, or larger datasets
- Use machine learning models



- E. Weed, (2021). Learning Statistics with Python, URL: [\[https://ethanweed.github.io/pythonbook/landingpage.html\]](https://ethanweed.github.io/pythonbook/landingpage.html)
- F. Quin, D. Weyns, M. Galster, and C. C. Silva, (2024). “A/B Testing: A Systematic Literature Review,” Journal of Systems and Software, vol. 211, p. 112011
- LaTeX Project, (2024). LaTeX Documentation, URL: [\[https://www.latex-project.org/help/documentation/\]](https://www.latex-project.org/help/documentation/)
- Python Software Foundation, (2024). Python 3 Documentation, URL: [\[https://docs.python.org/3/\]](https://docs.python.org/3/)
- Ziya07, (2025). Cloud-Enabled Marketing Strategy Dataset, Kaggle, URL: [\[https://www.kaggle.com/datasets/ziya07/cloud-enabled-marketing-strategy-dataset.\]](https://www.kaggle.com/datasets/ziya07/cloud-enabled-marketing-strategy-dataset.)