

Analysis of Call Center Data

Karissa D. Hems

Western Governors University

Contents

A1.	Research Question	Error! Bookmark not defined.
A2.	Context and Background	2
A3.	Summary of Published Works and Relation to Project	2
A4.	Summary of Data Analytics Solution	Error! Bookmark not defined.
A5.	Benefit and Support to Decision Making Process	Error! Bookmark not defined.
B1.	Goals Objectives and Deliverables	Error! Bookmark not defined.
B2.	Project Scope	Error! Bookmark not defined.
B3.	Methodology	Error! Bookmark not defined.
B4.	Timeline and Milestones	Error! Bookmark not defined.
B6.	Criteria for Success	Error! Bookmark not defined.
C1.	Hypothesis	Error! Bookmark not defined.
C2.	Analytical Method and Justification	Error! Bookmark not defined.
C3.	Tools and Environments	Error! Bookmark not defined.
C4.	Methods and Metrics and Justification	Error! Bookmark not defined.
C5.	Practical Significance	Error! Bookmark not defined.
C6.	Graphical Representations	Error! Bookmark not defined.
D1.	Source of Data	Error! Bookmark not defined.
D2.	Data Collection Methods	Error! Bookmark not defined.
D3.	Data Quality	Error! Bookmark not defined.
D4.	Data Governance and Precautions	Error! Bookmark not defined.

A. Project Highlights

Question

This project is an analysis of the helpdesk call center metrics at Sleuth Goose Shipping Services. It will identify any specific days of the week with longer hold times and to investigate if longer hold times lead to a rise in escalation incidents. The research question is whether there is a correlation between increased hold times and the occurrence of escalations.

Scope, Tools and Methodology

Python within a Jupyter Notebook was used to complete this project. The call center data was analyzed within the Jupyter Notebook using Python and produced answers to the research question. The fields I originally believed I would use from the Call Center Data file were date, incoming, answered, abandoned, answer speed, and waiting time. I did not utilize the incoming field as I was able to gather the information, I needed from the length of the datafile. I originally neglected to mention that I would use the escalations field, but in order to see if there is a correlation in the two metrics, I needed to use this field as well. I proceeded in a linear fashion, using the waterfall methodology which will be further discussed below.

B. Project Execution

Plan

The overall goal of this project is to reveal if the call center is meeting the internally set standards. The objective is to answer when call holding times surpass five minutes and whether this is correlated with an increase in escalations. The deliverables for the objectives will be reporting days with the longest hold times and whether there is a definitive correlation between the longer hold time and escalations.

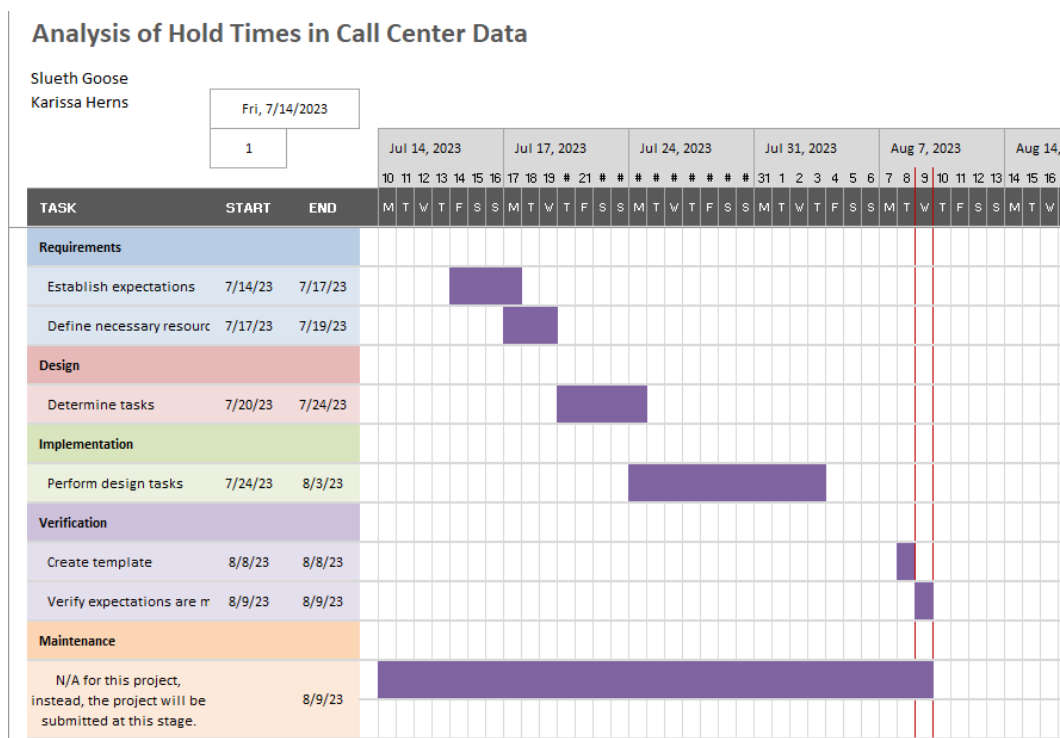
Planning Methodology

I completed this project in a linear fashion by using the waterfall method. First, in the requirements phase all customer requirements were gathered and the scope, expectations, and resources were defined. Next, in the design phase, tasks needed to complete the project were determined. Since the data was self-created, the data cleaning process was skipped and it was determined that descriptive statistics would be appropriate for this project. Next, during implementation, the tasks determined in the design phase were carried out. This included carrying out the analysis tasks using Python. For the fourth stage, which is the final stage for this project, verification was conducted. The results needed to be verified to be useful and the code is available to the customer to use at their discretion. The maintenance phase is not necessary for this project, but would be the fifth phase in other projects.

The methodology did not vary as there was no need to change the linear progression of the project.

Timeline and Milestones

My timeline was extended beyond the original estimation as it took longer than anticipated to complete the design tasks. The design task took five days longer than originally estimated. I made some time up during the template creation and verification processes as these each took one day instead of three and four days respectively.



C. Data Collection Process

Data Selection, Collection and Governance

The dataset was self-created, and it was specifically designed to resemble a typical call center dataset. Creating the dataset myself eliminated obstacles that could have been encountered such as omitting non-relevant information like the category of service. One thing that differed was that I intentionally thought I would need to use a count of calls in general, but I was able to use length of the data so it was not necessary. This dataset contains fictional data which I created myself, hence there are no security or privacy issues to be concerned with. I created a fictional company to accompany the fictional data. I stored the data locally and shut down my device when not in use.

C.1. Advantages and Limitations of the Data

An advantage of creating a dataset is knowing that the information needed for the desired analysis will be there. It is also complete and does not need to be wrangled and/or cleaned. The downside to this is that the type of variance and real-world metrics one would encounter in a call center are not present. Additionally, the volume needed for things like forecasting with any kind of confidence is not there.

D. Data Extraction and Preparation

Extraction (Creation) and Preparation

I created the data using <https://www.mockaroo.com/> and then saved the data in an Excel workbook. I had to create several sets of data and combine them in the spreadsheet as the website limits the amount of data that can be created at one time with a free account. I could not locate datasets that included the number of escalations or number of calls on a daily basis.

There were necessary fields in the dataset, and I included fields that would assist in future analysis if so desired.

E. Data Analysis Process

E.1 Data Analysis Methods

I used a descriptive analytical method for this project. Grouping the call center data by day of week and calculating averages for key metrics like hold times and escalations provides a descriptive summary of the dataset. Descriptive statistics and visualization provide a clear and concise summary of data patterns and distributions, facilitating a better understanding of the dataset and informing data-driven decision-making

E.2 Advantages and Limitations of Tools and Techniques

The advantages of Python and Jupyter Notebooks are that Python has extensive libraries and tools specialized for data analysis and visualization. Python can handle and process data efficiently and effectively. Jupyter Notebooks allows for interactive and iterative data analysis. This aids in easy exploration and visualization of data. A disadvantage to Python is that its power is limited to the knowledge of the user and the learning curve can be steep. The graphs are not necessarily as customizable using Python as another tool might be such as Power BI. One of the disadvantages to using Jupyter Notebooks is that others may not have this tool readily available. It can also perform slowly with larger datasets.

E.3. Application of Analytical Methods

This project applied descriptive statistics and visualizations to analyze call center performance through the following steps:

- Load the call center dataset

- Created a filtered dataset that only contains data with call hold time greater than five minutes
- Created a copy of the original dataframe to avoid modifying raw data and added a new column based on the condition of call time over 5 mins
- Extracting new datetime feature day_name and month from the date
- Converting boolean columns like escalation to numeric values for easier analysis
- Grouped by day of week and calculated the mean call duration for each day
- Created a bar chart visualization of average durations ordered by decreasing average
- Calculated overall percentage of calls with escalations across the entire dataset
- Utilized a pivot table to analyze escalations versus the 'over_five' column
- Generated a heatmap as a visualization see the correlation between escalations and call times based on the pivot table
- Created a contingency table containing the observed frequencies for the variables 'over_five' and 'escalation_num'
- Calculated the chi-square statistic, p-value, and degrees of freedom.
- Generated a histogram to visualize the density distribution of hold times, grouped by whether the call was escalated or not.

The analysis relied on several assumptions. First, that the acceptable threshold for hold times before customer dissatisfaction was 5 minutes based on company policy. Second, that sudden spikes in escalations were related to excessive hold times on a given day rather than other factors. Finally, that hold times and escalations would exhibit a measurable correlation. These assumptions were verified by depicting the average speed to answer for each day of the week in a bar chart. I was also able to verify the relationship between answer times and escalations by using a heat map.

F. Data Analysis Results

F.1 Statistical Significance

Null Hypothesis: There is no association between call wait times being over five minutes and the occurrence of escalations in the call center.

To investigate this, the statistical test used was a Chi-square test. The metrics generated were: Chi-square value: 38.164111759703005, p-value: 6.503767274372569e-10, and Degrees of freedom: 1. This is a hypothesis test and the results show that there is sufficient evidence to reject the null hypothesis and support the claim that there is an association between call wait times being over five minutes and the occurrence of escalations in the call center. The low p-value (6.503767274372569e-10) indicates that this association is statistically significant. The chi-square value of 38.164111759703005 and the degrees of freedom of 1 further support this conclusion. The low p-value indicates that the observed relationship is unlikely to have occurred by chance, providing evidence to support the alternative hypothesis.

Note: The original choice for testing in this instance was a regression analysis, however the models were not accurate. This is likely due to the dataset being highly imbalanced and small. In order to test the null hypothesis, I needed to change the type of analysis I used.

F.2 Practical Significance

The data analytics solution offers practical insights by identifying days with higher hold times and number of escalated events, enabling managers to make informed staffing decisions. Furthermore, visualizing these results can assist managers in evaluating the trade-off between addressing customer dissatisfaction due to long hold times or abandoned calls and the potential need to hire more staff to handle call volumes effectively. We can see in the visualizations that there are days of the week with longer hold times and this could be used to assist in the decision whether or not to adjust staffing.

F.3 Overall Success

The project has demonstrated its success in achieving the set goals. Mondays appear to have longer hold times. The heat map shows that longer hold times appear to be correlated with escalations, and the statistical test yielded significant metrics with a Chi-square value of 38.164111759703005, and a p-value of 6.503767274372569e-10. These results indicate that the null hypothesis can be rejected. The calculated low p-value demonstrates that this relationship is unlikely to have arisen by chance alone, providing strong evidence to support the alternative hypothesis.

G. Conclusion

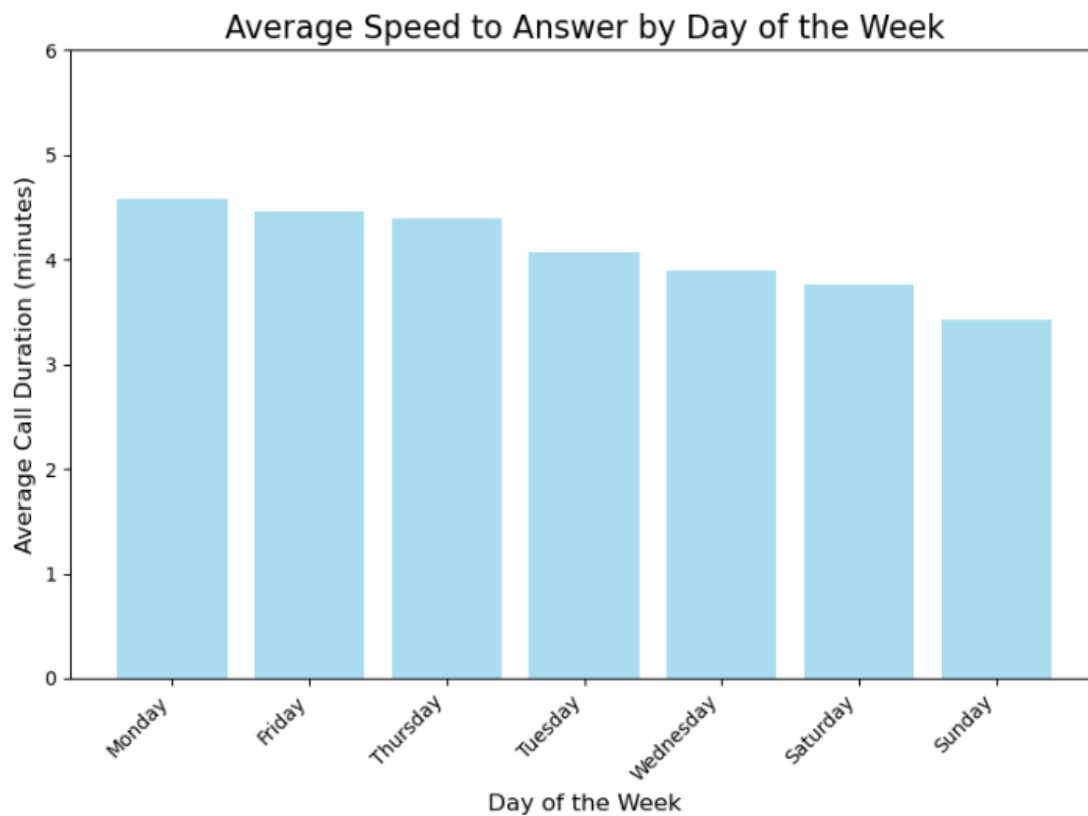
G.1 Summary of Conclusions

This project successfully analyzed call center metrics to uncover insights into the relationship between call wait times, escalations, and days of the week. Employing Python and Jupyter Notebooks, the project applied descriptive statistics and visualizations to reveal patterns and correlations. The Chi-square test established a statistically significant association between call wait times over five minutes and escalations. Practical implications were highlighted, including identifying peak hold times on certain days, specifically Mondays. The project's methodology, was supported by statistical evidence and validated the findings.

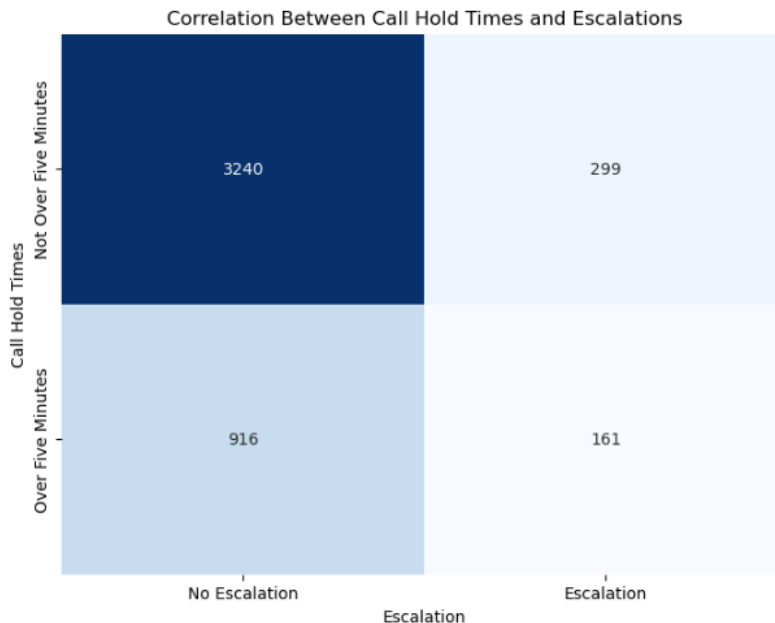
G.2

The first visualization was a bar chart that visually depicts the average speed to answer calls by day of the week. This allows readers to quickly see the patterns in the data - for example, that Mondays have the longest average speed to answer while Sundays have the shortest.

Ordering the bars from longest to shortest average duration creates a logical flow and highlights the day with the maximum value (Monday). The chart uses colors, labels, and formatting to make it easy to interpret. The visualization was generated using Python Pandas and Matplotlib.



I also used Python to generate a heatmap. The heatmap is used to reveal and quantify the relationship between call hold times and escalation rates. The heatmap's color gradient shows a distinctly higher proportion of escalations in the 'Over Five Minutes' row, visually highlighting the correlation. The total call counts provided through annotations show the large difference in sample sizes between the two hold time buckets. The heatmap conveys the takeaway that longer call holds have a higher likelihood of escalation. The title and axis labels provide context for interpreting the correlation shown in the heatmap.



G.3 Recommended Courses of Action

Recommendation 1: Maintain current staffing levels and scheduling, as average wait times are within acceptable thresholds. The analysis showed maximum average speed to answer across days of the week was 4.6 minutes, which is under the 5-minute threshold set by Sleuth Goose management. Since this metric is currently meeting organizational benchmarks, no changes need to be made to staff levels or shifts at this time.

Recommendation 2: Explore voluntary overtime shifts to provide additional capacity when experiencing high call volume. The analysis showed Mondays have the longest average speed to answer at 4.6 minutes, approaching the 5-minute threshold set by management. Implementing voluntary overtime shifts on Mondays or on an as needed basis could provide more capacity to prevent wait times exceeding the threshold.

H. Summary of Project

<https://wgu.hosted.panopto.com/Panopto/Pages/Viewer.aspx?id=f572c8af-c2a7-4502-abd8-b05a004fcd98>

A. Summary of Project

- Code provided in Notepad document
- Excel Spreadsheet contains data
- Data created at www.Mockaroo.com

References

No sources were cited.

