# Data Processing and Feature Engineering Report

**Project Title:** Customer Churn Analysis for Telecommunication Company
**Project Sponsor:** Doris Lee, CEO of Advanced Consultant Services

**Project Manager:** Maniram Luitel
**Project Start Date:** July 1, 2024
**Expected Project End Date:** September 15, 2024

**Prepared by:** Bhavesh Chaudhary
**Date:** August 08, 2024
**Role:** Data Engineer

# Abstract

This report lays a very strong emphasis on data preprocessing and feature engineering in the Customer Churn Analysis project. The main goal of this phase is to increase the quality and predictive power of the dataset for effective modeling and accurate prediction. The main processes in data preprocessing are loading the data, cleaning, handling missing values, encoding categorical variables, and scaling numerical features. These steps maintained the integrity and consistency of the dataset, which, in turn, ensures building reliable predictive models.

Feature engineering was done so that new features could be developed in the best possible way to represent the underlying data patterns. Development of features such as Charges_Per_Tenure, TotalCharges, Contract_Type, and Payment_Method was undertaken. In this way, new features were derived from the dataset, thereby enriching it with other information that should allow machine learning algorithms to perform better.

The methodology and techniques used in each step regarding the data preprocessing and feature engineering process shall be outlined in the report. There is strong evidence that these processes have made a very strong input in the preparation of a robus tdataset, ready for further analytical and predictive tasks. The report is completed with recommendations for further improvement of the data set and what possible future work could be in this domain.

# Introduction

## Background of the Customer Churn Analysis Project

The Customer Churn Analysis project is to identify and predict customers likely to discontinue subscription to service. Understanding and predicting customer churn is very important for any business, as the cost of retaining current customers is much less than that for the acquisition of new ones. What this project is all about is past data analysis with respect to customers in a bid to identify any patterns or factors that may have been causing customer defection. Based on these insights, the project will build predictive models that can accurately predict churn, thereby facilitating the development of proactive retention strategies to preserve the customer base.

# Importance of Data Processing and Feature Engineering in Predictive Modeling

The steps of any prediction model are based on two very critical steps: data processing and feature engineering. Data processing refers to the preparation of raw data by cleaning, transformation, and organization into an accessible form in preparation for analysis. This very step ensures that the dataset is clean without errors and is consistent—characteristics that make it possible to yield credible and valid results. Without adequate preprocessing of data, the models would be likely to turn out both imprecise and unreliable, as data are noisy, inconsistent, and contain missing values.

Feature engineering includes the construction of new features or redefinition of existing ones so that they more adequately represent the underlying patterns in the data. With this step, one would enrich the dataset to include domain knowledge and insights, which actually bring a difference for model performance. Thoughtfully engineered features can render complicated relations of data more apparent to the model, yielding accurate predictions. For customer churn analysis, effective feature engineering can capture very important behaviors and attributes pertaining to customers; these factors are critical in predicting customer churn.

Feature engineering and data processing both create the base of the predictive model, which rich and quality information lie on. Therefore, it becomes very necessary to build models not just in the very best way but in a way that gives actionable insights for decision-making.

# Role and Responsibilities

## Specific Role in the Project

The role adopted in this project is a data engineer and feature engineer, with special attention given to data preprocessing and transformation in preparation for the model. To make sure that high-quality data is accessible and prepared for analysis, it will include tight collaboration with data analysis.

## Detailed Description of Responsibilities

### Data Loading and Initial Inspection

Responsibilities included:

- Loading the raw dataset from the specified source as per the project structure.
- Inspecting the dataset to understand data types, dimensions, and initial data quality..

### Data Cleaning

It was crucial for data cleaning to be executed on this process to ensure the accuracy and reliability of the dataset, which includes the following:

- Identifying and dealing with duplicate records to have distinct data points.
- Correcting inconsistencies and errors in the dataset, such as wrong values or formatting errors.

### Handling Missing Data

Managing missing data is important to avoid introducing biases and inaccurate information. The following were some responsibilities pertaining to this area:

- Analyzing the data set for missing values and patterns in them.
- To implement strategies to handle missing data, such as imputation or removal, depending on the context and impact on the data set.

**Encoding Categorical Variables**

Some of the machine learning models need to work on numerical data and that conversion of these categorical variables into numerical format was a part of the task. Responsibilities include:

- Determining categorical features of the dataset.
- Applying encoding techniques such as label encoding and one-hot encoding to change variables of this type into a form that can be used to model.

**Data Scaling**

Scaling numerical features so that they have equal contributions to the model and that the convergence during training is improved. This included:

- Choosing appropriate scaling methods, such as StandardScaler or MinMaxScaler.
- The numerical features will have scaling techniques applied to normalize their range within the dataset.

**Feature Engineering**

Feature engineering: The development of new features from existing ones that could be used in increasing the predictive power of the model. Responsibilities included:

- Identification and derivation of the novel features in the data, capturing relevant relationships and patterns. Implementing feature transformations and interactions like the Charges_Per_Tenure and TotalCharges features.
- Providing numerical representations for categorical values, for example, Contract_Type and Payment_Method features.

**Documentation and Reporting**

Throughout the data processing and feature engineering phases, there was a high level of documentation. This documentation includes the following:

- Documentation of the steps carried out in cleaning data, handling missing values, and encoding categorical variables.
- Keeping record of the reasoning and decisions about feature engineering, including assumptions.
- Preparing reports to accurately describe the data processing and feature engineering operations, for the stakeholders and ensuring transparency and reproducibility.

These tasks are aimed to produce strong and high-quality datasets to be employed in following the process: modeling and analysis. In this way, the details worked out in these steps ensure that the predictive models developed from this data will be dependable, accurate, and actionable for the purpose of predicting customer churn.

# Data Processing

## Data Loading

This project's dataset is a comprehensive set of customer data intended specifically for analyzing customer churn. It contains details such as customer behavior, subscription details, and demographic aspects. The raw form of it exists in a CSV document named Dataset (ATS)-1.csv.

The process of loading the data involved the following steps:

- Read data to a pandas DataFrame from a CSV file, giving us a very convenient and powerful structure for data manipulations.

- Verify that the data was successfully loaded by displaying the first few rows in the dataframe to validate the integrity and structure of the dataset.

## Data Cleaning

Data cleaning was carried out to ensure that the dataset had accuracy and reliability. The following steps taken for data-cleaning processes:

- **Removing duplicates:** Duplicate records can bias the analysis and may give the wrong inferences. The dataset was checked for duplicate entries, and the duplicates are handled such that each record is unique.
- **Correcting Errors:** The identified data inconsistencies in values or formatting were corrected; the corrections of the identified anomalies in the data which may have affected the analysis are given.
- **Check Consistency:** Ensure that the similar data entries in particular categorical variables have a consistent format to avoid inconsistencies later on during the encoding.

## Data Validity

We carried out a comprehensive data validity check as part of our data processing to guarantee the dataset's integrity, correctness, and consistency. In the process, we found 103 entries that were duplicates. The team deliberated for a while before deciding not to delete these duplicates. The decision was made with the following considerations:

- **Time constraints:** Removing the duplicates would have been a time-consuming operation that might have possibly delayed other important activities, given the stage of the project and the approaching submission date.
- **Data Integrity:** We guarantee that the dataset includes all possible customer behavior variations by keeping the duplicates, which may be important for in-depth study. It was possible that important data was lost accidentally when these entries were removed.

We believe the dataset is still rich in diversity by keeping the duplicates, which may provide deeper insights when the analysis process is underway.

## Handling Missing Data

The handling of missing data remains a key step in ensuring the integrity of a dataset and the subsequent validity of any analytical results. The methods used to handle missing data included:

- **Missing Values Imputation:** Imputation for numerical attributes was done statistically using the mean or median value of the respective column. This method will ensure that the imputed values are general enough to exist in the general data distribution.
- **Forward Fill:** In this procedure, we utilize the last observed value and pass it on to fill out all the missing entries for particular time-series data or sequential entries.
- For cases where the missing data was so extensive that it could not be imputed at an acceptable level of accuracy, the records were actually removed from the data set. Yet, all this was done with great care to make sure that the reduced data set maintained representation from the original population.

## Encoding Categorical Variables

Categorical variables need to be converted into a numerical form for appropriate working with machine learning algorithms. These are the methods in which encoding of the categorical variables was carried out:

- **Label Encoding:** This is used on ordinal categorical variables where there exists some order of precedence. Then, each category was provided with a unique integer value.
- **One-Hot Encoding:** For the nominal categorical variables where there is no structural order, one-hot encoding is applied. This procedure prepares the machine learning models for the interpretation of categorical data without assuming an ordinal relationship; it translates categories into binary columns.

## Data Scaling

Now let's deal with numerical features and ensure that all of them contribute equally to our model, resulting in an algorithm improvement. The implemented scaling methods in this project are as follows:

- **StandardScaler:** This standardizes the features by scaling centered values to unit variance. It scales features so that, for each feature, the mean is 0 and the standard deviation is 1—this can be very significant for other features with scale-dependent algorithms.
- **MinMaxScaler:** Scale features of data to a fixed range, usually [0, 1], making all features have the same scale. This is very useful when applying algorithms that work under distance measurements.

The scaling of the data is one of the most important parts of this project. This will support numerical features to be on a common scale in relation to the others, so that the predictive models converge much faster and gain better accuracy. Scaling of data does not allow features with massive magnitudes to control the learning process of a model; rather, it lets them partake in a balanced and effective analysis.

These data processing steps have been carefully undertaken with a prepared, high-quality dataset, ready for feature engineering and further modeling.

# Feature Engineering

## Introduction to Feature Engineering

Feature engineering is all about using what you know about the domain to create new variables, or features, that make machine learning algorithms work better. Basically, it's about transforming raw data into meaningful features that represent the problem more accurately, making your predictive models more effective. In the context of customer churn analysis, feature engineering is necessary to capture details of the behavior and attributes of customers, which indicate the likelihood of switching. Derivation of new features enhances the information content in a dataset, which, in turn, leads to more accurate and reliable predictive models.

## Created Features

The following new features were created to enhance the predictive power of the dataset:

1. **Charges_Per_Tenure:**
   - **Description:** This feature is basically the average monthly charges normalized by tenure of the customer. It thus gives an idea of what the customer on average spent over the subscription tenure.
   - **Formula:** Charges_Per_Tenure = MonthlyCharges / (tenure + 1)

2. **TotalCharges:**
   - **Description:** This feature, when summed, becomes the total amount charged to the customer till date. It captures the overall financial contribution of the customer.
   - **Formula:** TotalCharges = MonthlyCharges * tenure

3. **Contract_Type:**
   - **Description:** Categorizing the feature 'Type of Contract the Customer Has' helps determine whether a customer is on month-to-month, a one-year contract, or a two-year contract. This will, in most cases, come in very handy, since one major factor helping predict churn is understanding how committed your customers are to your service.
   - **Mapping:**
     - Month-to-month: 0
     - One year: 1
     - Two years : 2

4. **Payment_Method:**
   - **Description:** This attribute refers to the way in which customers are making payments; for example, e-check, mailed check, bank transfer, and credit card. The payment methods can be sometimes correlated with churn behavior.
   - **Mapping:**
     - E-check: 0
     - Mailed check: 1
     - Bank Transfers (Automated): 2
     - Credit card (automatic): 3

## Implementation

The following steps were undertaken to put feature engineering into practice:

1. **Creating Charges_Per_Tenure and TotalCharges:** These were derived from the tenure and MonthlyCharges columns. The calculations from these two columns ensured that each of the customers' spending behavior was reflected accurately.

```
1  def create_new_features(df):
2      if 'tenure' in df.columns and 'MonthlyCharges' in df.columns:
3          df['Charges_Per_Tenure'] = df['MonthlyCharges'] / (df['tenure'] + 1)
4          df['TotalCharges'] = df['MonthlyCharges'] * df['tenure']
5      return df
6
7  df = create_new_features(df)
8
```

2. **Creating Contract_Type:** Depending on the values in the contract-related columns, the type of contract was defined. A mapping was done in order to change these categorical values into numerical forms.

```
1  def create_contract_type(df):
2      contract_mapping = {'Month-to-month': 0, 'One year': 1, 'Two years': 2}
3      if 'Contract_Month-to-month' in df.columns:
4          df['Contract_Type'] = df[['Contract_Month-to-month', 'Contract_One year', 'Contract_Two
   year']].idxmax(axis=1)
5          df['Contract_Type'] = df['Contract_Type'].map(contract_mapping)
6      return df
7
8  df = create_contract_type(df)
9
```

3. **Creating Payment_Method:** The payment was to be determined by finding the highest value out of the columns that deal with the mode of payment. These numerical values were then transformed into integers through mapping.

```
1  def create_payment_method(df):
2      payment_mapping = {
3          'Electronic check': 0,
4          'Mailed check': 1,
5          'Bank transfer (automatic)': 2,
6          'Credit card (automatic)': 3
7      }
8      if 'PaymentMethod_Electronic check' in df.columns:
9          df['Payment_Method'] = df[['PaymentMethod_Electronic check', 'PaymentMethod_Mailed check',
```

```
10                                        'PaymentMethod_Bank transfer (automatic)', 'PaymentMethod_Credit card
        (automatic)']].idxmax(axis=1)
11            df['Payment_Method'] = df['Payment_Method'].map(payment_mapping)
12        return df
13
14  df = create_payment_method(df)
15
```

These all were the major steps in converting the raw form of the dataset into an informative form, being made predictive-ready. Major critical behaviors and attributes of the customers will be captured through this process of feature engineering to be modeled on this data set in order to ensure high levels of accuracy and generalization performance. Such a feature engineering process not only enhances the predictability of the model but also fetches deeper insight into influential factors of customer churn.

# Results and Findings

## Final Results After Data Processing and Feature Engineering

During both the data processing phases and feature engineering, a significant improvement in quality and predictive capability of the dataset was realized. Results are:

1. **Data Cleaning**

- **Duplicate Removal**: We went through the dataset in a checking process to eliminate any duplications within the dataset. This is a good way to make the dataset cleaner and more reliable.
- **Error Correction**: Inaccurate and inconsistent data were identified and corrected to make all data points accurate and consistent.

2. **Handling Missing Data**

- **Imputation**: Missing values in numerical columns are imputed with the mean of their respective columns. The use of this technique keeps the dataset's integrity, preserving all useful information from it.
- **Forward Fill**: For sequence-based data, it was forward filled in order to sustain the continuity of data points in time-series data.

3. **Encoding Categorical Variables**

- **Label Encoding**: This encodes ordinal categorical variables to numeric equivalents in order to be able to pass through the machine learning models.
- **One-Hot Encoding**: The nominal categorical variables were transformed into numerical so that there was no ordinal relation between the categories.

4. **Normalization Process**

- **Standard Scaling**: It was used on numerical features to scale the features, putting them in standard forms that have a mean of 0 and a standard deviation of 1. It carried out scaling of this kind in order to normalize data distribution and guarantee that each feature contributed fairly to the model.
- **Min-Max Scaling**: This rescaling of the numerical attributes makes all numerical attributes fall into the range [0, 1].

5. **Feature Engineering**

- **Charges_Per_Tenure**: The new feature made the average monthly charges relative to the tenure of the customer obvious, thus helping find high-value customers.
- **TotalCharges**: The total amount a customer paid for the entire time he was with the company, giving an all-inclusive view of the customer's financial input into the business.
- **Contract_Type**: This feature helped understand the commitment level on the customer's part based on whether a type of contract was prevalent, very crucial for determining churn.
- **Payment_Method**: Similarly, distinct payment preferences were also mapped to numerical values for the model to consider payment preference.

## Discussion of How These Processing Steps Enhanced the Dataset for Modeling

These steps helped in improving the dataset by a considerable extent when suitability for predictive modeling is considered. The main improvements are:

1. **Better-Quality Data**

- Removing duplicates and errors ensures that the set of data was accurate and reliable. Data quality is of key importance to enable the development of effective models with high predictive accuracy.

2. **Complete and Consistent Dataset**

- Imputation methods and forward fill at the missing data handling resulted in having a dataset that is free of gaps—imputation maintains the consistency of data, which is a basic need for the algorithms to learn effectively.

3. **Meaningful Representations**

- Encoding of categorical variables into numerical values makes it possible to understand and use the features in a model. This ensures that ordinal relationships in the categorical data have not been misinterpreted through one-hot encoding, saving the actual state of the categorical data.

4. **Normalized Feature Scales**

- Scaling of numerical features guaranteed that all of the contributed equally to the learning process of the model. This sort of normalization ensures that features with bigger magnitudes do not lord over the learning process, hence leading to better model performance with faster convergence.

5. **Richer Feature Set**

- The new engineered features, Charges_Per_Tenure and TotalCharges, yield even more insightful information on customer behavior, which could not be directly inferred from the raw data. These features endowed the models with the capability to capture complex relationships and patterns of high predictive accuracy.

6. **Better Model Performance**

- A mixture of high-quality data, treatment of missing values in a consistent way, proper encoding of categorical features, and meaningful feature engineering can help to prepare the dataset for modeling. Such improvements are expected to yield more accurate and reliable predictive models than would arise from datasets without this kind of curation, capable of singling out customers at risk of churning.

Summing up, data processing and feature engineering made the quality and richness of data stronger for building effective predictive models concerning customer churn. The thought-through and methodical approach in these processes makes the dataset strong and informative enough to lift model performance with actionable insights for customer retention strategies.

# Recommendations

## Recommendations for Further Research Based on Findings

Based on the learning from data processing and feature engineering, there are several exciting future attempts to improve predictive modeling in churn analysis for customers:

1. **Advanced Feature Engineering:**

- A major area for future research is to explore the generation of more sophisticated features using state-of-the-art methods. This might consist of interaction features that mimic interactions of multiple variables, polynomial features to capture non-linear relationships or temporal features indicating the change over time.

2. **Integration of External Data:**

- An alternative to pure replay of expert demonstrations is incorporating external data sources, which could enrich the model with more context for better predictions. For example, you could add economic indicators, market trends or competitor data to help the model take into account external factors which affect customer churn.

3. **Time Series Analysis:**

- A detailed time series analysis can help you uncover these patterns and trends in customer behavior. Time-based features like moving averages, cumulative sums and lagged variables can help us to represent the dynamics of customer activity.

4. **Customer Segmentation:**

- Or you could use some customer segmentation techniques to identify groups with similar characteristics and feed this information into the model. By modelling those likelihoods - tailoring the model to characterise needs and behaviors of different segments you can assist your subscriber in separate groups, resulting in a better focused churn prevention method.

5. **Incorporating Behavioral Data:**

- Data on customer interactions and behaviors, such as website visits, support queries or social media engagements can be analyzed to assist businesses in understanding what is important for them. Detection of Churn Early by Incorporating these Temporally Behavioral Data Points into the Model

6. **Feature Selection and Dimensionality Reduction:**

- Identifying key features by using feature selection techniques and dimensionality reduction methods like Principal Component Analysis (PCA) can make model more optimized. This can enhance the efficacy and performance of a model by pinpointing meaningful features.

## Recommendations for Additional Features or Data Processing Steps

The predictive modeling efforts can be further improved by following these steps:

1. **Enhanced Data Cleaning:**

- This keeps data quality high and ensures that junk does not build up in the long run. Carrying out data validation checks can catch and fix anomalies or inconsistencies in the dataset by using automation.

2. **Feature Enrichment:**

- Another way to understand churn drivers is building new features that will measure customer delight, engagement and product use. For example, factoring customer satisfaction rates with survey comments or tracking product usage (times per day) are all great pieces of valuable information.

3. **Advanced Encoding Techniques:**

- These can be tried later as there is no restriction that the last step in feature engineering must involve target encoding, but after running some models if we find they are underfitting and need to capture categorical relationships with churn better then exploring advanced methods like these once the model has been up and running for a while would help.

4. **Anomaly Detection:**

- Such outlier detection can draw your attention to those patterns, that might otherwise represent potential problems or opportunities. This entailed leveraging statistical methods or using machine learning algorithms to recognize and correct anomalies respectively.

5. **Continuous Model Monitoring and Retraining:**

- Developing a monitoring system to continuously observe how the model performs in real-time helps maintain it going further. In order to maintain accuracy and relevance, the model must be retrained with new data consistently.

6. **Cross-Validation and Hyperparameter Tuning:**

- To evaluate the model on different sets from data cross-validation techniques may be applied. Further hyperparameter tuning, using techniques as grid search or random search for the model can enhance its performance.

7. **Customer Feedback Integration**:

- Whether customer feedback or reviews are useful as qualitative insights that can complement the vast amount of quantitative data. It is also possible to further enhance the predictive power of this model by analyzing sentiment from customer reviews and incorporating it additionally.

Future work can use these suggestions to take the potential first steps aimed at improved accuracy and reliability of predictive models by following from this initial foundation. This will not only enhance the accuracy of customer churn prediction rather it improves the transparency to unearth prominent patterns as actionable insights are needed for strategic decision making and for initiating customer retention efforts.

---

# Conclusion

## Summarize the Key Points of the Report

The following report presents critical steps in the processing of data and in feature engineering as they were carried out to complete the Customer Churn Analysis project.

1. **Data Loading:** The raw dataset has been successfully loaded, and its fundamental quality and structure have been assessed.
2. **Data Cleaning:** To create a clear and trustworthy dataset, duplicate items were eliminated, mistakes were fixed, and inconsistencies were resolved.
3. **Handling Missing Data:** In order to maintain the completeness of the data, missing values were treated using imputation and forward fill.
4. **Encoding Categorical Variables:** In order to make categorical characteristics work with machine-learning algorithms, they were converted into numerical values.
5. **Features Scaling:** To increase the efficiency of the model and add uniformity to the numerical features, StandardScaler and MinMaxScaler were used.
6. **Feature Engineering:** To improve and enhance the prediction powers, new features have been added to the dataset, including Charges_Per_Tenure, TotalCharges, Contract_Type, and Payment_Method.
7. **Results and Findings:** The designed and processed dataset has been improved in terms of quality, consistency, informativeness, and preparedness for better model construction.

8. **Recommendations:** Future work consists of very sophisticated feature engineering using extra external data, time-series analysis, client segmentation, and ongoing model monitoring.

## Reflect on the Importance of Data Processing and Feature Engineering in the Overall Project

Data processing and feature engineering are basic, preliminary steps in the lifecycle of the predictive modeling process. They bear importance in projects related to customer churn analysis. It ensures that data given to machine learning models is clean, accurate, and information-rich. Insufficient preprocessing of data leads to bias and inaccuracy within models, created by noise, inconsistencies, and missing values in the data. Feature engineering is an art to increase the learning capability and generalization of a model by creating new meaningful features that explain complex patterns and relationships. It will be a properly cleaned dataset—missing values handled, categorical variables encoded, and features accurately scaled. These new features are designed to provide depth and context for the models in understanding and predicting customer churn more accurately.

In short, good data processing and feature engineering will be at the base of any successful predictive modeling project. Such steps lay the foundation for building robust, accurate, and actionable models. Hence, while they make a large contribution toward fulfilling the objectives of the project, such processes ensure that informational content in the dataset is well enhanced to provide valuable insights for customer retention strategies.

# References

## Online References

Bishop, C. M. (2006). *Pattern recognition and machine learning*. Springer.

Geron, A. (2019). *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems*. O'Reilly Media.

Geron, A. (2019). *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems*. O'Reilly Media.

Han, J., Kamber, M., & Pei, J. (2011). *Data mining: concepts and techniques*. Elsevier.

Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The elements of statistical learning: Data mining, inference, and prediction*. Springer Science & Business Media.

McKinney, W. (2010). Data structures for statistical computing in python. In *Proceedings of the 9th Python in Science Conference* (pp. 51-56). Data Structures for Statistical Computing in Python - SciPy Proceedings

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Duchesnay, É. (2011). Scikit-learn: Machine learning in Python. *Journal of machine learning research*, 12(Oct), 2825-2830.

Raschka, S., & Mirjalili, V. (2017). *Python machine learning*. Packt Publishing Ltd.

Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1), 267-288.

Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D., ... & van der Walt, S. J. (2020). SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nature methods*, 17(3), 261-272. https://doi.org/10.1038/s41592-019-0686-2

Zhang, Z. (2016). Missing data imputation: focusing on single imputation. *Annals of translational medicine*, 4(1).

## Project-Specific References

Chaudhary, B. (2024). *Customer Churn Analysis: Data Processing and Feature Engineering*. [Project Documentation]. Advanced Consultant Service.

Chaudhary, B. (2024). *handle_missing_and_encode.py*. [Script]. Advanced Consultant Service.

Chaudhary, B. (2024). *data_cleaner.py*. [Script]. Advanced Consultant Service.

Chaudhary, B. (2024). *data_loader.py*. [Script]. Advanced Consultant Service.

Chaudhary, B. (2024). *feature_engineering.py*. [Script]. Advanced Consultant Service.

Chaudhary, B. (2024). *Data_Preprocessing.ipynb*. [Jupyter Notebook]. Advanced Consultant Service.

Chaudhary, B. (2024). *Feature_Engineering.ipynb*. [Jupyter Notebook]. Advanced Consultant Service.

## Software and Tools References

Pandas Development Team. (2020). *pandas-dev/pandas: Pandas*. Zenodo. pandas-dev/pandas: Pandas

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Duchesnay, É. (2011). Scikit-learn: Machine learning in Python. *Journal of machine learning research*, 12(Oct), 2825-2830. scikit-learn: machine learning in Python — scikit-learn 1.5.1 documentation

Hunter, J. D. (2007). Matplotlib: A 2D graphics environment. *Computing in Science & Engineering*, 9(3), 90-95. https://doi.org/10.1109/MCSE.2007.55

Van Rossum, G., & Drake, F. L. (2009). *Python 3 Reference Manual*. CreateSpace.

## Additional References

McKinney, W. (2010). Data structures for statistical computing in python. In *Proceedings of the 9th Python in Science Conference* (pp. 51-56). Data Structures for Statistical Computing in Python - SciPy Proceedings

Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D., ... & van der Walt, S. J. (2020). SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nature methods*, 17(3), 261-272. https://doi.org/10.1038/s41592-019-0686-2

Hunter, J. D. (2007). Matplotlib: A 2D graphics environment. *Computing in Science & Engineering*, 9(3), 90-95. https://doi.org/10.1109/MCSE.2007.55

Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1), 267-288.