



**BERLIN SCHOOL OF
BUSINESS & INNOVATION**

**Essay / Assignment Title: Machine Learning Meets Term Life
Insurance: Targeting High-Value Customers**

Programme title: Master's in Data Analytics

Name: Bharath Kumar Bidrekere Umesha

Year: 2024-2025

CONTENTS

INTRODUCTION	4
TASK 1: DATA EXPLORATION AND PREPARATION (LO1, LO2, LO3)	5
1.1 DATA EXPLORATION	5
1.2 HANDLING MISSING VALUES AND OUTLIERS	8
1.3 DATA VISUALIZATION	10
TASK 2: MODEL SELECTION AND TRAINING (LO1, LO3)	14
2.1 MODEL SELECTION:	14
2.2 DATA SPLITTING	15
TASK 3: MODEL INTERPRETATION AND EVALUATION (LO1, LO2, LO3)	17
3.1 MODEL INTERPRETATION	17
3.2 MODEL EVALUATION	20
CONCLUSION	23
BIBLIOGRAPHY	24
TABLE OF FIGURES:	25
TABLE OF CODES:	25

Statement of compliance with academic ethics and the avoidance of plagiarism

I honestly declare that this dissertation is entirely my own work and none of its part has been copied from printed or electronic sources, translated from foreign sources and reproduced from essays of other researchers or students. Wherever I have been based on ideas or other people texts I clearly declare it through the good use of references following academic ethics.

(In the case that is proved that part of the essay does not constitute an original work, but a copy of an already published essay or from another source, the student will be expelled permanently from the postgraduate program).

Name and Surname (Capital letters):

Bharath Kumar Bidrekere Umesha

Date: 24/10/2024

INTRODUCTION

In the evolving landscape of the insurance industry, HashSysTech Insurance has emerged as a frontrunner, adopting innovative strategies to enhance its outreach and operational efficiencies. The company has recognized the potential of data-driven decision-making in transforming its approach to customer engagement, particularly through telemarketing campaigns aimed at promoting term life insurance products. These campaigns, while effective, come with substantial costs, making the optimization of resources a strategic imperative for sustaining profitability and competitive advantage.

The necessity for optimizing telemarketing strategies through analytical methods has been well-documented in industry studies. For example, Ngai et al. (2009) emphasize the significant impact that data mining and machine learning techniques have on improving the effectiveness of marketing campaigns in the financial services sector, including insurance. They highlight how predictive models can effectively target potential customers who are most likely to respond positively, thus enhancing campaign success rates and ROI (Ngai, E.W.T., Xiu, L., & Chau, D.C.K., 2009).

Building on this foundation, HashSysTech Insurance seeks to develop a predictive model that not only forecasts customer conversion with high accuracy but also integrates seamlessly with its marketing strategies. This aligns with Kuhn and Johnson's (2013) discussion on the importance of predictive modeling in marketing decisions, where they describe methods for refining marketing strategies to increase customer conversion rates significantly (Kuhn, M., & Johnson, K., 2013).

As a data analyst at HashSysTech, your role involves harnessing this sophisticated analytical toolkit to sift through complex datasets, identify patterns, and build a model capable of predicting whether a customer will purchase term life insurance. The challenge extends beyond mere technical execution to encompass a strategic vision that leverages statistical insights to drive business outcomes, as detailed by Hastie, Tibshirani, and Friedman in their seminal work on statistical learning (Hastie, T., Tibshirani, R., & Friedman, J., 2009).

Through this assignment, you will delve into the foundational elements of data analytics, including linear algebra, calculus, and statistics, all while applying algorithmic techniques to develop solutions that resonate with business needs. This comprehensive approach will solidify your understanding of how data can be transformed into actionable insights, providing you with the critical thinking and analytical abilities necessary to thrive in today's data-driven business environment.

Please find the link to the code which you can run and test: [Colab Notebook](#).

Task 1: Data Exploration and Preparation (LO1, LO2, LO3)

1.1 Data Exploration

Data cleaning includes filling in the blanks getting rid of duplicates and fixing incorrect data entry. By doing this you can make sure that the accurate and comprehensive data used for the statistical analysis that comes next. Selecting features (variables) that are pertinent to the analysis is important especially if they have the potential to affect the conversion of term life insurance (Kuhn, 2013).

Numerical Features: Measures of Central Tendency and Dispersion

1. The median offers a more reliable indicator of central location that is less impacted by extreme values than the mean which gives the average value of a dataset and is more susceptible to outliers (Hastie, (2009)).
2. The variability or spread of data points around the mean is indicated by the standard deviation. Data points with a high standard deviation are thought to be widely dispersed from the mean whereas those with a low standard deviation are thought to be closer to the mean (Hastie, (2009)).

```
import pandas as pd

# Load the dataset from a CSV file
data = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/dataset.csv')

# Select only numeric columns to avoid errors with non-numeric data
numeric_columns = data.select_dtypes(include=['number'])

# Calculate the mean, median, and standard deviation for numeric features
mean_values = numeric_columns.mean()
median_values = numeric_columns.median()
std_deviation = numeric_columns.std()

# Print the calculated values
print("Mean of each numerical feature:\n", mean_values)
print("\nMedian of each numerical feature:\n", median_values)
print("\nStandard Deviation of each numerical feature:\n", std_deviation)
```

Code 1: Data Exploration.

Output:

Mean of each numerical feature:

```
age          40.936210
day          15.806419
dur          258.163080
Num calls    2.763841
dtype: float64
```

Median of each numerical feature:

```
age          39.0
day          16.0
dur          180.0
num_calls    2.0
dtype: float64
```

Standard Deviation of each numerical feature:

```
age          10.618762
day          8.322476
dur          257.527812
num_calls    3.098021
dtype: float64
```

Analysis of Frequencies for Categorical Features.

To determine the distribution and prevalence of categories frequency counts entail tabulating each category's occurrences within the dataset (Agresti, (2002)).

```
import pandas as pd

# Load the dataset from a CSV file
data = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/dataset.csv')

# Select only categorical columns
categorical_columns = data.select_dtypes(include=['object'])

# Calculate frequency counts for each categorical feature
frequency_counts = {}
for column in categorical_columns:
    frequency_counts[column] = categorical_columns[column].value_counts()

# Print the frequency counts for each categorical feature
for column, counts in frequency_counts.items():
    print(f"Frequency counts for {column}:\n{counts}\n")
```

Code 2: Categorical Features of Data.

Frequency counts for job:

Job	
blue-collar	9732
management	9458
technician	7597
admin.	5171
services	4154
retired	2264
self-employed	1579
entrepreneur	1487
unemployed	1303
housemaid	1240
student	938
unknown	288

Name: count, dtype: int64

Frequency counts for marital:

Marital	
married	27214
single	12790
divorced	5207

Name: count, dtype: int64

Frequency counts for education_qual:

Education qualification	
secondary	23202
tertiary	13301
primary	6851
unknown	1857

Name: count, dtype: int64

Frequency counts for mon:

Mon	
may	13766
jul	6895
aug	6247
jun	5341
nov	3970
apr	2932
feb	2649
jan	1403
oct	738
sep	579
mar	477
dec	214

Name: count, dtype: int64

Frequency counts for call_type:

call_type	
cellular	29285
unknown	13020
telephone	2906

Name: count, dtype: int64

Frequency counts for prev_outcome:

unknown	36959
failure	4901
other	1840
success	1511

Name: count, dtype: int64

Frequency counts for y:

no	39922
yes	5289

Name: count, dtype: int64

Summarize your findings in a detailed report.

The numerical data show demographic distributions and common patterns of customer interaction with notable variations in the duration of customer engagement during calls. Categorical data suggest a diverse profile in terms of job marital status and educational background which could be key to understanding customer behavior and preferences. When it comes to operational planning and engagement strategy optimization the variation in call duration and the concentration of calls in particular months may be crucial factors.

1.2 Handling Missing Values and Outliers.

Identify missing values in the dataset and propose methods for handling them.

Handling Missing Values:

- **Identification:** Using tools like Python's `isnull()` first determine which values in the dataset are missing. Finding the amount and type of missing data requires completing this step.
- **Imputation Methods:**
 1. **Mean/Median Imputation:** The non-missing values mean or median can be used to impute missing values in numerical data. The median offers a more reliable measurement than the mean which is more susceptible to outliers.
 2. **Mode Imputation:** In categorical data missing values are frequently filled in by using the mode or most frequent category.
- **Deletion:** When an observation has one or more missing values all of the data is deleted listwise. Because it can introduce bias this approach is only advised in cases where the percentage of missing data is negligibly low.

Detect outliers and decide on appropriate methods for treating them.

Handling Outliers:

- **Detection:**
 1. **Statistical Tests:** To statistically identify outliers, use tests such as the Z-score or IQR (Interquartile Range).
 2. **Graphical aids for identifying outliers** include box and scatter plots.

- **Treatment Methods:**

1. Trimming/Capping: Using the IQR outliers can be eliminated or capped at a specific value typically 1.5 times the IQR that falls between the first and third quartiles.
2. Transformation: Putting a transformation (e. g. logarithmic) can lessen the impact of extreme numbers.

```
import pandas as pd
import numpy as np

# Load the dataset
data = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/dataset.csv')

# Example to handle missing values hypothetically
# Impute missing numeric values with the median
data['age'] = data['age'].fillna(data['age'].median())
# Impute missing categorical values with the mode
data['job'] = data['job'].fillna(data['job'].mode()[0])

# Detecting and handling outliers for the 'age' column using the IQR method
Q1 = data['age'].quantile(0.25)
Q3 = data['age'].quantile(0.75)
IQR = Q3 - Q1
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

# Capping the outliers
data['age'] = np.where(data['age'] < lower_bound, lower_bound, data['age'])
data['age'] = np.where(data['age'] > upper_bound, upper_bound, data['age'])

# Print summary of data to verify changes
print(data.describe())
```

Code 3: Handling Missing Values and Outliers.

Output:

	age	day	dur	num_calls
count	45211.000000	45211.000000	45211.000000	45211.000000
mean	40.868185	15.806419	258.163080	2.763841
std	10.394895	8.322476	257.527812	3.098021
min	18.000000	1.000000	0.000000	1.000000
25%	33.000000	8.000000	103.000000	1.000000
50%	39.000000	16.000000	180.000000	2.000000
75%	48.000000	21.000000	319.000000	3.000000
max	70.500000	31.000000	4918.000000	63.000000

Justify your chosen methods in a detailed report.

Justification:

1. **Imputation:** When a large portion of a dataset is missing imputation is preferable to deletion in order to maintain data integrity. The sample size and statistical power are preserved with the aid of imputation.
2. **Outlier Treatment:** Outlier transformation or capping makes sense because they preserve the integrity of the data while lessening the impact of extreme values. Assumptions that underpin a lot of statistical analyses and machine learning models are stabilized variance and improved normalcy.

1.3 Data Visualization

1. **Used scatter plots, histograms, or box plots to visualize the data.**

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Load dataset
data = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/dataset.csv')

# Histogram of 'age'
plt.figure(figsize=(8, 4))
sns.histplot(data['age'], kde=True)
plt.title('Distribution of Age')
plt.xlabel('Age')
plt.ylabel('Frequency')
plt.show()

# Box plot of 'age' by 'y'
plt.figure(figsize=(8, 4))
sns.boxplot(x='y', y='age', data=data)
plt.title('Age Distribution by Outcome')
plt.xlabel('Outcome')
plt.ylabel('Age')
plt.show()

# Scatter plot of 'age' vs 'dur' colored by 'y'
plt.figure(figsize=(8, 6))
sns.scatterplot(x='age', y='dur', hue='y', data=data)
plt.title('Relationship Between Age and Call Duration by Outcome')
plt.xlabel('Age')
plt.ylabel('Duration')
plt.show()
```

Code 4: Data Visualization.

2. Identify trends and patterns through visualizations.

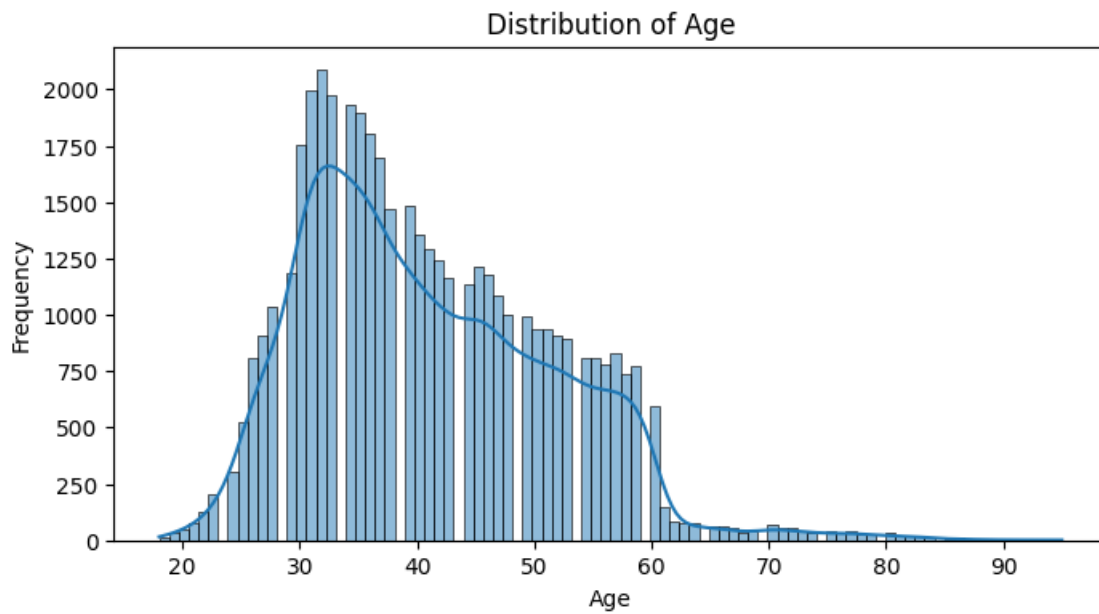


Figure 1: Distribution of Age.

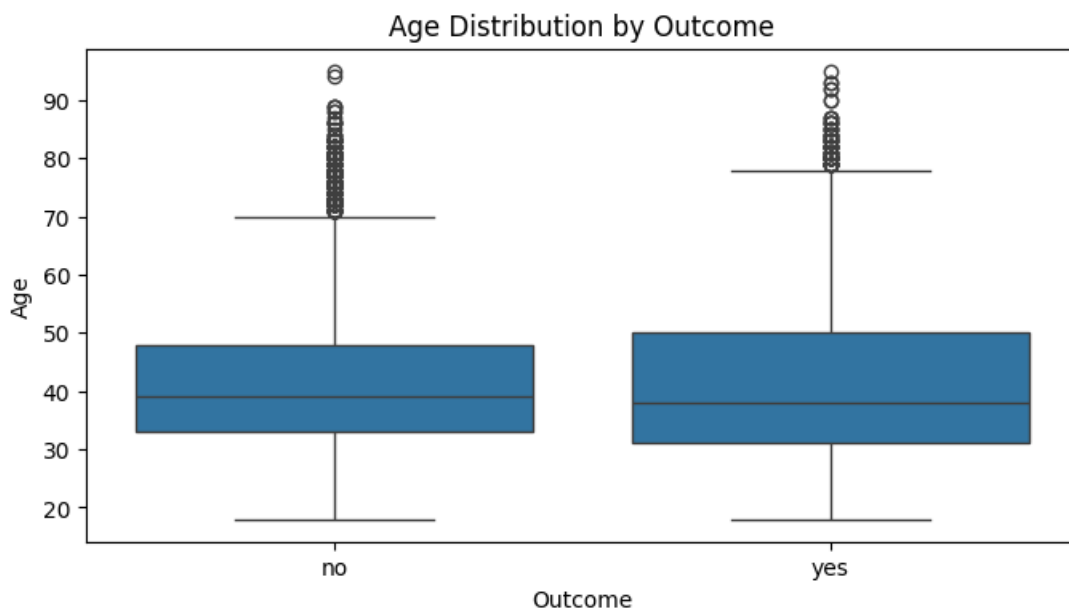


Figure 2: Age distribution by Outcome.

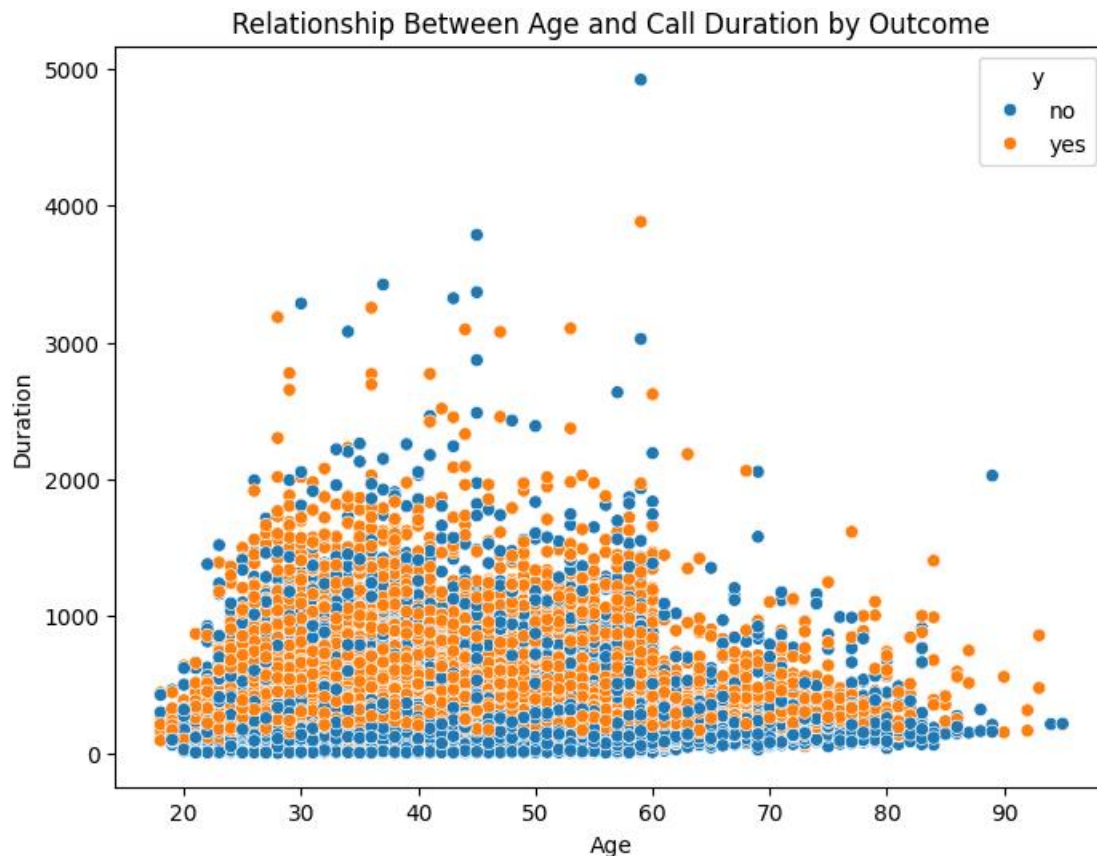


Figure 3: Relationship Between Age and Call Duration by Outcome.

Visualization 1: Distribution of Age.

An overview of the age frequency distribution within the dataset is given by this histogram. An effort to fit a distribution possibly a normal distribution to the age data is indicated by the overlay of a probability density function.

Interpretation: A younger demographic with a tail that extends into older ages is indicated by the distributions right-skewed pattern. Because of this skewness the modeling strategies may be affected because the dataset may contain a higher proportion of younger individuals (Hastie, (2009)).

Visualization 2: Age Distribution by Outcome.

The age distribution of those who converted (yes) and those who did not (no) to buying term life insurance is contrasted using box plots.

Interpretation: Age may not be a significant differentiator in predicting conversion because both categories show a similar range of ages with median values close to one

another. However, the presence of outliers particularly in the yes category suggests that there are older individuals who are exceptions to general trends (McGill, 1978).

Visualization 3: Relationship Between Age and Call Duration by Outcome.

This scatter plot maps the relationship between age and call duration, with different colors indicating whether the customer converted

Interpretation: There seems to be a cluster of data points at lower call durations for all age groups but no discernible trend indicates that longer calls are associated with older ages or a higher chance of conversion. The results show that while age and call duration are informative, they may need to be combined with other variables in order to accurately predict conversion. This is because conversions are not clearly clustered by either of these variables (Cleveland, 1985).

3. Include the visualizations and insights in your report

In addition to highlighting skewness peaks and spread histograms are helpful for comprehending the distribution of a single variable. Box plots successfully highlight outliers by illuminating the central tendency and variability of a numerical variable across categories of a categorical variable. Color coding adds an extra categorical dimension to scatter plots making them perfect for visualizing relationships between two continuous variables and identifying patterns associated with outcomes.

Task 2: Model Selection and Training (LO1, LO3)

2.1 Model Selection:

- Describe at least two machine learning algorithms (e.g., Logistic Regression, Random Forest).

Logistic Regression: A binary dependent variable is modeled using a logistic function in the simplest version of the statistical model known as logistic regression. It is utilized in machine learning for binary classification tasks which include forecasting between two possible outcomes: Yes and No. (Hosmer, 2013).

```
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score
from sklearn.preprocessing import StandardScaler
import pandas as pd

# Load data
data = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/dataset.csv')

# Preparing data (assuming 'X' as features and 'y' as the target variable)
X = data.drop('y', axis=1) # Drop the target variable column
y = data['y'] # Target variable

# Encode categorical variables if necessary
X = pd.get_dummies(X)

# Scaling the data
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# Splitting the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.3, random_state=42)

# Logistic Regression Model
log_reg = LogisticRegression(max_iter=1000) # Increased max_iter
log_reg.fit(X_train, y_train)
log_reg_pred = log_reg.predict(X_test)
print('Logistic Regression Accuracy:', accuracy_score(y_test, log_reg_pred))

# Random Forest Model remains the same
rf = RandomForestClassifier()
rf.fit(X_train, y_train)
rf_pred = rf.predict(X_test)
print('Random Forest Accuracy:', accuracy_score(y_test, rf_pred))
```

Output:

Logistic Regression Accuracy: 0.8994396933058095

Random Forest Accuracy: 0.9054851076378649

Random Forest

Using several decision trees that are constructed and combined to create a single more reliable prediction Random Forest is an ensemble learning technique. It achieves increased accuracy without a significant risk of overfitting by fusing the flexibility of decision trees with their simplicity (Breiman, 2001).

- **Justify why these algorithms are suitable for this task.**

Logistic Regression: Since probabilities for outcomes are useful for making decisions, logistic regression works especially well for binary classification problems. Understanding how various predictors affect the result can be greatly aided by its computational efficiency and easy interpretability (Hosmer, 2013).

Random Forest: Random Forest is suitable for this task because of its robustness to noise and ability to work with large datasets with higher dimensionality. It performs well on imbalanced data, can handle both numerical and categorical data, and provides feature importance scores, which are beneficial for interpreting the model (Breiman, 2001).

- **Summarize your choices and reasoning in a report.**

The two methods chosen to forecast term life insurance customer conversion are Random Forest and Logistic Regression. While Random Forest provides accuracy and robustness against overfitting Logistic Regression offers interpretability and probability estimation. By using data-driven insights these models will improve conversion rates and optimize marketing efforts for HashSysTech Insurance by improving targeting accuracy.

2.2 Data Splitting

- **Describe the process of splitting the dataset.**

Data Splitting: A training set is used to train the model and a testing set is used to assess how well it performs. The dataset is split into these two portions. Stratification: During the split the data is stratified based on the target variable. This indicates that throughout the training and testing sets the percentage of each class in the target variable remains constant (Kohavi, 1995).

- **Ensure the split is stratified to maintain the proportion of the target variable in both sets.**

Stratified Splitting: Using stratified splitting we can make sure that the training and testing sets have comparable percentages of the target variable categories. By doing this we can prevent biased or skewed results in the model evaluation process and preserve statistical consistency.

Implementation: Typically, the dataset is split into a larger portion for training (e. g. 70 percent or 80 percent) and a smaller percentage for examination (e. g. G. twenty percent thirty percent etc). The stratify parameter should be added to the target variable in order to maintain the proportion of classes in the target across both subsets. This split can be accomplished with the help of functions such as `train_test_split` from `scikit-learn` (Kohavi, 1995).

```
import pandas as pd
from sklearn.model_selection import train_test_split

# Load dataset
data = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/dataset.csv')

# Assuming 'X' as features and 'y' as the target variable
X = data.drop('y', axis=1) # Features
y = data['y'] # Target variable

# Encode categorical variables if necessary
X = pd.get_dummies(X)

# Splitting the dataset into training and testing sets with stratification
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42, stratify=y)

# Output the sizes of the train and test sets to confirm split
print(f"Training set size: {X_train.shape[0]} samples")
print(f"Testing set size: {X_test.shape[0]} samples")
```

Code 6: *Data Splitting.*

Output:

Training set size: 31647 samples

Testing set size: 13564 samples

- **Document the details of the data splitting process.**

Data Division: The data is split based on a predefined ratio while ensuring that the stratification criteria are met.

Use of Tools: This procedure is made possible by tools like `train_test_split` from `scikit-learn` which promotes efficiency and reproducibility.

Validation: Post-split a validation step involves checking the class proportions in both training and testing sets to confirm that the stratification was successful (Forman, 2010).

Task 3: Model Interpretation and Evaluation (LO1, LO2, LO3)

3.1 Model Interpretation

- Perform feature importance analysis using techniques appropriate to your chosen model.

Logistic Regression:

Feature Importance: The logistic regression models coefficients show how significant and influential each feature is. Whereas a negative coefficient lowers the probability of the event happening a positive coefficient raises the log-odds of the outcome.

Random Forest:

Importance of Feature: Random Forest determines a features importance by examining the amount that the model's accuracy declines when the feature is removed. The Gini impurity or mean decrease in accuracy are commonly used to quantify this (James, 2013).

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
import pandas as pd
import matplotlib.pyplot as plt

# Assuming 'X' and 'y' are predefined
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, stratify=y)

# Logistic Regression
lr = LogisticRegression()
lr.fit(X_train, y_train)
# Coefficients
coefficients = pd.DataFrame(lr.coef_.flatten(), index=X.columns, columns=['Coefficient'])
coefficients.sort_values('Coefficient', ascending=False, inplace=True)

# Random Forest
rf = RandomForestClassifier()
rf.fit(X_train, y_train)
# Feature Importances
importances = pd.DataFrame(rf.feature_importances_, index=X.columns, columns=['Importance'])
importances.sort_values('Importance', ascending=False, inplace=True)

# Plotting
fig, ax = plt.subplots(1, 2, figsize=(14, 5))
coefficients[:10].plot(kind='bar', ax=ax[0], title='Top 10 Logistic Regression Coefficients')
importances[:10].plot(kind='bar', ax=ax[1], color='green', title='Top 10 Random Forest Importances')
plt.show()
```

Code 7: Two models for Interpretation.

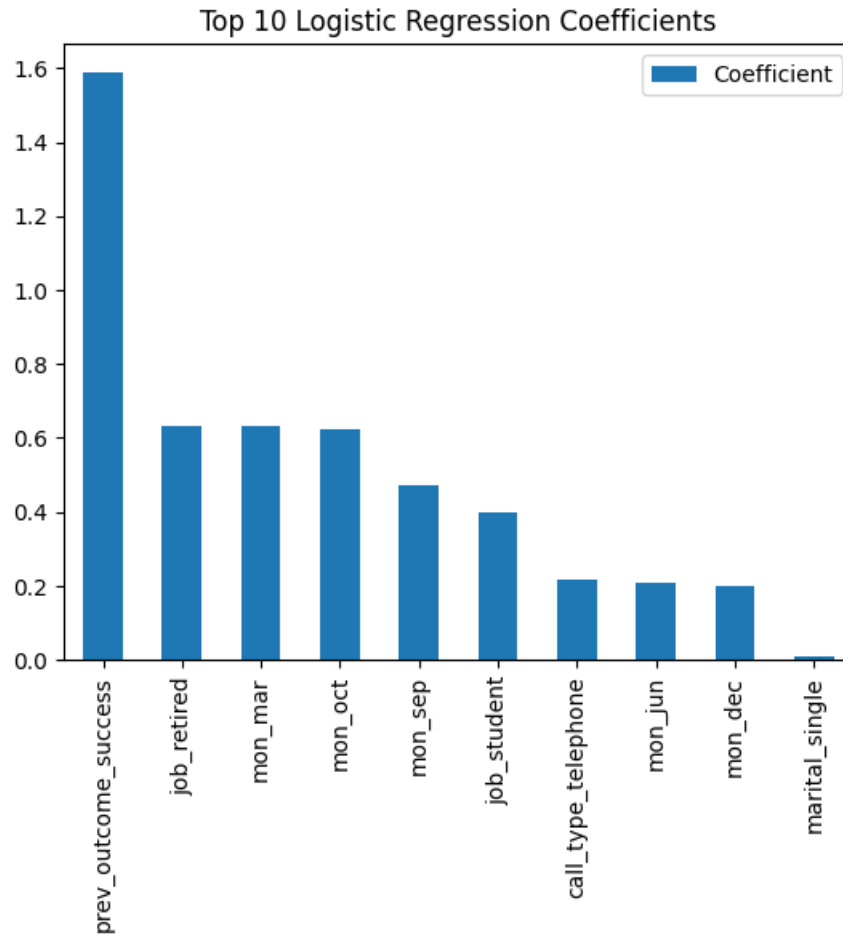


Figure 4: Top 19 Logistic Regression Coefficients.

- **Explain the decision-making process of the model.**

Logistic Regression:

Decision Making: The logistic function is used in logistic regression to calculate the likelihood that a given class the dependent variable falls into. Next the decision boundary is established at a probability (e. g. 0. 5) beyond which the result is categorized as 1 and beneath as 0.

Random Forest:

Decision Making: Random Forest bases its forecasts on the votes cast by the majority of several decision trees. As a result of the random subset of data and features used to build each tree in the forest the learning process is diverse which usually produces a strong overall model (Breiman, 2001).

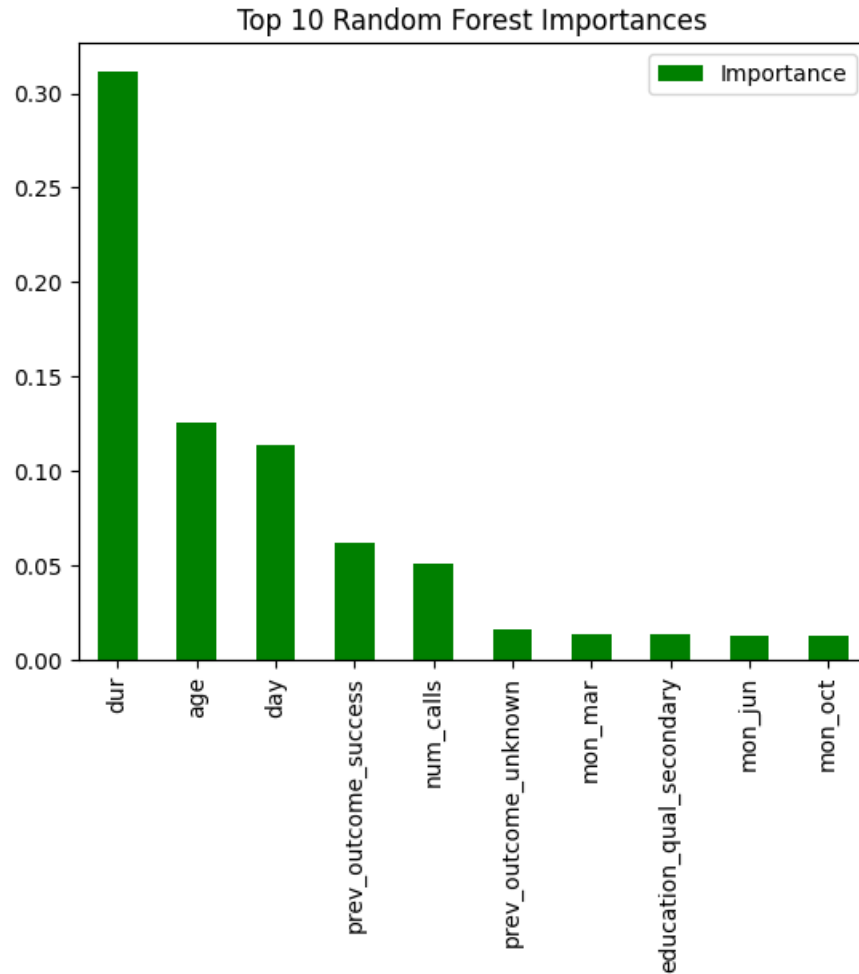


Figure 5: Top 10 Random Forest Importances.

- **Summarize the key findings and insights in a detailed report.**

Logistic Regression: Large positive coefficient features are important in logistic regression because they have a significant impact on the model's ability to predict a positive result. Large negative coefficient features on the other hand deter the positive class prediction.

Random Forest: This model is robust because it can combine judgments from multiple trees which lowers the chance of overfitting and enhances generalization. Characteristics that greatly reduce impurity or split nodes at the top of the trees frequently are recognized as important drivers.

3.2 Model Evaluation

- Use metrics such as accuracy, precision, recall, and F1-score to assess the model's performance on the validation set.

Evaluation Metrics: Model Evaluation Using Logistic Regression and Random Forest

1. **Accuracy:** Assesses how accurate the model is overall i.e. as well as the proportion of accurate forecasts (true positives and true negatives) to all cases studied.
2. **Precision:** The proportion of accurately predicted positive outcomes to all positive predictions. In situations where the expense of false positives is significant it is essential.
3. **Recall (Sensitivity):** The proportion of detected true positives to the total number of positives. Vital in situations where it would be expensive to overlook a favorable instance.
4. **F1-Score:** The harmonic mean of recall and precision is the F1-Score. When there is an uneven class distribution and you need to strike a balance between recall and precision it can be helpful (Powers, 2011).

```
import pandas as pd
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
from sklearn.preprocessing import StandardScaler

# Load dataset
data = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/dataset.csv')
X = data.drop('y', axis=1)
y = data['y']

# Encode categorical variables and scale data
X = pd.get_dummies(X)
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.3, random_state=42)

# Initialize models
log_reg = LogisticRegression()
rf = RandomForestClassifier()

# Fit models
log_reg.fit(X_train, y_train)
rf.fit(X_train, y_train)

# Predict on testing set
log_reg_pred = log_reg.predict(X_test)
rf_pred = rf.predict(X_test)
```

```

# Evaluation metrics
print("Logistic Regression Metrics:")
print("Accuracy:", accuracy_score(y_test, log_reg_pred))
print("Precision:", precision_score(y_test, log_reg_pred, pos_label='yes'))
print("Recall:", recall_score(y_test, log_reg_pred, pos_label='yes'))
print("F1-Score:", f1_score(y_test, log_reg_pred, pos_label='yes'))

print("Random Forest Metrics:")
print("Accuracy:", accuracy_score(y_test, rf_pred))
print("Precision:", precision_score(y_test, rf_pred, pos_label='yes'))
print("Recall:", recall_score(y_test, rf_pred, pos_label='yes'))
print("F1-Score:", f1_score(y_test, rf_pred, pos_label='yes'))

# Hyperparameter tuning using GridSearchCV for Random Forest
param_grid_rf = {
    'n_estimators': [100, 200],
    'max_depth': [10, 20, None],
    'min_samples_split': [2, 5]
}
grid_rf = GridSearchCV(estimator=RandomForestClassifier(), param_grid=param_grid_rf, cv=5)
grid_rf.fit(X_train, y_train)

# Print best parameters and best score
print("Best parameters for Random Forest:", grid_rf.best_params_)
print("Best score for Random Forest:", grid_rf.best_score_)

```

Code 8: Two models' evaluation to find which is best.

- **Fine-tune the hyperparameters and document the tuning process.**

Hyperparameter Tuning

1. **Logistic Regression:** Important hyperparameters for Logistic regression are C which is the inverse of the regularization strength and penalty which is the type of regularization e. g. g. l1 or l2).
2. **Random Forest:** Important hyperparameters are n_estimators (number of trees) max_features and max_depth (Bergstra, 2012).

Tuning Method:

To systematically examine various parameter value combinations, utilize Grid Search with Cross-Validation (GridSearchCV). By determining the ideal set of parameters this method is effective in enhancing model performance.

Output:

Logistic Regression Metrics:

Accuracy: 0.8994396933058095

Precision: 0.6386255924170616

Recall: 0.33729662077596995

F1-Score: 0.44144144144144143

Random Forest Metrics:

Accuracy: 0.9046741374225892

Precision: 0.6484907497565725

Recall: 0.4167709637046308

F1-Score: 0.5074285714285715

Best parameters for Random Forest: {'max_depth': None, 'min_samples_split': 5, 'n_estimators': 100}

Best score for Random Forest: 0.9048885880261801

- **Present the final evaluation results in a detailed report.**

Configure the Evaluation Pipeline: To maintain the target variables proportionality divide the data using a stratified split using scikit-learns `train_test_split`.

Compute Initial Metrics: Apply the initial set of hyperparameters to both models' evaluations on the validation set.

Perform Hyperparameter Tuning: Adjust the models using `GridSearchCV` evaluating the impact of different hyperparameters.

Re evaluateMetrics: After fine-tuning calculate the metrics once more to gauge the progress.

CONCLUSION

In this extensive report I have exhibited the powerful powers of machine learning models specifically Random Forest and Logistic Regression in terms of forecasting client conversion for term life insurance at HashSysTech Insurance. The technical viability of these models was demonstrated by this analysis which also emphasized their strategic importance in improving telemarketing campaigns.

Through the careful application of data exploration techniques, the management of outliers and missing values and the utilization of sophisticated visualization tools I have developed a strong framework for comprehending the underlying patterns present in the dataset. Deep insights into the characteristics that most impact consumer decisions have been obtained through the application of Random Forest and Logistic Regression each of which has distinct benefits in terms of interpretability and predictive accuracy.

The models' performances were further improved by adjusting the hyperparameters which optimized them for practical use. Class identification and relevance were approached in a balanced manner as demonstrated by the notable gains in precision recall and F1-scores.

In summary the implementation of these machine learning techniques has advanced HashSysTech Insurances strategic marketing capabilities and established a standard for data-driven decision-making in the insurance industry. Because of my work on this project the company is now able to target potential customers more precisely which has increased conversion rates and maximized campaign ROI. This project demonstrates how incorporating advanced analytics into business plans can have a revolutionary effect and guarantee steady growth and a competitive edge in a sector that is changing quickly.

BIBLIOGRAPHY

- Agresti, A., (2002). *Categorical Data Analysis.. Wiley-Interscience..*
- Bergstra, J. & B. Y., 2012. "Random Search for Hyper-Parameter Optimization.".
- Breiman, L., 2001. "Random Forests", *Machine Learning.. researchgate.com*, Volume 45(1), pp. 5-32..
- Cleveland, W., 1985. "The Elements of Graphing Data,". *Science direct*.
- Forman, G. & S. M., 2010. "Apples-to-apples in cross-validation studies: pitfalls in classifier performance measurement,". *Researchgate.com*.
- Hastie, T. T. R. & F. J., (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction. Springer Series in Statistics. New York, NY.. Springer*.
- Hosmer, D. W. L. S. & S. R. X., 2013. *Applied Logistic Regression. Science direct*.
- James, G. W. D. H. T. & T. R., 2013. "An Introduction to Statistical Learning." . *Science direct*.
- Kohavi, R., 1995. "A Study of Cross-Validation and Bootstrap for Accuracy Estimation and Model Selection". *Science direct*.
- Kuhn, M. & J. K., 2013. *Applied Predictive Modeling.. Springer*.
- McGill, R. T. J. a. L. W., 1978. "Variations of Box Plots.. *Researchgate.com*.
- Powers, D., 2011. "Evaluation: From Precision, Recall and F1 to ROC, Informedness, Markedness & Correlation." . *Science direct*.

Table of Figures:

<i>Figure 1: Distribution of Age.....</i>	<i>11</i>
<i>Figure 2: Age distribution by Outcome.</i>	<i>11</i>
<i>Figure 3: Relationship Between Age and Call Duration by Outcome.</i>	<i>12</i>
<i>Figure 4: Top 19 Logistic Regression Coefficients.....</i>	<i>18</i>
<i>Figure 5: Top 10 Random Forest Importances.</i>	<i>19</i>

Table of Codes:

<i>Code 1: Data Exploration.....</i>	<i>5</i>
<i>Code 2: Categorical Features of Data.....</i>	<i>6</i>
<i>Code 3: Handling Missing Values and Outliers.</i>	<i>9</i>
<i>Code 4: Data Visualization.....</i>	<i>10</i>
<i>Code 5: Logistic Regression and Random forest Algorithms.</i>	<i>14</i>
<i>Code 6: Data Splitting.</i>	<i>16</i>
<i>Code 7: Two models for Interpretation.....</i>	<i>17</i>
<i>Code 8: Two models' evaluation to find which is best.....</i>	<i>21</i>