

PG-DBDA Project

**Predict if a customer is likely to buy travel insurance
using Python**

By:

Swapnil Waghmare

Supervisor

Dr. Shantanu Pathak

IACSD Akurdi Pune

Contents:

1. Component 1

- 1.1. Introduction
- 1.2. Problem Statement
- 1.3. Objective
- 1.4. Available Data
- 1.5. Approach and Algorithms used
 - 1.5.1. GridSearch technique
 - 1.1.1. Logistic Regression
 - 1.1.2. Support Vector Machine
 - 1.1.3. Random Forest

2. Component 2

- 2.1. Basic data exploration
- 2.2. Outlier Treatment
- 2.3. Correlation matrix

3. Component 3

- 3.1. Splitting the data
- 3.2. GridSearch code
- 3.3. Algorithms/models code
 - 3.3.1. Logistic Regression
 - 3.3.2. Support Vector Machine
 - 3.3.3. Random Forest

4. Output

- 4.1. Observations
- 4.2. Conclusion

1. Component 1

1.1. Introduction:

In the contemporary time of data-driven decision-making, companies across various industries are leveraging predictive analytics to gain insights into customer behavior and preferences. One such sector where predictive modeling plays a crucial role is the insurance industry, particularly in the realm of travel insurance. Travel insurance provides coverage for unforeseen events during trips, including medical emergencies, trip cancellations, lost luggage, and more. Understanding the factors that influence customers' decisions to purchase travel insurance is paramount for insurance providers to tailor their marketing strategies effectively and improve their sales performance.

This project aims to utilize advanced machine learning techniques implemented in Python to predict whether a customer is likely to purchase travel insurance based on their information, travel history, and other relevant features. By analyzing historical data on customer purchases and their information, we seek to develop a predictive model that accurately identifies potential buyers of travel insurance.

The significance of this project lies in its potential to assist insurance companies in targeting their marketing efforts more efficiently, thereby maximizing their sales revenue while also enhancing customer satisfaction. Moreover, by identifying the key predictors of travel insurance purchase, insurers can refine their product offerings and pricing strategies to better meet the needs and preferences of their target market.

Throughout this report, we will outline the methodologies employed, including data preprocessing, model development, and evaluation metrics. Additionally, we will discuss which model or algorithms will be best for predicting the customer behavior. The implications of our findings will suggest which predictive model to be used by insurance companies to improve their business operations.

1.2. Problem Statement

In the fiercely competitive landscape of the insurance industry, particularly within the realm of travel insurance, understanding and predicting customer behavior is paramount for insurers to optimize their marketing strategies and improve sales performance. The problem revolves around the need to develop a robust predictive model that accurately forecasts whether a customer is likely to purchase travel insurance based on their demographic information, travel history, and other pertinent features.

The primary challenge lies in effectively leveraging available data to construct a predictive model that exhibits high accuracy, reliability, and interpretability. Additionally, the model must be scalable and adaptable to accommodate varying data sources and evolving customer preferences.

Addressing this problem necessitates a comprehensive approach encompassing data preprocessing, feature engineering, model selection, and evaluation. Moreover, it requires a deep understanding of the underlying factors influencing customer purchasing decisions in the context of travel insurance.

Ultimately, the goal is to equip insurance companies with a predictive tool that empowers them to tailor their marketing efforts, optimize resource allocation, and enhance customer engagement, thereby driving revenue growth and competitive advantage in the dynamic travel insurance market.

1.3. Objectives

The objective of this project is to develop a predictive model using Python that can accurately forecast whether a customer is likely to purchase travel insurance. This entails:

Data Collection and Preprocessing: Gather relevant data on customer demographics, travel history, and other pertinent features. Clean and preprocess the data to ensure accuracy and consistency.

Model Development: Utilize advanced machine learning algorithms, such as logistic regression, support vector machine, random forests, or gradient boosting, to build a predictive model. Train the model on historical data to learn patterns and associations between customer attributes and insurance purchase behavior.

Model Evaluation: Assess the performance of the predictive model using appropriate evaluation metrics, such as accuracy, precision, recall, and F1-score. Validate the model's robustness and generalization ability through cross-validation techniques.

By achieving these objectives, the project aims to provide insurance companies with a valuable predictive tool for optimizing marketing efforts, improving sales performance, and enhancing customer satisfaction in the competitive travel insurance market.

1.4. Available Data:

.csv file containing 9 features and 1967 rows

1df.shape

(1987, 9)

1df.head()

	Age	Employment Type	GraduateOrNot	AnnualIncome	FamilyMembers	ChronicDiseases	FrequentFlyer	EverTravelledAbroad	TravelInsurance
0	31	Government Sector	Yes	400000	6	1	No	No	0
1	31	Private Sector/Self Employed	Yes	1250000	7	0	No	No	0
2	34	Private Sector/Self Employed	Yes	500000	4	1	No	No	1
3	28	Private Sector/Self Employed	Yes	700000	3	1	No	No	0
4	28	Private Sector/Self Employed	Yes	700000	8	1	Yes	No	0

1.5. Approaches and Algorithms used.

1.5.1. GridSearch Technique

Grid Search is a technique used in machine learning to find the best combination of hyperparameters for a model. It works by exhaustively testing all possible combinations of hyperparameters from a predefined grid and selecting the set that yields the highest performance on a validation set. This process involves training the model multiple times with different hyperparameter configurations and evaluating its performance using cross-validation. It's widely used for its simplicity and effectiveness in optimizing model performance.

Since the problem is of classification type we have used classification algorithms as follows:

1.1.1. Logistic Regression

Binary Classification Method: Logistic regression serves as a statistical approach tailored for binary classification tasks, where the target variable comprises only two possible outcomes.

Probability Modeling: Its fundamental function involves modeling the probability of a binary outcome occurrence by fitting a logistic function to the observed data.

Linear Nature: Despite its title, logistic regression operates as a linear model, predicting the log-odds of the target variable belonging to a specific class.

Sigmoid Transformation: Employing the logistic function, also known as the sigmoid function, facilitates the transformation of the linear combination of features into probabilities, constrained within the range of 0 to 1.

Decision Boundary Establishment: Logistic regression delineates a decision boundary within the feature space to segregate classes effectively.

Parameter Estimation Mechanism: The model parameters undergo estimation via maximum likelihood estimation techniques, commonly utilizing optimization algorithms like gradient descent.

Interpretability Aspect: Logistic regression furnishes interpretable coefficients, elucidating the direction and strength of the association between input features and the target variable.

Assumptions Consideration: Key assumptions include assuming linearity between input features and the log-odds of the target variable, alongside the independence of observations and the absence of multicollinearity.

Integration of Regularization: To curb overfitting and enhance generalization, logistic regression can integrate regularization techniques such as L1 (Lasso) and L2 (Ridge) regularization.

Diverse Range of Applications: Logistic regression finds extensive utility across diverse fields such as finance, healthcare, marketing, and social sciences. Its application spectrum spans predicting binary outcomes like customer churn, disease diagnosis, and credit risk assessment.

1.1.2. Support Vector Machine

Supervised Learning Tool: SVM stands as a supervised learning technique suitable for tackling classification and regression tasks.

Linear and Non-linear Classification Ability: It showcases proficiency in both linear and non-linear classification by identifying the best-fit hyperplane or decision boundary for separating classes within the feature space.

Margin Maximization Objective: SVM strives to maximize the margin, representing the distance between the decision boundary and the closest data points of distinct classes, known as support vectors.

Kernel Technique Application: When facing non-linearly separable classes, SVM employs the kernel technique to map data into higher-dimensional spaces, facilitating linear separation.

Variety of Kernel Functions: SVM utilizes various kernel functions like linear, polynomial, radial basis function (RBF), and sigmoid, each suited to different data types and separation needs.

Incorporation of Regularization: SVM integrates a regularization parameter (C) to balance margin maximization and classification error minimization. Larger C values prioritize precise classification, potentially leading to overfitting, while smaller values emphasize a broader margin, possibly resulting in underfitting.

Significance of Support Vectors: Support vectors, being the data points nearest to the decision boundary, play a pivotal role in determining the optimal hyperplane.

Handling of Binary and Multi-class Classification: SVM demonstrates proficiency in both binary and multi-class classification tasks, employing strategies like one-vs-one or one-vs-all.

Resilience to Overfitting and High-dimensional Spaces: SVM exhibits robustness against overfitting, particularly in high-dimensional spaces, and effectively manages datasets with a small number of samples relative to features.

Diverse Range of Applications: SVM finds utility across numerous domains such as text classification, image recognition, bioinformatics, and finance, excelling in tasks necessitating clear decision boundaries and high-dimensional feature spaces.

1.1.3. Random Forest

Ensemble Technique: Random Forest serves as an ensemble learning method applicable to classification and regression tasks.

Construction of Decision Trees: This approach involves the creation of multiple decision trees during training, whose predictions are then combined to enhance accuracy and mitigate overfitting.

Utilization of Bagging: Employing bagging (Bootstrap Aggregating), Random Forest trains each tree on a random subset of the training data with replacement.

Random Feature Subset: At each decision split, only a random subset of features is considered, aiding in reducing variance and ensuring tree independence.

Consensus Voting: In classification, Random Forest aggregates tree predictions through majority voting; for regression, it computes predictions' averages.

Reliable Performance: Known for its accuracy and resistance to overfitting, Random Forest handles noisy data and missing values effectively.

Out-of-Bag (OOB) Error Estimation: The model estimates its performance using out-of-bag samples, streamlining the evaluation process without requiring a separate validation set.

Feature Importance Estimation: Random Forest provides insight into feature importance by assessing their contribution to reducing impurity across all trees in the forest.

Parallelizable: Its parallel nature allows for easy parallelization, as each tree can be trained independently.

Diverse Applications: Random Forest finds applications across various fields like finance, healthcare, and bioinformatics, excelling in tasks such as credit scoring, disease prediction, and gene expression analysis, owing to its ability to handle high-dimensional data and capture intricate feature-target relationships.

2. Component 2

2.1. Basic data exploration

`.shape` - returns the number of rows by the number of columns for dataset.

```
1 df.shape
```

```
(1987, 9)
```

`.head()` - returns the first 5 rows of dataset. This is useful if you want to see some example values for each variable.

```
In [8]: 1 df.head()
```

```
Out[8]:
```

	Age	Employment Type	GraduateOrNot	AnnualIncome	FamilyMembers	ChronicDiseases	FrequentFlyer	EverTravelledAbroad	TravellInsurance
0	31	Government Sector	Yes	400000	6	1	No	No	0
1	31	Private Sector/Self Employed	Yes	1250000	7	0	No	No	0
2	34	Private Sector/Self Employed	Yes	500000	4	1	No	No	1
3	28	Private Sector/Self Employed	Yes	700000	3	1	No	No	0
4	28	Private Sector/Self Employed	Yes	700000	8	1	Yes	No	0

`.columns` - returns the name of all of your columns in the dataset.

```
In [4]: 1 df.columns
```

```
Out[4]: Index(['Unnamed: 0', 'Age', 'Employment Type', 'GraduateOrNot', 'AnnualIncome',  
              'FamilyMembers', 'ChronicDiseases', 'FrequentFlyer',  
              'EverTravelledAbroad', 'TravelInsurance'],  
              dtype='object')
```

.nunique(axis=0) returns the number of unique values for each variable.

```
In [12]: 1 df.nunique()
```

```
Out[12]: Age                11
Employment Type            2
GraduateOrNot              2
AnnualIncome              30
FamilyMembers              8
ChronicDiseases            2
FrequentFlyer              2
EverTravelledAbroad        2
TravelInsurance            2
dtype: int64
```

.describe() summarizes the count, mean, standard deviation, min, and max for numeric variables.

```
In [10]: 1 df.describe()
```

```
Out[10]:
```

	Age	AnnualIncome	FamilyMembers	ChronicDiseases	TravelInsurance
count	1987.000000	1.987000e+03	1987.000000	1987.000000	1987.000000
mean	29.650226	9.327630e+05	4.752894	0.277806	0.357323
std	2.913308	3.768557e+05	1.609650	0.448030	0.479332
min	25.000000	3.000000e+05	2.000000	0.000000	0.000000
25%	28.000000	6.000000e+05	4.000000	0.000000	0.000000
50%	29.000000	9.000000e+05	5.000000	0.000000	0.000000
75%	32.000000	1.250000e+06	6.000000	1.000000	1.000000
max	35.000000	1.800000e+06	9.000000	1.000000	1.000000

.dropna(axis=0) to remove any rows with null values.

.isna().sum() to get the count of na values for each column

```
In [11]: 1 df.isna().sum()
```

```
Out[11]: Age                0
Employment Type            0
GraduateOrNot              0
AnnualIncome              0
FamilyMembers              0
ChronicDiseases            0
FrequentFlyer              0
EverTravelledAbroad        0
TravelInsurance            0
dtype: int64
```

2.2. Outlier Treatment

Identification: The first step in removing outliers involves identifying them within the dataset. Outliers are data points that significantly deviate from the rest of the observations.

Techniques for Detection: Common techniques for detecting outliers include statistical methods such as z-score, interquartile range (IQR), or visualization methods like scatter plots and box plots.

Decision Criteria: Determine the criteria for identifying outliers based on the chosen detection method, such as a z-score threshold or a specified range of values within the IQR.

Data Exclusion: Once outliers are identified, they can be excluded from the dataset. This can be achieved by removing the data points entirely or by applying transformations to reduce their impact.

```
In [21]: 1 def impute_IQR(df, col):
          2     q1,q3 = df[col].quantile([0.25, 0.75])
          3     iqr = q3-q1
          4     min_val_val = q1 - 1.5 * iqr
          5     max_val_val = q3 + 1.5 * iqr
          6     df.loc[df[col] < min_val_val , col] = min_val_val
          7     df.loc[df[col] > max_val_val , col] = max_val_val
          8     return df
```

```
In [22]: 1 for col in numeric_columns:
          2     df = impute_IQR(df,col)
```

2.3. Correlation Matrix:

Definition: A correlation matrix is a tabular representation of the correlation coefficients between variables in a dataset.

Correlation Coefficients: Each cell in the matrix displays the correlation coefficient, which quantifies the strength and direction of the linear relationship between two variables.

Range of Values: Correlation coefficients range from -1 to 1, where -1 indicates a perfect negative linear relationship, 0 denotes no linear relationship, and 1 signifies a perfect positive linear relationship.

Interpretation: Positive values suggest a positive association between variables, while negative values indicate a negative association. A correlation coefficient close to 0 implies little to no linear relationship.

Symmetry: The correlation matrix is symmetric across the diagonal, as the correlation between variable A and variable B is the same as the correlation between variable B and variable A.

Visualization: Correlation matrices are often visualized using heatmaps, with colors representing the strength and direction of correlations.

Diagonal Values: The diagonal of the correlation matrix typically displays correlation coefficients of 1, as each variable is perfectly correlated with itself.

Applications: Correlation matrices are useful for identifying patterns and relationships between variables in multivariate datasets. They aid in feature selection, identifying multicollinearity, and understanding the underlying structure of the data.

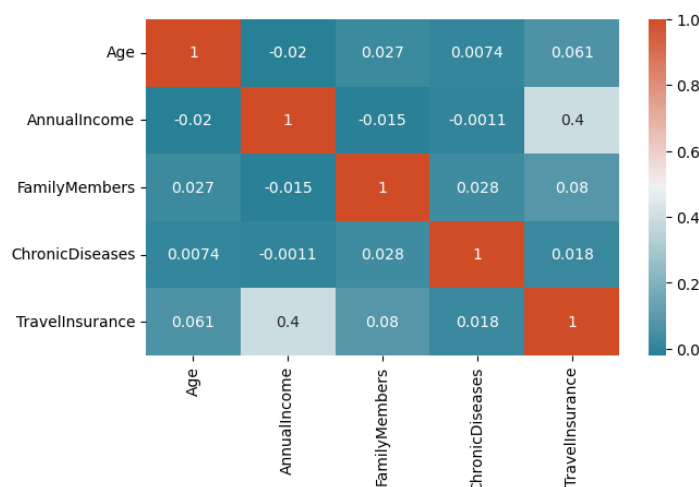
Interpretation Caveats: While correlation matrices provide insights into linear relationships, they may not capture non-linear associations or causality between variables. Additionally, outliers can influence correlation coefficients, necessitating careful interpretation.

calculate correlation matrix

```
In [23]: 1 corr = df.corr()# plot the heatmap
        2 plt.subplots(figsize=(7,4))
        3 sns.heatmap(corr, xticklabels=corr.columns, yticklabels=corr.columns, annot=True, cmap=sns.diverging_palette(220
        4

/var/folders/zv/jmryq7bx7fn1rxzkkmwnf87h0000gn/T/ipykernel_16122/1063527439.py:1: FutureWarning: The default value
of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid c
olumns or specify the value of numeric_only to silence this warning.
corr = df.corr()# plot the heatmap
```

Out[23]: <Axes: >



3. Component 3

3.1. Splitting the data

Split dataset into x and y

- Target variable y is TravelInsurance
- Except TravelInsurance column all remaining are independent features

```
In [24]: 1 x = df.drop('TravelInsurance',axis = 1)
          2 y = df['TravelInsurance']
```

```
In [25]: 1 x.shape,y.shape
```

```
Out[25]: ((1987, 8), (1987,))
```

Splitting the data into x_train, x_test, y_train, y_test

```
In [36]: 1 from sklearn.model_selection import train_test_split
```

```
In [37]: 1 x_train,x_test,y_train,y_test = train_test_split(x,y,test_size = 0.3,random_state = 7,sti
```

```
In [38]: 1 x_train.shape,x_test.shape,y_train.shape,y_test.shape
```

```
Out[38]: ((1390, 12), (597, 12), (1390,), (597,))
```

3.2. GridSearch technique

Here for example we are using GridSearch to fine tune the parameters of Logistic Regression algorithms

We are using Gridsearch technique to fine tune the parameters of the models

```
In [39]: 1 from sklearn.model_selection import GridSearchCV
          2 from sklearn.ensemble import RandomForestClassifier
          3 from sklearn.linear_model import LogisticRegression
          4
          5 param_grid = {
          6     'max_iter' : [x for x in range(50,401,50)],
          7     'C' : [x/10 for x in range(1,11,1)]
          8 }
```

```
In [40]: 1 gscv = GridSearchCV(LogisticRegression(penalty = 'l2',class_weight = 'balanced',
          2                                             random_state = 7, solver = 'saga'),
          3                        param_grid, cv=2, verbose=2)
```

```
In [41]: 1 gscv.fit(x_train,y_train)
```

```
In [42]: 1 gscv.best_score_
```

```
Out[42]: 0.733093525179856
```

```
In [43]: 1 print("The best parameters are: ")
          2
          3 print(gscv.best_params_)
```

```
The best parameters are:
{'C': 0.4, 'max_iter': 50}
```


Similarly for Support Vector Machine and Random Forest also we will use GridSearch technique to find out the best possible combination of parameter values.

3.3. Algorithms/models

3.3.1. Logistic Regression

```
In [44]: 1 lr = LogisticRegression(penalty = 'l2',class_weight = 'balanced',random_state = 7,  
2 solver = 'saga',  
3 max_iter = 50,C = 0.4)  
4
```

```
In [45]: 1 lr.fit(x_train,y_train)
```

```
/Users/swapnilsahebraowaghmare/anaconda3/lib/python3.11/site-packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was reached which means the coef_ did not converge  
warnings.warn(
```

```
Out[45]: LogisticRegression  
LogisticRegression(C=0.4, class_weight='balanced', max_iter=50, random_state=7, solver='saga')
```

```
In [46]: 1 y_pred_logr = lr.predict(x_test)
```

```
In [50]: 1 from sklearn.metrics import classification_report
```

```
In [51]: 1 print(classification_report(y_test, y_pred_logr))
```

	precision	recall	f1-score	support
0	0.78	0.80	0.79	384
1	0.62	0.59	0.61	213
accuracy			0.73	597
macro avg	0.70	0.70	0.70	597
weighted avg	0.72	0.73	0.72	597

3.3.2. Support Vector machine

```
In [54]: 1 svc = SVC(C=0.8, gamma = 0.4, random_state=7)
```

```
In [55]: 1 svc.fit(x_train, y_train)
```

```
Out[55]: SVC(C=0.8, gamma=0.4, random_state=7)
```

```
In [56]: 1 y_pred_svc = svc.predict(x_test)
```

```
In [57]: 1 print(classification_report(y_test, y_pred_svc))
```

	precision	recall	f1-score	support
0	0.80	0.95	0.87	384
1	0.86	0.58	0.69	213
accuracy			0.82	597
macro avg	0.83	0.77	0.78	597
weighted avg	0.82	0.82	0.81	597

3.3.3. Random Forest:

```
In [60]: 1 rfc = RandomForestClassifier(max_depth = 4, max_samples=0.8, n_estimators= 20,  
2 oob_score=True, class_weight='balanced', random_state=7)
```

```
In [61]: 1 rfc.fit(x_train, y_train)
```

```
Out[61]: RandomForestClassifier  
RandomForestClassifier(class_weight='balanced', max_depth=4, max_samples=0.8,  
n_estimators=20, oob_score=True, random_state=7)
```

```
In [62]: 1 y_pred_rfc = rfc.predict(x_test)
```

```
In [63]: 1 print(classification_report(y_test, y_pred_rfc))
```

	precision	recall	f1-score	support
0	0.81	0.96	0.88	384
1	0.89	0.61	0.72	213
accuracy			0.83	597
macro avg	0.85	0.78	0.80	597
weighted avg	0.84	0.83	0.82	597

4. Output:

4.1. Observations:

- 4.1.1. Accuracy of Logistic regression is 73%
- 4.1.2. Accuracy of SVM is 82%
- 4.1.3. Accuracy of random forest classifier is 83%

4.2. Conclusion:

- 4.2.1. Among the three classifiers i.e.

logistic regression

support vector machine

random forest

the maximum accuracy is of random forest model (83%)

hence the best model for classification for given dataset is Random Forest

Reasons: Random Forest is an ensemble technique

Ensemble models offer several benefits that make them advantageous in machine learning applications:

1. Improved Performance: Ensemble models typically yield higher predictive accuracy compared to individual models. By combining multiple base models, ensemble methods harness the collective wisdom of diverse algorithms, mitigating errors and biases present in any single model and leading to more accurate predictions.

2. Robustness: Ensemble models are more robust to noise and outliers in the data. By aggregating predictions from multiple models, ensemble methods can reduce the impact of individual model errors and variations, resulting in more stable and reliable predictions.

3. Generalization: Ensemble models tend to generalize well to unseen data. By leveraging the diversity of base models, ensemble methods can capture a broader range of patterns and relationships in the data, leading to more robust and generalizable models.

4. Reduction of Overfitting: Ensemble models help alleviate overfitting, a common problem in machine learning. By combining predictions from multiple models trained on different subsets of data or using different algorithms, ensemble methods can reduce the risk of overfitting and improve the model's ability to generalize to new data.

5. Flexibility: Ensemble methods are highly flexible and adaptable to different types of data and modeling tasks. They can be constructed using various base models and aggregation techniques, allowing them to be tailored to specific datasets and prediction objectives.

6. Interpretability: Ensemble models can provide insights into the underlying structure of the data. By analyzing the contributions of individual base models to the ensemble prediction, it is possible to gain a better understanding of the factors driving the predictions and identify important features or patterns in the data.

7. Scalability: Ensemble methods can scale to large datasets and complex modeling tasks. By parallelizing the training of base models or using scalable algorithms, ensemble methods can efficiently handle large volumes of data and compute-intensive tasks.

Overall, ensemble models offer a powerful framework for improving predictive performance, robustness, and generalization across a wide range of machine learning applications. Their ability to combine multiple models and leverage diverse viewpoints makes them a valuable tool in the machine learning toolkit.