# Final Project - Travel Insurance Predictions Team 7

October 8, 2021

---

**ADS-505 Final Project** - Travel Insurance Predictions

**Team:** #7

**Team Members:** Jimmy Nguyen, Christopher Robinson, Nima Amin Taghavi

**Date:** 09/20/2021

**Programmin Language:** Python Code

---

Table of Contents

---

# 1   Problem statement

**About the Client**

The client in this data mining project is a tour & travels company that is offering travel insurance package to their customers. This new insurance package also includes COVID-19 coverage for their flights. However, the client wants to know which customers based on their data base history are potential purchasers who may be interested in buying this new insurance package. Previously, the insurance package was offered to some of the customers in 2019 and data was collected from the performance and sales of the package during that period. The sample data given has close to 2000 customers from that period. The client is requesting information on which customer are most likely going to buy the travel insurance given their information such as employment type, income level, etc.

**Business Problem**

The client may use the solutions presented to them for customer-targeted advertising of the new travel insurance package. Also, data visualizations provided will help derive interesting insights about their potential buyers in order to optimize marketing strategies.

**Data Mining Problem**

- A supervised classification task, where the outcome variable of interest is *TravelInsurance* that indicates whether the customer will buy the travel insurance. Performance metrics should take in consideration the positive class of buyers/purchasers.
- Find out interesting patterns and trends for better customer segmentations through data exploration and visualizations.
- An unsupervised task, where the goal is to cluster customers.

---

# 2   Packages

**Python code:**

```
[1]: %%javascript
IPython.OutputArea.prototype._should_scroll = function(lines) {
    return false;
}
```

<IPython.core.display.Javascript object>

```python
[2]: from pathlib import Path
     import numpy as np
     import pandas as pd
     import matplotlib.pylab as plt
     import seaborn as sns
     from sklearn.linear_model import LinearRegression, LogisticRegression,
      ↪LogisticRegressionCV
     from sklearn.model_selection import train_test_split
     import statsmodels.api as sm
     import scikitplot as skplt
     from mord import LogisticIT
     from sklearn import preprocessing
     from sklearn.metrics import accuracy_score
     from sklearn.neighbors import KNeighborsClassifier
     from sklearn.tree import DecisionTreeRegressor, DecisionTreeClassifier
     from sklearn.neural_network import MLPClassifier
     from sklearn.model_selection import cross_val_score, GridSearchCV
     from dmba import regressionSummary, stepwise_selection, plotDecisionTree
     from dmba import regressionSummary, stepwise_selection
     from dmba import forward_selection, backward_elimination, exhaustive_search
     from dmba import classificationSummary, gainsChart, liftChart
     from dmba.metric import AIC_score
     from tabulate import tabulate
     import matplotlib.patches as mpatches
     import warnings
     sns.set_theme()
     plt.rcParams['figure.figsize'] = [11, 9]


     warnings.filterwarnings('ignore')
```

## 3 Data Set

**Data Dictionary**

1. **Age** - Age Of The Customer

2. **Employment Type** - The Sector In Which Customer Is Employed

3. **GraduateOrNot** - Whether The Customer Is College Graduate Or Not

4. **AnnualIncome** - The Yearly Income Of The Customer In Indian Rupees[Rounded To Nearest 50 Thousand Rupees]

5. **FamilyMembers** - Number Of Members In Customer's Family

6. **ChronicDisease** - Whether The Customer Suffers From Any Major Disease Or Conditions Like Diabetes/High BP or Asthama,etc.

7. **FrequentFlyer** - Derived Data Based On Customer's History Of Booking Air Tickets On Atleast 4 Different Instances In The Last 2 Years[2017-2019].

8. **EverTravelledAbroad** - Has The Customer Ever Travelled To A Foreign Country[Not Necessarily Using The Company's Services]

9. **TravelInsurance** - Did The Customer Buy Travel Insurance Package During Introductory Offering Held In The Year 2019.

**Python code:**

```python
[3]: # Load data set
df = pd.read_csv("../../Data/TravelInsurancePrediction.csv")

# First few rows of data set
df.head(3)
```

```
[3]:    Age            Employment Type GraduateOrNot  AnnualIncome  \
    0   31           Government Sector           Yes        400000
    1   31  Private Sector/Self Employed         Yes       1250000
    2   34  Private Sector/Self Employed         Yes        500000

       FamilyMembers  ChronicDiseases FrequentFlyer EverTravelledAbroad  \
    0              6                1            No                  No
    1              7                0            No                  No
    2              4                1            No                  No

       TravelInsurance
    0                0
    1                0
    2                1
```

---

# 4 Exploratory Data Analysis (EDA)

- Graphical and non-graphical representations of relationships between the response variable and predictor variables

---

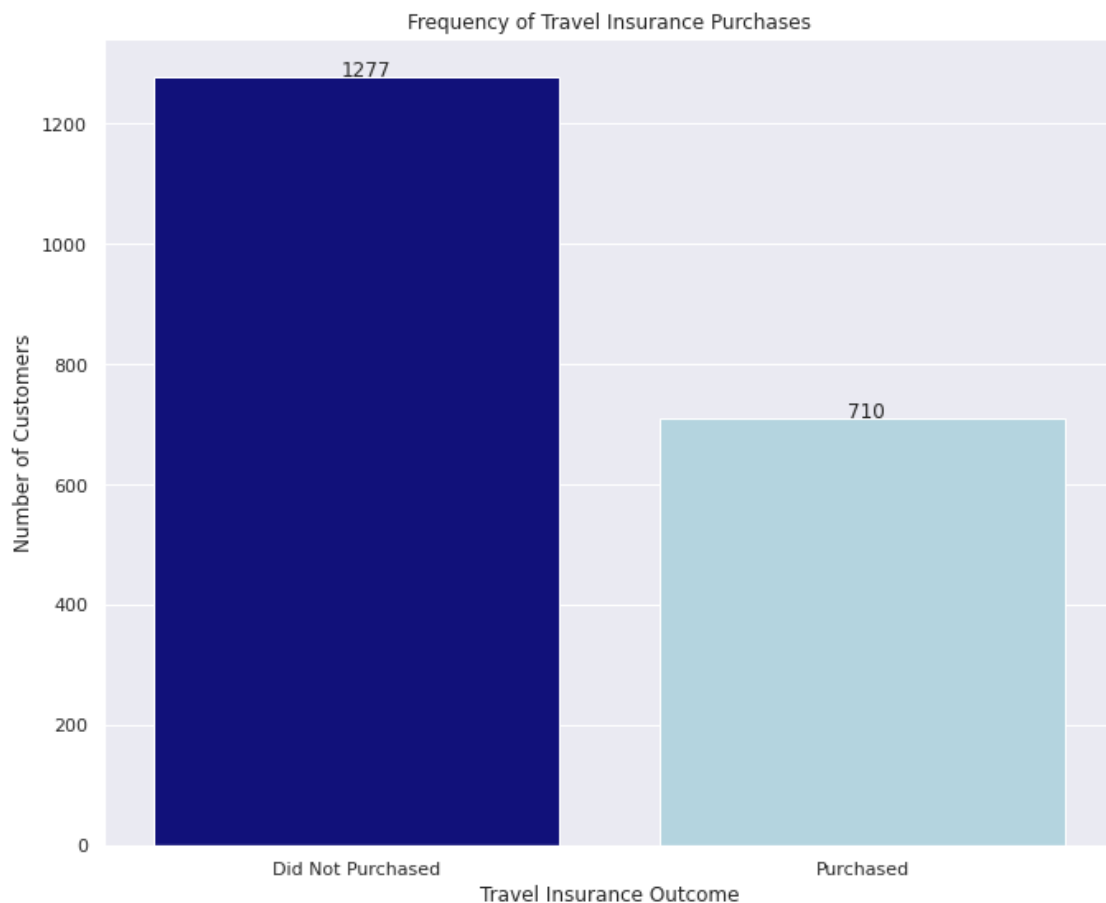## 4.1 Exploring Response Variable - TravelInsurance

- Examine the frequency of travel insurance purchases
- 0 = No
- 1 = Yes

**Python code:**

```
[4]: # Graph count plot of Non-Purchases vs. Purchases
     ax = sns.countplot(data=df, x="TravelInsurance", palette = ["darkblue",␣
      ↪"lightblue"])

     # Add count labels
     for p, label in zip(ax.patches, df['TravelInsurance'].value_counts()):
         ax.annotate(label, (p.get_x()+0.37, p.get_height()+0.15))

     # Graph Properties
     plt.title("Frequency of Travel Insurance Purchases")
     plt.xticks([0, 1], ["Did Not Purchased", "Purchased"])
     plt.xlabel("Travel Insurance Outcome")
     plt.ylabel('Number of Customers')
     plt.show()
```



**Summary**

- There is a higher number of customers who flew without buying travel insurance in this sample data set.

## 4.2 Examing customers' Age

- Age distributions
- Age with Target Variable Overlaid
- Normalized Histogram with Target Overlaid on Age
- Age Group Comparisons (20s vs. 30s)
- Percentage of Purchases between Age groups (20s vs. 30s)

### 4.2.1 Age Distribution

**Python code:**

```python
# Get a range of customer ages
age = pd.DataFrame({'Age': df['Age'].value_counts().sort_index()})
age
```

```
     Age
25   146
26   148
27   131
28   506
29   192
30    65
31   234
32    91
33   148
34   266
35    60
```

```python
# Histogram of Age and set the range of bins from 25-35
bins = np.arange(25, 36)
ax = df['Age'].plot.hist(bins=bins)

# add labels
for p, label in zip(ax.patches, df['Age'].value_counts().sort_index()):
    ax.annotate(label, (p.get_x() + 0.37, p.get_height() + 0.15))

# title and axis
plt.title("Histogram of Customer Ages")
plt.xlabel("Age Range of Customers")
plt.ylabel("Number of Customers")
plt.show()
```

Histogram of Customer Ages

**Summary**

- There are 506 customers who are 28 years-old which is visualized as the most in this data set
- While customers who are 30 years-old are the least in this data set.

### 4.2.2 Age with Target Variable Overlaid

**Python code:**

```
[7]: # Set up plot with response overlaid
     n, bins, patches = plt.hist(
         [
             df[df['TravelInsurance'] == 1]['Age'],
             df[df['TravelInsurance'] == 0]['Age']
         ],
         bins=10,
         stacked=True,
         color=["lightblue", "darkblue"],
     )
```

7

```python
# title and axis
labels = ["Travel Insurance Purchased", "Non-purchase"]
plt.legend(labels)
plt.title("Histogram of Customer Ages Travel Insurance Overlaid")
plt.xlabel("Age Range of Customers")
plt.ylabel('Number of Customers')
plt.show()
```



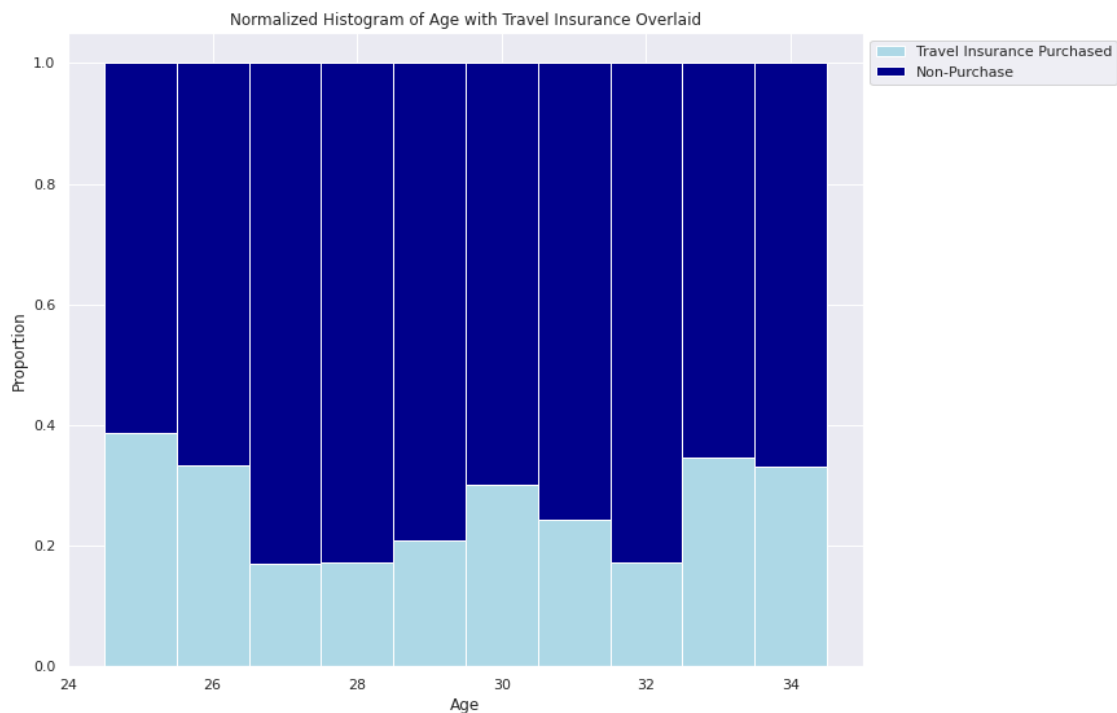Histogram of Customer Ages Travel Insurance Overlaid

**Summary**

- It is difficult to compare between age groups with target variable overlaid
- Therefore, it is better to focus on one class from the target variable and analyze age in a normalized histogram.
- The following is visualized below.

### 4.2.3  Normalized Histogram with Target Variable Overlaid on Age

**Python code:**

```
[8]: # Create normalized histogram for age groups by target overlay
     n_table = np.column_stack((n[0], n[1]))  # stack the tables
     n_norm = n_table / n_table.sum(axis=1)[:,
                                             None]  # create normalized tables by sum
     ourbins = np.column_stack((bins[0:10], bins[1:11]))  # create table bins

     p1 = plt.bar(x=ourbins[:, 0],
                  height=n_norm[:, 0],
                  width=ourbins[:, 1] - ourbins[:, 0],color = "lightblue")  # first␣
      ↪bar chart
     p2 = plt.bar(
         x=ourbins[:, 0],
         height=n_norm[:, 1],
         width=ourbins[:, 1] - ourbins[:, 0],  # second bar chart
         bottom=n_norm[:, 0], color = "darkblue")
     # Annotate legend, title with x and y labels
     plt.legend(['Travel Insurance Purchased', 'Non-Purchase'],
                bbox_to_anchor=(1, 1))
     plt.title('Normalized Histogram of Age with Travel Insurance Overlaid')
     plt.xlabel('Age')
     plt.ylabel('Proportion')
     plt.show()
```



**Summary**

9

- Insights for this graph show that it may be better to compare age classified into 2 groups instead such as customers who are in their twenties (20s) vs customers in their thirties (30s).

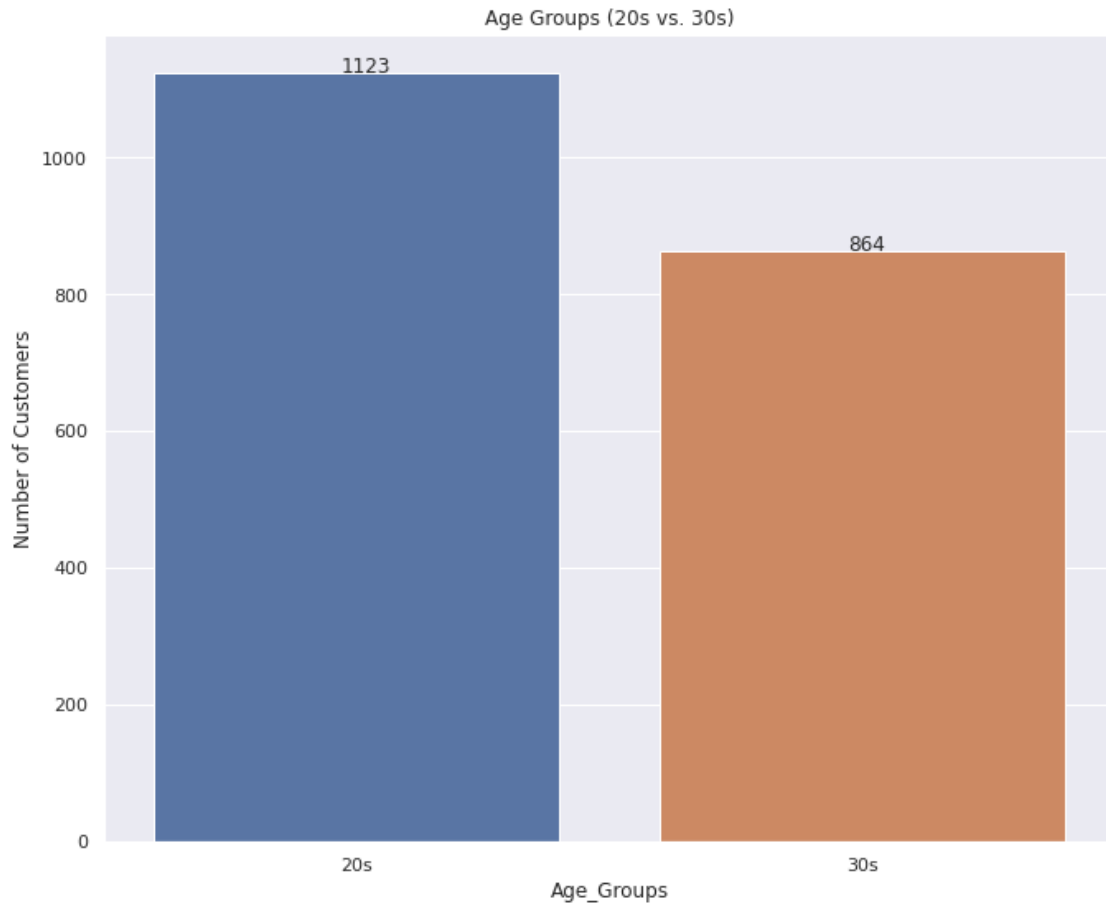### 4.2.4 Age Groups Comparison (20s vs. 30s)

**Python code:**

```
[9]: # Create function to categorize age groups
     def age_groups(x):
         '''
         x: This is a value from df['Age']
         returns each as a new categorical value of 20s or 30s
         '''
         if x < 30:
             return '20s'
         else:
             return '30s'


     # Apply age_groups function on each value
     age_groups = pd.DataFrame(
         {'Age_Groups': df['Age'].apply(lambda x: age_groups(x)),
          'AnnualIncome':df['AnnualIncome'],
          'TravelInsurance':df['TravelInsurance']})

     # Graph count plot of age groups (20s vs. 30s)
     ax = sns.countplot(data=age_groups, x="Age_Groups", order=['20s', '30s'])

     # add labels
     for p, label in zip(ax.patches, age_groups['Age_Groups'].value_counts()):
         ax.annotate(label, (p.get_x()+0.37, p.get_height()+0.15))
     plt.title("Age Groups (20s vs. 30s)")
     plt.ylabel('Number of Customers')
     plt.show()
```

Age Groups (20s vs. 30s)

**Summary**

- Before visualizing the target variable overlaid, we can see here that after binning age into two groups, there are more customers who are in their 20s than customers in their 30s in this data set.

### 4.2.5 Percentage of Travel Insurance Purchases on Various Features

**Python code:**

```
[10]: def make_stacked_barcharts(df, x):
          '''
          Takes in a data frame 'df' and a column 'x'
          and returns a stacked bar chart of the column
          with the percentage of purchases overlaid

          '''

          # Calculate total counts from both groups
          total = df.groupby(x)['TravelInsurance'].count().reset_index()
```

```python
# Calculate total counts from only purchases
purchase = df[df.TravelInsurance == 1].groupby(
    x)['TravelInsurance'].count().reset_index()

# get percentages for purchases
purchase['TravelInsurance'] = [
    i / j * 100
    for i, j in zip(purchase['TravelInsurance'], total['TravelInsurance'])
]

# get percentages
total['TravelInsurance'] = [
    i / j * 100
    for i, j in zip(total['TravelInsurance'], total['TravelInsurance'])
]

# bar chart 1 -> top bars (group of 'TravelInsurance=0')
bar1 = sns.barplot(x, y="TravelInsurance", data=total, color='darkblue')

# bar chart 2 -> bottom bars (group of 'TravelInsurance=1')
bar2 = sns.barplot(x,
                   y="TravelInsurance",
                   data=purchase,
                   color='lightblue')

# add legend
top_bar = mpatches.Patch(color='darkblue', label='Non-Purchase')
bottom_bar = mpatches.Patch(color='lightblue',
                            label='Travel Insurance Purchased')
plt.legend(handles=[top_bar, bottom_bar],
           bbox_to_anchor=(1.05, 1),
           loc=2,
           borderaxespad=0.)
# Aesthetics
plt.title("Percentage of Travel Insurance Purchases by " + x)
plt.xlabel(x)
plt.ylabel("% of Purchases")

# Change ticks on x-axis for ChronicDiseases column
if x == "ChronicDiseases":
    chronicdiease = [0, 1]
    labels = ['No', 'Yes']
    plt.xticks(chronicdiease, labels)

# show the graph
plt.show()
```

```
[11]: # Call stacked bar chart function
      make_stacked_barcharts(age_groups, 'Age_Groups')

      # Make stacked bar charts on various columns
      plot_columns = [
          'Employment Type', 'GraduateOrNot', 'ChronicDiseases', 'FrequentFlyer',
          'EverTravelledAbroad'
      ]

      # Plot each column
      for i in plot_columns:
          make_stacked_barcharts(df, i)
```



Percentage of Travel Insurance Purchases by Age_Groups

Percentage of Travel Insurance Purchases by Employment Type



Percentage of Travel Insurance Purchases by GraduateOrNot

Percentage of Travel Insurance Purchases by ChronicDiseases



Percentage of Travel Insurance Purchases by FrequentFlyer

Percentage of Travel Insurance Purchases by EverTravelledAbroad

**Summary**

- There is a higher proportion of customers in their 30s that purchased travel insurance.
- There is a higher proportion of customers who works in a private sector or is self-employed that purchased travel insurance.
- There is no significant difference in proportion between customers who are a college graduate or not that purchased travel insurance.
- This also applies to customers with or without chronic diseases that purchased travel insurance.
- However, there is a higher proportion of customers who are frequent flyers and/or have traveled abroad that purchased travel insurance.

## 4.3   Monthly Income (Indian Rupee) Range

```
[12]: # Counts by Monthly Rupee
n = np.array([[20.,  63.,  60.,  67.,  41.,  67.,  74., 191.,  84.,  43.],
        [208., 239., 223., 273., 153., 268., 274., 204.,  87.,  58.]])

# Bins for Montly Rupee
bins= np.array([ 25000.,  37500.,  50000.,  62500.,  75000.,  87500., 100000.,
        112500., 125000., 137500., 150000.])

            # Create normalized histogram for groups by target overlay
n_table = np.column_stack((n[0], n[1]))  # stack the tables
n_norm = n_table / n_table.sum(
```

16

```
        axis=1)[:, None]  # create normalized tables by sum
ourbins = np.column_stack((bins[0:10], bins[1:11]))  # create table bins

p1 = plt.bar(x=ourbins[:, 0],
             height=n_norm[:, 0],
             width=ourbins[:, 1] - ourbins[:, 0],
             color="lightblue")  # first bar chart
p2 = plt.bar(
        x=ourbins[:, 0],
        height=n_norm[:, 1],
        width=ourbins[:, 1] - ourbins[:, 0],  # second bar chart
        bottom=n_norm[:, 0],
        color="darkblue")
    # Annotate legend, title with x and y labels
plt.legend(['Travel Insurance Purchased', 'Non-Purchase'],
             bbox_to_anchor=(1, 1))
plt.title('Normalized Histogram of Montly Income with Travel Insurance␣
 ↪Overlaid')
plt.xlabel('$ Monthly Income (Indian Rupee)')
plt.ylabel('Proportion')
plt.show()
```



Normalized Histogram of Montly Income with Travel Insurance Overlaid

## 4.4 Side-by-side Box-plots between Annual Income and Different Attributes

**Python code:**

```python
[13]: def make_boxplots(df, x):
          '''
          Takes in 'x' as a column from data frame 'df'
          and returns a side by side box-plot of
          x on the x-axis and AnnualIncome on the y-axis
          seperated by different colors noted by TravelInsurance


          '''

          # Palatte to color the target variable
          palatte = {0: "darkblue", 1: "lightblue"}

          # Change x-axis labels if age_groups or GraduatedOrNot
          order = None
          if x == "Age_Groups":
              order = ["20s", "30s"]
          if x == "GraduateOrNot":
              order = ["No", "Yes"]

          #Convert AnnualIncome to Monthly
          df['Monthly_Income'] = round(df['AnnualIncome']/12,2)

          # Boxplot
          sns.boxplot(x=x,
                      y="Monthly_Income",
                      hue="TravelInsurance",
                      data=df,
                      order=order,
                      palette=palatte)

          # Legend properties
          top_bar = mpatches.Patch(color='darkblue', label='Non-purchase')
          bottom_bar = mpatches.Patch(color='lightblue',
                                      label='Travel Insurance Purchased')
          plt.legend(handles=[top_bar, bottom_bar],
                     bbox_to_anchor=(1.05, 1),
                     loc=2,
                     borderaxespad=0.)

          # Graph Properties
          plt.title(x + " vs. Montly Income with Travel Insurance Overlaid ")
          plt.xlabel(x)
          plt.ylabel("$ Monthly Income (Indian Rupee)")
```

```
      # Change ticks on x-axis for ChronicDiseases column
      if x == "ChronicDiseases":
          chronicdiease = [0, 1]
          labels = ['No', 'Yes']
          plt.xticks(chronicdiease, labels)

      # show the graph
      plt.show()
```

[14]:
```
# call stacked bar chart function
make_boxplots(age_groups, 'Age_Groups')

# call boxplot function on various columns
plot_columns = ['Employment Type', 'GraduateOrNot','ChronicDiseases',␣
 ↪'FrequentFlyer',
        'EverTravelledAbroad']

# plot each column
for i in plot_columns:
    # boxplot function
    make_boxplots(df, i)
```

Employment Type vs. Montly Income with Travel Insurance Overlaid



GraduateOrNot vs. Montly Income with Travel Insurance Overlaid

ChronicDiseases vs. Montly Income with Travel Insurance Overlaid



FrequentFlyer vs. Montly Income with Travel Insurance Overlaid

EverTravelledAbroad vs. Montly Income with Travel Insurance Overlaid

# 5 Feature Engineering and Pre-Processing

- Invent new columns
- Handle missing values, outliers, correlated features, etc.

**Python code:**

- Bin age into age groups

```
[15]:  # Create function to categorize age groups
       def age_groups(x):
           '''
           x: This is a value from df['Age']
           returns each as a new categorical value of 20s or 30s
           '''
           if x < 30:
               return '20s'
           else:
               return '30s'

       # group age into 20s and 30s
       df['Age'] = df['Age'].apply(lambda x: age_groups(x))
       df = df.rename(columns = {'Age':'Age_Groups'})
       df.head()
```

```
[15]:    Age_Groups              Employment Type GraduateOrNot  AnnualIncome  \
      0       30s              Government Sector           Yes        400000
      1       30s  Private Sector/Self Employed           Yes       1250000
      2       30s  Private Sector/Self Employed           Yes        500000
      3       20s  Private Sector/Self Employed           Yes        700000
      4       20s  Private Sector/Self Employed           Yes        700000

         FamilyMembers  ChronicDiseases FrequentFlyer EverTravelledAbroad  \
      0              6                1            No                  No
      1              7                0            No                  No
      2              4                1            No                  No
      3              3                1            No                  No
      4              8                1           Yes                  No

         TravelInsurance  Monthly_Income
      0                0        33333.33
      1                0       104166.67
      2                1        41666.67
      3                0        58333.33
      4                0        58333.33
```

- Create poor- lower - middle - high income classes

```python
[16]: # Create function to categorize wealth groups
      def wealth_groups(x):
          '''
          x: This is a value from df['Monthly_Income']
          returns each as a new categorical value of 20s or 30s

          The range for categorizing income class

          - Poor Class:  2500- 6500 per month

          - Lower Class:  6500- 15000 per month

          - Middle Class:  15000- 100000 per month

          - Upper Class:  100000- 350000 per month

          '''

          if x <= 6500:
              return 'poor'
          elif x <= 15000:
              return 'lower'
          elif x <= 100000:
              return 'middle'
```

```
    else:
        return 'upper'

# group Monthly Income into wealth groups
df['AnnualIncome'] = df['Monthly_Income'].apply(lambda x: wealth_groups(x))
df = df.drop('Monthly_Income', axis = 1)
df = df.rename(columns = {'AnnualIncome':'Income_Class'})
df.head()
```

[16]:
```
   Age_Groups              Employment Type GraduateOrNot Income_Class  \
0        30s              Government Sector           Yes       middle
1        30s  Private Sector/Self Employed           Yes        upper
2        30s  Private Sector/Self Employed           Yes       middle
3        20s  Private Sector/Self Employed           Yes       middle
4        20s  Private Sector/Self Employed           Yes       middle

   FamilyMembers  ChronicDiseases FrequentFlyer EverTravelledAbroad  \
0              6                1            No                  No
1              7                0            No                  No
2              4                1            No                  No
3              3                1            No                  No
4              8                1           Yes                  No

   TravelInsurance
0                0
1                0
2                1
3                0
4                0
```

- Convert Family Members to household size categories

[17]:
```
# Create function to categorize age groups
def household_groups(x):
    '''
    x: This is a value from df['FamilyMembers']
    returns each as a new categorical value of 20s or 30s

    The range for categorizing income class

    - 1 = single

    - 2-4 = small

    - 5-10 = medium
```

```
        - >10 = large

        '''

        if x == 1:
            return 'single'
        elif x <= 4:
            return 'small'
        elif x <= 10:
            return 'medium'
        else:
            return 'large'


    # group Family Members into household groups
    df['FamilyMembers'] = df['FamilyMembers'].apply(lambda x: household_groups(x))
    df = df.rename(columns = {'FamilyMembers':'Household_Size'})
    df.head()
```

[17]:

| | Age_Groups | Employment Type | GraduateOrNot | Income_Class |
|---|---|---|---|---|
| 0 | 30s | Government Sector | Yes | middle |
| 1 | 30s | Private Sector/Self Employed | Yes | upper |
| 2 | 30s | Private Sector/Self Employed | Yes | middle |
| 3 | 20s | Private Sector/Self Employed | Yes | middle |
| 4 | 20s | Private Sector/Self Employed | Yes | middle |

| | Household_Size | ChronicDiseases | FrequentFlyer | EverTravelledAbroad |
|---|---|---|---|---|
| 0 | medium | 1 | No | No |
| 1 | medium | 0 | No | No |
| 2 | small | 1 | No | No |
| 3 | small | 1 | No | No |
| 4 | medium | 1 | Yes | No |

| | TravelInsurance |
|---|---|
| 0 | 0 |
| 1 | 0 |
| 2 | 1 |
| 3 | 0 |
| 4 | 0 |

- Convert Binary Categorical Variable values to 0/1

[18]:
```
# Convert Frequent Flyer to 0/1 binary values
df['FrequentFlyer'] = np.where((df['FrequentFlyer'] == 'No'), 0, 1)
```

```
# Convert Ever Traveled Abroad to 0/1 binary values
df['EverTravelledAbroad'] = np.where((df['EverTravelledAbroad'] == 'No'), 0, 1)


# Convert GraduateOrNot to 0/1 binary values
df['GraduateOrNot'] = np.where((df['GraduateOrNot'] == 'No'),0,1)

# first few rows
df.head()
```

[18]: 
| | Age_Groups | Employment Type | GraduateOrNot | Income_Class |
|---|---|---|---|---|
| 0 | 30s | Government Sector | 1 | middle |
| 1 | 30s | Private Sector/Self Employed | 1 | upper |
| 2 | 30s | Private Sector/Self Employed | 1 | middle |
| 3 | 20s | Private Sector/Self Employed | 1 | middle |
| 4 | 20s | Private Sector/Self Employed | 1 | middle |

| | Household_Size | ChronicDiseases | FrequentFlyer | EverTravelledAbroad |
|---|---|---|---|---|
| 0 | medium | 1 | 0 | 0 |
| 1 | medium | 0 | 0 | 0 |
| 2 | small | 1 | 0 | 0 |
| 3 | small | 1 | 0 | 0 |
| 4 | medium | 1 | 1 | 0 |

| | TravelInsurance |
|---|---|
| 0 | 0 |
| 1 | 0 |
| 2 | 1 |
| 3 | 0 |
| 4 | 0 |

- One hot encode categorical variables to dummy variables

[19]: 
```
# one-hot encoding the categorical variables
df = pd.get_dummies(df)
df.head()
```

[19]: 
| | GraduateOrNot | ChronicDiseases | FrequentFlyer | EverTravelledAbroad |
|---|---|---|---|---|
| 0 | 1 | 1 | 0 | 0 |
| 1 | 1 | 0 | 0 | 0 |
| 2 | 1 | 1 | 0 | 0 |
| 3 | 1 | 1 | 0 | 0 |
| 4 | 1 | 1 | 1 | 0 |

| | TravelInsurance | Age_Groups_20s | Age_Groups_30s |
|---|---|---|---|
| 0 | 0 | 0 | 1 |
| 1 | 0 | 0 | 1 |

```
2                      1                  0                  1
3                      0                  1                  0
4                      0                  1                  0

    Employment Type_Government Sector  \
0                                    1
1                                    0
2                                    0
3                                    0
4                                    0

    Employment Type_Private Sector/Self Employed  Income_Class_middle  \
0                                              0                    1
1                                              1                    0
2                                              1                    1
3                                              1                    1
4                                              1                    1

    Income_Class_upper  Household_Size_medium  Household_Size_small
0                    0                      1                     0
1                    1                      1                     0
2                    0                      0                     1
3                    0                      0                     1
4                    0                      1                     0
```

---

# 6  Data splitting

- Training, validation, and test sets
- Since there does not exist a class imbalance problem, we split the data set into 75% training and 25% validation.

**Python code:**

```python
[20]: # Response Variable
      outcome = 'TravelInsurance'
      y = df[outcome]

      # features - Do not use Target_B or Target_D
      predictors = [
          'GraduateOrNot', 'ChronicDiseases', 'FrequentFlyer', 'EverTravelledAbroad',
          'Age_Groups_20s', 'Age_Groups_30s', 'Employment Type_Government Sector',
          'Employment Type_Private Sector/Self Employed', 'Income_Class_middle',
          'Income_Class_upper', 'Household_Size_medium', 'Household_Size_small'
      ]
      X = df[predictors]
```

```python
# Set seed to 1 and split on 30% validation
train_X, valid_X, train_y, valid_y = train_test_split(X,
                                                      y,
                                                      test_size=0.25,
                                                      random_state=1)


# Check dimensions
train_X.shape, valid_X.shape
```

[20]: ((1490, 12), (497, 12))

---

# 7 Model building strategies

- Describing main research questions and appropriate analytics methods

**Python code:**

[ ]:

---

# 8 Model performance and hyper-parameter tuning

- Model tuning, comparison, and evaluations

## 8.1 Decision Tree

**Python code:**

[21]:
```python
# user grid search to find optimized tree
param_grid = {
    'max_depth': [5, 10, 15, 20, 25],
    'min_impurity_decrease': [0, 0.001, 0.005, 0.01],
    'min_samples_split': [10, 20, 30, 40, 50],
}

# Run Exhaustive Search
gridSearch = GridSearchCV(DecisionTreeClassifier(random_state=1), param_grid,␣
  ↪cv=5, n_jobs=-1)
gridSearch.fit(X=train_X, y=train_y)

# Initial Parameters
print('Initial Score: ', gridSearch.best_score_)
print('Initial parameters: ', gridSearch.best_params_)
```

```python
# Improving the parameters
param_grid = {
    'max_depth': [3, 4, 5, 6, 7, 8],
    'min_impurity_decrease': [0, 0.001, 0.002, 0.003, 0.005, 0.006, 0.007, 0.
 →008],
    'min_samples_split': [6,7,8,9,10,11,12]
}

# Run Exhaustive Search with fine-tuned parameters
gridSearch = GridSearchCV(DecisionTreeClassifier(random_state=1), param_grid,
 →cv=5, n_jobs=-1)
gridSearch.fit(train_X, train_y)

# Final parameters
print('Improved Score: ', gridSearch.best_score_)
print('Improved parameters: ', gridSearch.best_params_)

# Final Decision Tree
tree_model = gridSearch.best_estimator_

# Fit to Training Data
tree_model.fit(train_X, train_y)
```

```
Initial Score:  0.7785234899328859
Initial parameters:  {'max_depth': 5, 'min_impurity_decrease': 0.005,
'min_samples_split': 10}
Improved Score:  0.7785234899328859
Improved parameters:  {'max_depth': 3, 'min_impurity_decrease': 0,
'min_samples_split': 6}
```
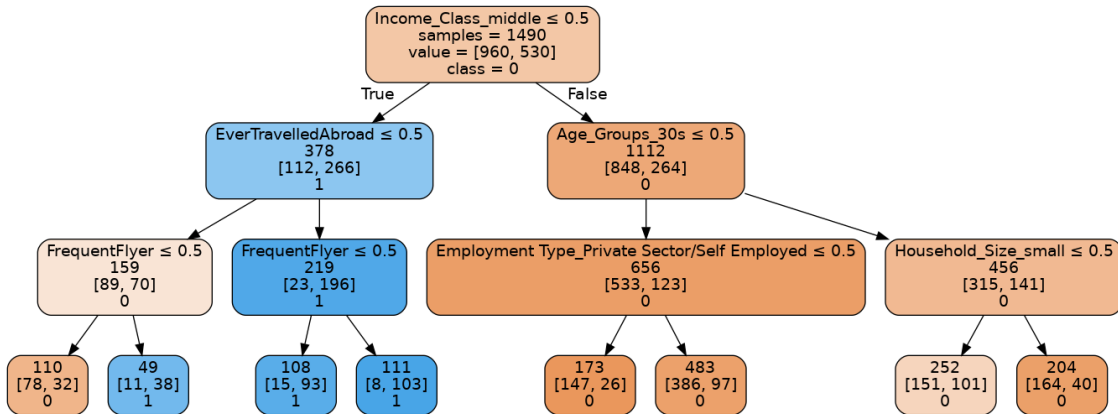
```
[21]: DecisionTreeClassifier(max_depth=3, min_impurity_decrease=0,
                             min_samples_split=6, random_state=1)
```

```python
[22]: # Plot Decision tree
plotDecisionTree(tree_model,
                 feature_names=train_X.columns,
                 class_names=tree_model.classes_)
```

```
[22]:
```

## 8.2 Logistic Regression

```
[23]: param_grid = {
          'penalty': ['l1', 'l2'],
          'C': np.logspace(-4, 4, 20),
          'solver': ['liblinear', 'saga']
      }

      # Create grid search object
      gridSearch = GridSearchCV(LogisticRegression(random_state=1, max_iter=5000),
                              param_grid=param_grid,
                              cv=5,
                              n_jobs=-1)

      # Fit to training data
      gridSearch.fit(train_X, train_y)

      # Final parameters
      print('Final Score: ', gridSearch.best_score_)
      print('Final parameters: ', gridSearch.best_params_)

      # Final logistic regression model
      logit_model = gridSearch.best_estimator_
```

```
Final Score:  0.7657718120805369
Final parameters:  {'C': 0.004832930238571752, 'penalty': 'l2', 'solver':
'liblinear'}
```

```
[24]: # Train logistic regression model to find the best predictors
      def train_model(variables):
          if len(variables) == 0:
              return None
```

```python
    model = LogisticRegressionCV(penalty="l2",
                                 solver='liblinear',
                                 cv=5,
                                 random_state=1,
                                 max_iter=5000)
    return model.fit(train_X[variables], train_y)


# Return the accuracy score in the validation set over each predictor
def score_model(model, variables):
    if len(variables) == 0:
        return 0
    logit_reg_valid = model.predict(valid_X[variables])
    return -accuracy_score(valid_y,
                           [1 if p > 0.5 else 0 for p in logit_reg_valid])

# Use step-wise regression to select the best subset of features
logit_model, best_variables = stepwise_selection(predictors,
                                                 train_model,
                                                 score_model,
                                                 direction='forward',
                                                 verbose=True)
print(best_variables)

# Use the previous columns
columns = best_variables

# Fit to Training Data with previous columns
logit_model.fit(train_X[columns], train_y)
```

Variables: GraduateOrNot, ChronicDiseases, FrequentFlyer, EverTravelledAbroad,
Age_Groups_20s, Age_Groups_30s, Employment Type_Government Sector, Employment
Type_Private Sector/Self Employed, Income_Class_middle, Income_Class_upper,
Household_Size_medium, Household_Size_small
Start: score=0.00, constant
Step: score=-0.77, add EverTravelledAbroad
Step: score=-0.79, add Income_Class_middle
Step: score=-0.80, add FrequentFlyer
Step: score=-0.80, unchanged None
['EverTravelledAbroad', 'Income_Class_middle', 'FrequentFlyer']

[24]: LogisticRegressionCV(cv=5, max_iter=5000, random_state=1, solver='liblinear')

## 8.3 Multi-Layered Neural Network

```
[25]: # user grid search to find optimized hidden layers
      param_grid = {
          'hidden_layer_sizes': [(1), (2), (3), (4), (5)],
      }

      # Run Exhaustive search for neural networks hyper-parameters
      gridSearch = GridSearchCV(MLPClassifier(activation = 'logistic',
                                          solver='lbfgs', random_state=1,
       ↪max_iter=5000),
                              param_grid,
                              cv=5,
                              n_jobs=-1,
                              return_train_score=True)
      # Fit to training set
      gridSearch.fit(train_X, train_y)

      # Initial Scores and hyper-parameters
      print('Initial score: ', gridSearch.best_score_)
      print('Initial parameters: ', gridSearch.best_params_)

      # Look at Initial Scores with averages
      display=['param_hidden_layer_sizes', 'mean_test_score', 'std_test_score']
      pd.DataFrame(gridSearch.cv_results_)[display]
```

```
Initial score:  0.7718120805369127
Initial parameters:  {'hidden_layer_sizes': 1}
```

```
[25]:    param_hidden_layer_sizes  mean_test_score  std_test_score
      0                         1         0.771812        0.018009
      1                         2         0.771812        0.008751
      2                         3         0.771812        0.020022
      3                         4         0.755705        0.019954
      4                         5         0.766443        0.020312
```

```
[26]: # user grid search to fine-tune hyper-parameters
      param_grid = {
          'hidden_layer_sizes': [(1, 2), (1, 3), (1, 4), (1, 5), (1, 6)],
      }

      # Run Exhaustive search for fine-tuning the hyper-parameters
      gridSearch = GridSearchCV(MLPClassifier(activation='logistic',
                                          solver='lbfgs',
                                          random_state=1,
                                          max_iter=5000),
                              param_grid,
```

```
                          cv=5,
                          n_jobs=-1,
                          return_train_score=True)
# Fit to training set
gridSearch.fit(train_X, train_y)

# Improved Scores and hyper-parameters
print('Improved score: ', gridSearch.best_score_)
print('Final parameters: ', gridSearch.best_params_)

# Look at fine-tuned hyper-parameters with averages
display = ['param_hidden_layer_sizes', 'mean_test_score', 'std_test_score']
pd.DataFrame(gridSearch.cv_results_)[display]
```

```
Improved score:  0.7751677852348993
Final parameters:  {'hidden_layer_sizes': (1, 3)}
```

[26]:
| | param_hidden_layer_sizes | mean_test_score | std_test_score |
|---|---|---|---|
| 0 | (1, 2) | 0.768456 | 0.016576 |
| 1 | (1, 3) | 0.775168 | 0.012006 |
| 2 | (1, 4) | 0.768456 | 0.018133 |
| 3 | (1, 5) | 0.770470 | 0.011547 |
| 4 | (1, 6) | 0.771141 | 0.012979 |

[27]:
```
# Final Network
network_model = gridSearch.best_estimator_

# Fit to training data
network_model.fit(train_X, train_y)
```

[27]:
```
MLPClassifier(activation='logistic', hidden_layer_sizes=(1, 3), max_iter=5000,
              random_state=1, solver='lbfgs')
```

[28]:
```
# Helper functions
def confusionMatrices(model, title):
    '''
    Takes in a model and the title to return classification
    summary in accuracy and confusion matrix of the model
    '''

    if model == logit_model:
        print(title + ' - Training results')
        classificationSummary(train_y, model.predict(train_X[columns]))
        print(title + ' - Validation results')
        valid_pred = model.predict(valid_X[columns])
        classificationSummary(valid_y, valid_pred)
    else:
```

```
        print(title + ' - Training results')
        classificationSummary(train_y, model.predict(train_X))
        print(title + ' - Validation results')
        valid_pred = model.predict(valid_X)
        classificationSummary(valid_y, valid_pred)

# Confusion Matrix - Decision Tree
tree_confusion = confusionMatrices(tree_model, '\n\tDecision Tree')

# Confusion Matrix - Logistic Regression
logit_confusion = confusionMatrices(logit_model, '\n\tLogistic regression')

# Confusion Matrix - Multi-layered Neural Network
network_confusion = confusionMatrices(network_model, '\n\tNeural Network')
```

```
        Decision Tree - Training results
Confusion Matrix (Accuracy 0.7785)

      Prediction
Actual   0    1
     0 926   34
     1 296  234

        Decision Tree - Validation results
Confusion Matrix (Accuracy 0.8089)

      Prediction
Actual   0    1
     0 310    7
     1  88   92

        Logistic regression - Training results
Confusion Matrix (Accuracy 0.7772)

      Prediction
Actual   0    1
     0 919   41
     1 291  239

        Logistic regression - Validation results
Confusion Matrix (Accuracy 0.8048)

      Prediction
Actual   0    1
     0 308    9
     1  88   92
```

```
        Neural Network - Training results
Confusion Matrix (Accuracy 0.7805)

        Prediction
Actual    0    1
     0  928   32
     1  295  235

        Neural Network - Validation results
Confusion Matrix (Accuracy 0.8068)

        Prediction
Actual    0    1
     0  309    8
     1   88   92
```

---

# 9   Results and final model selection

- Performance measures on test Set

**Python code:**

[ ]:

---

# 10   Discussion and conclusion

- Address the problem statement and suggestions that could go beyond the scope of the course

[ ]: