

# CUSTOMER PURCHASE PREDICTION

*Using Machine Learning models to predict user behavior*

**ME781 DATA MINING/GROUP 3**

**203100068 MUDIT SAND**

**170040106 BAVISH K**

**193100008 JAYKUMAR GUPTA**

**180100029 BURRA SRIHITH BHARADWAJ**

**150040010 TANYA GUPTA**

## Table of Content

<b>1</b>	<b>Problem definition</b>
<b>2</b>	<b>Technology landscape assessment</b>
<b>3</b>	<b>Project planning report</b>
<b>4</b>	<b>Conceptual design document</b>
<b>5</b>	<b>Auto-generated code document</b>
<b>6</b>	<b>Screenshots of user interface and output</b>
<b>7</b>	<b>Model training and testing report</b>
<b>8</b>	<b>Product brochure and video presentation</b>
<b>9</b>	<b>User manual</b>

**NOTE : All the required documents are also submitted on  
Moodle for better visualization.**

## Problem Definition

COVID-19 has accelerated the transition from visiting physical stores to online shopping. Predicting customer behavior in the context of e-commerce is gaining importance. It can increase customer satisfaction and sales, resulting in higher conversion rates and competitive advantage, by facilitating a more personalized shopping process.

At PREDICT.AI (<http://www.predictai.design/>), we aim to HELP GROWING STARTUPS AND BUSINESSES utilize their customer data and build models for PREDICTING CUSTOMER BEHAVIOUR. Comparing models will give further insight into the performance differences in static customer data. Conducting descriptive data analysis visualization will help our clients extract more value from data and make decisions to boost their customer satisfaction.

<b>CUSTOMER REQUIREMENTS</b>	<b>BUSINESS CASE</b>	<b>BARRIER TO ENTRY &amp; EXISTING PRODUCTS/ SERVICES</b>	<b>UNIQUE SELLING PROPOSITION &amp; PROTECTION OF USP</b>
Accurate models	Target startups and small businesses	Companies not wanting to share data	Easy to use Interface
User satisfaction	Publish conclusions from publicly available data	Companies building their own AI Teams	High accuracy models
Increasing revenue	Subscription model like Bloomberg for companies	Google Analytics	Data protection and privacy
24/7 Help and support	Testimonial from clients	NTENT	Branding of USP

## Technology landscape assessment

### Patents

---

Jivox Kairos™ is the market's first purchase prediction engine for eCommerce marketing, with a patent granted to Jivox for SCORING USERS BASED ON INTENT FOR ONLINE ADVERTISING. As a global brand marketer, you can confidently use this cutting-edge technology to drive sales, by personalizing messaging in real-time to individual consumers based on their interests and in-the-moment purchase intent.

Built on the Jivox Neuron™ AI and machine learning technology, Kairos “learns” how a product is relevant to a specific consumer based on their purchase intent. Kairos algorithms use these purchase intent signals to score individual users' likelihood and immediacy to make a purchase, and rank products relative to the consumer's interest. The pairing of user scoring and product ranking creates for you the opportune moment to serve the right message with the right product offer at the right time.

Source:

(<https://info.jivox.com/kairos-purchase-prediction-ecommerce> )

---

Amex Advance is a data-driven business that partners with companies across the advertising, travel, and service industries to deliver curated personalization services optimized for their customers. Leveraging best-in-class predictive machine learning, deep consumer insights, connectivity capabilities, and an integrated platform, Amex Advance transforms its deterministic data insights into customized solutions to solve partners' key business challenges.

Source :

([www.acxiom.com/news/acxiom-amex-advance-launch-new-data-driven-offering-predict-consumer-purchase-intent/](http://www.acxiom.com/news/acxiom-amex-advance-launch-new-data-driven-offering-predict-consumer-purchase-intent/) )

---

## Libraries Used

NumPy	Pandas	SciKit-Learn	Matplotlib
Plotly	XGBoost	Unittest	Seaborn

## Published Literature

---

Kumar, A., Kabra, G., Mussada, E.K. et al. "Combined artificial bee colony algorithm and machine learning techniques for prediction of online consumer repurchase intention." *Neural Comput & Applic* 31, 877–890 (2019)

This Literature Focuses on predicting online consumer repurchase intention within the context of shopping malls and consumer characteristics using intelligent techniques. This study integrates characteristics of both consumers as well as shopping malls to predict consumer repurchase intention in the online platform. The experimental results have been analyzed through five machine learning classification models, i.e., decision trees, AdaBoost, random forest (RF), support vector machine (SVM) and neural network (NN) in different settings. In the partitions of data into 70–30 training–testing, among all models, the performance of AdaBoost has the highest sensitivity (0.95%) and accuracy (97.58%).

---

Chen, Zhen-Yu, and Zhi-Ping Fan. "Distributed customer behavior prediction using multiplex data: a collaborative MK-SVM approach." *Knowledge-Based Systems* 35 (2012): 111-119.

This paper presents the understanding of customer behavior is a critical success factor. The big databases in an organization usually involve multiplex data such as static, time series, symbolic sequential and textual data which are separately stored in different databases of different sections. It poses a challenge to traditional centralized customer behavior prediction. In this study, a

novel approach called collaborative multiple kernel support vector machine (C-MK-SVM) is developed for distributed customer behavior prediction using multiplex data. The alternating direction method of multipliers (ADMM) is used for the global optimization of the distributed sub-models in C-MK-SVM. Computational experiments on a practical retail dataset are reported. Computational results show that C-MK-SVM exhibits better customer behavior prediction performance and higher computational speed than support vector machine and multiple kernel support vector machine.

---



# Project planning report

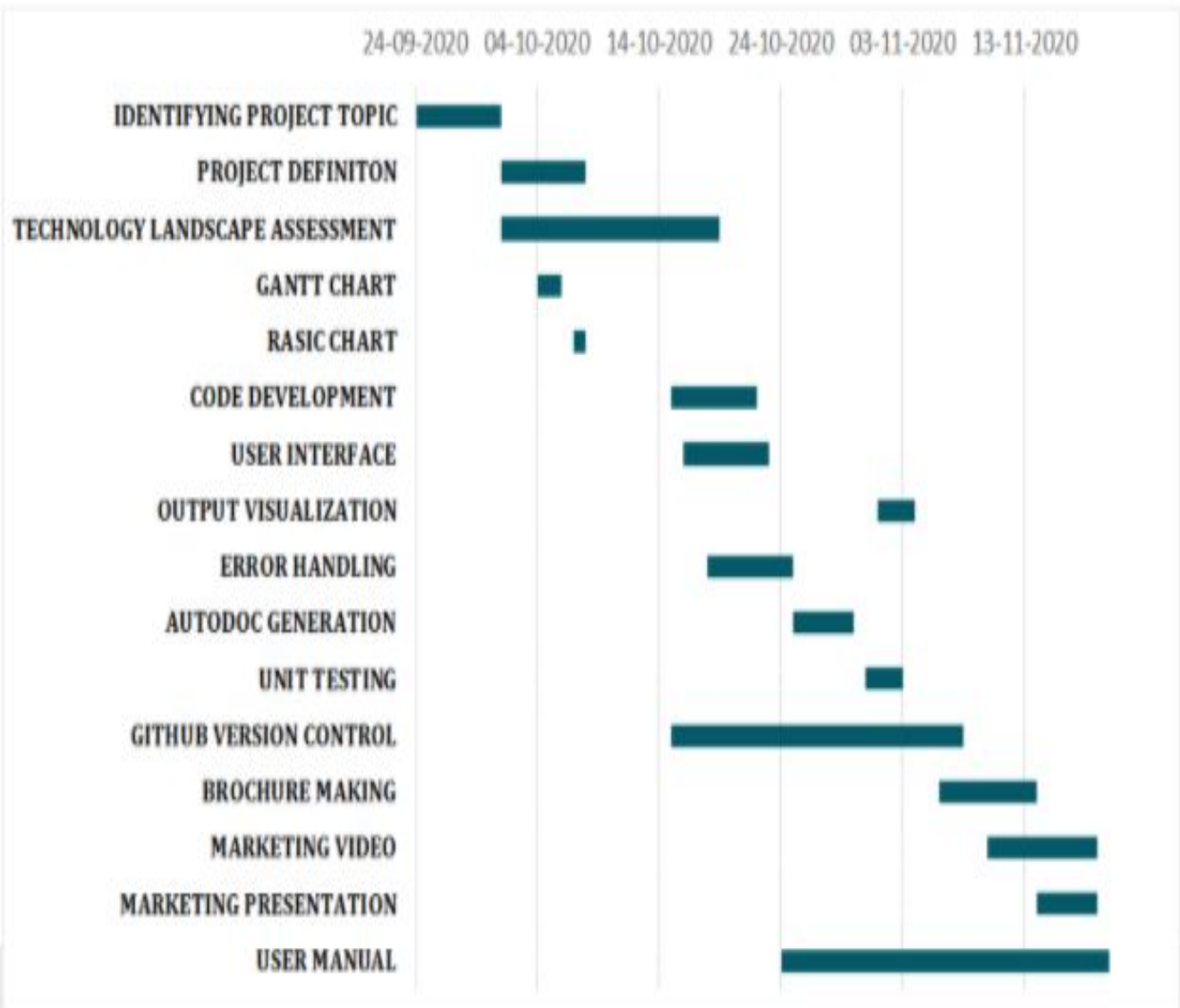
## RASIC CHART

*predict.ai*

Tasks \ People	MUDIT	BAVISH	JAY	SRIHITH	TANYA
OBJECTIVE AND DEFINITION	A	A	S	S	R
TECHNOLOGY LANDSCAPE ASSESSMENT	R	S			S
PLANNING - TIMELINE, GANTT, RASIC CHARTS	A	A	A	A	R
CONCEPTUAL DESIGN - MODEL/DATASET SELECT	A	A	R	S	A
CODE DEVELOPMENT PHASE 1		R	S	S	
CODE DEVELOPMENT PHASE 2	S	R		S	S
MARKETING BROCHURE, PRESENTATION, VIDEO	S		S		R
USER MANUAL & PROJECT REPORT	S	S	S	R	S

# GANTT CHART

*predict.ai*



## Conceptual design document

Class Containing The Information about DataSet

---

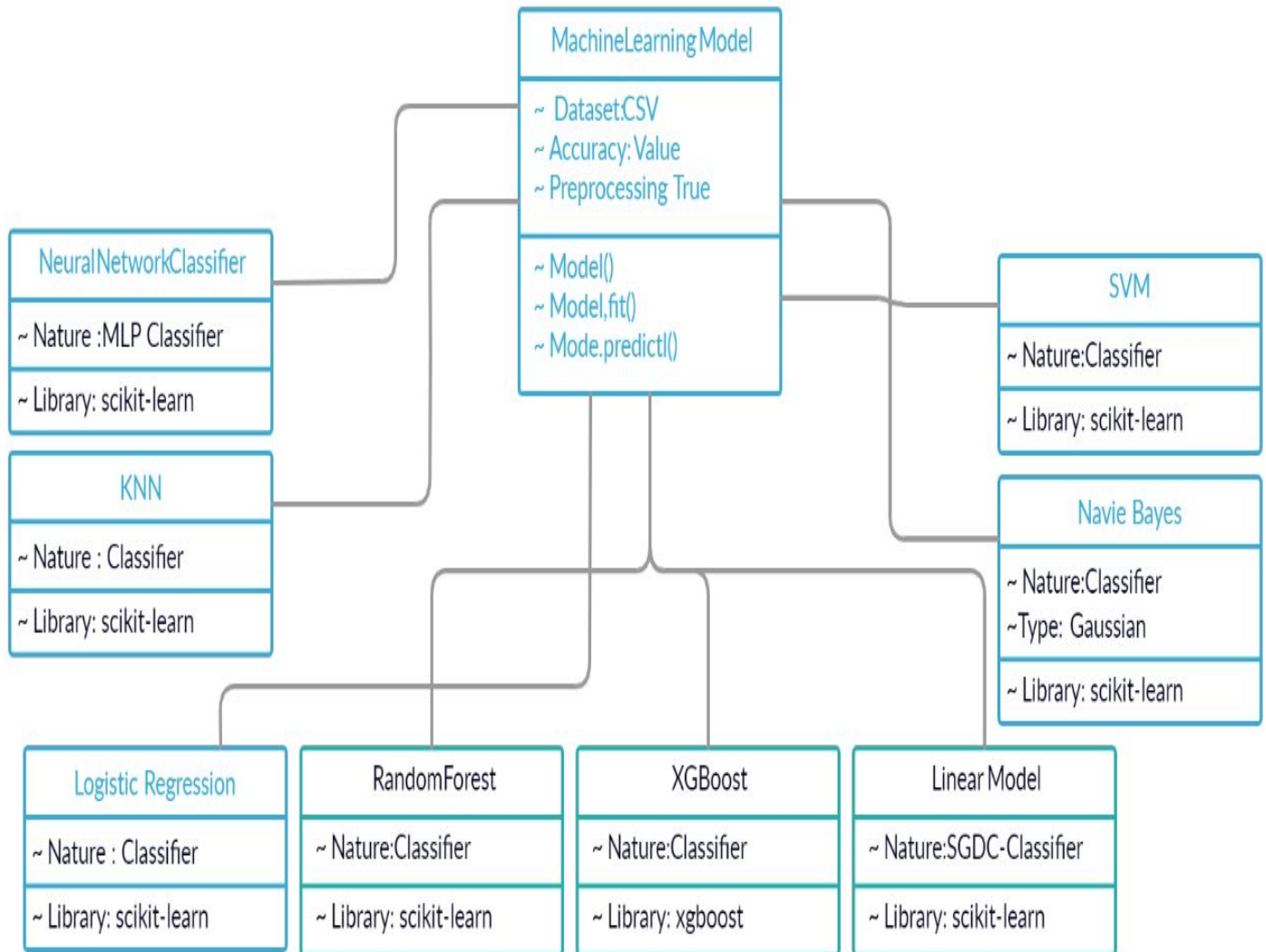
### Shopping mall DataSet Features

- + Administrative: Numerical
- + Administrative\_Duration: Numerical
- + Informational: Numerical
- + Informational\_Duration: Numerical
- + ProductRelated: Numerical
- + ProductRelated\_Duration: Numerical
- + BounceRates: Numerical
- + ExitRates: Numerical
- + PageValues: Numerical
- + SpecialDay: Numerical
- + Month: Categorical
- + OperatingSystems: Categorical
- + Browser: Categorical
- + Region: Categorical
- + TrafficType: Categorical
- + VisitorType: Categorical
- + Weekend: Categorical
- + Revenue: Categorical

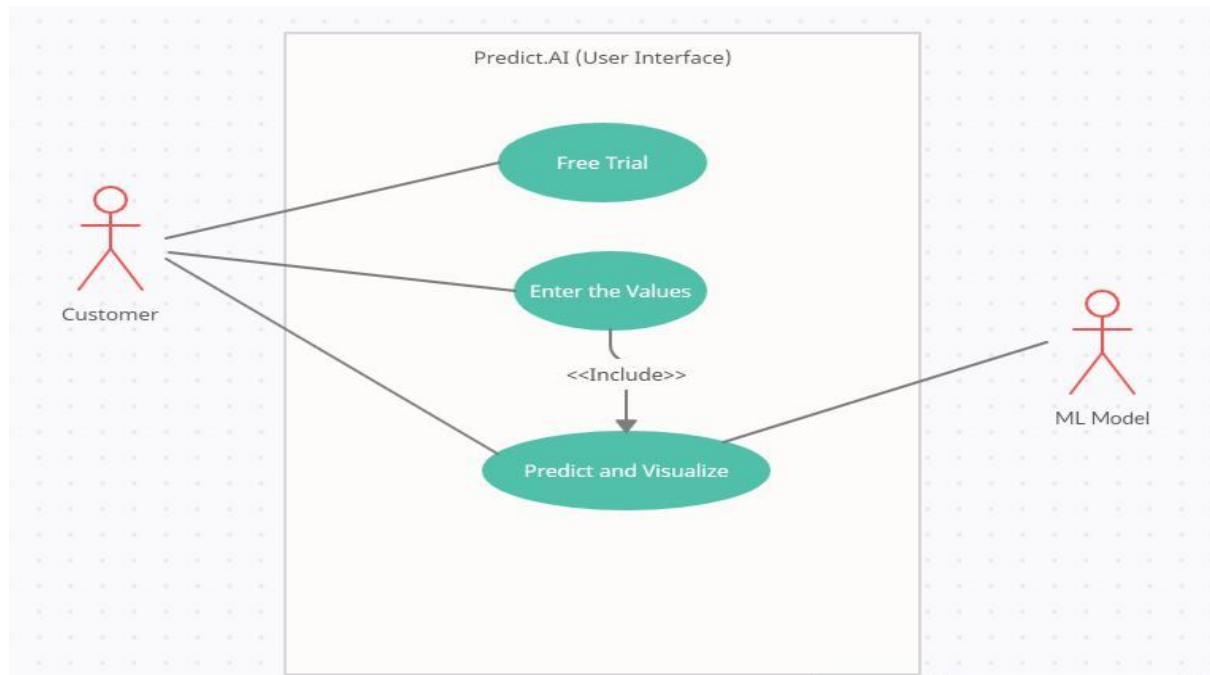
`pandas.read_csv()`

## Machine learning Model Class

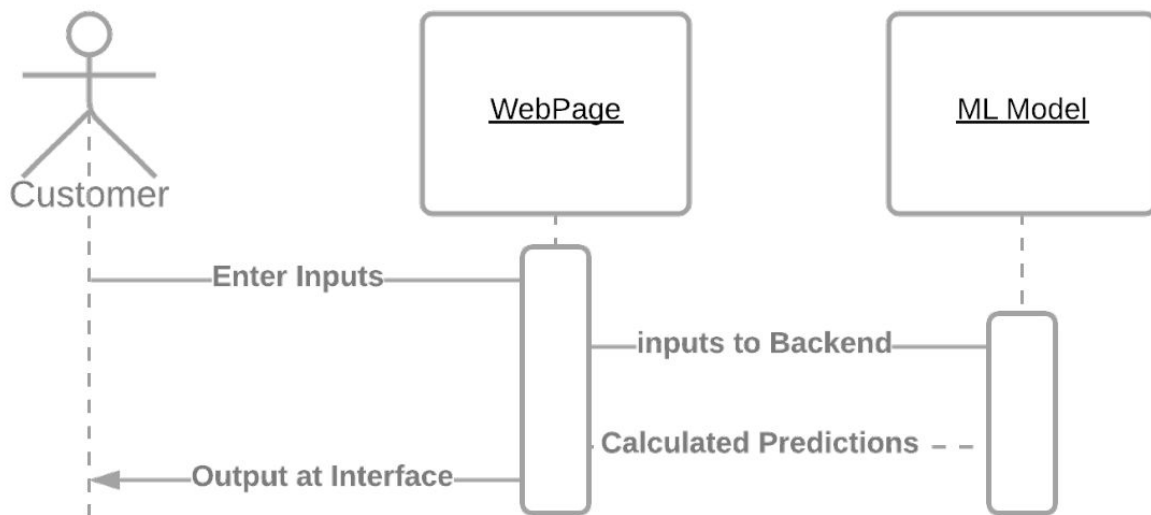
---



## Use Case



## UML Sequence Diagram:





# Auto generated code document

## final\_models.py

Co-authored by Tanya, Bhavik, Mudit, Srihit and Jaykumar as a result of our ME781 Data Mining final project.

Based on Shopping Data set from UCI's Machine Learning Repository

1. Predict.ai Website <http://www.predictai.design/>
2. GitHub Repo <https://github.com/Tannybuoy/predictai>
3. Demo Video <https://www.youtube.com/watch?v=xPt4c4daKM/>

Importing necessary libraries for running Machine Learning models

### Google Drive Mount in Google Colab

Go to directory containing the dataset

Reading shopping data

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.preprocessing import RobustScaler, StandardScaler, MinMaxScaler
from sklearn.model_selection import cross_val_score, train_test_split, cv
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import roc_auc_score
import xgboost as xgb
from xgboost import XGBClassifier
```

```
from google.colab import drive
drive.mount('/content/drive')
```

```
cd "/content/drive/My Drive/Colab Notebooks"
```

```
X_train = pd.read_csv('ShoppingData.csv')
df = X_train.copy()
df.head()
```

## Producing dummy variables for categorical data and cleaning data

A dummy dataframe is created to clean the dataset and preprocess

1. **Visitor** - Columns based on New, Returning or Other
2. **Month** - New columns for each month
3. **Class** - Column after changing data type to int

```
dummiesdf = pd.get_dummies(df['VisitorType'])
df.drop('VisitorType', inplace = True, axis = 1)
df['New_Visitor'] = dummiesdf['New_Visitor']
df['Other'] = dummiesdf['Other']
df['Returning_Visitor'] = dummiesdf['Returning_Visitor']

dfmonth = pd.get_dummies(df['Month'])
df.drop('Month', inplace = True, axis = 1)
dfwithdummies = pd.concat([df, dfmonth], axis = 1, sort = False)
```

```
dfwithdummies['Class'] = df['Revenue'].astype(int)
dfwithdummies.drop('Revenue', axis = 1, inplace = True)
dfwithdummies['Weekend'] = df['Weekend'].astype(int)
dfwithdummies.drop('Returning_Visitor', axis = 1, inplace = True)
dfcleaned = dfwithdummies.copy()
```

```
X = dfcleaned.drop('Class', axis = 1)
Y = dfcleaned['Class'].copy()
```

```
cor = X.corr()
sns.heatmap(cor, xticklabels=cor.columns, yticklabels=cor.columns)
```

## Checking for Collinearity Between Features and Creating Reducing Feature Size

The cor and heatmap help in visualising correlation between various features. Accordingly we do remove the columns/ pre-process.

```
def AvgMinutes(Count, Duration):
    if Duration == 0:
```

AvgMinutes function is used to calculate the average time spent by a customer on the given page. It is obtained by dividing the "Count" by "Duration"

Three new column features hence get added and six columns can now be dropped

Correlation matrix is plotted again using sns heatmap to check if the correlation between the above dropped six features has been dealt with

## Quick overview of features

Histogram of all features

Checking for NA values

Visualising no of unique values and the unique values in each column of the training dataset

Scaling to normalize data

Plotting the histogram obtained post above processing functions

## Linear Model with All Features

Linear model

Accuracy score imported to calculate accuracy

roc\_auc\_score imported to calculate accuracy

It illustrates in a binary classifier system the discrimination threshold created by plotting the true positive rate vs false positive rate

## Random Forest with all Features

```
        output = 0
    elif Duration != 0:
        output = float(Duration)/float(Count)
    return output

Columns = [['Administrative', 'Administrative_Duration'], ['Informational', 'Informational_Duration'], ['ProductRelated', 'ProductRelated_Duration']]

X['AvgAdministrative'] = X.apply(lambda x: AvgMinutes(Count = x['Administrative'], Duration = x['Administrative_Duration']), axis=1)
X['AvgInformational'] = X.apply(lambda x: AvgMinutes(Count = x['Informational'], Duration = x['Informational_Duration']), axis=1)
X['AvgProductRelated'] = X.apply(lambda x: AvgMinutes(Count = x['ProductRelated'], Duration = x['ProductRelated_Duration']), axis=1)
X.drop(['Administrative', 'Administrative_Duration', 'Informational', 'Informational_Duration', 'ProductRelated', 'ProductRelated_Duration'], axis=1)

cor = X.corr()
sns.heatmap(cor, xticklabels=cor.columns, yticklabels=cor.columns)
```

```
for idx, column in enumerate(X.columns):
    plt.figure(idx)
    X.hist(column=column, grid=False)

for i in X.columns:
    print('Feature:', i)
    print('# of N/A:', X[i].isna().sum())

for i in X_train.columns:
    print('*****')
    print('COLUMN TITLE:', i)
    print('# UNIQUE VALUES:', len(X_train[i].unique()))
    print('UNIQUE VALUES:', X_train[i].unique())
    print('*****')
    print()
```

```
X_copy = X.copy()
rc = RobustScaler()
X_rc = rc.fit_transform(X_copy)
X_rc = pd.DataFrame(X_rc, columns=X.columns)

for idx, column in enumerate(X_rc.columns):
    plt.figure(idx)
    X_rc.hist(column=column, grid=False)

from sklearn import linear_model
from sklearn import metrics
X_train, X_test, y_train, y_test = train_test_split(X_rc, Y, test_size=.2)

model = linear_model.SGDClassifier()
model.fit(X_train, y_train)
y_pred = model.predict(X_test)

from sklearn.metrics import accuracy_score
accuracy_score(y_test, y_pred)

from sklearn.metrics import roc_auc_score
roc_auc_score(y_test, y_pred)
```

```
from sklearn.ensemble import RandomForestClassifier
clf = RandomForestClassifier(max_depth=17, random_state=0)
clf.fit(X_train, y_train)
y_pred1 = clf.predict(X_test)

accuracy_score(y_test, y_pred1)

roc_auc_score(y_test, y_pred1)
```

## Finding Important Features then Removing from Dataframe

SelectKBest to obtain a list of importance of each feature column

On seeing the list, we drop the ones which have a very low weightage and less importance

## Random Forest Classifier with Feature Selection Dataframe

Now once again we run Random Forest Classifier, but after retaining only the important features as determined by SelectKBest

## XGBoost Classifier with Feature Selection Dataframe

```
from sklearn import svm
from sklearn.feature_selection import SelectKBest
from sklearn.feature_selection import chi2
list_one = []

feature_ranking = SelectKBest(chi2, k=5)
fit = feature_ranking.fit(X, Y)

fmt = 'X-8s%-20s%'

for i, (score, feature) in enumerate(zip(feature_ranking.scores_, X.columns)):
    list_one.append((score, feature))

dfObj = pd.DataFrame(list_one)
dfObj.sort_values(by=[0], ascending = False)

X_rc.drop(['Aug', 'TrafficType', 'OperatingSystems', 'Other', 'Jul'], axis=1,
          inplace=True)

X_train1, X_test1, y_train1, y_test1 = train_test_split(X_rc, Y, test_size=0.2)

clf1 = RandomForestClassifier(n_estimators=200, max_depth=30)
clf1.fit(X_train1, y_train1)
y_pred2 = clf1.predict(X_test1)

accuracy_score(y_test1, y_pred2)

roc_auc_score(y_test1, y_pred2)

model = XGBClassifier(learning_rate=0.1, n_estimators=150, min_child_weight=5)
model.fit(X_train1, y_train1)
y_pred3 = model.predict(X_test1)

accuracy_score(y_test1, y_pred3)

roc_auc_score(y_test1, y_pred3)
```

## LogisticRegression with Feature Selection Dataframe

## Gaussian Naive Bayes with Feature Selection Dataframe

## KNN classifier with Feature Selection Dataframe

## SVM Classification with PCA feature reduction technique

```
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
log_reg = LogisticRegression(solver='lbfgs', multi_class='multinomial', max_iter=1000)
log_reg.fit(X_train1, y_train1)
y_pred4 = log_reg.predict(X_test1)
print(accuracy_score(y_test1, y_pred4))
print(roc_auc_score(y_test1, y_pred4))

from sklearn.naive_bayes import GaussianNB
GNB = GaussianNB()
GNB.fit(X_train1, y_train1)
y_pred5 = GNB.predict(X_test1)
print(accuracy_score(y_test1, y_pred5))
print(roc_auc_score(y_test1, y_pred5))

from sklearn.neighbors import KNeighborsClassifier
knn = KNeighborsClassifier(n_neighbors=6)
knn.fit(X_train1, y_train1)
y_pred6 = knn.predict(X_test1)
print(accuracy_score(y_test1, y_pred6))
print(roc_auc_score(y_test1, y_pred6))

from sklearn.decomposition import PCA
pca = PCA(n_components=15)
d=pca.fit_transform(X_train1)
e=pca.fit_transform(X_test1)
print(pca.explained_variance_ratio_.sum())

from sklearn.svm import SVC
svm = SVC()
svm.fit(d, y_train1)
y_pred7 = svm.predict(e)
print(accuracy_score(y_test1, y_pred7))
print(roc_auc_score(y_test1, y_pred7))

from sklearn.svm import SVC
svm = SVC()
```



## SVM Classification with Feature Selection Dataframe

## Neural Network Classifier With Feature Selection Dataframe

```
svm.fit(X_train1,y_train1)
y_pred8 = svm.predict(X_test1)
print(accuracy_score(y_pred8,y_test1))
print(roc_auc_score(y_test1, y_pred8))
```

```
from sklearn.neural_network import MLPClassifier

mlp = MLPClassifier(hidden_layer_sizes=(19,19,19), activation='relu', sol
mlp.fit(X_train1,y_train1)
y_pred9= mlp.predict(X_test1)
print(accuracy_score(y_pred9,y_test1))
print(roc_auc_score(y_test1, y_pred9))
```

# Screenshots of user interface and output visualization

USER INTERFACE and OUTPUT:

link(<http://www.predictai.design/trial>)

---

## CUSTOMER PREDICTION TRIAL

---

**Name:**

**PageValues:**

**AvgInformational:**

**AvgAdministrative:**

**AvgProductRelated:**

**Visitor Type:**

**SpecialDay:**

**BounceRates:**

**ExitRates:**

---


## PREDICTION

---


Hi Guest

Unfortunately, this customer will not complete the transaction

what are other words for unfortunately?




unluckily, regrettably, alas, sadly, unhappily, inopportunistly, lamentably, worse luck, unsuccessfully, badly



Thesaurus.plus

---

Made with  by Predict.ai

# CUSTOMER PREDICTION TRIAL

---

Name:

PageValues:

AvgInformational:

AvgAdministrative:

AvgProductRelated:

Visitor Type:

SpecialDay:

BounceRates:

ExitRates:

## PREDICTION

---

Hi Guest

This customer will complete the transaction



Made with  by Predict.ai

## Model training and testing report

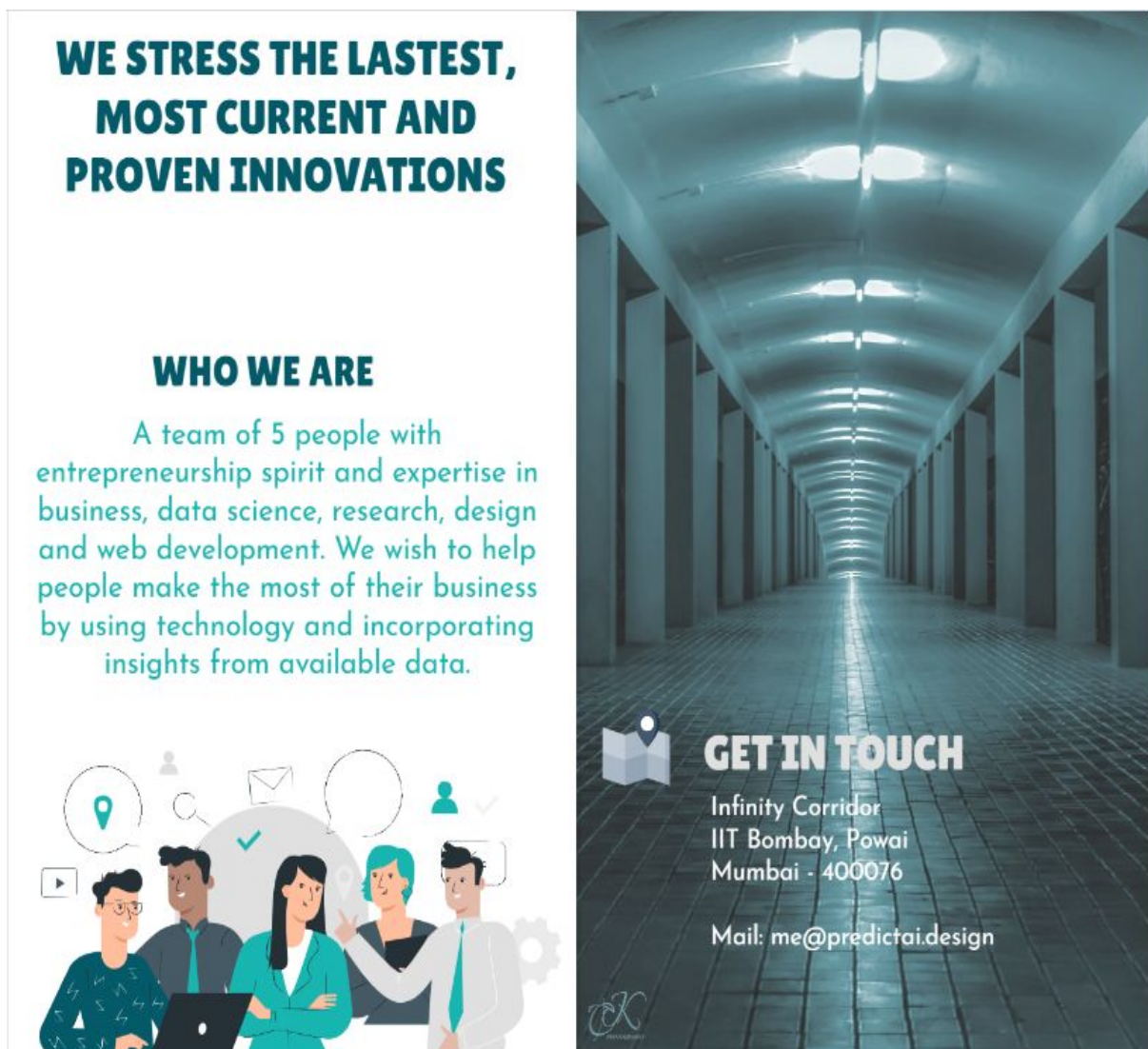
Classifier	Accuracy	R2 accuracy
Logistic regression	0.8892	0.6859
KNN Classifier	0.8819	0.6933
Gaussian Naïve Bayes Classifier	0.6597	0.7475
Neural Network classifier	0.8921	0.729
<i>SVM classifier with Pca Reduction</i>	<i>0.9026</i>	<i>0.7618</i>
<i>Non linear SVM classifier</i>	<i>0.9038</i>	<i>0.7636</i>
<i>Random Forest classifier</i>	<i>0.9042</i>	<i>0.7617</i>
SGDC-classifier	0.8892	0.7465
<i>XGBoost Classifier</i>	<i>0.9022</i>	<i>0.7711</i>

## Product brochure and video presentation

Link of the Video presentation:

<https://www.youtube.com/watch?v=xFt4cl4daKM>

Broucher Pages:



## OUR VALUES

### Innovation

We value the trust and respect of our community and coworkers. We commit to becoming a place where we do what's right because we love what's right.

### Integrity

We share our clients' aspirations, work vigorously to understand their reality, and align our incentives with their objectives – so they know we're in this for the long haul.

### Expertise

We celebrate the power of technology to transform lives. We commit to helping people use technology creatively and ethically.



**We deliver transformational results that address your needs.**



## OUR PRODUCTS

### Data Analytics

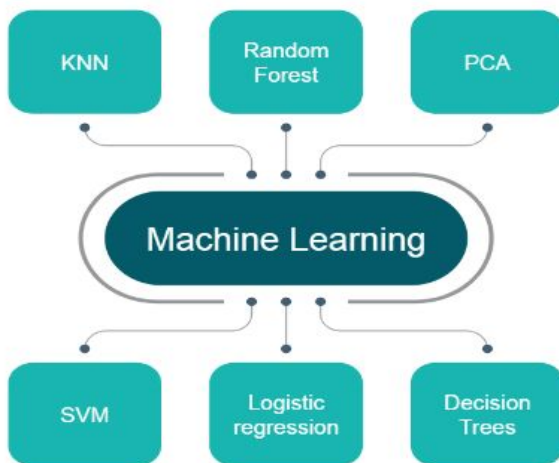
Dashboard to get customer insights from data. View activity across products, contact us for support

### Customer Purchase Prediction

Use exit rates, bounce rates, time spent on a page, month and more to predict whether customer will make a purchase. Plan your marketing strategies accordingly.

## OUR ALGORITHMS

Our state of the art technology and analytics ensures the highest accuracy and best extraction of information from data. After all, data is the new oil. Make decisions based on data analytics and customer behaviour prediction.



## PRICING

Free for hobby project enthusiasts and individuals  
Can be customised based on company size and requirements

*predict.ai*

**Driving business development and innovation through technology**





# User manual

## USER MANUAL

### Steps:

- 1) Visit <http://www.predictai.design/>
- 2) Click on 'FREE TRIAL' at Navigation
- 3) Enter The Values in the Input Fields
- 4) Click on Predict
- 5) Visualize the result Predicted by Machine Learning Models.

### Contact Us:

Drop a message in contact section of the website mentioned in step 1.

### Meaning of Input Fields:

- 1) Name : Enter your Name
- 2) PageValues: Numerical value
- 3) AvgInformational: Numerical value
- 4) AvgAdministrative: Numerical value
- 5) AvgProductRelated: Numerical value
- 6) Visitor Type: Binary (0 or 1)
- 7) SpecialDay: Between 0 and 1
- 8) BounceRates: Between 0 and 1
- 9) ExitRates: Between 0 and 1

Thanks for your association with Predict.AI