

INSURANCE FRAUD DETECTION

#Data Collection

- we have insurance claims data (insurance_claim.csv). However, this data is not perfect.
- It may contain outliers, missing values and irrelevant (not related to insurance claims) information

```
In [1]: import warnings
warnings.filterwarnings(action='ignore', category=FutureWarning)
```

```
In [2]: # importing required libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

# To get all the rows and columns
pd.set_option('display.max_rows', None)
pd.set_option('display.max_columns', None)
```

```
In [3]: # Reading the Dataset File

df=pd.read_csv('insurance_claims.csv' )
```

```
In [4]: df.head()
```

```
Out[4]:
```

	months_as_customer	age	policy_number	policy_bind_date	policy_state	policy_csl	policy_deductable	policy_annu
0	328	48	521585	2014-10-17	OH	250/500	1000	
1	228	42	342868	2006-06-27	IN	250/500	2000	
2	134	29	687698	2000-09-06	OH	100/300	2000	
3	256	41	227811	1990-05-25	IL	250/500	2000	
4	228	44	367455	2014-06-06	IL	500/1000	1000	

```
In [5]: # Check the columns
df.columns

# _c39 is extra columns Let's remove it
df.drop('_c39', axis = 1, inplace=True)
```

```
In [6]: df.head()
```

```
Out[6]:
```

	months_as_customer	age	policy_number	policy_bind_date	policy_state	policy_csl	policy_deductable	policy_annu
0	328	48	521585	2014-10-17	OH	250/500	1000	
1	228	42	342868	2006-06-27	IN	250/500	2000	
2	134	29	687698	2000-09-06	OH	100/300	2000	
3	256	41	227811	1990-05-25	IL	250/500	2000	
4	228	44	367455	2014-06-06	IL	500/1000	1000	

```
In [7]: # Data description
def display_basic_information(df, column_name):
    df.describe()
```

```
Out[7]:
```

Function is responsible for printing out basic information about the dataset such as Null values, Data types, etc..

	months_as_customer	age	policy_number	policy_deductable	policy_annual_premium	umbrella_limit
count	1000.000000	1000.000000	1000.000000	1000.000000	1000.000000	1.000000e+03
data_type	float64	int64	int64	int64	int64	int64
null_values = df[column_name].isna().sum()	0	0	0	0	0	0
mean	200.350000	38.940000	546238.648000	1136.000000	1256.406150	1.101000e+06
std	157.766000	10.287000	257063.005276	611.864673	244.167395	2.297407e+06
min	0.000000	18.000000	100804.000000	500.000000	433.330000	-1.000000e+06
25%	157.766000	28.000000	335980.250000	500.000000	1089.607500	0.000000e+00
50%	199.500000	38.000000	533135.000000	1000.000000	1257.200000	0.000000e+00
75%	276.250000	44.000000	759099.750000	2000.000000	1415.695000	0.000000e+00
max	328.000000	64.000000	999435.000000	2000.000000	2047.590000	1.000000e+07

display_basic_information(df,i)

```
In [8]: # Number of rows and columns of data
#####
Basic Information: rows: ", len(df))
#####
Basic Information: columns: ", len(df.columns))
```

```
Number of rows: 1000
Number of columns: 39
total number of null values: 0
```

Data Preparation

```
#####
Basic Information:
#####
• We prepare the data column wise, as every column would require different preprocessing

column name: age
data type: int64
total number of null values: 0
```

Out[7]:

In [8]:

In [10]:

In [13]:

Out[13]:

In [11]:

In [12]:

Out[12]:

```

authorities_contacted      0
incident_month             0
incident_month_as_customer, 'age', 'policy_state', 'policy_deductible',
incident_provided_education_level, 'umbrella_limit', 'insured_zip', 'insured_sex',
number_of_vehicles_involved, 'insured_occupation', 'insured_relationship',
property_claimed_gains', 'capital_loss', 'incident_type', 'collision_type',
bodily_injury_claim_severity', 'authorities_contacted', 'incident_state',
witnesses', 'incident_hour_of_the_day', 'number_of_vehicles_involved',
police_report_available', 'bodily_injuries', 'witnesses',
total_claim_amount_report_available', 'total_claim_amount', 'injury_claim',
injury_claim_property_claim', 'vehicle_claim', 'auto_make', 'auto_year',
property_claim_reported'],
vehicle_claim(object')
auto make
0

```

Checking the quality of data

```
# Saving the output in HTML format to view it
profile.to_file(output_file = 'Insurance_Prediction_Report.html')
```

[illegible]

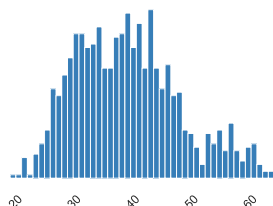
```
In [17]: fraud_reported['Y'].N
from pandas_profiling import ProfileReport

# Exploratory Data Analysis for the dataset
profile = ProfileReport(df, title= 'Insurance Fraud Claims Report', html={ 'style': { 'full_width': True } })
profile.to_notebook_iframe()

# Saving the output in HTML format to view it
profile.to_file(output_file = 'Insurance Prediction Report.html')
```

```
Summarize dataset: 0%|          | 0/5 [00:00<?, ?it/s]
Generate report structure: 0%|          | 0/1 [00:00<?, ?it/s]
Render HTML: 0%|          | 0/1 [00:00<?, ?it/s]
```

Missing	0
Missing (%)	0.0%
Infinite	0
Infinite (%)	0.0%
Mean	38.948
Minimum	19
Maximum	64
Zeros	0
Zeros (%)	0.0%
Negative	0
Negative (%)	0.0%
Memory size	7.9 KiB



```
In [20]: eda_df = df
eda_df.head(10)
```

months_as_customer	age	policy_state	policy_deductable	policy_annual_premium	umbrella_limit	insured_zip	insured_age
12	25	CA	0	1200	1000000	90210	25
18	30	TX	500	1500	500000	75001	30
24	35	FL	1000	1800	200000	33133	35
30	40	NY	1500	2200	100000	10001	40
36	45	IL	2000	2500	50000	60601	45
42	50	WA	2500	2800	20000	98101	50
48	55	OH	3000	3000	10000	43201	55
54	60	PA	3500	3200	5000	19101	60
60	65	MD	4000	3500	2000	21201	65
66	70	DE	4500	3800	1000	19701	70
72	75	NC	5000	4000	500	27601	75
78	80	SC	5500	4200	200	29401	80
84	85	GA	6000	4500	100	30301	85
90	90	LA	6500	4800	50	70110	90

	months_as_customer	age	policy_state	policy_deductable	policy_annual_premium	umbrella_limit	insured_zip	insured
0	Quantile statistics			OH	1000	1406.91	0	466132
1	Minimum	228	42	IN	2000	1197.22	5000000	468176
2	5-th percentile	234	29	OH	2000	1413.14	5000000	430632
3	Q1	256	41	IL	2000	1415.74	6000000	608117
4		228	44	IL	1000	1583.91	6000000	610706
5	Export report to file:	256	39	OH	1000	1351.10	0	478456
6	!pip install pyyaml==5.4.1	137	34	IN	1000	1333.35	0	441716
	Requirement already satisfied: pyyaml==5.4.1 in c:\users\shekhar\anaconda3\lib\site-packages (5.4.1)	165	37	IL	1000	1137.03	0	603195
	import plotly.graph_objects as go				500	1442.99	0	601734
	import plotly.express as px							
9		212	42	IL	500	1315.68	0	600983

Data Distribution using Histogram

```
In [21]: def count_plot(df):
```

```

In [20]: eda_df = df
eda_df.head(10)

Out[20]:
  months_as_customer  age  policy_state  policy_deductable  policy_annual_premium  umbrella_limit  insured_zip  insured_...
0  Quantile statistics
1  Minimum  228  42  IN  2000  1197.22  5000000  468176
2  5-th percentile  134  29  OH  2000  1413.14  5000000  430632
3  Q1  256  41  IL  2000  1415.74  6000000  608117
4  228  44  IL  1000  1583.91  6000000  610706
Export report to file: 0%| | 0/1 [00:00<?, ?it/s]
5  256  39  OH  1000  1351.10  0  478456
6  137  34  IN  1000  1333.35  0  441716
Requirement already satisfied: pyyaml==5.4.1 in c:\users\shekhar\anaconda3\lib\site-packages
(p.4.1) 165 37 IL 1000 1137.03 0 603195
In [18]: !pip install pyyaml==5.4.1
6 137 34 IN 1000 1333.35 0 441716
Requirement already satisfied: pyyaml==5.4.1 in c:\users\shekhar\anaconda3\lib\site-packages
(p.4.1) 165 37 IL 1000 1137.03 0 603195
In [19]: import plotly.graph_objects as go
import plotly.express as px
9 212 42 IL 500 1315.68 0 600983

```

Data Distribution using Histogram

```

In [21]: def count_plot(df):
        """
        Function is responsible for printing out the count plot on sys.out
        """
        for col in df.columns:
            fig = px.histogram(df, col, nbins=20, width=500, height=400)
            fig.show()
        count_plot(eda_df)

```



###Correlation

```

In [22]: # Correlation heatmap for all numeric values
plt.figure(figsize = (18,10))
sns.heatmap(eda_df.corr(), annot=True)

Out[22]: <AxesSubplot:~>

```

```

In [25]: Following observations were made
from matplotlib import pyplot as plt
import seaborn as sns
1. months_as_customer is highly correlated with age
2. total_claim_amount is highly correlated with injury_claim and vehicle_claim and property_claim
plt.figure(figsize=(15,8))
sns.countplot(data=eda_df, x='age', hue='fraud_reported')

```

```

In [23]: # Let's remove months_as_customer and total_claim_amount column as they provides high correlation
Out[25]: eda_df.drop(['months_as_customer', 'total_claim_amount'], axis=1, inplace=True)

```

```

In [24]: # Can right education level prevent fraud on high scale ?
plt.figure(figsize=(15,8))
sns.heatmap(eda_df.corr(), annot=True)
sns.countplot(data=eda_df, x='insured_education_level', hue='fraud_reported')
Out[24]: <AxesSubplot:~>
Out[26]: <AxesSubplot:xlabel='insured_education_level', ylabel='count'~>

```

```

###Hypothesis Based on EDA 1.Is the age of person have any prominent relationship with insurance fraud ?
In [27]: #which incident type have more probability of fraud?
sns.countplot(data=eda_df, x='incident_type', hue='fraud_reported')
Out[27]: <AxesSubplot:xlabel='incident_type', ylabel='count'~>

```

```

5.which city incident are leads to fraud ?
In [28]: #which incident_severity have more impact on fraud?
plt.figure(figsize=(15,8))
sns.countplot(data=eda_df, x='incident_severity', hue='fraud_reported')
Out[28]: <AxesSubplot:xlabel='incident_severity', ylabel='count'~>

```

```

In [29]: #which state incident are Leads to fraud ?
plt.figure(figsize=(15,8))

```

```

In [25]: #Following observations were made
from matplotlib import pyplot as plt
import seaborn as sns
#important variables customer is highly correlated with age
#total claim amount is highly correlated with injury_claim and vehicle_claim and property_claim
plt.figure(figsize=(15,8))
sns.countplot(data=eda_df, x='age', hue='fraud_reported')

In [23]: #Let's remove months_as_customer and total_claim_amount column as they provides high correlation
eda_df.drop(['months_as_customer', 'total_claim_amount'], axis=1, inplace=True)

Out[25]: <AxesSubplot: xlabel='age', ylabel='count'>

In [26]: #Correlation heatmap to see which variables are involved in fraud on high scale ?
plt.figure(figsize=(16,10))
sns.heatmap(eda_df.corr(), annot=True)
sns.countplot(data=eda_df, x='insured_education_level', hue='fraud_reported')

Out[24]: <AxesSubplot: >
Out[26]: <AxesSubplot: xlabel='insured_education_level', ylabel='count'>

###Hypothesis Based on EDA 1.Is the age of person have any prominent relationship with insurance fraud ?
In [27]: #which incident type have more probability of fraud?
sns.countplot(data=eda_df, x='incident_type', hue='fraud_reported')

Out[27]: <AxesSubplot: xlabel='incident_type', ylabel='count'>

5.which city incident are leads to fraud ?
In [28]: #which incident severity have more impact on fraud?
plt.figure(figsize=(15,8))
sns.countplot(data=eda_df, x='incident_severity', hue='fraud_reported')

Out[28]: <AxesSubplot: xlabel='incident_severity', ylabel='count'>

In [29]: #which state incident are Leads to fraud ?
plt.figure(figsize=(15,8))
sns.countplot(data=eda_df, x='incident_state', hue='fraud_reported')

Out[29]: <AxesSubplot: xlabel='incident_state', ylabel='count'>

In [30]: #does more old the auto/vehicle more the probability to get lead to fraud ?
plt.figure(figsize=(15,8))
sns.countplot(data=eda_df, x='auto_year', hue='fraud_reported')

Out[30]: <AxesSubplot: xlabel='auto_year', ylabel='count'>

```

#Data Preprocessing (One Hot)

Since there are multiple categorical columns, we need to processed it using multiple methods.

For the current problem statement we would be using:

1. [One hot encoding \(https://www.google.com/url?sa=t&rc=1&q=&esrc=s&source=web&cd=&cad=rja&uact=8&ved=2ahUKEwjhybaXiYD8AhUqS2wGHXrXC48QFnoqne-hot-encode-data-in-machine-learning%2F&usq=AOvVaw0wM6DDmQLcewKfleh-O-v\)](https://www.google.com/url?sa=t&rc=1&q=&esrc=s&source=web&cd=&cad=rja&uact=8&ved=2ahUKEwjhybaXiYD8AhUqS2wGHXrXC48QFnoqne-hot-encode-data-in-machine-learning%2F&usq=AOvVaw0wM6DDmQLcewKfleh-O-v)
2. [Weight of evidence encoding \(https://www.google.com/url?sa=t&rc=1&q=&esrc=s&source=web&cd=&cad=rja&uact=8&ved=2ahUKEwii9eekiYD8AhVWSWwGHRiJCdoQFnof-of-evidence-woe-and-information-value-iv%2F&usq=AOvVaw0L9j12BFxR75K46dbAi_VR\)](https://www.google.com/url?sa=t&rc=1&q=&esrc=s&source=web&cd=&cad=rja&uact=8&ved=2ahUKEwii9eekiYD8AhVWSWwGHRiJCdoQFnof-of-evidence-woe-and-information-value-iv%2F&usq=AOvVaw0L9j12BFxR75K46dbAi_VR)

Handling columns

1.Creating copy of eda_df i.e eda data frame into new data frame for data processing

```

In [31]: processed_df=eda_df.copy()
processed_df.head(10)

Out[31]:
age  policy_state  policy_deductable  policy_annual_premium  umbrella_limit  insured_zip  insured_sex  insured_educ
0      48          OH                1000                1406.91             0          466132      MALE
1      42          IN                2000                1197.22             5000000      468176      MALE
2      29          OH                2000                1413.14             5000000      430632      FEMALE

policy_state ['OH' 'IN' 'IL']
insured_sex ['MALE' 'FEMALE']
insured_education_level ['MD' 'PhD' 'Associate' 'Masters' 'High School' 'College' 'JD']
insured_occupation ['craft-repair' 'machine-oper' 'sales' 'armed-forces' 'technical-support'
'prof-specialty' 'other-service' 'priv-house-serv' 'exec-managerial'
'protective-serv' 'transport-moving' 'handler-cleaners' 'admclerical' 'farming-fishing']
insured_relationship ['husband' 'other-relative' 'own-child' 'unmarried-wife' 'MALE-in-famil
y']
incident_type ['Single Vehicle Collision' 'Vehicle Theft' 'Multi-vehicle Collision'
'Parked Car']
collision_type ['Side Collision' 'Not Known' 'Rear Collision' 'Front Collision']
incident_severity ['Major Damage' 'Minor Damage' 'Total Loss' 'Trivial Damage']
authorities_contacted ['Police' 'None' 'Fire' 'Other' 'Ambulance']
incident_state ['SC' 'VA' 'NY' 'OH' 'WV' 'NC']
property_damage ['YES' 'Not Known' 'NO']
police_report_available ['YES' 'Not Known' 'NO']
auto_make ['Saab' 'Mercedes' 'Dodge' 'Chevrolet' 'Accura' 'Nissan' 'Audi' 'Toyota'
'Ford' 'Suburu' 'BMW' 'Jeep' 'Honda' 'Volkswagen']
fraud_reported ['Y' 'N']

```

1. fraud_reported is our target column in string format such as "Yes","No" we will convert it to numeric indication "Yes"= 1 and "No"= 0
2. insured_sex also have gender column in format such as "male" and "female" we will convert that to numeric indicating "male"=1 "female"= 0

```
processed_df.head(10)
```

Out[31]:

Handling rows

In [32]:

```
age policy_state policy_deductable policy_annual_premium umbrella_limit insured_zip insured_sex insured_educ
0 48 OH 1000 1406.91 0 466132 MALE
1 42 IN 2000 1197.22 5000000 468176 MALE
2 29 OH 2000 1413.14 5000000 430632 FEMALE
3 41 IL 2000 1415.74 6000000 608117 FEMALE
4 44 IL 1000 1583.91 6000000 610706 MALE
5 39 OH 1000 1351.10 0 478456 FEMALE
6 34 IN 1000 1333.35 0 441716 MALE
7 37 IL 1000 1137.03 0 603195 MALE
8 33 IL 500 1442.99 0 601734 FEMALE
9 42 IL 500 1315.68 0 600983 MALE

policy_state ['OH' 'IN' 'IL']
insured_sex ['MALE' 'FEMALE']
insured_education_level ['MD' 'PhD' 'Associate' 'Masters' 'High School' 'College' 'JD']
insured_occupation ['craft-repair' 'machine-op-asst' 'sales' 'armed-forces' 'tech-support'
'prof-specialty' 'other-service' 'priv-house-serv' 'exec-managerial'
'protective-serv' 'transport-moving' 'handler-cleaners' 'admin-clerical'
'farming-fishing']
insured_relationship ['husband' 'other-relatives' 'own-child' 'unmarried-wife' 'in-family']
incident_type ['Single Vehicle Collision' 'Vehicle Theft' 'Multi-vehicle Collision'
'Parkd Car']
collision_type ['Side Collision' 'Not Known' 'Rear Collision' 'Front Collision']
incident_severity ['Major Damage' 'Minor Damage' 'Total Loss' 'Trivial Damage']
authorities_contacted ['Police' 'None' 'Fire' 'Other' 'Ambulance']
incident_state ['SC' 'VA' 'NY' 'OH' 'WV' 'NC']
property_damage ['YES' 'Not Known' 'NO']
police_report_available ['YES' 'Not Known' 'NO']
auto_make ['Saab' 'Mercedes' 'Dodge' 'Chevrolet' 'Accura' 'Nissan' 'Audi' 'Toyota'
'Ford' 'Suburu' 'BMW' 'Jeep' 'Honda' 'Volkswagen']
fraud_reported ['Y' 'N']
```

1. fraud_reported is our target column in string format such as "Yes","No" we will convert it to numeric indication "Yes"= 1 and "No"= 0
2. insured_sex also have gender column in format such as "male" and "female" we will convert that to numeric indicating "male"=1 "female"= 0
3. education_level and incident_severity is replaced with Ordinal Encoding

In [33]:

```
processed_df.fraud_reported.replace(('Y','N'),(1,0),inplace=True)
processed_df.insured_sex.replace(('MALE','FEMALE'),(1,0),inplace=True)
processed_df.insured_education_level.replace(('MD','PhD','Associate','Masters','High School',
processed_df.incident_severity.replace(('Trivial Damage','Minor Damage','Major Damage','Total Loss'))
processed_df.head(10)
```

Out[33]:

```
age policy_state policy_deductable policy_annual_premium umbrella_limit insured_zip insured_sex insured_educ
0 48 OH 1000 1406.91 0 466132 1
1 42 IN 2000 1197.22 5000000 468176 1
2 29 OH 2000 1413.14 5000000 430632 0
3 41 IL 2000 1415.74 6000000 608117 0
4 44 IL 1000 1583.91 6000000 610706 1
5 39 OH 1000 1351.10 0 478456 0
6 34 IN 1000 1333.35 0 441716 1
7 37 IL 1000 1137.03 0 603195 1
8 33 IL 500 1442.99 0 601734 0
9 42 IL 500 1315.68 0 600983 1
```

In [35]:

```
from sklearn.preprocessing import StandardScaler
```

In [34]:

```
# Applying one hot encoding using get dummies method
# Importing the standard scaler class
one_hot_df = pd.get_dummies(processed_df, drop_first=True)
scaler = StandardScaler()
one_hot_df.head(10)
```

Out[34]:

```
# Dropping the target columns as scaling is not required on that
scaled_policy_deductable, policy_annual_premium, umbrella_limit, insured_zip, insured_sex, insured_education_level, c
```

In [37]:

```
# Scaling the dataframe
one_hot_scaled_df = pd.DataFrame(scaler.fit_transform(scaled_df), columns=scaled_df.columns)
```

In [38]:

```
# Before applying feature selection, let's just save the dependent column
y = one_hot_df['fraud_reported']
x = one_hot_scaled_df
```

In [39]:

```
from sklearn.decomposition import PCA
pca = PCA()
pca.fit(x)
# Get the principle components for x dataframe
pca.components_
```

In [40]:

```
# Variance explained by each principle component
pd.DataFrame(np.cumsum(pca.explained_variance_ratio_))
```

Out[40]:

```
Due to one hot, we got large set of columns. To select the best columns among them we need to perform Principle Component analysis (PCA).
1 0.102400
Before PCA we need to standardize and scale all the columns using Standard Scaler.
2 0.126466
```



```

In [35]: from sklearn.preprocessing import StandardScaler

In [34]: # Applying one hot encoding using get_dummies method
# Importing the standard scaler class
one_hot_df = pd.get_dummies(processed_df, drop_first=True)
scaler = StandardScaler()
one_hot_df.head(10)

Out[34]: # Dropping the target columns as scaling is not required on that
scaled_policy_onehot_df = pd.DataFrame(scaled_df.drop('fraud_reported', axis=1), columns=scaled_df.columns)

In [37]: # Scaling the dataframe
one_hot_scaled_df = pd.DataFrame(scaler.fit_transform(scaled_df), columns=scaled_df.columns)

In [38]: # Before applying feature selection, let's just save the dependent column
y = one_hot_df['fraud_reported']
x4 = one_hot_scaled_df

In [39]: from sklearn.decomposition import PCA
pca = PCA()
pca.fit(x4)
# Get the principle components for wef dataframe
principle_components_wef = pca.fit_transform(x)

In [40]: # Variance explained by each principle component
pd.DataFrame(np.cumsum(pca.explained_variance_ratio_))

```

Feature Selection (One Hot)

Due to one hot, we got large set of columns. To select the best columns among them we need to perform Principle Component Analysis (PCA).

Before PCA we need to standardize and scale all the columns using Standard Scaler.

```

1 0.192400
2 0.126466
3 0.149562
4 0.172372
5 0.194346
6 0.216174
7 0.237454
8 0.258261
9 0.278391
10 0.298387

```

```

In [41]: plt.figure()
plt.plot(np.cumsum(pca.explained_variance_ratio_))
plt.xlabel('Number of components')
plt.ylabel('Variance')
plt.title('explained_variance_ratio')
plt.show()

```

C:\Users\SHEKHAR\AppData\Local\Temp\ipykernel_10448\2189574362.py:6: UserWarning:

Matplotlib is currently using module://matplotlib_inline.backend_inline, which is a non-GUI backend, so cannot show the figure.

We can see that around 99% of variance is being explained by 66 components. So instead of giving all columns as input in our algorithm let's use these principle components instead

```

In [42]: pca = PCA(n_components=66)
pca_onehot_df = pca.fit_transform(x)

column_name = [f'PC-{i}' for i in range(1, 67)]
print(column_name)

In [43]: principal_onehot_df = pd.DataFrame(pca_onehot_df, columns=column_name)
principal_onehot_df.head(5)

from sklearn.model_selection import KFold
['PC-1', 'PC-2', 'PC-3', 'PC-4', 'PC-5', 'PC-6', 'PC-7', 'PC-8', 'PC-9', 'PC-10', 'PC-11', 'PC-12', 'PC-13', 'PC-14', 'PC-15', 'PC-16', 'PC-17', 'PC-18', 'PC-19', 'PC-20', 'PC-21', 'PC-22', 'PC-23', 'PC-24', 'PC-25', 'PC-26', 'PC-27', 'PC-28', 'PC-29', 'PC-30', 'PC-31', 'PC-32', 'PC-33', 'PC-34', 'PC-35', 'PC-36', 'PC-37', 'PC-38', 'PC-39', 'PC-40', 'PC-41', 'PC-42', 'PC-43', 'PC-44', 'PC-45', 'PC-46', 'PC-47', 'PC-48', 'PC-49', 'PC-50', 'PC-51', 'PC-52', 'PC-53', 'PC-54', 'PC-55', 'PC-56', 'PC-57', 'PC-58', 'PC-59', 'PC-60', 'PC-61', 'PC-62', 'PC-63', 'PC-64', 'PC-65', 'PC-66']

Out[42]: kf = KFold(n_splits=10, random_state=0, shuffle=True)
for train_index, test_index in kf.split(x):
    x_train, x_test = x.values[train_index], x.values[test_index]
    y_train, y_test = y.values[train_index], y.values[test_index]

# Training model
model = LogisticRegression()
model.fit(x_train, y_train)

# Evaluating model
y_pred = model.predict(x_test)

accuracy_scores.append(sm.accuracy_score(y_test, y_pred))
# Data Modeling (One Hot)
precision_scores.append(sm.precision_score(y_test, y_pred))
recall_scores.append(sm.recall_score(y_test, y_pred))
f1_scores.append(sm.f1_score(y_test, y_pred))

#displaying average results
print("#####")
print("Results")
print("#####\n")
print("Average accuracy score: ", (sum(accuracy_scores)/len(accuracy_scores)))

results = [(sum(accuracy_scores)/len(accuracy_scores))]
return results

```



```
In [43]: principal_onehot_df = pd.DataFrame(pca_onehot_df, columns=column_name)
principal_onehot_df.head()
from sklearn.model_selection import KFold
['PC-1', 'PC-2', 'PC-3', 'PC-4', 'PC-5', 'PC-6', 'PC-7', 'PC-8', 'PC-9', 'PC-10', 'PC-11', 'PC-12', 'PC-13', 'PC-14', 'PC-15', 'PC-16', 'PC-17', 'PC-18', 'PC-19', 'PC-20', 'PC-21', 'PC-22', 'PC-23', 'PC-24', 'PC-25', 'PC-26', 'PC-27', 'PC-28', 'PC-29', 'PC-30', 'PC-31', 'PC-32', 'PC-33', 'PC-34', 'PC-35', 'PC-36', 'PC-37', 'PC-38', 'PC-39', 'PC-40', 'PC-41', 'PC-42', 'PC-43', 'PC-44', 'PC-45', 'PC-46', 'PC-47', 'PC-48', 'PC-49', 'PC-50', 'PC-51', 'PC-52', 'PC-53', 'PC-54', 'PC-55', 'PC-56', 'PC-57', 'PC-58', 'PC-59', 'PC-60', 'PC-61', 'PC-62', 'PC-63', 'PC-64', 'PC-65', 'PC-66']

Out[42]:
kf = KFold(n_splits=10, random_state=0, shuffle=True)
for PC in PC_list:
    for train_index, test_index in kf.split(X):
        X_train, X_test = X.values[train_index], X.values[test_index]
        y_train, y_test = y.values[train_index], y.values[test_index]
        # Training model
        model = RandomForestClassifier()
        model.fit(X_train, y_train)
        # Evaluating model
        y_pred = model.predict(X_test)

        accuracy_scores.append(sm.accuracy_score(y_test, y_pred))
        precision_scores.append(sm.precision_score(y_test, y_pred))
        recall_scores.append(sm.recall_score(y_test, y_pred))
        f1_scores.append(sm.f1_score(y_test, y_pred))

    #displaying average results
    print("#####")
    print("Results")
    print("#####\n")
    print("Average accuracy score: ", (sum(accuracy_scores)/len(accuracy_scores)))

    results = [(sum(accuracy_scores)/len(accuracy_scores))]
    return results
```

```
##Random Forest
```

```
In [44]: import numpy as np # Linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
from sklearn.preprocessing import LabelEncoder
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
import lightgbm as lgb
```

```
In [45]: y=one_hot_df['fraud_reported']
X= principal_onehot_df
```

```
In [46]: columns = ["Algorithm", "Encoding", "Accuracy"]
data = []
```

```
In [47]: from sklearn.ensemble import RandomForestClassifier
model = RandomForestClassifier()
result= ["Random Forest", "OneHot"]
result.extend(evaluate_model(model, X, y))
data.append(result)
```

C:\Users\SHEKHAR\anaconda3\lib\site-packages\sklearn\metrics_classification.py:1318: UndefinedMetricWarning:

Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division` parameter to control this behavior.

C:\Users\SHEKHAR\anaconda3\lib\site-packages\sklearn\metrics_classification.py:1318: UndefinedMetricWarning:

Recall is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division` parameter to control this behavior.

```
In [48]: from sklearn.tree import DecisionTreeClassifier
model = DecisionTreeClassifier()
result= ["Decision tree", "OneHot"]
result.extend(evaluate_model(model, X, y))
data.append(result)
```

C:\Users\SHEKHAR\anaconda3\lib\site-packages\sklearn\metrics_classification.py:1318: UndefinedMetricWarning:

Results
Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division` parameter to control this behavior.

Average accuracy score: 0.641
C:\Users\SHEKHAR\anaconda3\lib\site-packages\sklearn\metrics_classification.py:1318: UndefinedMetricWarning:

```
In [49]: from sklearn.ensemble import XGBClassifier
model = XGBClassifier()
result= ["XGBoost", "OneHot"]
# on WoE data
model.fit(X_train, y_train)
result.extend(evaluate_model(model, X, y))
data.append(result)
```

#####accuracy score: 0.7489999999999999
Results

#####C:\Users\SHEKHAR\anaconda3\lib\site-packages\sklearn\metrics_classification.py:1318: UndefinedMetricWarning:

Average accuracy score: 0.7270000000000001
Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division` parameter to control this behavior.
##AdaBoost

```
In [50]: from sklearn.ensemble import AdaBoostClassifier
```

```

#Decision Tree will-defined and being set to 0.0 due to no predicted samples. Use `zero_division`
parameter to control this behavior.

In [48]: from sklearn.tree import DecisionTreeClassifier
model = DecisionTreeClassifier()
result = ["Decision tree", "OneHot"]
result.extend(evaluate_model(model, X, y))
Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division`
parameter to control this behavior.

C:\Users\SHEKHAR\anaconda3\lib\site-packages\sklearn\metrics\_classification.py:1318: UndefinedMetricWarning:
Results
#####
Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division`
parameter to control this behavior.
Average accuracy score: 0.641

C:\Users\SHEKHAR\anaconda3\lib\site-packages\sklearn\metrics\_classification.py:1318: UndefinedMetricWarning:
#####

```

```

Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division`
parameter to control this behavior.

In [49]: from xgboost import XGBClassifier
result = ["XGBoost", "OneHot"]
# on WoE data
#####XGBClassifier()
result.extend(evaluate_model(model, X, y))
#####append(result)

#####accuracy score: 0.7489999999999999
Results
#####C:\Users\SHEKHAR\anaconda3\lib\site-packages\sklearn\metrics\_classification.py:1318: UndefinedMetricWarning:
Average accuracy score: 0.7270000000000001
Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division`
parameter to control this behavior.
#####AdaBoost

```

```

In [50]: from sklearn.ensemble import AdaBoostClassifier
model = AdaBoostClassifier(random_state=0)
result = ["Adaboost", "OneHot"]
result.extend(evaluate_model(model, X, y))
data.append(result)

#####
Results
#####

Average accuracy score: 0.6990000000000001

```

Logistic Regression

```

In [51]: from sklearn.linear_model import LogisticRegression
model = LogisticRegression()
result = ["LogisticRegression", "OneHot"]
result.extend(evaluate_model(model, X, y))
data.append(result)

#####
Results
#####

Average accuracy score: 0.732

```

SVM

```

In [52]: from sklearn.svm import SVC

model = SVC()
result = ["SVM", "OneHot"]

result.extend(evaluate_model(model, X, y))
data.append(result)

C:\Users\SHEKHAR\anaconda3\lib\site-packages\sklearn\metrics\_classification.py:1318: UndefinedMetricWarning:

Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division`
parameter to control this behavior.

C:\Users\SHEKHAR\anaconda3\lib\site-packages\sklearn\metrics\_classification.py:1318: UndefinedMetricWarning:

Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division`
parameter to control this behavior.

C:\Users\SHEKHAR\anaconda3\lib\site-packages\sklearn\metrics\_classification.py:1318: UndefinedMetricWarning:

Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division`
parameter to control this behavior.

C:\Users\SHEKHAR\anaconda3\lib\site-packages\sklearn\metrics\_classification.py:1318: UndefinedMetricWarning:

Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division`
parameter to control this behavior.

```

SVM

In [52]: `from sklearn.svm import SVC`

```
model = SVC()
result = ["SVM", "OneHot"]

result.extend(evaluate_model(model, X, y))
data.append(result)
```

C:\Users\SHEKHAR\anaconda3\lib\site-packages\sklearn\metrics_classification.py:1318: UndefinedMetricWarning:

Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division` parameter to control this behavior.

C:\Users\SHEKHAR\anaconda3\lib\site-packages\sklearn\metrics_classification.py:1318: UndefinedMetricWarning:

Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division` parameter to control this behavior.

C:\Users\SHEKHAR\anaconda3\lib\site-packages\sklearn\metrics_classification.py:1318: UndefinedMetricWarning:

Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division` parameter to control this behavior.

C:\Users\SHEKHAR\anaconda3\lib\site-packages\sklearn\metrics_classification.py:1318: UndefinedMetricWarning:

Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division` parameter to control this behavior.

C:\Users\SHEKHAR\anaconda3\lib\site-packages\sklearn\metrics_classification.py:1318: UndefinedMetricWarning:

Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division` parameter to control this behavior.

C:\Users\SHEKHAR\anaconda3\lib\site-packages\sklearn\metrics_classification.py:1318: UndefinedMetricWarning:

Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division` parameter to control this behavior.

C:\Users\SHEKHAR\anaconda3\lib\site-packages\sklearn\metrics_classification.py:1318: UndefinedMetricWarning:

Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division` parameter to control this behavior.

C:\Users\SHEKHAR\anaconda3\lib\site-packages\sklearn\metrics_classification.py:1318: UndefinedMetricWarning:

Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division` parameter to control this behavior.

C:\Users\SHEKHAR\anaconda3\lib\site-packages\sklearn\metrics_classification.py:1318: UndefinedMetricWarning:

Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division` parameter to control this behavior.

```
#####
Results
#####
```

Average accuracy score: 0.753

C:\Users\SHEKHAR\anaconda3\lib\site-packages\sklearn\metrics_classification.py:1318: UndefinedMetricWarning:

In [53]: `from sklearn.preprocessing import LabelEncoder`

Requirement already satisfied: numpy in c:\users\shekhar\anaconda3\lib\site-packages (from category_encoders) (1.16.0)

Requirement already satisfied: statsmodels>=0.9.0 in c:\users\shekhar\anaconda3\lib\site-packages (from category_encoders) (0.13.2)

Requirement already satisfied: pandas>=1.0.5 in c:\users\shekhar\anaconda3\lib\site-packages (from category_encoders) (1.4.2)

Requirement already satisfied: scikit-learn>=0.20.0 in c:\users\shekhar\anaconda3\lib\site-packages (from category_encoders) (1.0.2)

Requirement already satisfied: python-dateutil>=2.8.1 in c:\users\shekhar\appdata\roaming\python\python39\site-packages (from pandas>=1.0.5->category_encoders) (2.8.2)

Requirement already satisfied: pytz>=2020.1 in c:\users\shekhar\anaconda3\lib\site-packages (from pandas>=1.0.5->category_encoders) (2021.3)

Requirement already satisfied: scipy>=1.0.0 in c:\users\shekhar\anaconda3\lib\site-packages (from pandas>=1.0.5->category_encoders) (1.8.1)

Requirement already satisfied: python-dateutil>=2.8.1 in c:\users\shekhar\appdata\roaming\python\python39\site-packages (from pandas>=1.0.5->category_encoders) (2.8.2)

Requirement already satisfied: pytz>=2020.1 in c:\users\shekhar\anaconda3\lib\site-packages (from pandas>=1.0.5->category_encoders) (2021.3)

Requirement already satisfied: statsmodels>=0.9.0 in c:\users\shekhar\anaconda3\lib\site-packages (from pandas>=1.0.5->category_encoders) (0.13.2)

Requirement already satisfied: pandas>=1.0.5 in c:\users\shekhar\anaconda3\lib\site-packages (from pandas>=1.0.5->category_encoders) (1.4.2)

Requirement already satisfied: scikit-learn>=0.20.0 in c:\users\shekhar\anaconda3\lib\site-packages (from pandas>=1.0.5->category_encoders) (1.0.2)

Requirement already satisfied: python-dateutil>=2.8.1 in c:\users\shekhar\appdata\roaming\python\python39\site-packages (from pandas>=1.0.5->category_encoders) (2.8.2)

Requirement already satisfied: pytz>=2020.1 in c:\users\shekhar\anaconda3\lib\site-packages (from pandas>=1.0.5->category_encoders) (2021.3)

Requirement already satisfied: statsmodels>=0.9.0 in c:\users\shekhar\anaconda3\lib\site-packages (from pandas>=1.0.5->category_encoders) (0.13.2)

In [54]: `progressed_df.fraud_reported.replace(('Y', 'N'), (1, 0), inplace=True)`


```
In [56]: for col in processed_df.columns:
        if processed_df[col].dtype == 'object':
            print(processed_df[col].value_counts())
```

```
OH      352
IL      338
IN      310
Name: policy_state, dtype: int64
FEMALE  537
MALE    463
Name: insured_sex, dtype: int64
JD       161
High School  160
Associate   145
MD          144
Masters     143
PhD         125
College    122
Name: insured_education_level, dtype: int64
machine-op-inspct  93
prof-specialty    85
tech-support      78
sales             76
exec-managerial   76
craft-repair      74
transport-moving  72
other-service     71
priv-house-serv   71
armed-forces      69
adm-clerical      65
protective-serv   63
handlers-cleaners  54
farming-fishing   53
Name: insured_occupation, dtype: int64
own-child  183
other-relative  177
not-in-family  174
husband     170
wife        155
unmarried   141
Name: insured_relationship, dtype: int64
Multi-vehicle Collision  419
Single Vehicle Collision  403
Vehicle Theft           94
Parked Car              84
Name: incident_type, dtype: int64
Rear Collision  292
Side Collision  276
Front Collision  254
Not Known      178
Name: collision_type, dtype: int64
Minor Damage  354
Total Loss    280
Major Damage  276
Trivial Damage  90
Name: incident_severity, dtype: int64
Police  292
Fire    223
Other   198
Ambulance  196
None       91
Name: authorities_contacted, dtype: int64
NY  262
SC  248
WV  217
VA  110
NC  110
PA  30
OH  23
Name: incident_state, dtype: int64
Not Known  360
NO         338
YES        302
Name: property_damage, dtype: int64
Not Known  343
NO         343
YES        314
Name: policy_state, policy_deductible, policy_annual_premium, umbrella_limit, insured_zip, insured_sex, insured_educ
```

```
In [58]: woe_df.head()
```

```
Out[58]:
```

age	policy_state	policy_deductible	policy_annual_premium	umbrella_limit	insured_zip	insured_sex	insured_educ
80	OH	1000	1406.91	0	466132	0.075555	
48	Dodge	800	1197.22	5000000	468176	0.075555	
80	Subaru	0.044884	2000	1413.14	430632	-0.067571	
78	Nissan	0.062770	2000	1415.74	608117	-0.067571	
78	Chevrolet	-0.102374	2000	6000000	610706	0.075555	
41	Ford	-0.102374	1000	1583.91	610706	0.075555	
72	BMW	-0.102374	1000	1583.91	610706	0.075555	
70	Toyota	-0.102374	1000	1583.91	610706	0.075555	
60	Acura	0.062770	1000	1406.91	466132	0.075555	
68	Volkswagen	0.044884	2000	1413.14	430632	-0.067571	
67	Jeep	0.062770	2000	1415.74	608117	-0.067571	
55	Honda	-0.102374	1000	1583.91	610706	0.075555	

Feature Selection (WOE)

```
In [59]: from sklearn.preprocessing import StandardScaler
```

```
In [57]: # Dropping the target variable from the Dataframe
```

```
scaler = StandardScaler()
from category_encoders.woe import WOEEncoder
```

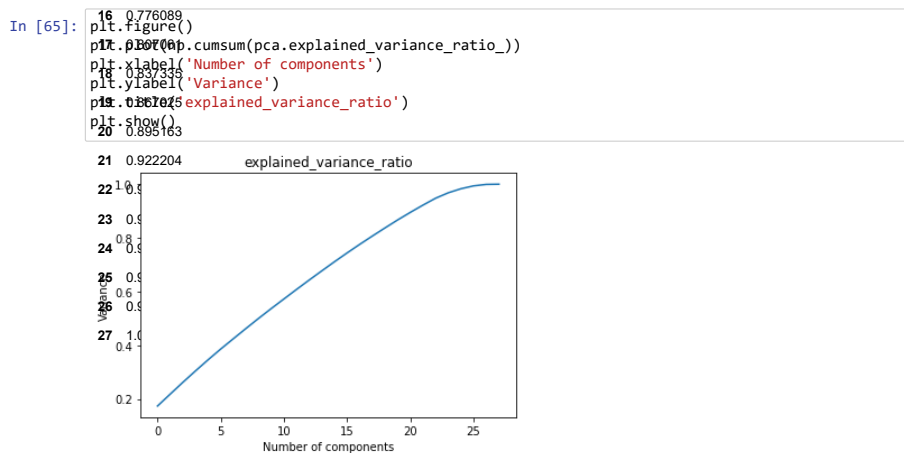
```
In [60]: # Dropping the target column as scaling is not required on it.
```

```
woe_encoder = WOEEncoder()
processed_df.drop('fraud_reported', axis=1, inplace=True)
for column in categorical_cols:
```

```
In [61]: processed_df[column] = woe_encoder.fit_transform(processed_df[column], processed_df['fraud_reported'])
woe_df_scaled_df = pd.DataFrame(scaler.fit_transform(scaled_df), columns=scaled_df.columns)
```

```
woe_df = processed_df.copy()
```

```
In [62]: y = one_not_at_fraud_reported
x = woe_df_scaled_df
```

We can see that around 99% of variance is being explained by 25 components. So instead of giving all columns as input in our algorithm let's use these principle components instead

```
In [66]: pca = PCA(n_components=25)
pca_woe_df = pca.fit_transform(x)

column_name = [f'PC-{i}' for i in range(1, 26)]

principal_woe_df = pd.DataFrame(pca_woe_df, columns=column_name)
principal_woe_df.head()
```

Out[66]:

	PC-1	PC-2	PC-3	PC-4	PC-5	PC-6	PC-7	PC-8	PC-9	PC-10	PC-11
0	-0.997570	0.780926	-0.533583	0.184909	-0.679628	-1.819256	-1.283710	-0.288213	0.726652	-1.200167	0.265421
1	4.026343	0.842083	1.913565	0.716060	-0.067698	1.320849	0.535631	1.558847	0.755646	0.518444	1.281701
2	0.253120	-0.283566	-1.111903	-0.477507	1.440599	0.566111	-1.542508	-1.292452	-0.435717	0.977837	0.053611
3	-0.529702	2.001171	0.996700	-1.318997	1.149643	1.095464	-1.880315	-1.191650	0.908376	0.964282	1.682311
4	4.738184	0.456561	-0.062727	1.323874	2.119771	0.866995	-0.387641	0.888265	0.054812	-0.839057	-0.126781

Data Modelling (WoE)

```
In [67]: import sklearn.metrics as sm
from sklearn.model_selection import KFold

# Function to evaluate model
def evaluate_model(model, x, y):
    accuracy_scores = []
    precision_scores = []
    recall_scores = []
    f1_scores = []

    kf = KFold(n_splits=10, random_state=0, shuffle=True)
    for train_index, test_index in kf.split(x):
        #setting up the data
        x_train, x_test = x.values[train_index], x.values[test_index]
        y_train, y_test = y.values[train_index], y.values[test_index]

        #Training model
        model.fit(x_train,y_train)

        #Evaluating model
        y_pred = model.predict(x_test)
        y_woe_df = y_woe_df.add_reported(y_pred)
        X = principal_woe_df
        accuracy_scores.append(sm.accuracy_score(y_test,y_pred))
        precision_scores.append(sm.precision_score(y_test,y_pred))
        recall_scores.append(sm.recall_score(y_test,y_pred))
        f1_scores.append(sm.f1_score(y_test, y_pred))

    result = ["Random Forest", "WoE"]
    result.extend(evaluate_model(model, X, y))
    data.append(result)
    print("#####")
    #####("Results")
    Result = int("#####\n")
    #####("Average accuracy score: ", (sum(accuracy_scores)/len(accuracy_scores)))
    results = [(sum(accuracy_scores)/len(accuracy_scores))]
    Average accuracy score: 0.7730000000000001
```

Decision Tree

Random Forest

```
In [71]: from sklearn.tree import DecisionTreeClassifier
model = DecisionTreeClassifier()

In [68]: import numpy as np #Linear Algebra
import pandas as pd #Data processing
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.preprocessing import LabelEncoder
import matplotlib.pyplot as plt
import lightgbm as lgb

Average accuracy score: 0.706
```

XGBoost


```

# Evaluating model
In [69]: y=woe_df.y_pred = model.predict(x_test)
X= principal_woe_df
accuracy_scores.append(sm.accuracy_score(y_test,y_pred))
precision_scores.append(sm.precision_score(y_test,y_pred))
In [70]: from sklearn.ensemble import RandomForestClassifier
model = RandomForestClassifier()
result= ["Random Forest", "WoE"]
result.extend(evaluate_model(model, X, y))
data.append(result)
print("#####")
#####
Results
#####
Average accuracy score: 0.7730000000000001

```

Decision Tree

```

# Random Forest
In [71]: from sklearn.tree import DecisionTreeClassifier
model = DecisionTreeClassifier()
In [68]: result= ["Decision Tree", "WoE"]
result.extend(evaluate_model(model,X, y))
from sklearn.preprocessing import LabelEncoder
import matplotlib.pyplot as plt
#####
seaborn as sns
matplotlib inline
#####
lightgbm as lgb
Average accuracy score: 0.706

```

XGBoost

```

In [72]: from xgboost import XGBClassifier
result= ["XGBoost", "WoE"]
# on WoE data
model = XGBClassifier()
result.extend(evaluate_model(model,X, y))
data.append(result)

#####
Results
#####

Average accuracy score: 0.787

```

Adaboost

```

In [73]: from sklearn.ensemble import AdaBoostClassifier
model = AdaBoostClassifier(random_state=0)
result= ["Adaboost", "WoE"]
result.extend(evaluate_model(model, X, y))
data.append(result)

#####
Results
#####

Average accuracy score: 0.781

```

Logistic Regression

```

In [74]: from sklearn.linear_model import LogisticRegression
model = LogisticRegression()
result= ["LogisticRegression", "WoE"]
result.extend(evaluate_model(model, X, y))
data.append(result)

```

KNN

```

#####
Results
In [76]: from sklearn.naive_bayes import GaussianNB
model = GaussianNB()
result= ["GaussianNB", "WoE"]
result.extend(evaluate_model(model, X, y))
data.append(result)

```

SVM

```

#####
Results
In [75]: from sklearn.svm import SVC
model = SVC()
result= ["SVM", "WoE"]
Average accuracy score: 0.806

```

Saving Results for Insurance Fraud Detection

```

In [77]: #####df = pd.DataFrame(data, columns=columns)
Results_df = results_df[['Algorithm', 'Encoding', 'Accuracy']]
#####df.sort_values(by=['Accuracy', 'Encoding'], ascending=False)

```

Out[77]: Average accuracy score: 0.8019999999999999

	Algorithm	Encoding	Accuracy
13	GaussianNB	WoE	0.806
11	LogisticRegression	WoE	0.804
12	SVM	WoE	0.802
9	XGBoost	WoE	0.787
10	Adaboost	WoE	0.781
7	Random Forest	WoE	0.773

```
result.extend(evaluate_model(model, X, y))
data.append(result)
```

KNN

```
#####
```

```
Results
In [76]: from sklearn.naive_bayes import GaussianNB
model = GaussianNB()
result = ["GaussianNB", "WoE"]
Average accuracy score: 0.804
result.extend(evaluate_model(model, X, y))
data.append(result)
```

SVM

```
#####
```

```
Results
In [75]: from sklearn.svm import SVC
model = SVC()
Average accuracy score: 0.806
result = ["SVM", "WoE"]
```

Saving Results for Insurance Fraud Detection

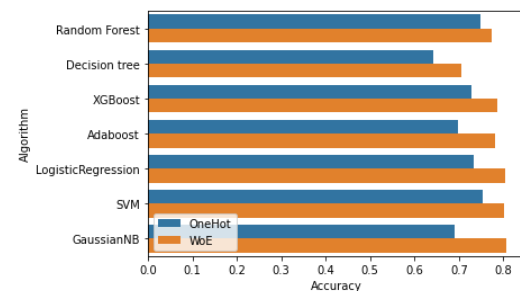
```
#####df = pd.DataFrame(data, columns=columns)
Results_df = results_df[['Algorithm', 'Encoding', 'Accuracy']]
#####df.sort_values(by=['Accuracy', 'Encoding'], ascending=False)
```

```
Out[77]: Average accuracy score: 0.8019999999999999
```

	Algorithm	Encoding	Accuracy
13	GaussianNB	WoE	0.806
11	LogisticRegression	WoE	0.804
12	SVM	WoE	0.802
9	XGBoost	WoE	0.787
10	Adaboost	WoE	0.781
7	Random Forest	WoE	0.773
5	SVM	OneHot	0.753
0	Random Forest	OneHot	0.749
4	LogisticRegression	OneHot	0.732
2	XGBoost	OneHot	0.727
8	Decision tree	WoE	0.706
3	Adaboost	OneHot	0.699
6	GaussianNB	OneHot	0.691
1	Decision tree	OneHot	0.641

```
In [78]: sns.barplot(data=results_df, x='Accuracy', y='Algorithm', hue='Encoding')
plt.legend(loc='lower left')
```

```
Out[78]: <matplotlib.legend.Legend at 0x15e50c899d0>
```



```
In [81]: results_df.to_csv("Insurance_Claims.csv")
```

```
In [85]: # Reading the Dataset
df=pd.read_csv('insurance_data_final.csv' )

# df = df.drop(df.columns[0], axis = 1) # To drop the first index column
```

INCOME PREDICTION

```
In [86]: df.head()
```

```
Out[86]: #Data Collection
Work Marital Status Hours_Week Country Above_Below 50K Education Gender Occupation CustomerID Age months_in_service
0 emp-not-inc divorced 13 United-States <=50K Bachelors Male Handlers-cleaners 5506 38
```

```
In [83]: # Mounting Google Drive
1 Private Divorced 40 United-States <=50K HS-grad Male Handlers-cleaners 5506 38
```

```
In [84]: # Importing required Libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

2 Private Divorced 40 United-States <=50K 11th Male Handlers-cleaners 3869 53
3 Private Divorced 40 United-States <=50K Bachelors Female Prof-specialty 7106 28
4 Private Married-civ-spouse 40 United-States <=50K Masters Female Exec-managerial 1760 37

# To get all rows and columns
pd.set_option('display.max_rows', None)
pd.set_option('display.max_columns', None)
```

Describing some common statistics on numeric columns

```
In [87]: # Data description:
df.describe()
```

```
Out[87]:
```

```
In [86]: df.head()
```

Out[86]: #Data Collection

In [83]:

In [84]:

Describing some Data common statistics on numeric columns

Out[87]:

Checking the null values present in the dataset

Out[88]:

In [90]:

In [89]:

```
total number of null values: 0
# Number of rows and columns of data
#####
Basic Information:
Basic Information:
#####
Number of rows: 9000
Number of columns: 7
Status
data type: object
total number of null values: 0
```


In [94]:

	Work	Marital Status	Hours_Week	Country	Above 50K	Education	Gender	Occupation	Age
0	Self-emp-not-inc	Married-civ-spouse	13	United-States	0	Bachelors	Male	Exec-managerial	50
1	Private	Divorced	40	United-States	0	HS-grad	Male	Handlers-cleaners	38
2	Private	Married-civ-spouse	40	United-States	0	11th	Male	Handlers-cleaners	53
3	Private	Married-civ-spouse	40	Cuba	0	Bachelors	Female	Prof-specialty	28
4	Private	Married-civ-spouse	40	United-States	0	Masters	Female	Exec-managerial	37
5	Private	Married-spouse-absent	16	Jamaica	0	9th	Female	Other-service	49
6	Self-emp-not-inc	Married-civ-spouse	45	United-States	1	HS-grad	Male	Exec-managerial	52

In [96]: `import plotly.graph_objects as go`
`import plotly.express as px`

Reprocessing the Dataframe for Exploratory Data Analysis

In [97]: `eda_df = income_df`

`# Changing the column name of Above_Below 50k and it's values to 1 and 0`
`eda_df.rename(columns = {'Above_Below 50K': 'Above 50K'}, inplace=True)`

`# Replacing the less than equal to values with 0 and greater than equal`
`#to values with 1`

`eda_df['Above 50K'].replace(['<=50K', '<=50K.'], 0, inplace=True)`
`eda_df['Above 50K'].replace(['>50K', '>50K.'], 1, inplace=True)`

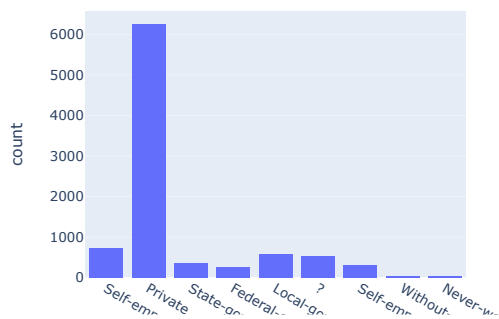
`eda_df.head(10)`

Out[97]:

	Work	Marital Status	Hours_Week	Country	Above 50K	Education	Gender	Occupation	Age
0	Self-emp-not-inc	Married-civ-spouse	13	United-States	0	Bachelors	Male	Exec-managerial	50
1	Private	Divorced	40	United-States	0	HS-grad	Male	Handlers-cleaners	38
2	Private	Married-civ-spouse	40	United-States	0	11th	Male	Handlers-cleaners	53
3	Private	Married-civ-spouse	40	Cuba	0	Bachelors	Female	Prof-specialty	28
4	Private	Married-civ-spouse	40	United-States	0	Masters	Female	Exec-managerial	37
5	Private	Married-spouse-absent	16	Jamaica	0	9th	Female	Other-service	49
6	Self-emp-not-inc	Married-civ-spouse	45	United-States	1	HS-grad	Male	Exec-managerial	52
7	Private	Never-married	50	United-States	1	Masters	Female	Prof-specialty	31
8	Private	Married-civ-spouse	40	United-States	1	Bachelors	Male	Exec-managerial	42
9	Private	Married-civ-spouse	80	United-States	1	Some-college	Male	Exec-managerial	37

Data Distribution using Histogram

In [98]: `def count_plot(df):`
 `for col in df.columns:`
 `fig = px.histogram(df, col, nbins=20, width=500, height=400)`
 `fig.show()`
`count_plot(eda_df)`



Box plot for outliers

In [99]: `# Correlation heatmap for all numeric values`
`corr = eda_df.corr()`
`# We can see from the histogram that there are not much significant outliers present in any of the columns and`
`# hence outliers removal can be skipped=True)`

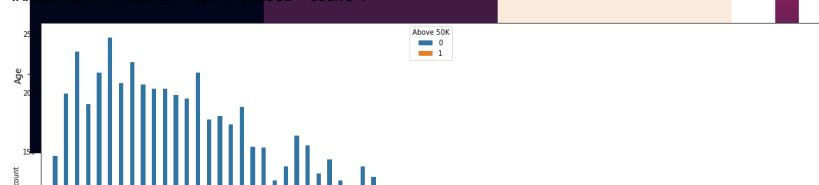
Out[99]: <AxesSubplot: >

Hypothesis Based on EDA

In [100]: `# Is the age of person have any prominent relationship with insurance fraud ?`
`from matplotlib import pyplot as plt`
`import seaborn as sns`

`plt.figure(figsize=(10, 4))`
`sns.countplot(data=eda_df, x='Age', hue='Above 50K')`

Out[100]: <AxesSubplot: xlabel='Age', ylabel='count'>



Correlation

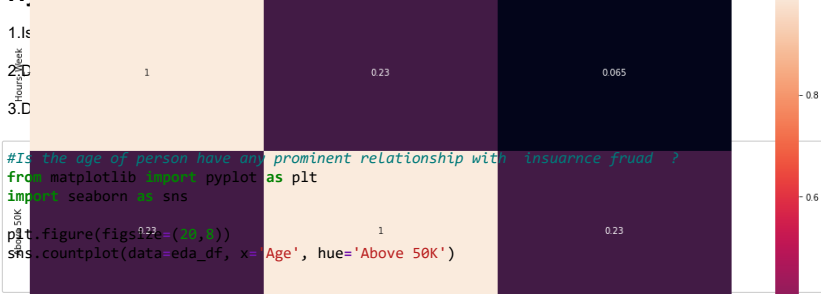
Box plot for outliers

```
In [99]: # Correlation heatmap for all numeric values
```

```
As we can see from the histogram that there are not much significant outliers present in any of the columns and hence outliers removal can be skipped=True)
```

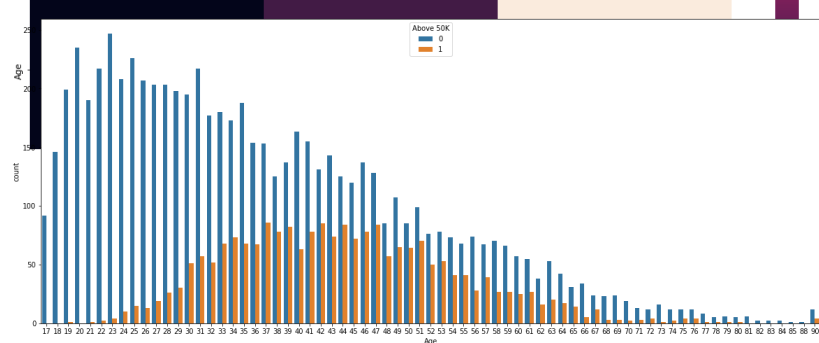
```
Out[99]: <AxesSubplot:>
```

Hypothesis Based on EDA



```
In [100]: #Is the age of person have any prominent relationship with insurance fraud ?
from matplotlib import pyplot as plt
import seaborn as sns
plt.figure(figsize=(20,8))
sns.countplot(data_eda_df, x= 'Age', hue='Above 50K')
```

```
Out[100]: <AxesSubplot:xlabel='Age', ylabel='count'>
```

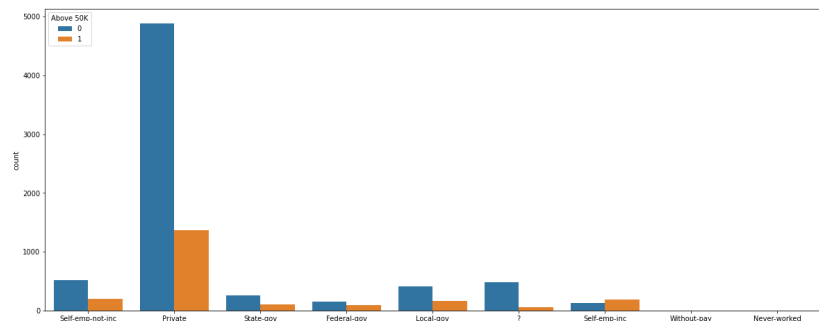


```
In [101]: # Does WorkClass affects the income range
```

```
from matplotlib import pyplot as plt
import seaborn as sns

plt.figure(figsize=(20,8))
sns.countplot(data=eda_df, x='Work', hue='Above 50K')
```

```
Out[101]: <AxesSubplot:xlabel='Work', ylabel='count'>
```



```
In [102]: # Does Education affects the income range
```

```
# Creating copy of eda_df to eda_data frame into new data frame for data processing
from matplotlib import pyplot as plt
import seaborn as sns
```

```
In [103]: processed_df=eda_df.copy()
plt.figure(figsize=(20,8))
sns.countplot(data=eda_df, x='Occupation', hue='Above 50K')
```

```
Out[102]: <AxesSubplot:xlabel='Occupation', ylabel='count'>
```

Handling rows

```
In [104]: for col in processed_df.columns:
if processed_df[col].dtype == 'object':
print(col, processed_df[col].unique())

Work ['Self-emp-not-inc' 'Private' 'State-gov' 'Federal-gov' 'Local-gov' '?']
Marital Status ['Married-civ-spouse' 'Never-married' 'Married-spouse-absent' 'Divorced' 'Separated' 'Married-AF-spouse' 'Widowed']
Country ['United-States' 'Cuba' 'Jamaica' 'India' '?' 'Mexico' 'South' 'Puerto-Rico' 'Honduras' 'England' 'Canada' 'Germany' 'Iran' 'Philippines' 'Italy' 'Poland' 'Columbia' 'Cambodia' 'Thailand' 'Ecuador' 'Laos' 'Taiwan' 'Haiti' 'Portugal' 'Dominican-Republic' 'El-Salvador' 'France' 'Guatemala' 'China' 'Japan' 'Yugoslavia' 'Peru' 'Outlying-US(Guam-USVI-etc)' 'Scotland' 'Trinidad&Tobago' 'Greece' 'Nicaragua' 'Vietnam' 'Hong' 'Ireland' 'Hungary']
Education ['Bachelors' 'HS-grad' 'Masters' '11th' 'Some-college' 'Assoc-acdm' 'Assoc-voc' '7th-8th' 'Doctorate' 'Prof-school' '5th-6th' '10th' '1st-4th' 'Preschool' '12th']
Occupation ['Exec-managerial' 'Handlers-cleaners' 'Prof-specialty' 'Other-service' 'Adm-clerical' 'Sales' 'Craft-repair' 'Transport-moving' 'Farming-fishing' 'Machine-op-inspct' 'Tech-support' '?' 'Protective-serv' 'Armed-Forces' 'Priv-house-serv']
```


In [102]: **Handling columns** *Drop the Education column as the income range*

```
# Creating copy of eda df to new data frame for data processing
import seaborn as sns
```

In [103]: `processed_df = eda_df.copy()`
`plt.figure(figsize=(20,8))`
`sns.countplot(data=eda_df, x='Occupation', hue='Above 50K')`

Out[102]: **Handling rows** `processed_df['Occupation', ylabel='count']`

In [104]: `for col in processed_df.columns:`
`if processed_df[col].dtype == 'object':`
`print(col, processed_df[col].unique())`

Work ['Self-emp-not-inc' 'Private' 'State-gov' 'Federal-gov' 'Local-gov' '?']
 Self-emp-inc ['Without-pay' 'Never-worked']
 Marital Status ['Married-civ-spouse' 'Divorced' 'Married-spouse-absent' 'Never-married'
 'Separated' 'Married-AF-spouse' 'Widowed']
 Country ['United-States' 'Cuba' 'Jamaica' 'India' '?' 'Mexico' 'South'
 'Puerto-Rico' 'Honduras' 'England' 'Canada' 'Germany' 'Iran'
 'Philippines' 'Italy' 'Poland' 'Columbia' 'Cambodia' 'Thailand' 'Ecuador'
 'Laos' 'Taiwan' 'Haiti' 'Portugal' 'Dominican-Republic' 'El-Salvador'
 'France' 'Guatemala' 'China' 'Japan' 'Yugoslavia' 'Peru'
 'Outlying-US(Guam-USVI-etc)' 'Scotland' 'Trinidad&Tobago' 'Greece'
 'Nicaragua' 'Vietnam' 'Hong' 'Ireland' 'Hungary']
 Education ['Bachelors' 'HS-grad' '11th' 'Masters' '9th' 'Some-college' 'Assoc-acdm'
 'Assoc-voc' '7th-8th' 'Doctorate' 'Prof-school' '5th-6th' '10th'
 '1st-4th' 'Preschool' '12th']
 # Data Pre-Processing (One-hot Encoding)
 Occupation ['Exec-managerial' 'Handlers-cleaners' 'Prof-specialty' 'Other-service'
 'Adm-clerical' 'Sales' 'Craft-repair' 'Transport-moving'
 'Farming-fishing' 'Machine-op-inspct' 'Tech-support' '?'
 'Protective-serv' 'Armed-Forces' 'Priv-house-serv']

1. Dropping the Education column as EducationNum column is provided.
2. Adding new category as 'Not Known' for the '?' values in Work , Country , and Occupation

In [105]: *# Drop the Education columns*

```
processed_df.drop('Education', axis=1, inplace=True)

# Work, Country, and Occupation column have '?' values. We will introduce the new category for t
df['Work'].replace('?', 'Not Known', inplace=True)
df['Country'].replace('?', 'Not Known', inplace=True)
df['Occupation'].replace('?', 'Not Known', inplace=True)

# Applying One Hot Encoding to the Dataframe
one_hot_df = pd.get_dummies(processed_df, drop_first=True)

one_hot_df.head(5)
```

Out[105]:

	Hours_Week	Above 50K	Age	Work_Federal-gov	Work_Local-gov	Work_Never-worked	Work_Private	Work_Self-emp-inc	Work_Self-emp-not-inc	Work_?
0	13	0	50	0	0	0	0	0	1	
1	40	0	38	0	0	0	1	0	0	
2	40	0	53	0	0	0	1	0	0	
3	40	0	28	0	0	0	1	0	0	
4	40	0	37	0	0	0	1	0	0	

In [106]: *# processed_df.to_csv('drive/MyDrive/Inventory Project/processed_df.csv')*

In [112]: *# Variance values for all the principle components*
Feature Selection (One-hot) `processed_df.to_csv('drive/MyDrive/Inventory Project/processed_df.csv')`

Out[112]: Due to one hot, we got large set of columns. Before PCA we need to standardize and scale all the columns using Standard Scaler.

In [107]: `from sklearn.preprocessing import StandardScaler`
Using Standard Scaler Class
`scaler = StandardScaler()`

In [108]: *# Dropping the target columns as scaling don't require it*
`scaled_df = one_hot_df.drop('Above 50K', axis=1)`

In [109]: `one_hot_scaled_df = pd.DataFrame(scaler.fit_transform(scaled_df), columns=scaled_df.columns)`

In [110]: *# Before applying feature selection Let's just save the Dependent column*
`y10 = one_hot_scaled_df['Above 50K']`
`x = one_hot_scaled_df`

In [113]: `plt.figure()`
In [111]: `from sklearn.decomposition import PCA`
`plt.xlabel('Number of components')`
`plt.ylabel('Variance')`
`plt.title('explained_variance_ratio')`
`plt.show()`
`principle_components_wef = pca.fit_transform(x)`



Feature Selection (One Hot)

```
In [112]: # Variance values for all the principle components
pca = PCA(n_components=70, explained_variance_ratio=0.99)

Out[112]: Due to one hot, we got large set of columns. Before PCA we need to standardize and scale all the columns using Standard Scaler.

In [107]: from sklearn.preprocessing import StandardScaler

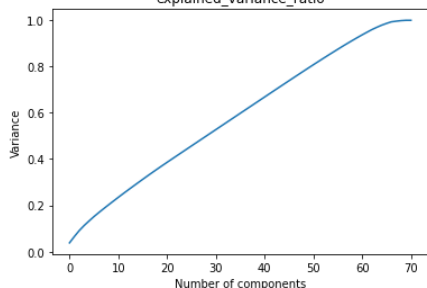
# Using Standard Scaler Class
scaler = StandardScaler()

In [108]: # Dropping the target columns as scaling don't require it
scaled_df = one_hot_df.drop('Above 50K', axis=1)

In [109]: one_hot_scaled_df = pd.DataFrame(scaler.fit_transform(scaled_df), columns=scaled_df.columns)

In [110]: # Before applying feature selection Let's just save the Dependent column
y10 = one_hot_scaled_df['Above 50K']
x = one_hot_scaled_df

In [113]: plt.figure()
In [111]: # Complete PCA decomposition explained variance_ratio))
plt.xlabel('Number of components')
plt.ylabel('Variance')
plt.title('explained_variance_ratio')
plt.show()
principle_components_wef = pca.fit_transform(x)
```



We can see that around 99% of variance is being explained by 66 components. So instead of giving all columns as input in our algorithm let's use these principle components instead

```
In [114]: pca = PCA(n_components=66)
pca_onehot_df = pca.fit_transform(x)

column_name = [f'PC-{i}' for i in range(1, 67)]

principal_onehot_df = pd.DataFrame(pca_onehot_df, columns=column_name)
principal_onehot_df.head()
```

```
Out[114]:
```

	PC-1	PC-2	PC-3	PC-4	PC-5	PC-6	PC-7	PC-8	PC-9	PC-10	PC-11
0	2.305275	-0.862720	0.473330	-0.061877	1.888662	-0.025479	0.124982	0.001428	-1.134315	-0.176133	1.247617
1	-0.599388	0.670837	-1.293432	-0.333856	-0.519609	-0.589518	-0.193658	-0.892236	0.153488	0.218329	0.316471
2	0.703507	0.908241	-1.343380	0.301603	-0.773007	-0.768850	-0.164568	-0.743841	0.033581	0.216421	0.216561
3	-0.112300	2.812372	3.201290	1.246660	-0.947925	0.631300	-0.809534	-0.447217	-0.972917	2.486118	-1.994861
4	0.473289	-0.030681	-0.613164	1.221729	-0.146936	1.434805	-0.103488	-0.387173	-0.433728	-0.303914	1.072821

Data Modeling (One Hot)

```
In [115]: import sklearn.metrics as sm
from sklearn.model_selection import KFold

# Function to evaluate model
def evaluate_model(model, x, y):
    accuracy_scores = []
    precision_scores = []
    recall_scores = []
    f1_scores = []

    kf = KFold(n_splits=10, random_state=0, shuffle=True)
    for train_index, test_index in kf.split(x):
        #setting up the data
        x_train, x_test = x.values[train_index], x.values[test_index]
        y_train, y_test = y.values[train_index], y.values[test_index]

        #Training model
        model.fit(x_train, y_train)

        #Evaluating model
        y_pred = model.predict(x_test)

        accuracy_scores.append(sm.accuracy_score(y_test, y_pred))
        precision_scores.append(sm.precision_score(y_test, y_pred))
        recall_scores.append(sm.recall_score(y_test, y_pred))
        f1_scores.append(sm.f1_score(y_test, y_pred))

    #displaying average results
    print("#####")
    print("Results")
    print("#####\n")
```

Data Modeling (One Hot)

```
In [115]: import sklearn.metrics as sm
from sklearn.model_selection import KFold

# Function to evaluate model
def evaluate_model(model, x, y):
    accuracy_scores = []
    precision_scores = []
    recall_scores = []
    f1_scores = []

    kf = KFold(n_splits=10, random_state=0, shuffle=True)
    for train_index, test_index in kf.split(x):
        #setting up the data
        x_train, x_test = x.values[train_index], x.values[test_index]
        y_train, y_test = y.values[train_index], y.values[test_index]

        #Training model
        model.fit(x_train,y_train)

        #Evaluating model
        y_pred = model.predict(x_test)

        accuracy_scores.append(sm.accuracy_score(y_test,y_pred))
        precision_scores.append(sm.precision_score(y_test,y_pred))
        recall_scores.append(sm.recall_score(y_test,y_pred))
        f1_scores.append(sm.f1_score(y_test, y_pred))

    #displaying average results
    print("#####")
    print("Results")
    print("#####\n")
    print("Average accuracy score: ", (sum(accuracy_scores)/len(accuracy_scores)))
    results = [(sum(accuracy_scores)/len(accuracy_scores))]
    return results
```

###Random Forest

```
In [116]: import numpy as np # Linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
from sklearn.preprocessing import LabelEncoder
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
import lightgbm as lgb
```

```
In [117]: y=one_hot_df['Above 50K']
X= principal_onehot_df
```

```
In [118]: columns = ["Algorithm", "Encoding", "Accuracy"]
data = []
```

```
In [119]: from sklearn.ensemble import RandomForestClassifier
model = RandomForestClassifier()
result= ["Random Forest", "OneHot"]
result.extend(evaluate_model(model, X, y))
data.append(result)
```

```
#####
Results
#####
```

Average accuracy score: 0.7946666666666666

###XGBoost

```
In [120]: from sklearn.ensemble import XGBClassifier
model = XGBClassifier()
result= ["XGBoost", "OneHot"]
result= ["DecisionTree", "OneHot"]
model = DecisionTreeClassifier()
result= XGBClassifier()
result.extend(evaluate_model(model, X, y))
data.append(result)
```

```
#####
Results
#####
```

Average accuracy score: 0.9029999999999999

###AdaBoost

SVM

```
In [123]: from sklearn.ensemble import AdaBoostClassifier
In [121]: model = AdaBoostClassifier(random_state=0)
result= ["AdaBoost", "OneHot"]
model.extend(evaluate_model(model, X, y))
data.append(result)
```

```
#####
Results
#####
```

Average accuracy score: 0.8112222222222222

Logistic Regression: 0.8187777777777777

```
In [124]: from sklearn.linear_model import LogisticRegression
```

```
###XGBost Tree
```

```
In [120]: from sklearn.ensemble import XGBClassifier
model = XGBClassifier()
result = ["Decision tree", "OneHot"]
result.extend(evaluate_model(model, X, y))
data.append(result)
```

```
#####
Results
#####
```

```
Average accuracy score: 0.9079999999999999
```

```
###Adaboost
```

```
SVM
```

```
In [123]: from sklearn.ensemble import AdaBoostClassifier
In [121]: model = AdaBoostClassifier(random_state=0)
result = ["Adaboost", "OneHot"]
model.score(X_train, y_train)
data.append(result)
```

```
#####
data.append(result)
#####
```

```
Average accuracy score: 0.8112222222222222
```

```
Average accuracy score: 0.8187777777777777
```

```
Logistic Regression
```

```
In [124]: from sklearn.linear_model import LogisticRegression
model = LogisticRegression()
result = ["LogisticRegression", "OneHot"]
result.extend(evaluate_model(model, X, y))
data.append(result)
```

```
C:\Users\SHEKHAR\anaconda3\lib\site-packages\sklearn\linear_model\_logistic.py:814: ConvergenceWarning:
```

```
lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
```

Increase the number of iterations (max_iter) or scale the data as shown in:
<https://scikit-learn.org/stable/modules/preprocessing.html> (<https://scikit-learn.org/stable/modules/preprocessing.html>)
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression (https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)

```
C:\Users\SHEKHAR\anaconda3\lib\site-packages\sklearn\linear_model\_logistic.py:814: ConvergenceWarning:
```

```
lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
```

Increase the number of iterations (max_iter) or scale the data as shown in:
<https://scikit-learn.org/stable/modules/preprocessing.html> (<https://scikit-learn.org/stable/modules/preprocessing.html>)
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression (https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)

```
C:\Users\SHEKHAR\anaconda3\lib\site-packages\sklearn\linear_model\_logistic.py:814: ConvergenceWarning:
```

```
lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
```

Increase the number of iterations (max_iter) or scale the data as shown in:
<https://scikit-learn.org/stable/modules/preprocessing.html> (<https://scikit-learn.org/stable/modules/preprocessing.html>)

```
In [125]: from sklearn.naive_bayes import GaussianNB
model = GaussianNB()
result = ["GaussianNB", "OneHot"]
result.extend(evaluate_model(model, X, y))
data.append(result)
```

```
#####
lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
#####
```

Increase the number of iterations (max_iter) or scale the data as shown in:
<https://scikit-learn.org/stable/modules/preprocessing.html> (<https://scikit-learn.org/stable/modules/preprocessing.html>)
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression (https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)

Handling columns

```
#####
Results
Creating copy of eda_df i.e eda data frame into new data frame for data processing
#####
```

```
In [126]: Average accuracy score: 0.8195555555555556
```

Handling rows

Convert Categorical columns into numeric using Weight Of Evidence encoding

```
In [127]: ! pip install category_encoders
```

```
Requirement already satisfied: category_encoders in c:\users\shakhar\anaconda3\lib\site-packages
```

KNN <https://scikit-learn.org/stable/modules/preprocessing.html> (<https://scikit-learn.org/stable/modules/preprocessing.html>)

```
In [125]: from sklearn.naive_bayes import GaussianNB
model = GaussianNB()
result = ["GaussianNB", "OneHot"]
CesUserException: C:\Users\user\AppData\Local\Programs\Python\Python38-32\Scripts\sklearn\linear_model.py:814: ConvergenceWarning:


warn(msg)


databend.append(result)
```

Increase the number of iterations (`max_iter`) or scale the data as shown in:
Average accuracy score: 0.3347777777777778/[preprocessing.html](https://scikit-learn.org/stable/modules/preprocessing.html) (<https://scikit-learn.org/stable/modules/preprocessing.html>)

#Data Preprocessing (WOE Encoding)
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression ([https://sci](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)
[kit-learn.org/stable/modules/linear_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression))

#####

```
In [126]: averaged_df = df.copy()
```

Convert Categorical columns into numeric using Weight Of Evidence encoding

```
In [127]: ! pip install category_encoders
```

```
In [128]: # Drop the Education columns
```

```
# Work, Country, and Occupation column have '?' values. We will introduce the new category for them
df['Work'].replace('?', 'Not Known', inplace=True)
df['Country'].replace('?', 'Not Known', inplace=True)
df['Occupation'].replace('?', 'Not Known', inplace=True)
from sklearn.preprocessing import StandardScaler
```

```
from category_encoders.woe import WOEEncoder
```

```
categorical_cols = ['Work', 'Marital Status', 'Occupation', 'Gender', 'Country']
```

```
In [135]: We'll defer processing of copy() selection let's just save the Dependent column
```

```
In [130]: wef_df.head()
```

Out[130]:

```
In [137]: pd.DataFrame(np.cumsum(pca.explained_variance_ratio_))
```

100

1 0.447477

Feature

3 0.724503

4 0.829041

5 0.930157

Feature Selection (WOE)

```

In [132]: df['Occupation'].replace('?', 'Not Known', inplace=True)
from sklearn.preprocessing import StandardScaler

In [129]: scaler = StandardScaler()
# Apply the standard scaling to the Dataframe
from category_encoders.woe import WOEEncoder

In [133]: WOE_encoder = WOEEncoder()
scaled_df = woe_df.drop('Above 50K', axis=1)

categorical_cols = ['Work', 'Marital Status', 'Occupation', 'Gender', 'Country']

In [134]: woe_df = pd.DataFrame(scaler.fit_transform(scaled_df), columns=scaled_df.columns)
processed_df[column] = WOE_encoder.fit_transform(processed_df[column], processed_df['Above 50K'])

In [135]: # We define processed_df as a copy of the collection Let's just save the Dependent column
y = one_hot_df['Above 50K']
x = woe_df_scaled_df

In [130]: woe_df.head()
In [136]: from sklearn.decomposition import PCA
Out[130]:
pca = PCA(n_components=5)
# Get the principle components for woe dataframe
principal_components = pca.fit_transform(x)
pd.DataFrame(np.cumsum(pca.explained_variance_ratio_))

In [137]:
Out[137]:
0 0.447477
1 0.596197
2 0.724505
3 0.829041
4 0.930157
5 1.000000

In [131]: # Save the processed_df to csv('drive/MyDrive/Inventory Project/woe_processed_df.csv')

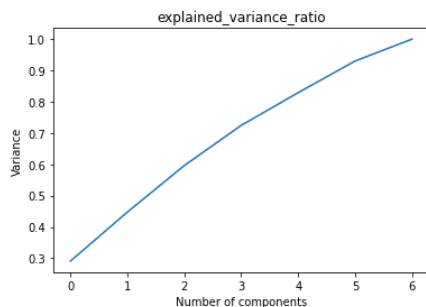
```

Feature Selection (WOE)

```

In [138]: plt.figure()
plt.plot(np.cumsum(pca.explained_variance_ratio_))
plt.xlabel('Number of components')
plt.ylabel('Variance')
plt.title('explained_variance_ratio')
plt.show()

```



We can see that around 95% of variance is being explained by 5 components. So instead of giving all columns as input in our algorithm let's use these principle components instead

```

In [139]: pca = PCA(n_components=5)
pca_woe_df = pca.fit_transform(x)

column_name = [f'PC-{i}' for i in range(1, 6)]

principal_woe_df = pd.DataFrame(pca_woe_df, columns=column_name)

In [140]: principal_woe_df.head()
Out[139]:
# Function to evaluate model
def evaluate_model(model, x, y):
    accuracy_scores = []
    precision_scores = []
    recall_scores = []
    f1_scores = []

    kf = KFold(n_splits=10, random_state=0, shuffle=True)
    for train_index, test_index in kf.split(x):
        #setting up the data
        x_train, x_test = x.values[train_index], x.values[test_index]
        y_train, y_test = y.values[train_index], y.values[test_index]

        ##Data Modeling

        #Training model
        model.fit(x_train, y_train)

        #Evaluating model
        y_pred = model.predict(x_test)

        accuracy_scores.append(sm.accuracy_score(y_test, y_pred))
        precision_scores.append(sm.precision_score(y_test, y_pred))
        recall_scores.append(sm.recall_score(y_test, y_pred))
        f1_scores.append(sm.f1_score(y_test, y_pred))

    #displaying average results
    print("#####")
    print("Results")
    print("#####\n")
    print("Average accuracy score: ", (sum(accuracy_scores)/len(accuracy_scores)))

    results = [(sum(accuracy_scores)/len(accuracy_scores))]
    return results

```

```

In [140]: principal_woe_df = pd.DataFrame(pca_woe_df, columns=column_name)
principal_woe_df.head()
Out[139]:
# Function to evaluate model
def evaluate_model(model, X, y):
    accuracy_scores = []
    precision_scores = []
    recall_scores = []
    f1_scores = []
    kf = KFold(n_splits=10, random_state=0, shuffle=True)
    for train_index, test_index in kf.split(X):
        #setting up the data
        x_train, x_test = X.values[train_index], X.values[test_index]
        y_train, y_test = y.values[train_index], y.values[test_index]

        #Training model
        model.fit(x_train,y_train)

        #Evaluating model
        y_pred = model.predict(x_test)

        accuracy_scores.append(sm.accuracy_score(y_test,y_pred))
        precision_scores.append(sm.precision_score(y_test,y_pred))
        recall_scores.append(sm.recall_score(y_test,y_pred))
        f1_scores.append(sm.f1_score(y_test, y_pred))

    #displaying average results
    print("#####")
    print("Results")
    print("#####\n")
    print("Average accuracy score: ", (sum(accuracy_scores)/len(accuracy_scores)))

    results = [(sum(accuracy_scores)/len(accuracy_scores))]
    return results

```

###Random Forest

```

In [141]: import numpy as np # Linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
from sklearn.preprocessing import LabelEncoder
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
import lightgbm as lgb

```

```

In [142]: y=wef_df['Above 50K']
X= principal_woe_df

```

```

In [143]: from sklearn.ensemble import RandomForestClassifier
model = RandomForestClassifier()
result= ["Random Forest", "WOE"]
result.extend(evaluate_model(model, X, y))
data.append(result)

```


Results
#####

Average accuracy score: 0.7963333333333333

###Decision Tree

```

In [144]: from sklearn.tree import DecisionTreeClassifier
model = DecisionTreeClassifier()
result= ["Decision tree", "WOE"]

```

```

In [146]: result.extend(evaluate_model(model, X, y))
from sklearn.ensemble import XGBClassifier
result.extend(evaluate_model(model, X, y))
data.append(result)

```

```

model = XGBClassifier()
result.extend(evaluate_model(model,X, y))
data.append(result)
#####
Results
#####
Average accuracy score: 0.7628888888888888

```

Average accuracy score: 0.7971111111111111
SVM

###Adaboost

```

In [145]: from sklearn.svm import SVC

```

```

In [147]: from sklearn.ensemble import AdaBoostClassifier
model = AdaBoostClassifier(random_state=0)
result= ["Adaboost", "WOE"]
result.extend(evaluate_model(model, X, y))
data.append(result)

```


Results
#####

Average accuracy score: 0.7981111111111111

###XGBoost
Logistic Regression

```

In [148]: from sklearn.linear_model import LogisticRegression
model = LogisticRegression()
result= ["LogisticRegression", "WOE"]
result.extend(evaluate_model(model, X, y))

```



```

result = [ "Decision tree", "WOE" ]
from sklearn.metrics import accuracy_score
result.extend(evaluate_model(model, X, y))
model = SVC()
result = [ "SVM", "WOE" ]
result.extend(evaluate_model(model, X, y))
data.append(result)

model = XGBClassifier()
result.extend(evaluate_model(model, X, y))
data.append(result)
Results
#####
Average accuracy score: 0.7628888888888888

```

SVM

```

###AdaBoost
from sklearn.svm import SVC

In [147]: from sklearn.ensemble import AdaBoostClassifier
model = AdaBoostClassifier(random_state=0)
result = [ "AdaBoost", "WOE" ]
result.extend(evaluate_model(model, X, y))
data.append(result)

#####
Results
#####

Average accuracy score: 0.7911111111111111

```

Logistic Regression

```

In [148]: from sklearn.linear_model import LogisticRegression
model = LogisticRegression()
result = [ "LogisticRegression", "WOE" ]
result.extend(evaluate_model(model, X, y))
data.append(result)

#####
Results
#####

Average accuracy score: 0.7936666666666667

```

KNN

```

In [149]: from sklearn.naive_bayes import GaussianNB
model = GaussianNB()
result = [ "GaussianNB", "WOE" ]
result.extend(evaluate_model(model, X, y))
data.append(result)

#####
Results
#####

Average accuracy score: 0.7947777777777778

```

Saving Results for Income Prediction

```

In [150]: results_df = pd.DataFrame(data, columns=columns)
results_df = results_df[['Algorithm', 'Encoding', 'Accuracy']]
results_df.sort_values(by=['Accuracy', 'Encoding'], ascending=False)

```

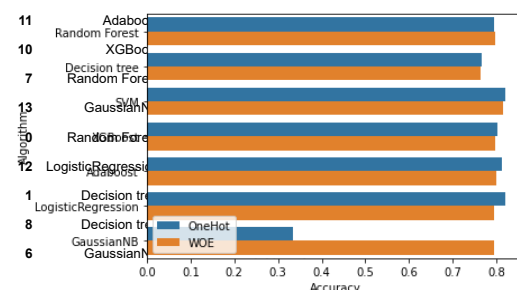
```

Out[150]:
   Algorithm Encoding Accuracy
5  LogisticRegression  OneHot  0.819556
2             SVM      OneHot  0.818778
4  Adaboost      OneHot  0.811222
3  XGBoost      OneHot  0.803000

In [151]: sns.barplot(data=results_df, x='Accuracy', y='Algorithm', hue='Encoding')
plt.legend(loc='lower left')

Out[151]: <matplotlib.legend.Legend at 0x15e5742b0a0>

```



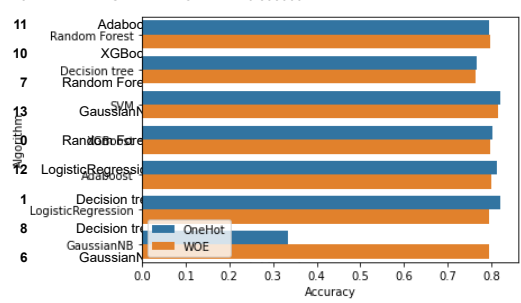
```

In [152]: results_df.to_csv("income_prediction_results.csv")

```

```
In [151]: sns.barplot(data=results_df, x='Accuracy', y='Algorithm', hue='Encoding')
plt.legend(loc='lower left')
```

```
Out[151]: <matplotlib.legend.Legend at 0x15e5742b0a0>
```



A horizontal bar chart comparing the accuracy of various machine learning algorithms using two different encoding methods: OneHot (blue bars) and WOE (orange bars). The x-axis represents Accuracy, ranging from 0.0 to 0.8. The y-axis lists the algorithms. The chart shows that for most algorithms, the OneHot encoding performs slightly better than the WOE encoding, with accuracies generally ranging between 0.75 and 0.80. The legend is located in the lower-left corner of the plot area.

Algorithm	OneHot Accuracy	WOE Accuracy
Adaboost	0.811222	0.803000
Random Forest	0.79	0.78
XGBoost	0.78	0.77
Decision tree	0.77	0.76
Random Forest	0.78	0.77
SVM	0.78	0.77
GaussianNB	0.78	0.77
Random Forest	0.78	0.77
Logistic Regression	0.78	0.77
Adaboost	0.78	0.77
Decision tree	0.78	0.77
Logistic Regression	0.78	0.77
Decision tree	0.78	0.77
GaussianNB	0.78	0.77
GaussianNB	0.78	0.77

```
In [152]: results_df.to_csv("income_prediction_results.csv")
```