

▼ Problem Statement

To identify potential defaulters from a pool of new customers based on learning based tree metho

Dataset

Source of Data

We have a dataset of 1319 records of previous credit card applicants to the bank.

Goal: Predict whether a credit card application will be accepted based upon various data about the

▼ Importing necessary packages

```
1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 import sklearn
6 from sklearn.preprocessing import StandardScaler
7 from sklearn.tree import DecisionTreeClassifier
8 %matplotlib inline
```

```
↳ /usr/local/lib/python3.6/dist-packages/statsmodels/tools/_testing.py:19: FutureWarning:
    import pandas.util.testing as tm
```

```
1 print(np.__version__)
2 print(pd.__version__)
3 print(sns.__version__)
4 print(sklearn.__version__)
```

```
↳ 1.18.4
   1.0.3
   0.10.1
   0.22.2.post1
```

▼ load the dataset from the csv file using pandas

```
1 !wget https://gist.githubusercontent.com/ucalyptus/905b6e56e5db73b2953624be72
2 data=pd.read_csv("/content/dataset.csv")
3 data.shape
```

```
↳
```

```
--2020-05-24 14:30:29-- https://gist.githubusercontent.com/ucalyptus/905b6e56
Resolving gist.githubusercontent.com (gist.githubusercontent.com)... 151.101.0
Connecting to gist.githubusercontent.com (gist.githubusercontent.com)|151.101.
HTTP request sent, awaiting response... 200 OK
Length: 73250 (72K) [text/plain]
Saving to: 'dataset.csv'
```

▼ displaying the data

```
2020-05-24 14:30:30 (2.50 MB/s) (dataset.csv) saved [73250/73250]
```

```
1 data.head()
```

	card	reports	age	income	share	expenditure	owner	selfemp	depende
0	yes	0	37.66667	4.5200	0.033270	124.983300	yes	no	
1	yes	0	33.25000	2.4200	0.005217	9.854167	no	no	
2	yes	0	33.66667	4.5000	0.004156	15.000000	yes	no	
3	yes	0	30.50000	2.5400	0.065214	137.869200	no	no	
4	yes	0	32.16667	9.7867	0.067051	546.503300	yes	no	

▼ Exploratory Data Analysis

▼ Column Description

- card: Dummy variable, 1 if application for credit card accepted, 0 if not
- reports: Number of major derogatory reports
- age: Age n years plus twelfths of a year
- income: Yearly income (divided by 10,000)
- share: Ratio of monthly credit card expenditure to yearly income
- expenditure: Average monthly credit card expenditure
- owner: 1 if owns their home, 0 if rent
- selfempl: 1 if self employed, 0 if not.
- dependents: 1 + number of dependents
- months: Months living at current address
- majorcards: Number of major credit cards held
- active: Number of active credit accounts

▼ What are the attributes of our dataset?

```
1 # printing the column names
2 print(data.columns)
```



```
Index(['card', 'reports', 'age', 'income', 'share', 'expenditure', 'owner',
```

```
1 # displaying the dimension of dataset
2 print(data.shape)
```

```
(1319, 12)
```

▼ Let's check the statistics of our dataset and look for missing values, other disci

```
1 #describing the numerical attributes
2 print(data.describe())
```

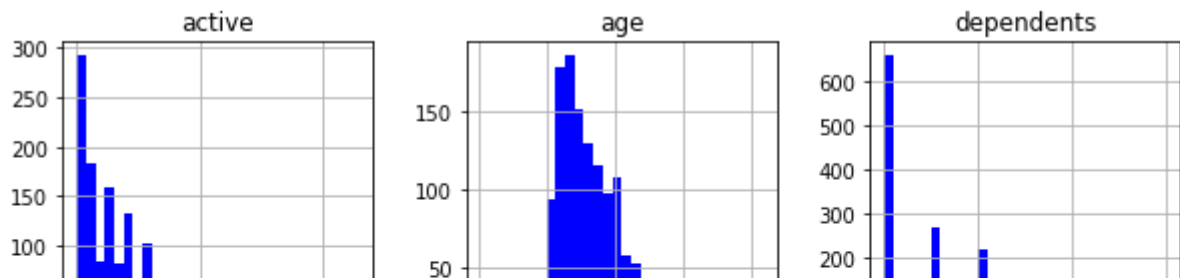
```
(1319, 12)
```

	reports	age	...	majorcards	active
count	1319.000000	1319.000000	...	1319.000000	1319.000000
mean	0.456406	33.213103	...	0.817286	6.996967
std	1.345267	10.142783	...	0.386579	6.305812
min	0.000000	0.166667	...	0.000000	0.000000
25%	0.000000	25.416670	...	1.000000	2.000000
50%	0.000000	31.250000	...	1.000000	6.000000
75%	0.000000	39.416670	...	1.000000	11.000000
max	14.000000	83.500000	...	1.000000	46.000000

```
[8 rows x 9 columns]
```

```
1 #plot the histogram of each parameter
2 import warnings
3 warnings.filterwarnings('ignore')
4 data.hist(color='blue',bins=30,figsize=(10,10))
5 plt.show()
```

```
(1319, 12)
```



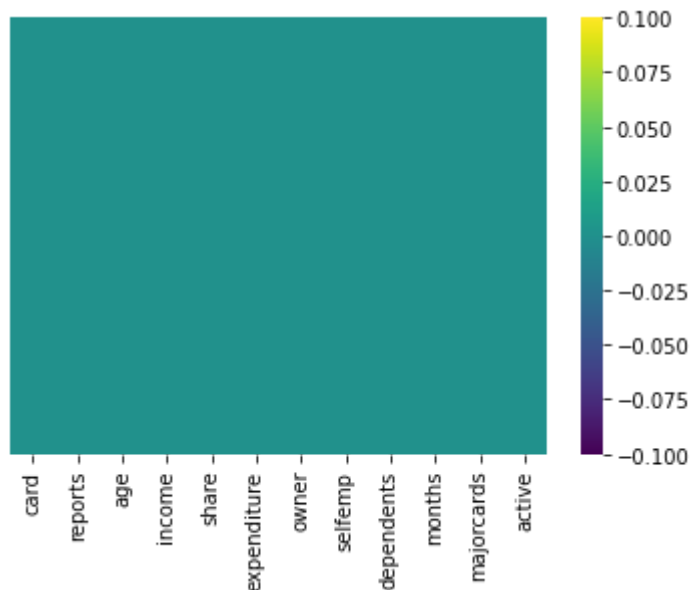
▼ Does my data have missing values?

```

1 sns.heatmap(data.isnull(),yticklabels=False,cbar=True,cmap='viridis')
2 #no null values exists in any of the attributes

```

↳ <matplotlib.axes._subplots.AxesSubplot at 0x7ff0b9738b38>

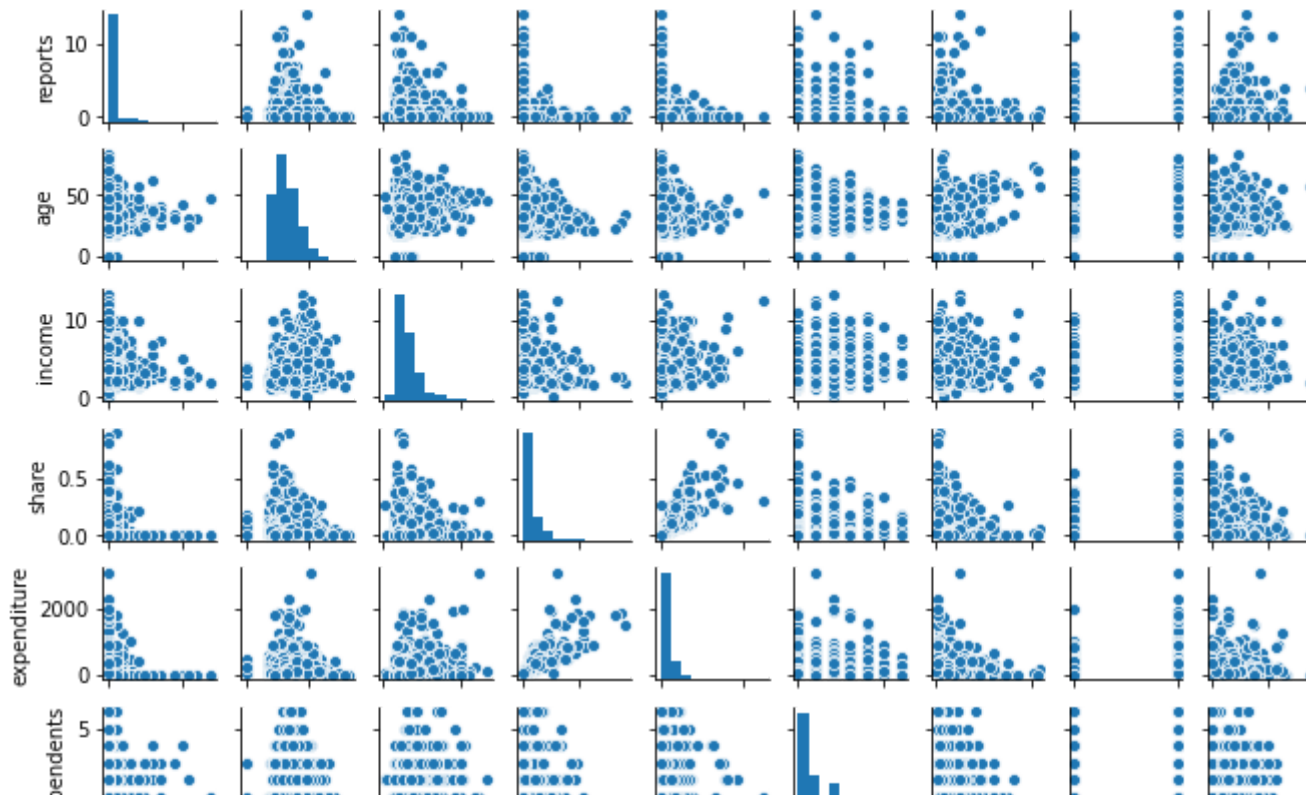


▼ No missing values found. Can we check out for any correlations between dataset?

```
1 sns.pairplot(data,height=1)
```

↳

<seaborn.axisgrid.PairGrid at 0x7ff0b5c5e860>



▼ As you can see

There is no linear relation between pairs of age, income, expenditure, reports. These seem to be the manually removed the features which have correlations with these features for dimensionality reduction essential for small datasets to avoid overfitting. We do not want our model to memorize the dataset a new input during inference phase. Common algorithms for Dimensionality Reduction are :

- Principal Component Analysis
- t-SNE
- UMAP
- Linear Discriminant Analysis (Gaussian)

```
1 #displaying the dataset
2 data.head()
```



	card	reports	age	income	share	expenditure	owner	selfemp	dependents
0	yes	0	37.66667	4.5200	0.033270	124.983300	yes	no	
1	yes	0	33.25000	2.4200	0.005217	9.854167	no	no	
2	yes	0	33.66667	4.5000	0.004156	15.000000	yes	no	
3	yes	0	30.50000	2.5400	0.065214	137.869200	no	no	
4	yes	0	32.16667	9.7867	0.067051	546.503300	yes	no	

▼ Converting Categorical features to numerical variables.

This can be done using One-Hot Encoding. However, one-hot encoding is recommended only if you

```
1 # replacing category to numerical value
2 data.card.replace(['yes','no'], ['1', '0'], inplace=True)
3 data.owner.replace(['yes','no'], ['1', '0'], inplace=True)
4 data.selfemp.replace(['yes','no'], ['1', '0'], inplace=True)
5 data.head()
```

	card	reports	age	income	share	expenditure	owner	selfemp	depende
0	1	0	37.66667	4.5200	0.033270	124.983300	1	0	
1	1	0	33.25000	2.4200	0.005217	9.854167	0	0	
2	1	0	33.66667	4.5000	0.004156	15.000000	1	0	
3	1	0	30.50000	2.5400	0.065214	137.869200	0	0	
4	1	0	32.16667	9.7867	0.067051	546.503300	1	0	

▼ Classwise distribution

```
1 # Count Target Variable Values
2 data.card.value_counts()
3
```



▼ We see that this is an imbalanced dataset.

Below are the percentages of the two classes.

```
1 # count target variable percentage
2 round(data.card.value_counts()*100/len(data.axes[0]),2)
```



▼ selecting the feature attributes

```
1 X=data[["reports","expenditure","age","income"]]
2 y=data["card"]
```

```
1 X.head()
```



```
1 y.head()
```

```
0    1
1    1
2    1
3    1
4    1
Name: card, dtype: object
```

▼ splitting the train and test set

```
1 from sklearn.model_selection import train_test_split
```

```
1 X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.3,random_sta
```

```
1 #length of train set
2 len(X_train)
```

```
923
```

```
1 #length of test set
2 len(X_test)
```

```
396
```

```
1 # displaying train set
2 X_train.head()
```

```

reports  expenditure    age  income
831      0      105.32500  39.91667    4.60
1239     0       0.00000  53.00000   11.00
280      0      80.45084  29.83333    2.80
82       0     108.60750  19.58333    1.65
1117     0     128.07250  22.91667    1.56
```

```
1 # displaying test set
2 X_test.head()
```

	reports	expenditure	age	income
56	0	1898.03300	34.33333	4.800
711	0	20.33333	36.41667	2.200
912	0	0.00000	29.33333	3.500
1053	1	0.00000	28.50000	2.000

```

1  #Load library
2  #from sklearn import tree
3
4  # Building Decision Tree - CART Algorithm (gini criteria)
5  dt_train_gini_decision = DecisionTreeClassifier(criterion = "gini", random_state=100,
6                                                  max_depth=5, min_samples_leaf=5)
7  # Train
8  dt_train_gini_decision.fit(X_train,y_train)

```

```

DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None, criterion='gini',
                        max_depth=5, max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=5, min_samples_split=2,
                        min_weight_fraction_leaf=0.0, presort='deprecated',
                        random_state=100, splitter='best')

```

```

1  #apply the classifier we trained to the test data(which,remember,it has never
2  dt_train_gini_decision.predict(X_test) #test matching
3  #passing x and getting y in predict func

```

```


```



```

0.9870550161812297
      precision    recall  f1-score   support

     0       0.92      1.00      0.96         87
     1       1.00      0.97      0.99        309

 accuracy          0.98         396
 macro avg          0.96      0.99      0.97         396
 weighted avg       0.98      0.98      0.98         396

```

```

1 #view a list of the features and theirv imoportant score
2 list(zip(X_train,dt_train_gini_decision.feature_importances_))
3 #which attributes are more decisive...used in overfitting to drop columns
4

```

```

0.9870550161812297
      precision    recall  f1-score   support

     0       0.92      1.00      0.96         87
     1       1.00      0.97      0.99        309

 accuracy          0.98         396
 macro avg          0.96      0.99      0.97         396
 weighted avg       0.98      0.98      0.98         396

[('reports', 0.004721763304312109),
 ('expenditure', 0.9783798350814344),
 ('age', 0.010235979525498706),
 ('income', 0.0066624220887548545)]

```

▼ As you can see

Expenditure has a feature importance out 0.97 in a scale of 0 to 1 (continuous). This shows that on the basis of expenditure which is logically meaningful as **Repayment Styles** and **"How we spend money"** we can decide if they deserve credit approval or not.

```

1 eff=dt_train_gini_decision.score(X_test,y_test)

```

```

1 print(eff*100)

```

```

0.9870550161812297
      precision    recall  f1-score   support

     0       0.92      1.00      0.96         87
     1       1.00      0.97      0.99        309

 accuracy          0.98         396
 macro avg          0.96      0.99      0.97         396
 weighted avg       0.98      0.98      0.98         396

97.97979797979798

```

```

1 import pydotplus
2 import collections
3 from sklearn import tree
4 ## sudo apt-get install graphviz
5 def visualize_graph(clf,features,name):
6     # Visualize data
7     dot_data = tree.export_graphviz(clf,
8                                     feature_names=features,
9                                     out_file=None,
10                                    filled=True,
11                                    rounded=True)
12     graph = pydotplus.graph_from_dot_data(dot_data)
13
14     colors = ('turquoise', 'orange')
15     edges = collections.defaultdict(list)
16
17     for edge in graph.get_edge_list():
18         edges[edge.get_source()].append(int(edge.get_destination()))
19

```

```

20     for edge in edges:
21         edges[edge].sort()
22         for i in range(2):
23             dest = graph.get_node(str(edges[edge][i]))[0]
24             dest.set_fillcolor(colors[i])
25
26     graph.write_png(name+".png")

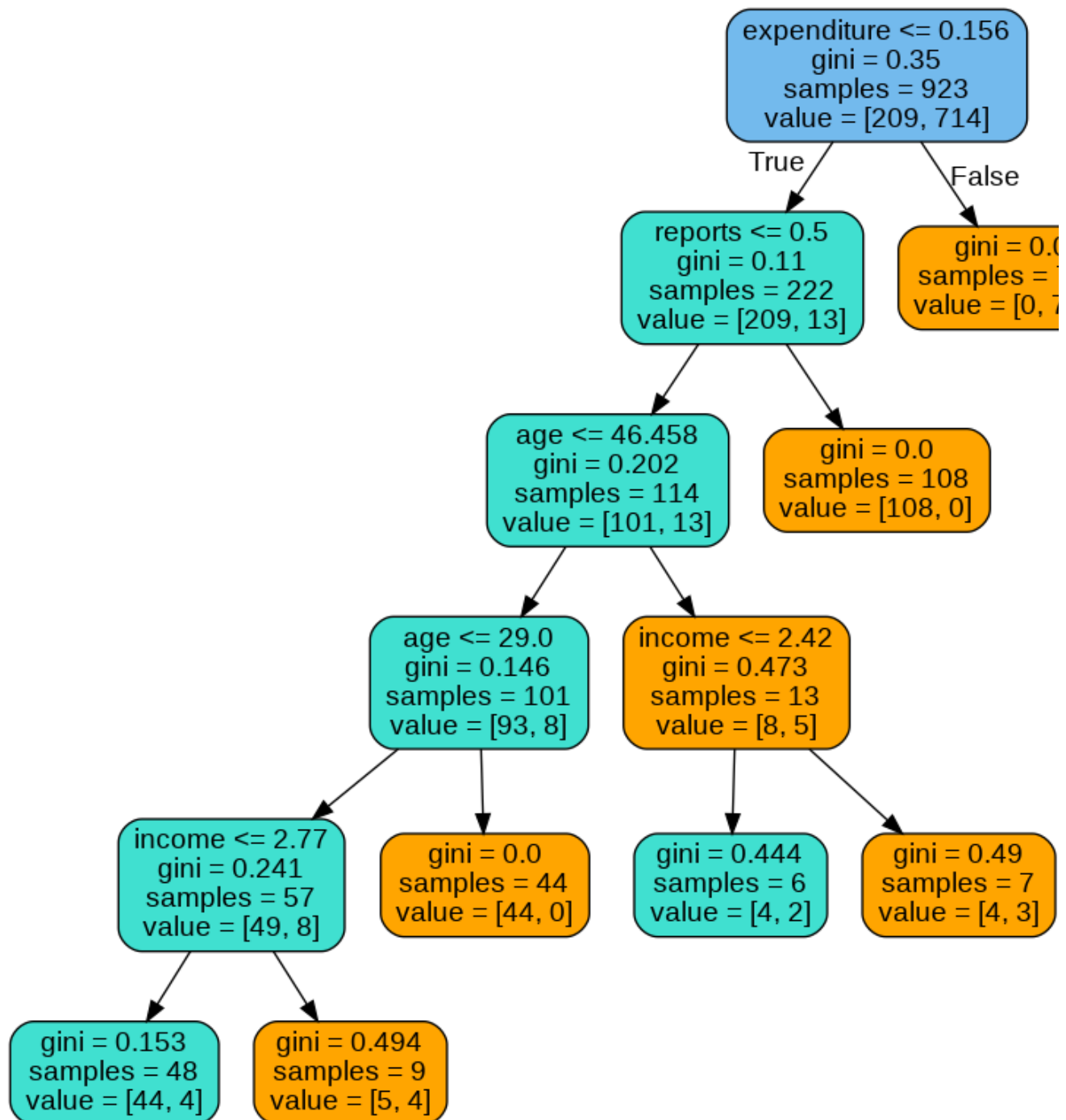
```

```
1 visualize_graph(dt_train_gini_decision,X.columns,"dt_en")
```

```

1 from IPython.display import Image
2 Image('/content/dt_en.png',width=700, height=700)

```



▼ Random Forest Classifier

```

1  from sklearn.ensemble import RandomForestClassifier
2  clf=RandomForestClassifier(n_jobs=-1,random_state=0)
3  clf.fit(X_train,y_train)
4  clf.predict(X_test)
5  clf.predict_proba(X_test)[0:10]

```

```

↳ array([[0.   , 1.   ],
         [0.   , 1.   ],
         [0.97, 0.03],
         [1.   , 0.   ],
         [0.95, 0.05],
         [0.95, 0.05],
         [0.29, 0.71],
         [0.91, 0.09],
         [0.   , 1.   ],
         [0.   , 1.   ]])

```

```

1  list(zip(X_train,clf.feature_importances_))

```

```

↳ [('reports', 0.11000605280019297),
    ('expenditure', 0.7961741426706236),
    ('age', 0.04579285219474255),
    ('income', 0.048026952334440993)]

```

```

1  eff=clf.score(X_test,y_test)*100
2  print(eff)

```

```

↳ 97.72727272727273

```

```

1  import joblib
2  from joblib import dump, load
3  dump(clf, 'model.joblib')

```

```

↳ ['model.joblib']

```

```

1  data.income.min()
2  #data.income.max()

```

```

↳ 0.21

```

```

1  from sklearn.model_selection import cross_val_score
2  scores = cross_val_score(clf, X_train, y_train, cv=5)

```

```

1  scores

```

```

↳ array([0.97837838, 0.98918919, 0.98378378, 0.97282609, 0.97282609])

```

```

1  print("Accuracy: %0.2f (+/- %0.2f)" % (scores.mean(), scores.std() * 2))

```

```

↳ Accuracy: 0.98 (+/- 0.01)

```

```

1  scores = cross_val_score(clf, X_train, y_train, cv=5, scoring='f1_macro')
2

```

```
1 print("Accuracy: %0.2f (+/- %0.2f)" % (scores.mean(), scores.std() * 2))
```

```
↳ Accuracy: 0.97 (+/- 0.02)
```

```
1 from sklearn.model_selection import ShuffleSplit
2 n_samples = X.shape[0]
3 cv = ShuffleSplit(n_splits=5, test_size=0.3, random_state=0)
4 cross_val_score(clf, X, y, cv=cv)
```

```
↳ array([0.97979798, 0.98232323, 0.97222222, 0.97474747, 0.97222222])
```

```
1 X_test
```

```
↳
```

	reports	expenditure	age	income
56	0	1898.03300	34.33333	4.800
711	0	20.33333	36.41667	2.200
912	0	0.00000	29.33333	3.500
1053	1	0.00000	28.50000	2.000
1142	1	0.00000	20.91667	2.625
...
481	0	3.75000	28.16667	3.500
378	0	398.85080	26.58333	3.200
1038	0	66.73250	44.08333	4.000
980	1	42.31083	30.66667	2.740
1299	0	38.99750	38.75000	3.200

396 rows × 4 columns

```
1 len(clf.predict(X_test))
```

```
↳
```

```
1
```