

Tsne (Credit card fraud detection)

Information about Credit Card Fraud Detection:

- The datasets contains transactions made by credit cards in September 2013 by european cardholders.
- This dataset presents transactions that occurred in two days, where we have 492 frauds out of 284,807 transactions.
- It contains only numerical input variables which are the result of a PCA transformation.
- Features V1, V2, ... V28 are the principal components obtained with PCA, the only features which have not been transformed with PCA are 'Time' and 'Amount'.

Features Information:

- (i). Amount: is the transaction Amount
- (ii). Time : contains the seconds elapsed between each transaction and the first transaction in the dataset.
- (iii). V1,V2,...V28: are the principal components obtained with PCA.
- (IV). Class: fraud =1,otherwise =0

Link: <https://www.kaggle.com/mlg-ulb/creditcardfraud> (<https://www.kaggle.com/mlg-ulb/creditcardfraud>)

Task 1: Visualize the T-sne plot of the credit card fraud detection data

In [12]:

```
1 # import the required library
2
3 import numpy as np
4 import pandas as pd
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7 sns.set_style('whitegrid')
8
```

In [13]:

```

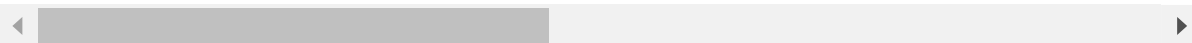
1 # Load the Credit Card Fraud Detection dataset
2 # recap the the dataset
3
4 credit = pd.read_csv('creditcard.csv')
5 credit.head() # show the 5 datapoints of dataset

```

Out[13]:

	Time	V1	V2	V3	V4	V5	V6	V7	V8
0	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698
1	0.0	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	0.085102
2	1.0	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	0.247676
3	1.0	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0.377436
4	2.0	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	-0.270533

5 rows × 31 columns



In [54]:

```

1 # split the features and lables
2 X = credit.iloc[:, :-1]
3 y = credit.iloc[:, -1]
4

```

In [55]:

```

1 # Dataset is very large, so we should standardize the data
2 # StandardScaler(): Standardize features by removing the mean and scaling to unit variance
3
4 from sklearn.preprocessing import StandardScaler
5 scaler = StandardScaler()
6 X_scaled = scaler.fit_transform(X)
7

```

Principal Component Analysis (PCA)

In [64]:

```

1 # IMPORT the PCA scikit-learn library
2 from sklearn.decomposition import PCA
3 pca = PCA() # pca() function
4 pca.fit(X_scaled)
5

```

Out[64]:

```

PCA(copy=True, iterated_power='auto', n_components=None, random_state=None,
     svd_solver='auto', tol=0.0, whiten=False)

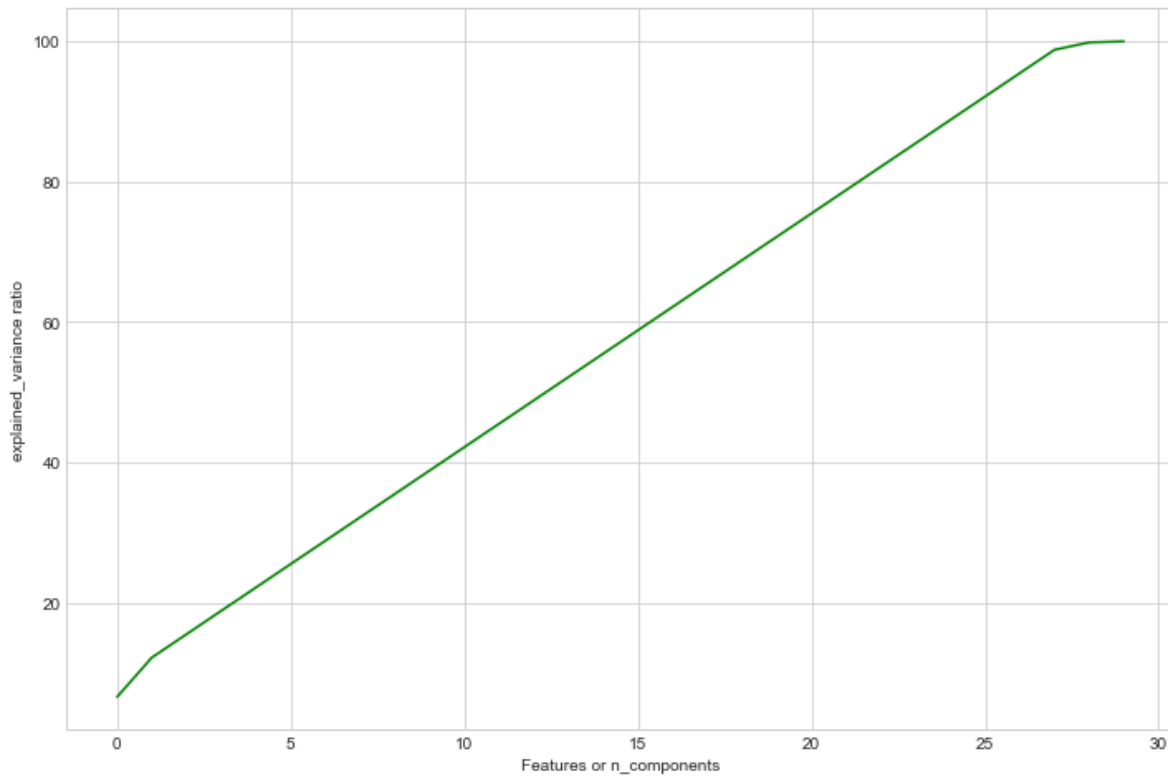
```

In [66]:

```
1 # plot cummulative variance of PCA()
2 plt.figure(figsize=(12,8))
3 plt.plot(np.cumsum(pca.explained_variance_ratio_)*100,color='green')
4 plt.xlabel("Features or n_components")
5 plt.ylabel("explained_variance ratio")
6
```

Out[66]:

Text(0,0.5,'explained_variance ratio')



Observation:

- We require greater than or equal 25 features for getting 95% experience variance ratio

Visualization of t-sne (t-distributed stochastic neighbor embedding):

In [76]:

```

1 # we are plotting the t-sne
2 # but t-sne take more time on full data set
3 # so we are taking 20K data points
4
5 X_Scaled_Fraud = X_scaled[y == 1]
6 X_Scaled_NonFraud = X_scaled[y == 0]
7
8 total_sample_size = 20000 # we are taking 20000 data points
9
10 #non-fraud sample calculation
11 NonFraud_sample_size = total_sample_size - len(X_Scaled_Fraud)
12
13 # create the final sample for fraud and non-fraud transaction.
14 NonFraud_sample = np.random.choice(len(X_Scaled_NonFraud),size = NonFraud_sample_size,
15
16 # X_sample for dataset for fraud and non fraud
17 X_sample = np.concatenate([X_Scaled_Fraud,X_Scaled_NonFraud[NonFraud_sample]])
18
19 # fraud transaction are in the starting
20 Y_sample = np.zeros((total_sample_size), dtype=np.int)
21 Y_sample[:len(X_Scaled_Fraud)] = 1
22
23
24
25
26
27
28
29
30
31

```

In [75]:

```

1 print("sample shape:",X_sample.shape)
2 print("Fraud tranasaction in sample:",y_sample.sum())
3

```

sample shape: (20000, 30)

Fraud tranasaction in sample: 492

In [77]:

```

1 colors = {0: 'blue', 1: 'darkred'}

```

In [80]:

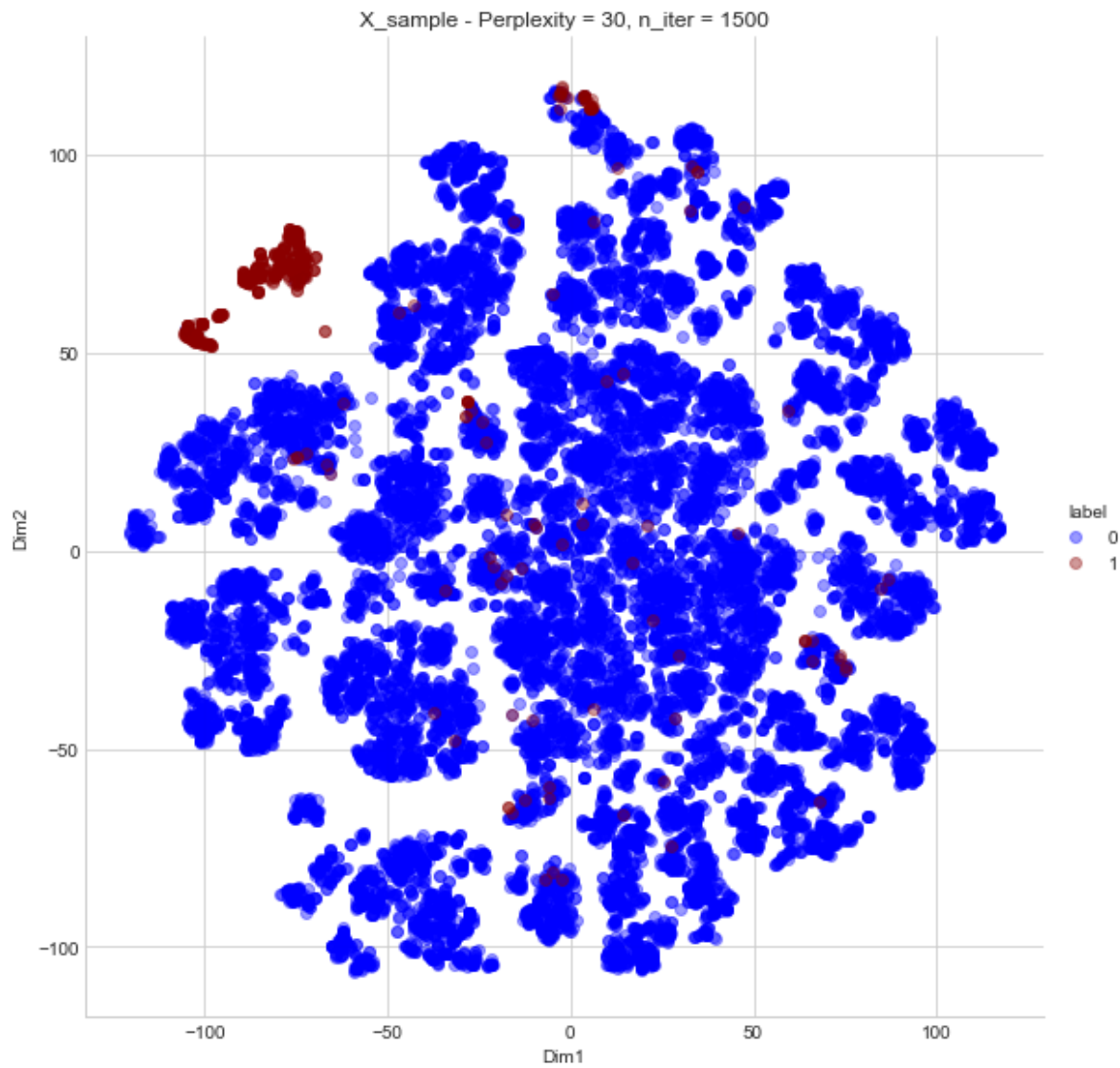
```

1 from sklearn.manifold import TSNE
2
3 def tsne_plot(perplexity=30,n_iter=1000,verbose=0):
4     tsne = TSNE(n_components=2, perplexity=perplexity, n_iter=n_iter, verbose=verbose)
5     tsne_df = pd.DataFrame(data={'Dim1': tsne[:, 0], 'Dim2': tsne[:, 1], 'label': y_sample})
6     sns.FacetGrid(tsne_df, hue="label", size=8, palette=colors).map(plt.scatter, "Dim1", "Dim2")
7     plt.title(f'X_sample - Perplexity = {perplexity}, n_iter = {n_iter}')
8

```

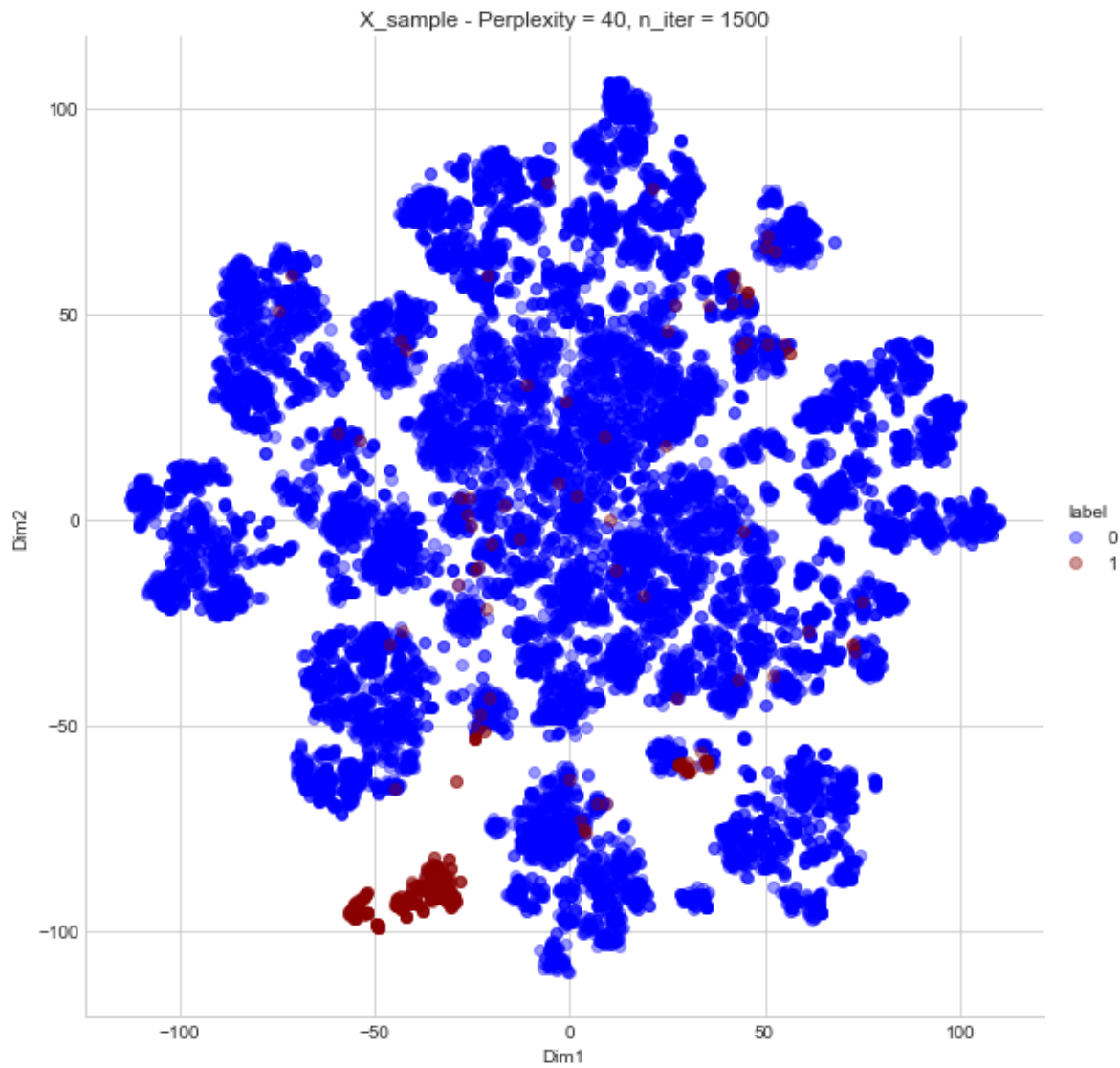
In [81]:

```
1 tsne_plot(perplexity=30, n_iter=1500)
```



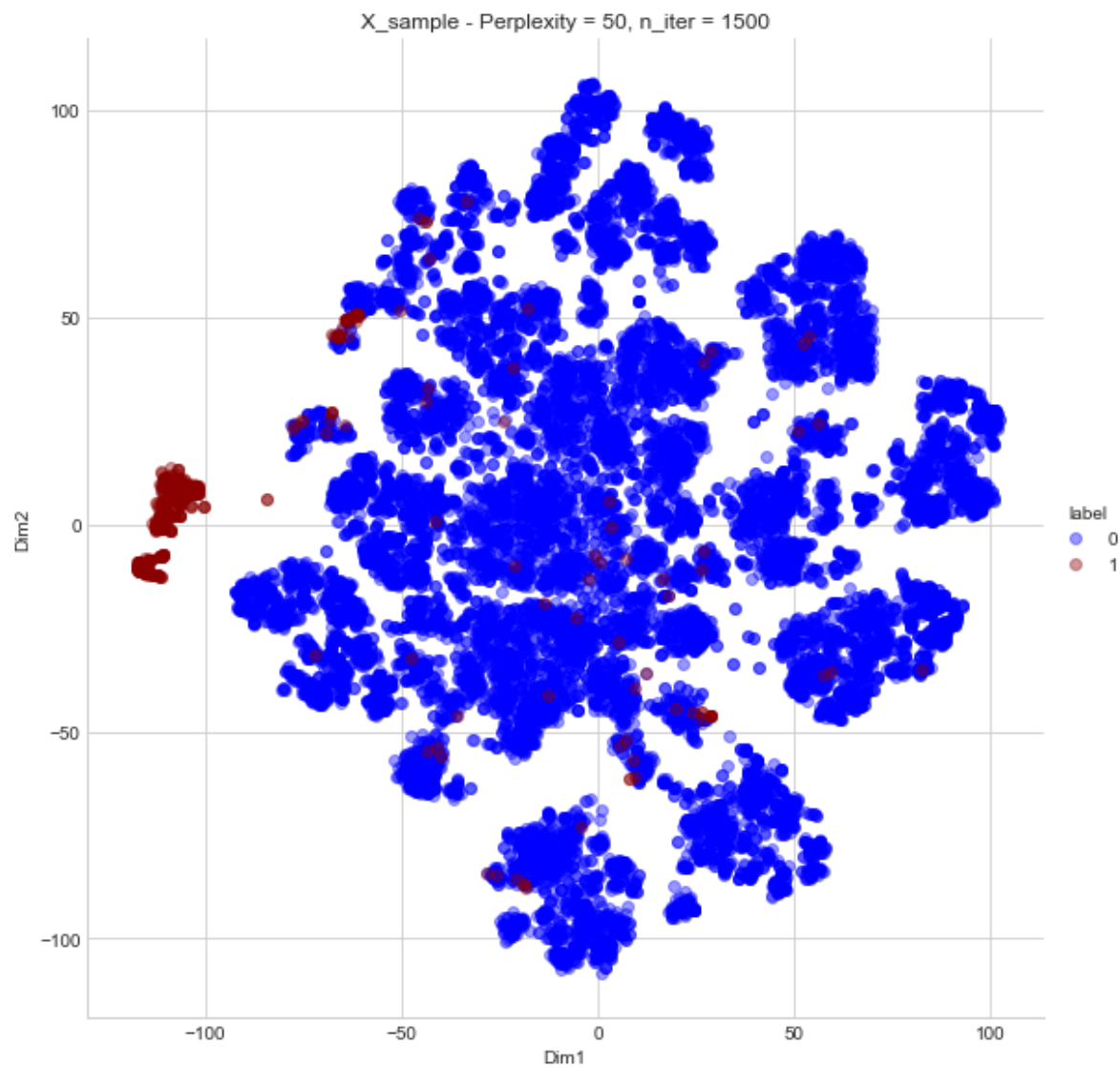
In [82]:

```
1 tsne_plot(perplexity=40, n_iter=1500)
```



In [83]:

```
1 tsne_plot(perplexity=50, n_iter=1500)
```



Observation:

- Most of fraud transaction data are well separated from non- fraud transaction data

In []:

```
1
```