

In [1]:

```

1 # Import our relevant libraries
2 import numpy as np
3 import pandas as pd
4 import seaborn as sns
5 from matplotlib import pyplot as plt
6 %matplotlib inline

```

In [2]:

```

1 data = pd.read_csv('Breast_cancer_wisconsin_diagnosis.csv')
2 data

```

Out[2]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothnes
0	842302	M	17.99	10.38	122.80	1001.0	
1	842517	M	20.57	17.77	132.90	1326.0	
2	84300903	M	19.69	21.25	130.00	1203.0	
3	84348301	M	11.42	20.38	77.58	386.1	
4	84358402	M	20.29	14.34	135.10	1297.0	
...	
564	926424	M	21.56	22.39	142.00	1479.0	
565	926682	M	20.13	28.25	131.20	1261.0	
566	926954	M	16.60	28.08	108.30	858.1	
567	927241	M	20.60	29.33	140.10	1265.0	
568	92751	B	7.76	24.54	47.92	181.0	

569 rows × 33 columns

In [3]:

```

1 # Drop the id column
2 data = data.drop('id', axis=1)
3 data = data.drop('Unnamed: 32', axis=1)
4 # Convert the diagnosis column to numeric format
5 data['diagnosis'] = data['diagnosis'].factorize()[0]
6 # Fill all Null values with zero
7 data = data.fillna(value=0)
8 # Store the diagnosis column in a target object and then drop it
9 X = data.drop('diagnosis', axis=1)
10 y = data['diagnosis']

```

In [4]:

```

1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X,y,test_size = 0.3, random

```

In [5]:

```

1 print('X_train.shape:', X_train.shape)
2 print('X_test.shape:', X_test.shape)
3 print('y_train.shape:', y_train.shape)
4 print('y_test.shape:', y_test.shape)

```

```

X_train.shape: (398, 30)
X_test.shape: (171, 30)
y_train.shape: (398,)
y_test.shape: (171,)

```

VISUALISING PCA AND TSNE PLOTS¶

Let's get to the meat of this notebook which is to produce high-level PCA and TSNE visuals

In [6]:

```

1 from sklearn.decomposition import PCA # Principal Component Analysis module
2 from sklearn.manifold import TSNE # TSNE module

```

In [7]:

```

1 # Turn dataframe into arrays
2 X_train = X_train.values
3 X_test = X_test.values
4
5 # Invoke the PCA method. Since this is a binary classification problem
6 # let's call n_components = 2
7 pca = PCA(n_components=2)
8 train_pca_2d = pca.fit_transform(X_train)
9 test_pca_2d = pca.fit_transform(X_test)
10
11 # Invoke the TSNE method
12 tsne = TSNE(n_components=2, verbose=1, perplexity=40, n_iter=2000)
13 train_tsne_results = tsne.fit_transform(X_train)
14 test_tsne_results = tsne.fit_transform(X_test)

```

```

[t-SNE] Computing 121 nearest neighbors...
[t-SNE] Indexed 398 samples in 0.001s...
[t-SNE] Computed neighbors for 398 samples in 0.015s...
[t-SNE] Computed conditional probabilities for sample 398 / 398
[t-SNE] Mean sigma: 17.977994
[t-SNE] KL divergence after 250 iterations with early exaggeration: 5
0.973656
[t-SNE] KL divergence after 1800 iterations: 0.173023
[t-SNE] Computing 121 nearest neighbors...
[t-SNE] Indexed 171 samples in 0.000s...
[t-SNE] Computed neighbors for 171 samples in 0.007s...
[t-SNE] Computed conditional probabilities for sample 171 / 171
[t-SNE] Mean sigma: 11.093358
[t-SNE] KL divergence after 250 iterations with early exaggeration: 5
0.035530
[t-SNE] KL divergence after 900 iterations: 0.130900

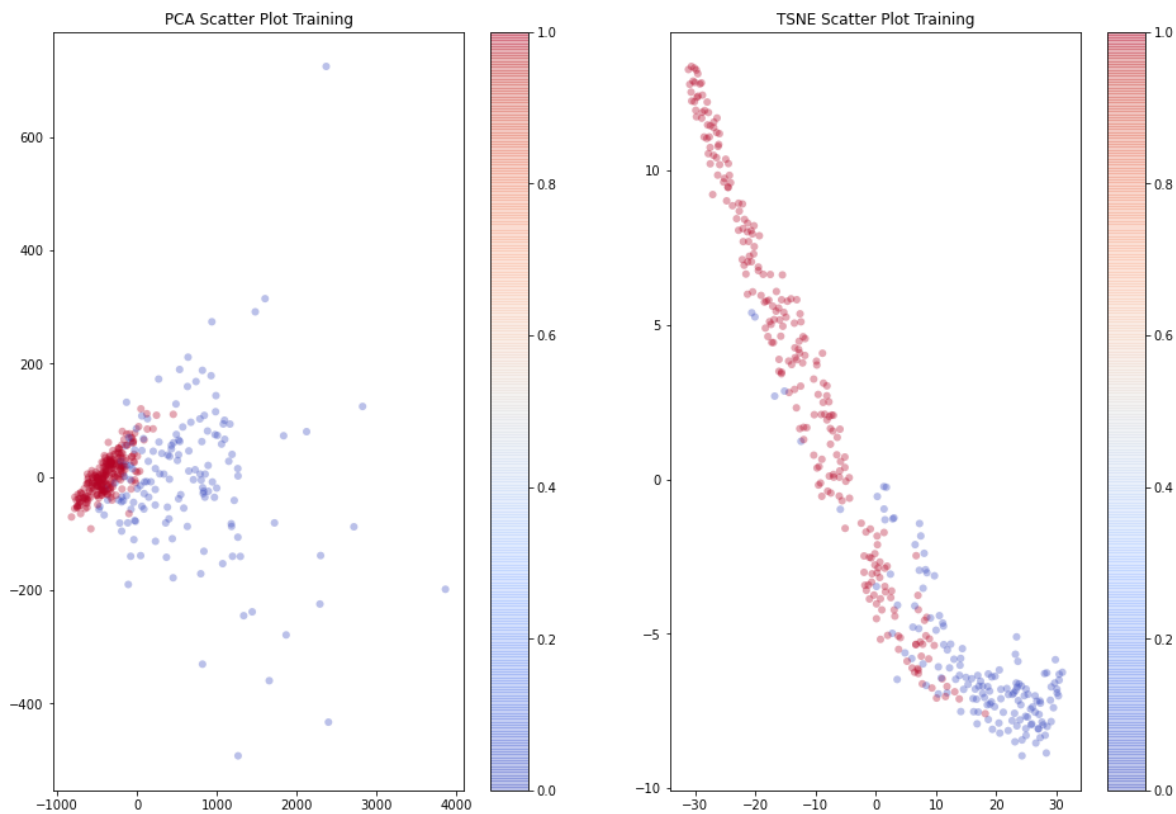
```

In [8]:

```

1 # Plot the TSNE and PCA visuals side-by-side
2 plt.figure(figsize = (16,11))
3 plt.subplot(121)
4 plt.scatter(train_pca_2d[:,0],train_pca_2d[:,1], c = y_train,
5             cmap = "coolwarm", edgecolor = "None", alpha=0.35)
6 plt.colorbar()
7 plt.title('PCA Scatter Plot Training')
8 plt.subplot(122)
9 plt.scatter(train_tsne_results[:,0],train_tsne_results[:,1], c = y_train,
10            cmap = "coolwarm", edgecolor = "None", alpha=0.35)
11 plt.colorbar()
12 plt.title('TSNE Scatter Plot Training')
13 plt.show()

```

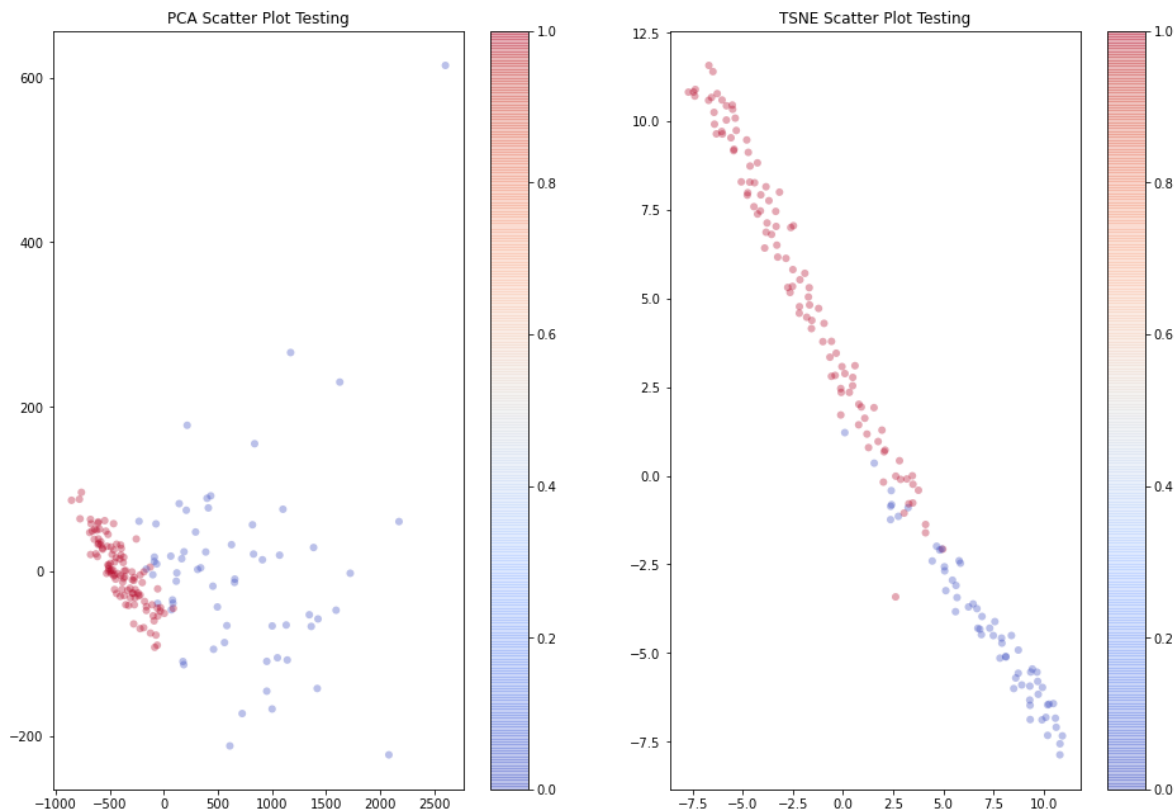


In [9]:

```

1 # Plot the TSNE and PCA visuals side-by-side
2 plt.figure(figsize = (16,11))
3 plt.subplot(121)
4 plt.scatter(test_pca_2d[:,0],test_pca_2d[:,1], c = y_test,
5             cmap = "coolwarm", edgecolor = "None", alpha=0.35)
6 plt.colorbar()
7 plt.title('PCA Scatter Plot Testing')
8 plt.subplot(122)
9 plt.scatter(test_tsne_results[:,0],test_tsne_results[:,1], c = y_test,
10           cmap = "coolwarm", edgecolor = "None", alpha=0.35)
11 plt.colorbar()
12 plt.title('TSNE Scatter Plot Testing')
13 plt.show()

```



As one can see from these high-level plots, even though PCA does quite a decent job of visualising our two target clusters (M for Malignant and B for Benign - cheating a bit here with the labels), the visuals in TSNE is much more obvious in terms of the demarcation in the target.

STANDARDISATION AND VISUALISATION

Let's now try scaling (or standardising) our features and see if we can get even more obvious/intuitive clusters in our plots.

In [10]:

```
1 # Calling Sklearn scaling method
2 from sklearn.preprocessing import StandardScaler
3 X_train_std = StandardScaler().fit_transform(X_train)
4 X_test_std = StandardScaler().fit_transform(X_test)
```

In [11]:

```
1 # Invoke the PCA method on the standardised data
2 pca = PCA(n_components=2)
3 train_pca_2d_std = pca.fit_transform(X_train_std)
4 test_pca_2d_std = pca.fit_transform(X_test_std)
5 # Invoke the TSNE method
6 tsne = TSNE(n_components=2, verbose=1, perplexity=40, n_iter=2000)
7 train_tsne_results_std = tsne.fit_transform(X_train_std)
8 test_tsne_results_std = tsne.fit_transform(X_test_std)
```

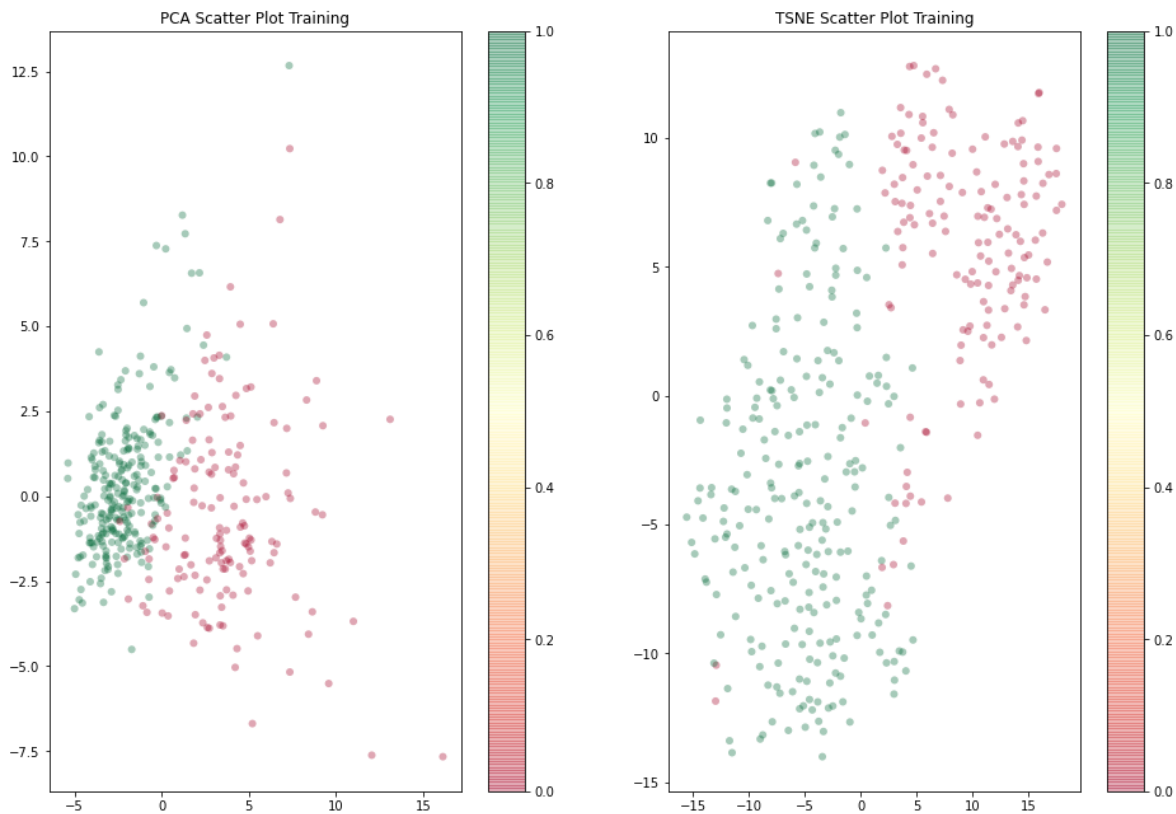
```
[t-SNE] Computing 121 nearest neighbors...
[t-SNE] Indexed 398 samples in 0.001s...
[t-SNE] Computed neighbors for 398 samples in 0.021s...
[t-SNE] Computed conditional probabilities for sample 398 / 398
[t-SNE] Mean sigma: 1.678350
[t-SNE] KL divergence after 250 iterations with early exaggeration: 6
3.064606
[t-SNE] KL divergence after 900 iterations: 0.727239
[t-SNE] Computing 121 nearest neighbors...
[t-SNE] Indexed 171 samples in 0.000s...
[t-SNE] Computed neighbors for 171 samples in 0.003s...
[t-SNE] Computed conditional probabilities for sample 171 / 171
[t-SNE] Mean sigma: 2.245331
[t-SNE] KL divergence after 250 iterations with early exaggeration: 5
7.437195
[t-SNE] KL divergence after 1000 iterations: 0.436906
```

In [12]:

```

1 # Plot the TSNE and PCA visuals side-by-side
2 plt.figure(figsize = (16,11))
3 plt.subplot(121)
4 plt.scatter(train_pca_2d_std[:,0],train_pca_2d_std[:,1], c = y_train,
5             cmap = "RdYlGn", edgecolor = "None", alpha=0.35)
6 plt.colorbar()
7 plt.title('PCA Scatter Plot Training')
8 plt.subplot(122)
9 plt.scatter(train_tsne_results_std[:,0],train_tsne_results_std[:,1], c = y_train,
10            cmap = "RdYlGn", edgecolor = "None", alpha=0.35)
11 plt.colorbar()
12 plt.title('TSNE Scatter Plot Training')
13 plt.show()

```

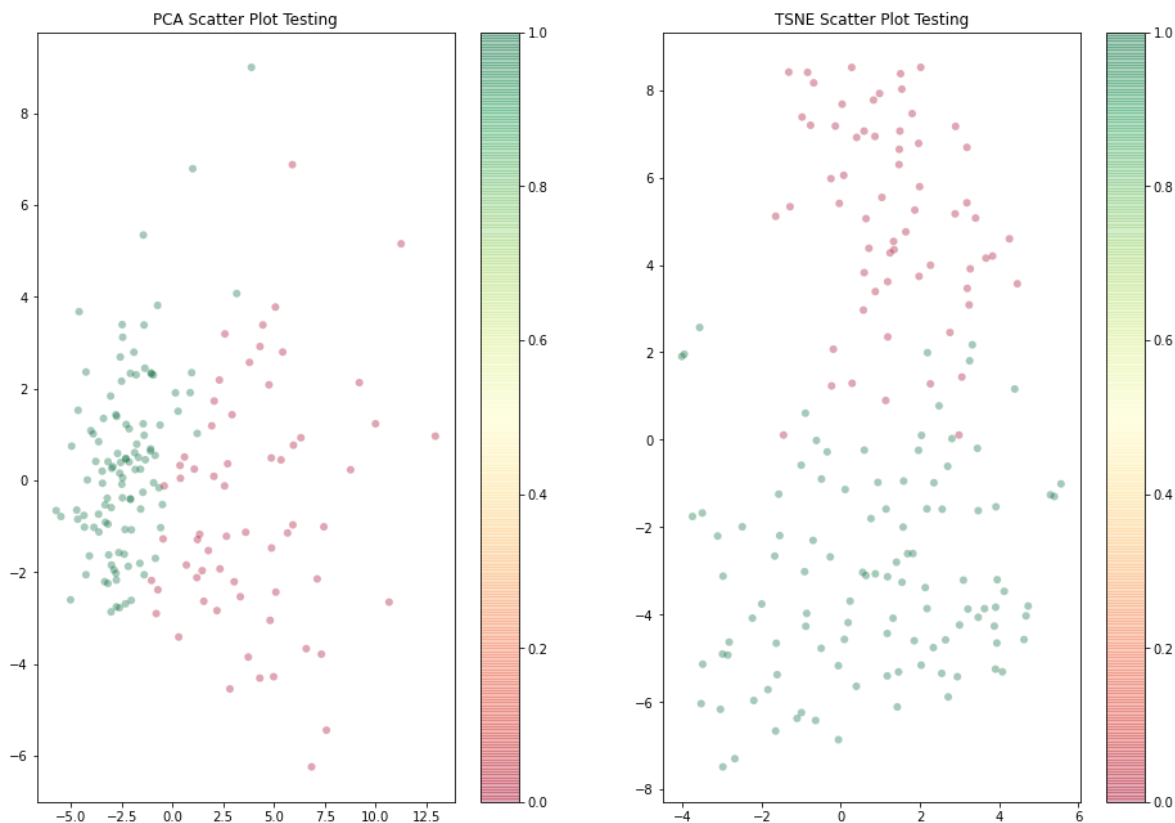


In [13]:

```

1 # Plot the TSNE and PCA visuals side-by-side
2 plt.figure(figsize = (16,11))
3 plt.subplot(121)
4 plt.scatter(test_pca_2d_std[:,0],test_pca_2d_std[:,1], c = y_test,
5             cmap = "RdYlGn", edgecolor = "None", alpha=0.35)
6 plt.colorbar()
7 plt.title('PCA Scatter Plot Testing')
8 plt.subplot(122)
9 plt.scatter(test_tsne_results_std[:,0],test_tsne_results_std[:,1], c = y_test,
10            cmap = "RdYlGn", edgecolor = "None", alpha=0.35)
11 plt.colorbar()
12 plt.title('TSNE Scatter Plot Testing')
13 plt.show()

```



In [14]:

```

1 X_train = train_tsne_results_std
2 X_test = test_tsne_results_std

```

Normalisation

In [15]:

```

1 X_train = (X_train - np.min(X_train))/(np.max(X_train) - np.min(X_train))
2 X_test = (X_test - np.min(X_test))/(np.max(X_test) - np.min(X_test))

```

LogisticRegression

In [16]:

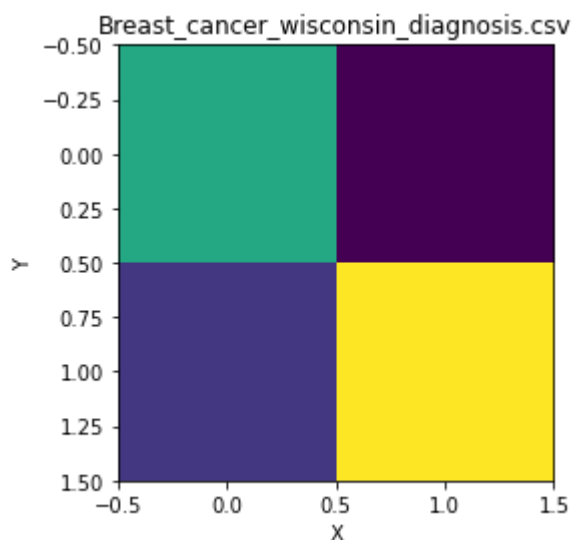
```
1 from sklearn.linear_model import LogisticRegression
2
3 model = LogisticRegression(random_state=0)
4 model.fit(X_train, y_train)
5
6 y_pred = model.predict(X_test)
```

In [17]:

```
1 from sklearn.metrics import confusion_matrix
2
3
4 conf_matrix = confusion_matrix(y_test, y_pred)
5 print('\n\nThe Confusion Matrix for our TSNE Logistic Regression is:\n')
6 print(conf_matrix)
7
8 plt.imshow(conf_matrix)
9 plt.title('Breast_cancer_wisconsin_diagnosis.csv')
10 plt.xlabel('X')
11 plt.ylabel('Y')
12 plt.show()
```

The Confusion Matrix for our TSNE Logistic Regression is:

```
[[56  7]
 [20 88]]
```



In [18]:

```

1 # Visualising the Training set results
2 from matplotlib.colors import ListedColormap
3 x_set, y_set = X_train, y_train
4 X1, X2 = np.meshgrid(np.arange(start = x_set[:, 0].min() - 1, stop = x_set[:, 0].max() + 1, step = 0.5),
5                       np.arange(start = x_set[:, 1].min() - 1, stop = x_set[:, 1].max() + 1, step = 0.5))
6 plt.contourf(X1, X2, model.predict(np.array([X1.ravel(), X2.ravel()]).T).reshape(X1.shape),
7              alpha = 0.75, cmap = ListedColormap(('green', 'green')))
8 plt.xlim(X1.min(), X1.max())
9 plt.ylim(X2.min(), X2.max())
10 for i, j in enumerate(np.unique(y_set)):
11     plt.scatter(x_set[y_set == j, 0], x_set[y_set == j, 1],
12                c = ListedColormap(('red', 'blue'))(i), label = j)
13 plt.title('Logistic Regression (Training set)')
14 plt.xlabel('x')
15 plt.ylabel('y')
16 plt.legend()
17 plt.show()

```

c argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with *x* & *y*. Please use the *color* keyword-argument or provide a 2-D array with a single row if you intend to specify the same RGB or RGBA value for all points.

c argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with *x* & *y*. Please use the *color* keyword-argument or provide a 2-D array with a single row if you intend to specify the same RGB or RGBA value for all points.



In [19]:

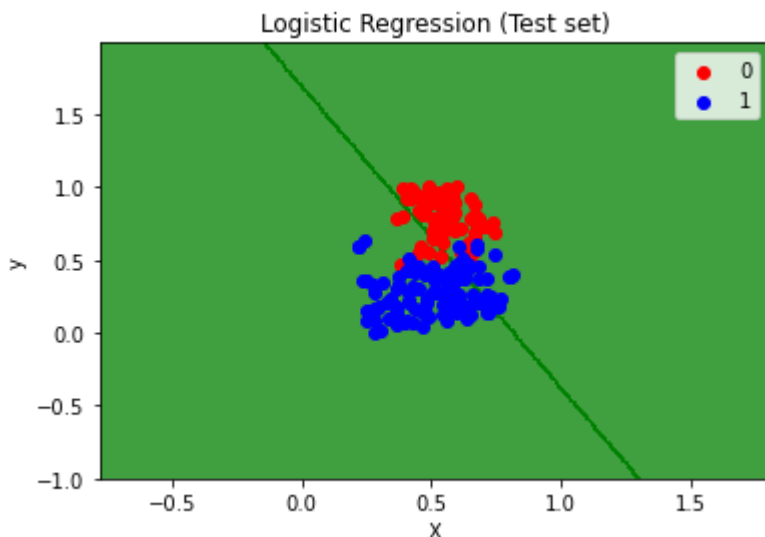
```

1 # Visualising the Test set results
2 from matplotlib.colors import ListedColormap
3 X_set, y_set = X_test, y_test
4 X1, X2 = np.meshgrid(np.arange(start = X_set[:, 0].min() - 1, stop = X_set[:, 0].max() + 1, step = 0.5),
5                       np.arange(start = X_set[:, 1].min() - 1, stop = X_set[:, 1].max() + 1, step = 0.5))
6 plt.contourf(X1, X2, model.predict(np.array([X1.ravel(), X2.ravel()]).T).reshape(X1.shape),
7              alpha = 0.75, cmap = ListedColormap(('green', 'green')))
8 plt.xlim(X1.min(), X1.max())
9 plt.ylim(X2.min(), X2.max())
10 for i, j in enumerate(np.unique(y_set)):
11     plt.scatter(X_set[y_set == j, 0], X_set[y_set == j, 1],
12                c = ListedColormap(('red', 'blue'))(i), label = j)
13 plt.title('Logistic Regression (Test set)')
14 plt.xlabel('X')
15 plt.ylabel('y')
16 plt.legend()
17 plt.show()
18
19 from sklearn.metrics import accuracy_score
20 accuracy = accuracy_score(y_test, y_pred)
21 print('\n\n Hence the accuracy of the t-sne for Logistic Regression is:', accuracy)
22 print('\n\n Done :)')

```

c argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with *x* & *y*. Please use the *color* keyword-argument or provide a 2-D array with a single row if you intend to specify the same RGB or RGBA value for all points.

c argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with *x* & *y*. Please use the *color* keyword-argument or provide a 2-D array with a single row if you intend to specify the same RGB or RGBA value for all points.



Hence the accuracy of the t-sne for Logistic Regression is: 0.8421052631578947

Done :)

