

**NOTE: All the TSNE plots are drawn by using 8k samples.**

## Import necessary libraries

In [1]:

```
1 import warnings  
2 warnings.filterwarnings("ignore")
```

In [2]:

```
1 from MulticoreTSNE import MulticoreTSNE as TSNE  
2 from sklearn.preprocessing import StandardScaler  
3 import pickle  
4 import numpy as np  
5 import pandas as pd  
6 import matplotlib.pyplot as plt  
7 import seaborn as sns  
8 from sklearn.externals import joblib  
9 %matplotlib inline
```

## Load preprocessed data

```
In [3]: 1 # Load data from .pkl files
2 #Functions to save objects for later use and retireve it
3 def savetofile(obj,filename):
4     pickle.dump(obj,open(filename+".pkl","wb"))
5 def openfromfile(filename):
6     temp = pickle.load(open(filename+".pkl","rb"))
7     return temp
8 y=openfromfile('y')
9
10 count_vect =openfromfile('count_vect')
11 X_bigram = openfromfile('X_bigram')
12
13 tf_idf_vect =openfromfile('tf_idf_vect')
14 X_tf_idf = openfromfile('X_tfidf')
15
16 avg_sent_vectors=openfromfile('avg_sent_vectors')
17
18 tfidf_sent_vectors=openfromfile('tfidf_sent_vectors')
```

## Save and Load Model:

```
In [4]: 1 from sklearn.externals import joblib
2 def saveModeltofile(obj,filename):
3     joblib.dump(obj,open(filename+".pkl","wb"))
4 def openModelfromfile(filename):
5     temp = joblib.load(open(filename+".pkl","rb"))
6     return temp
```

## Standardizing data

```
In [5]: 1 def std_data(data,mean):
2     scaler=StandardScaler(with_mean=mean)
3     std_data=scaler.fit_transform(data)
4     return std_data
```

## Function for hyperparameter tunning and draw TSNE Plot:

```
In [6]: 1 def tsne_plots(data, label, params, vect):
2     j=1;
3     for perp in params['perplexity']:
4         i=1;
5         for itr in params['n_iter']:
6             model=TSNE(n_components=2,random_state=0,perplexity=perp,n_iter=itr,n_jobs=-1)
7             tsne_data=model.fit_transform(data)
8             #saveModelToFile(model,vect+str(itr)+'_'+str(perp))
9             #saveModelToFile(tsne_data,'tsnedata'+str(itr)+'_'+str(perp))
10            tsne_data_label=np.vstack((tsne_data.T, label)).T
11            tsne_df=pd.DataFrame(data=tsne_data_label, columns=('X-axis','Y-axis','label'))
12            plt.figure(j,figsize=(18,7))
13            plt.subplot(int('12'+str(i)))
14            plt.title('TSNE (%s, Perplexity=%d, iteration=%d)' %(vect,perp,itr))
15            sns.scatterplot(data=tsne_df,x='X-axis',y='Y-axis',hue='label',palette=['blue','red'])
16            i+=1
17        plt.show()
18    j+=1
```

## Initialization of common objects required for all vectorization:

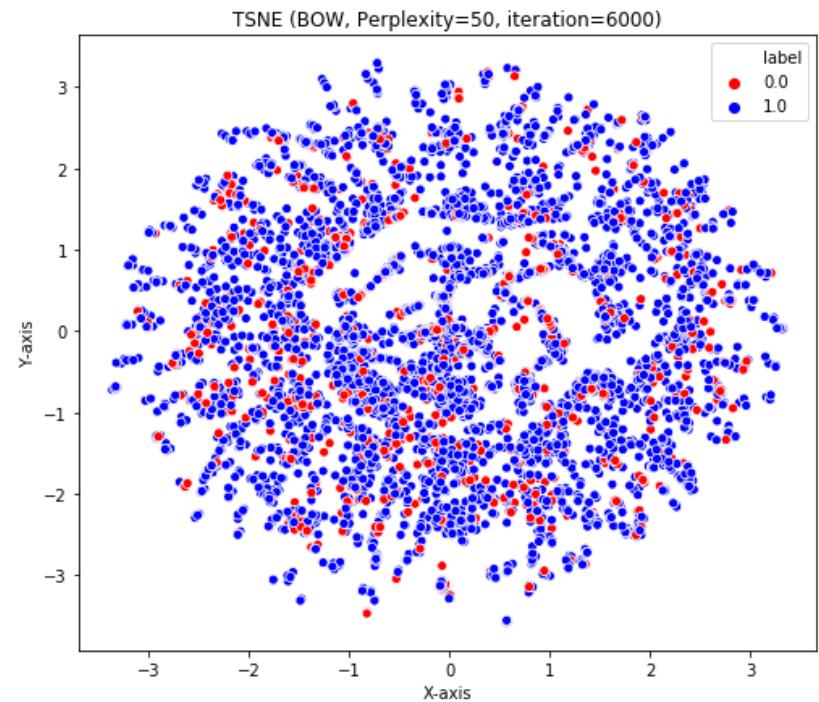
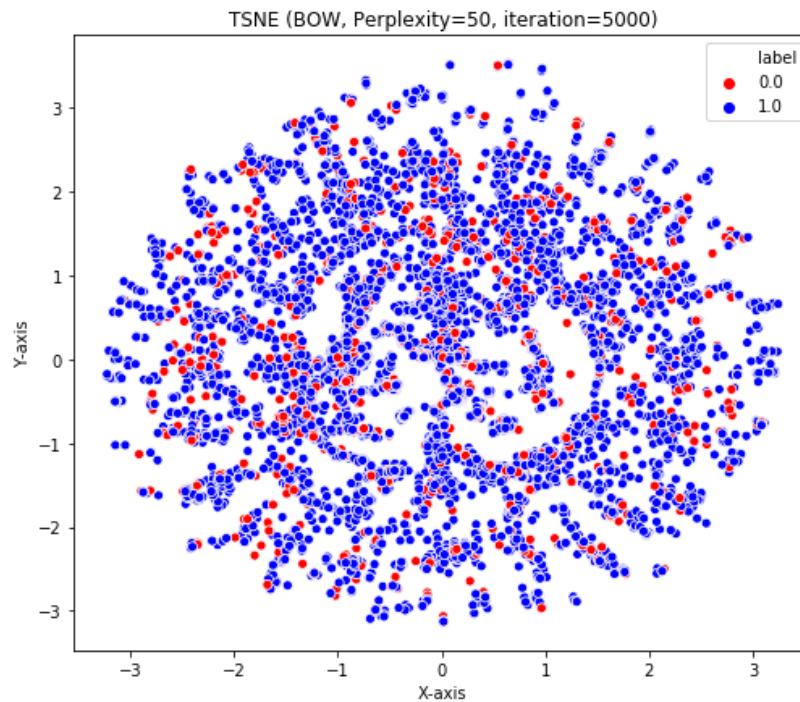
```
In [7]: 1 perplexity=[50,70,80,100]
2 iterations=[5000,6000]
3 params={'perplexity':perplexity,'n_iter':iterations}
4 vect=['BOW','TF-IDF','AVG-W2V','TFIDF-W2V']
```

## BoW

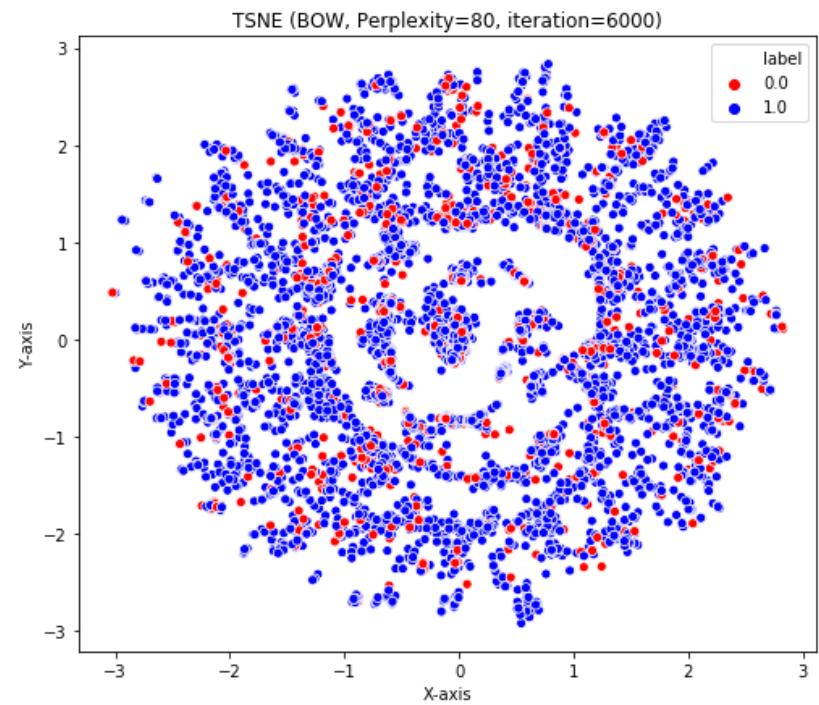
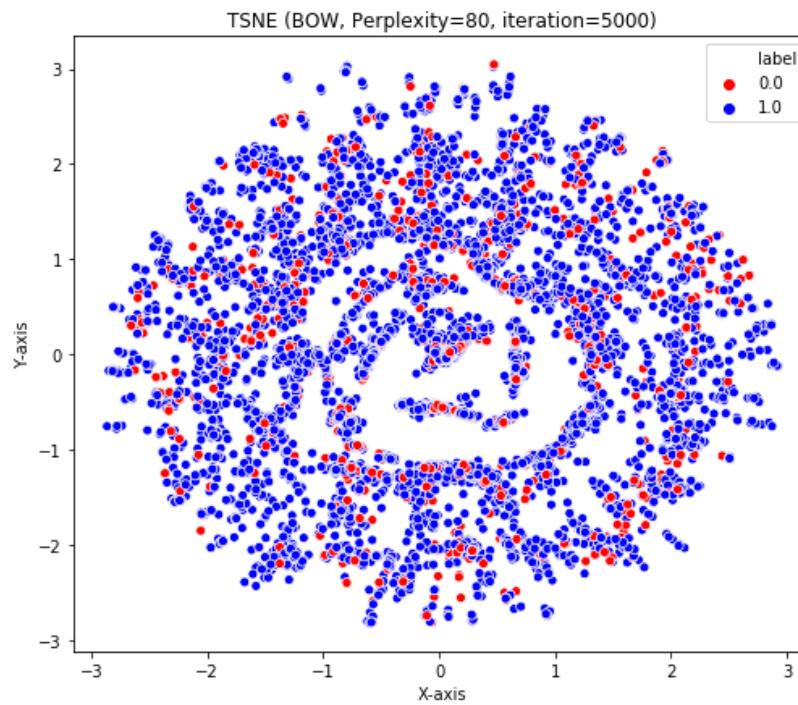
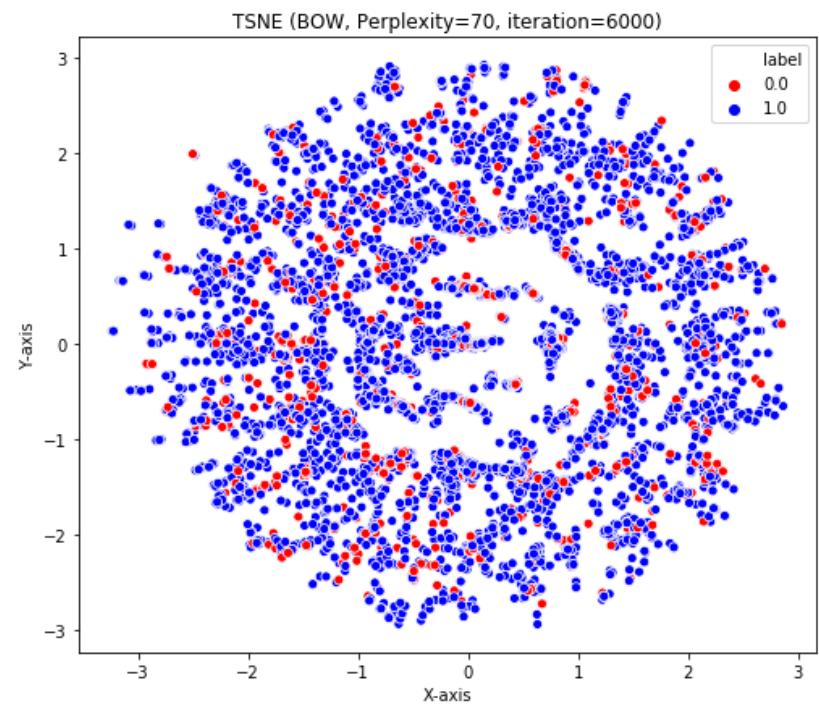
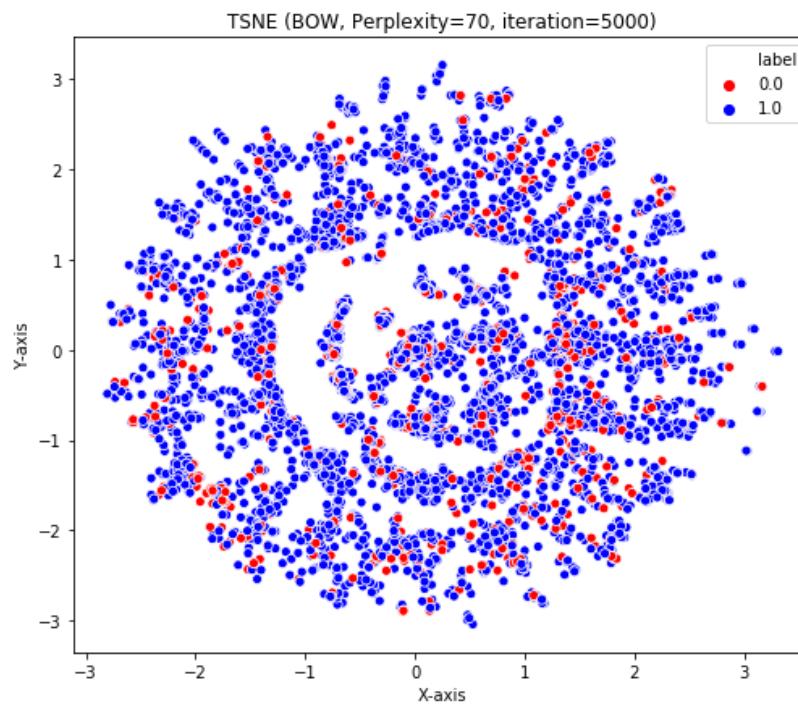
In [8]:

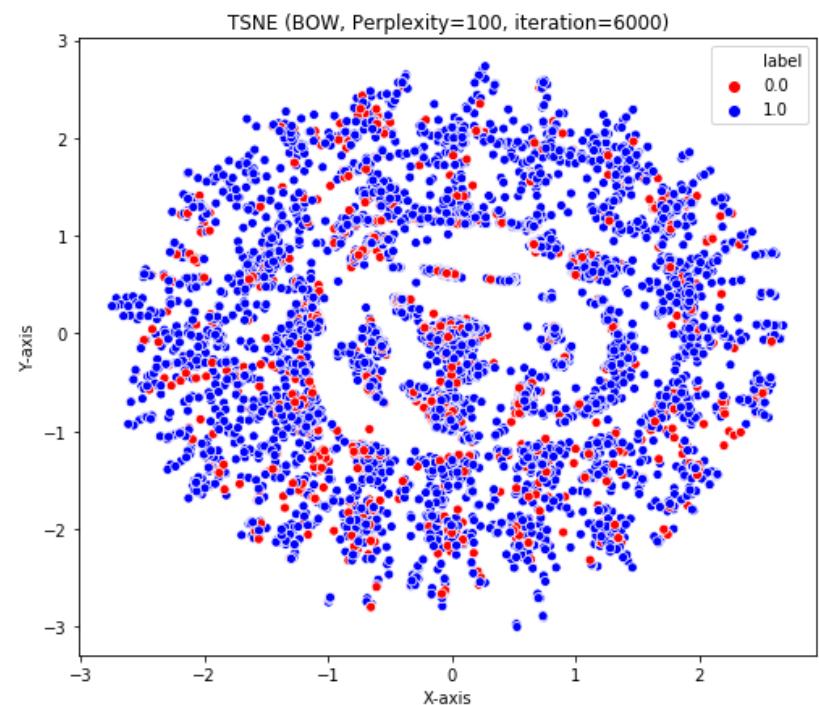
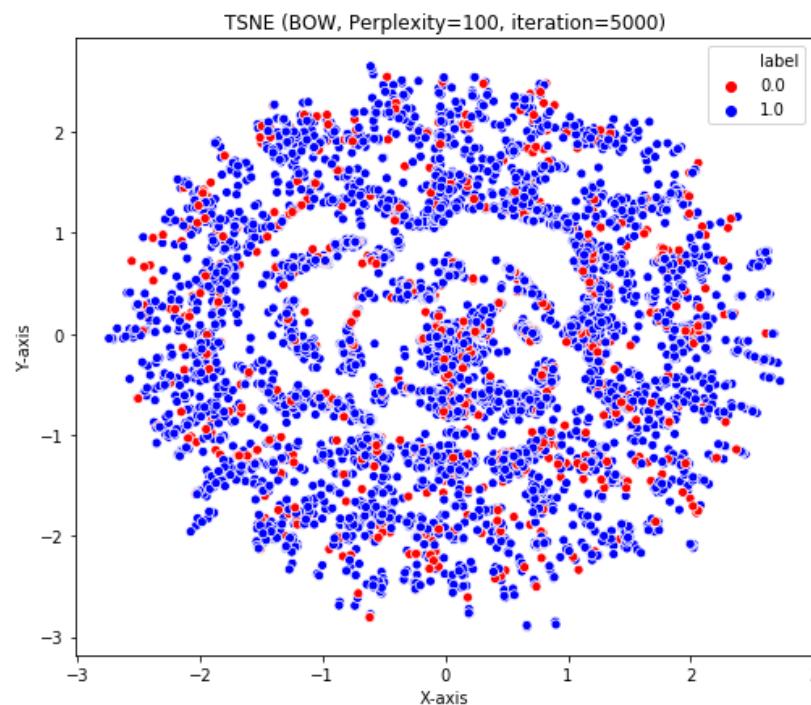
```
1 print('1st run :\n')
2 data=std_data(X_bigram,mean=False).toarray()
3 %time tsne_plots(data, y, params, vect[0])
```

1st run :



## TSNE\_AFFR\_8K\_POINTS





CPU times: user 5h 43min 14s, sys: 2h 53min 37s, total: 8h 36min 51s  
Wall time: 33min 27s

**Observation:**

Note: For clear tuning of perplexity saw AVGW2V and TFIDFW2V section.

1. from the above TSNE plots we observe that for :

[a.] perplexity=70, the plot almost stabilized when the iterations increased from 5000 to 6000

2. from the above TSNE plots we observe that for :

[a.] perplexity=50, 80, 100, the plot is still changing its shape(not stable), when the iterations increased from 5000 to 6000 .

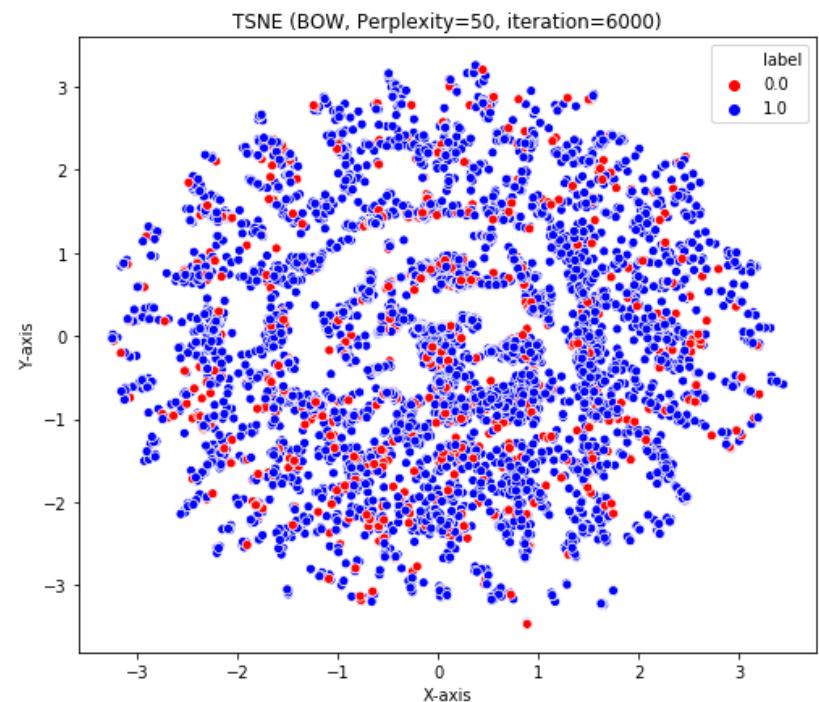
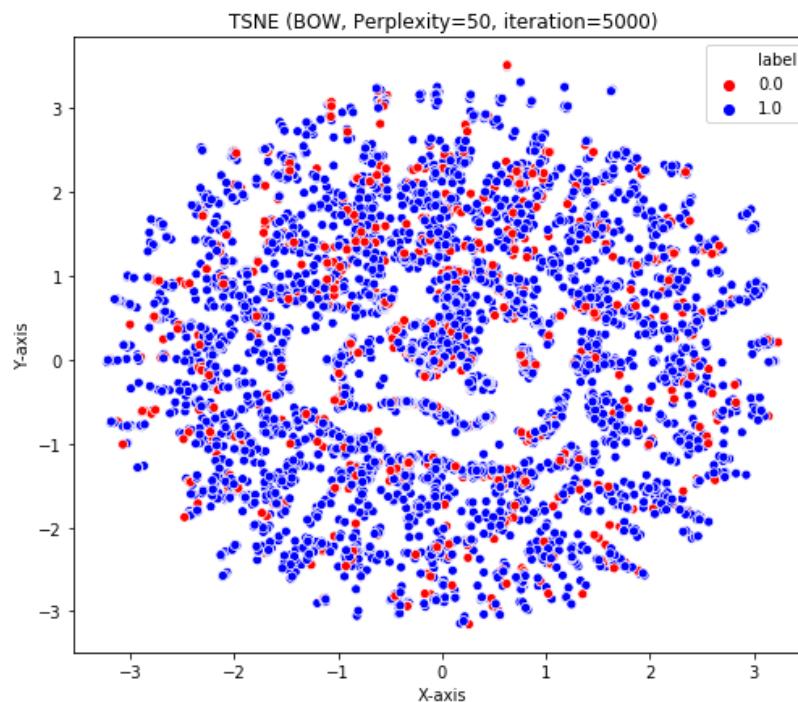
3. so we choose optimal perplexity=70

**Note:**

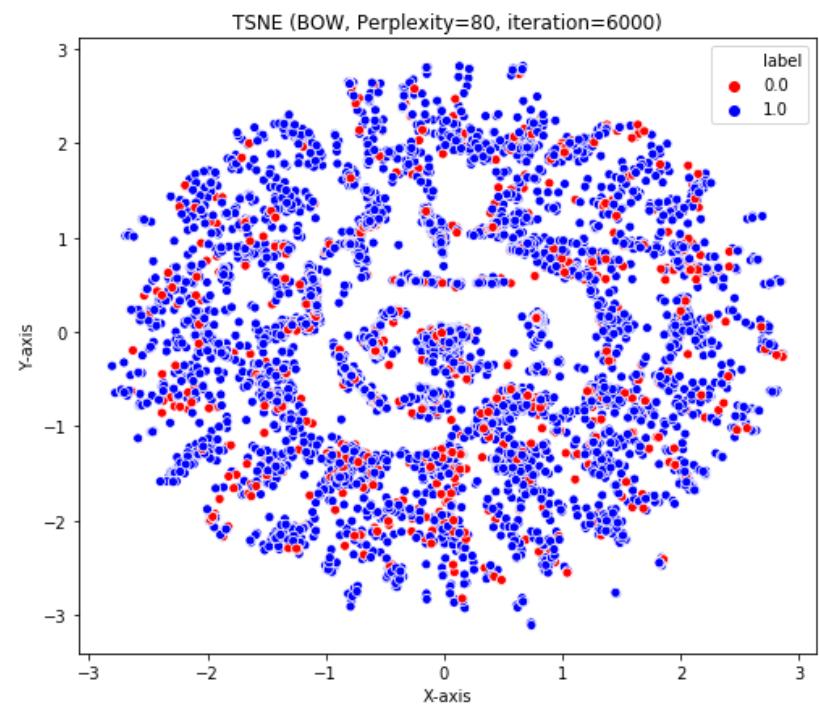
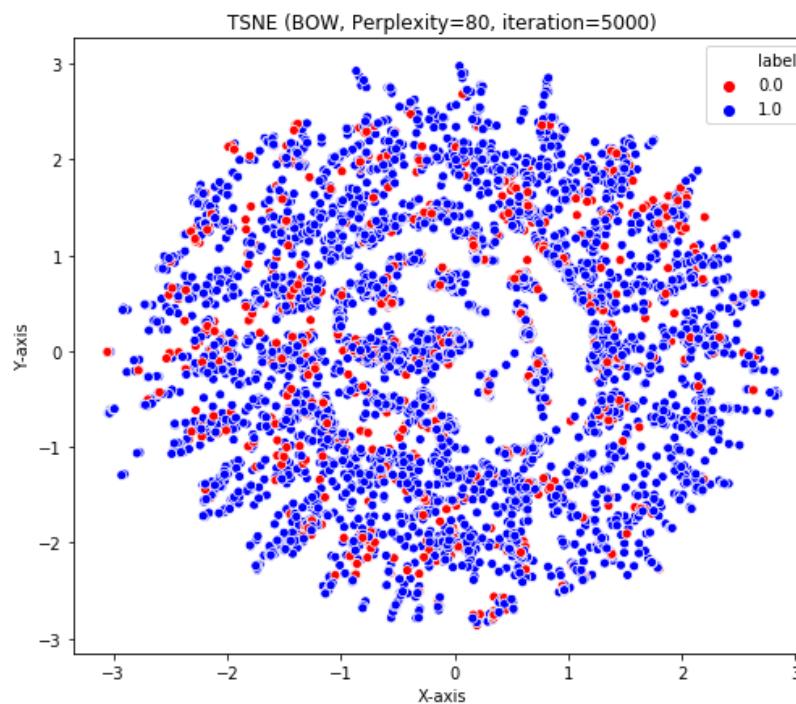
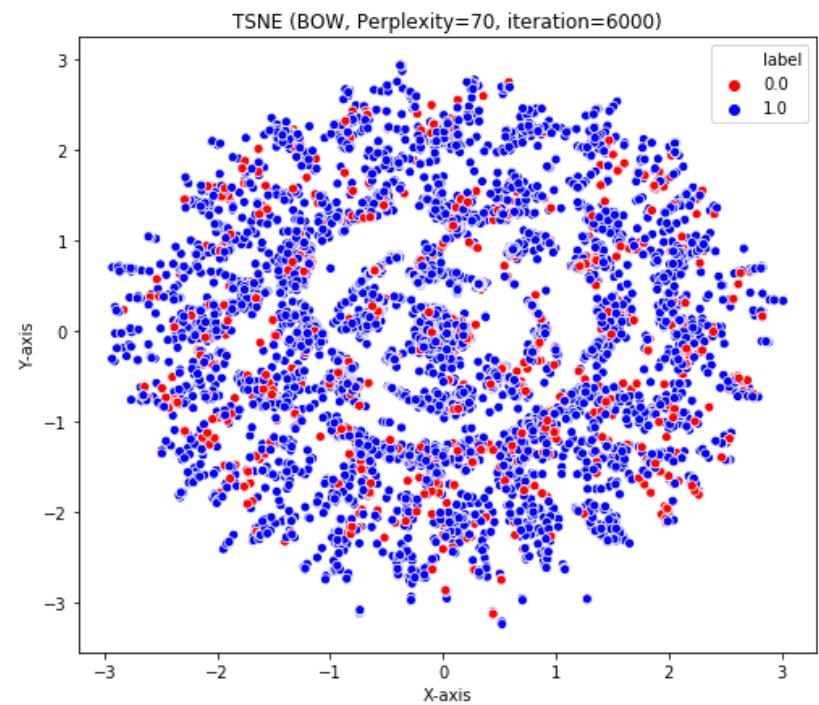
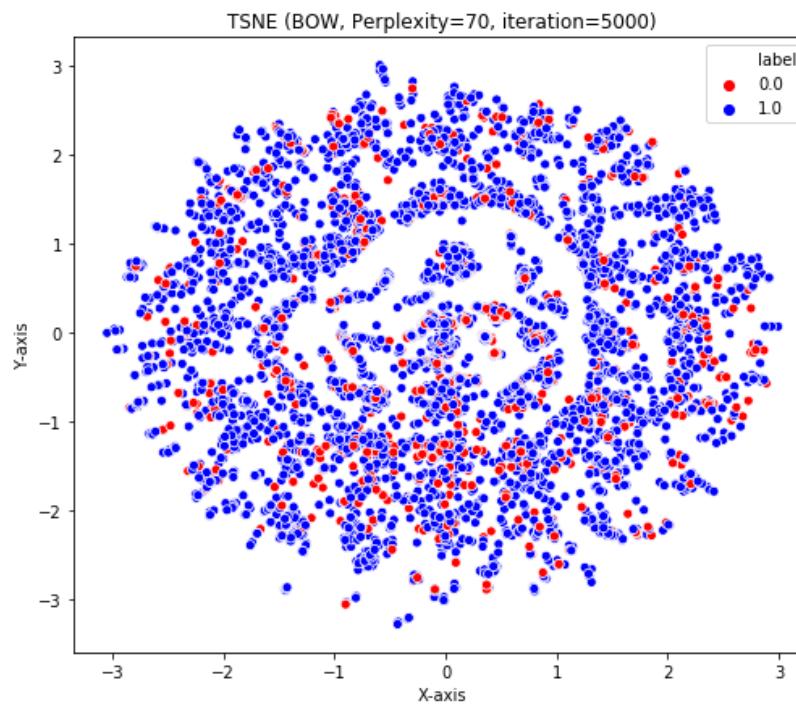
1. To check whether our result is correct or not we re-run TSNE.

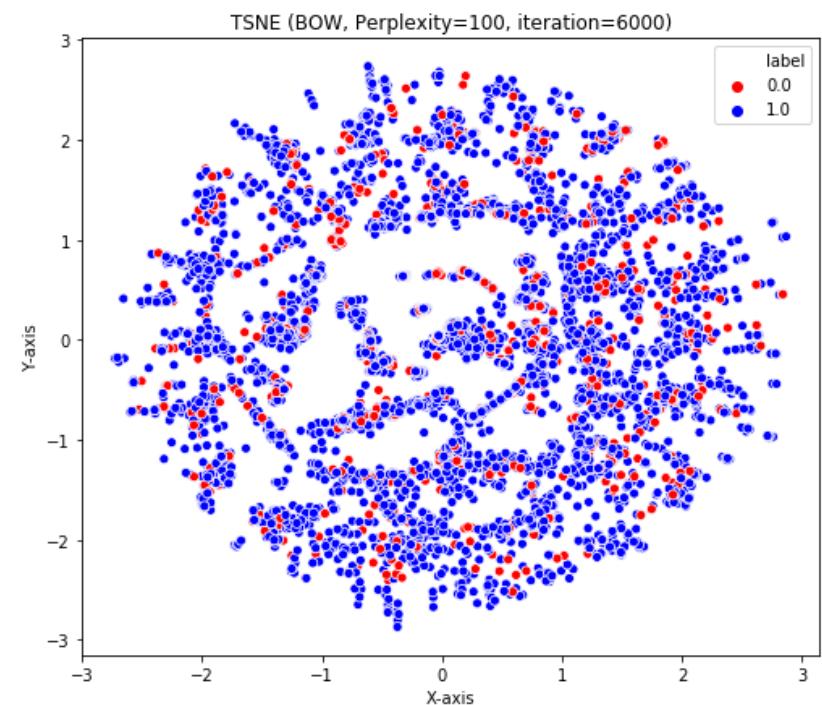
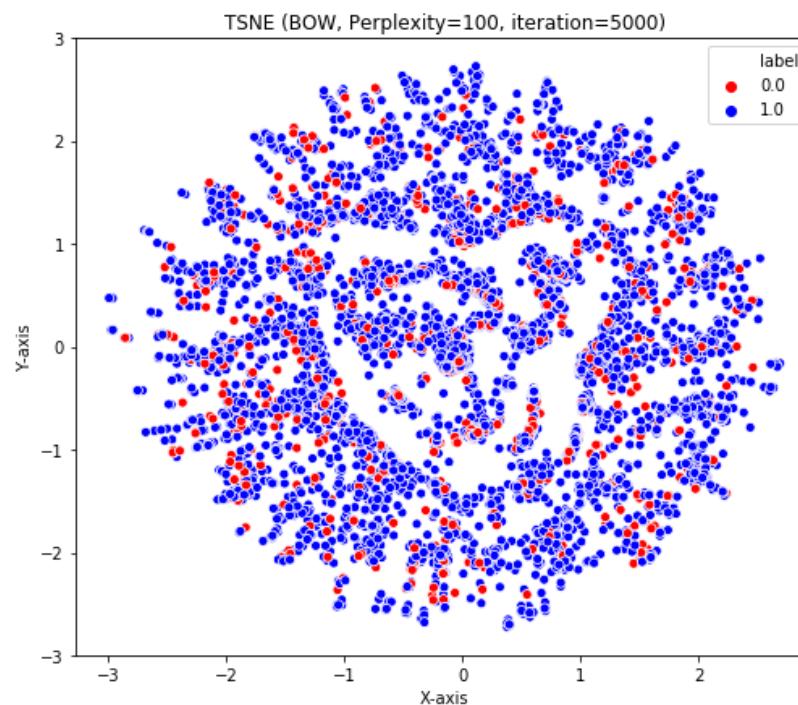
```
In [9]: 1 print('2nd run :\n')
2 %time tsne_plots(data, y, params, vect[0])
```

2nd run :



## TSNE\_AFFR\_8K\_POINTS





CPU times: user 5h 42min 41s, sys: 2h 48min 14s, total: 8h 30min 55s  
Wall time: 33min 5s

**Observation:**

1. from the above TSNE plots we observe that for multiple run of TSNE we get same plot so finally we choose optimal perplexity as:
  - [a.] perplexity=70

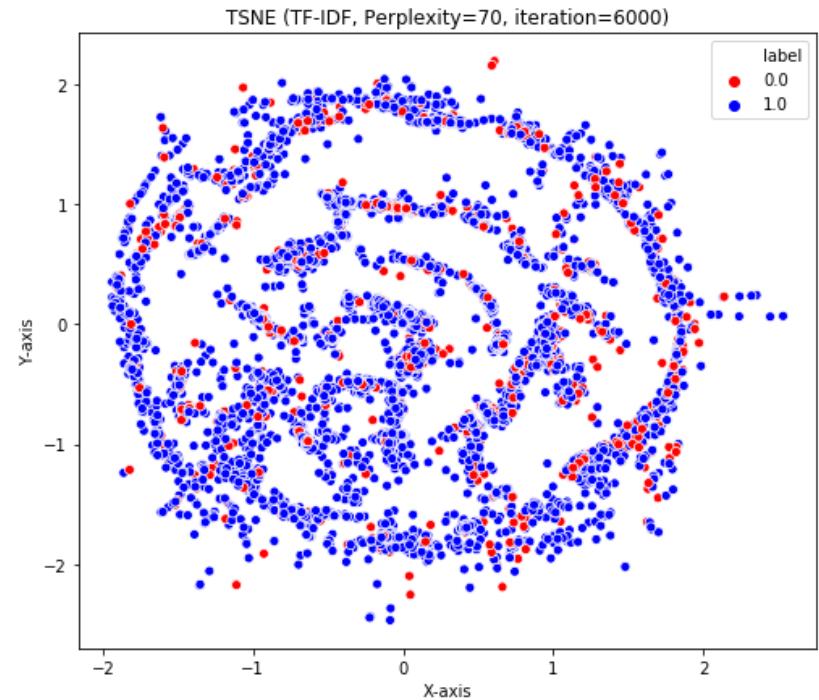
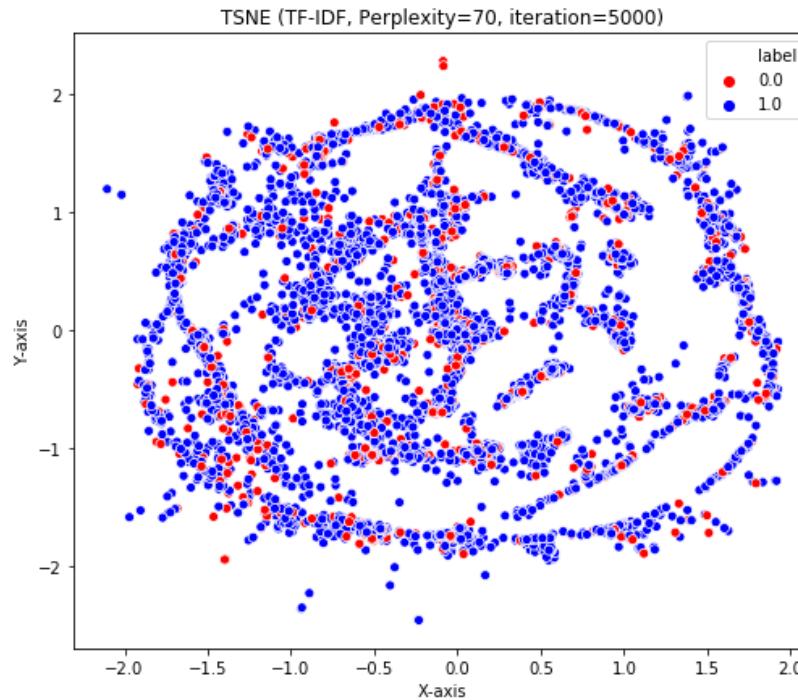
**TFIDF**

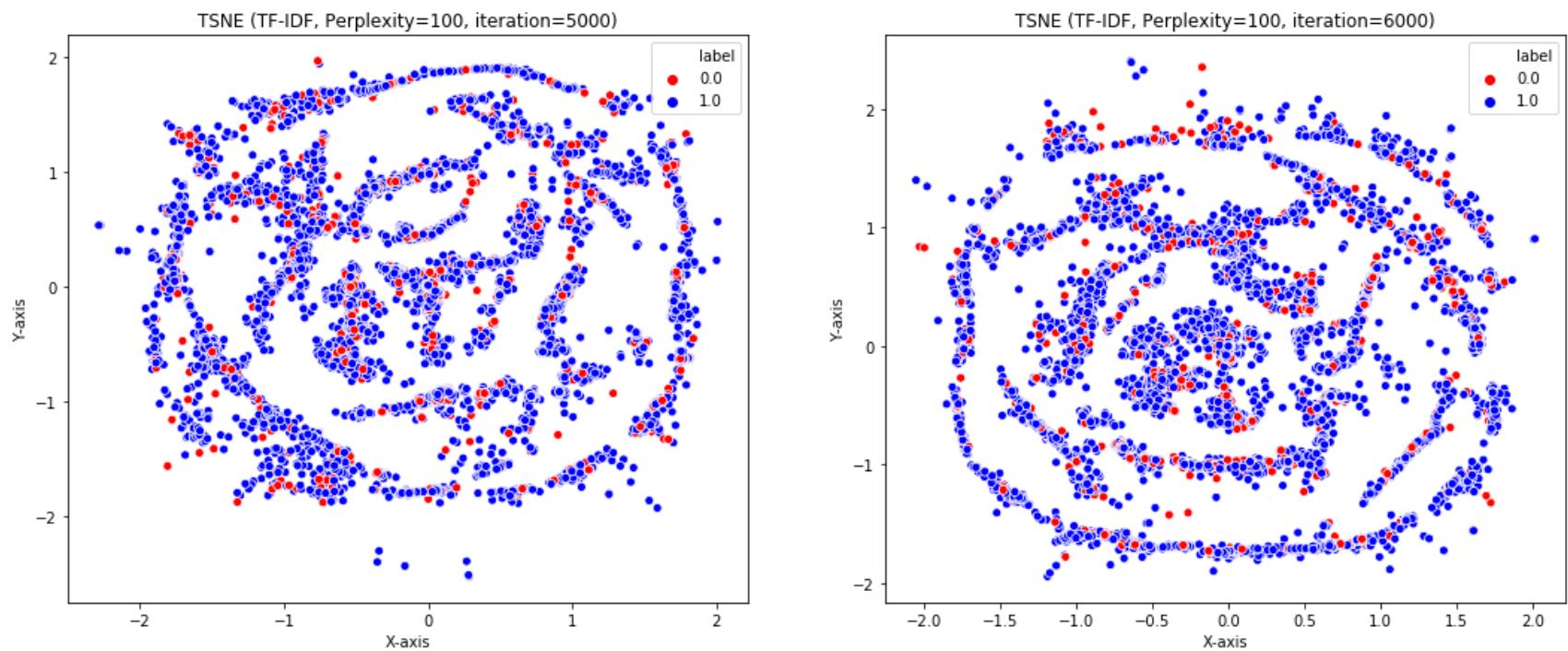
```
In [10]: 1 params['perplexity']=[70,100]
```

In [11]:

```
1 print('1st run :\n')
2 data=std_data(X_tf_idf,mean=False).toarray()
3 %time tsne_plots(data, y, params, vect[1])
```

1st run :

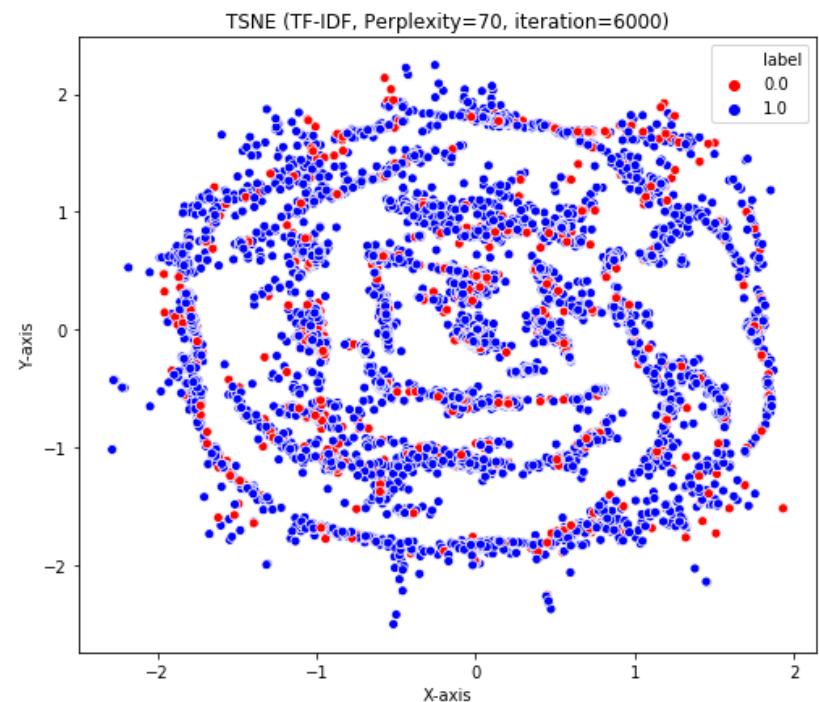
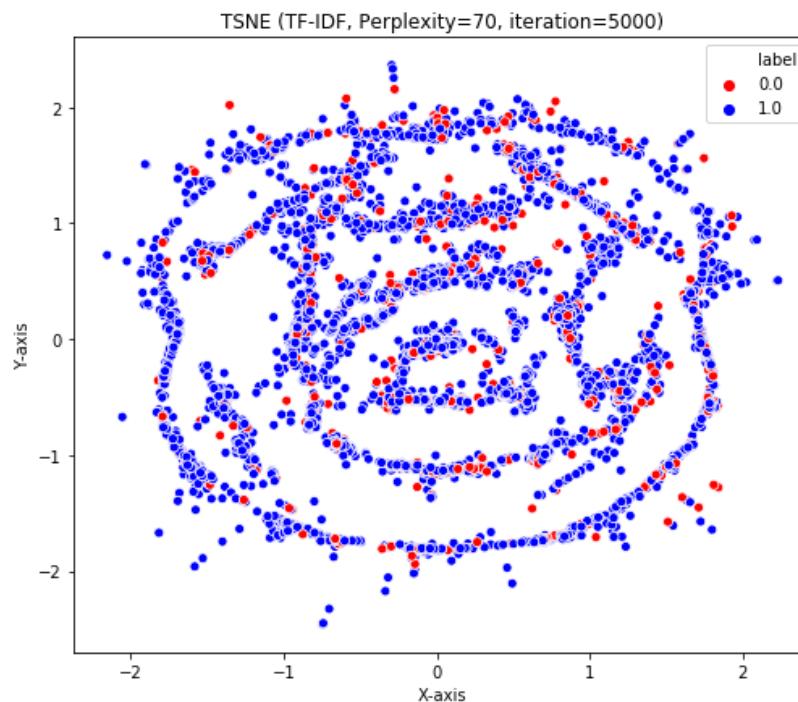


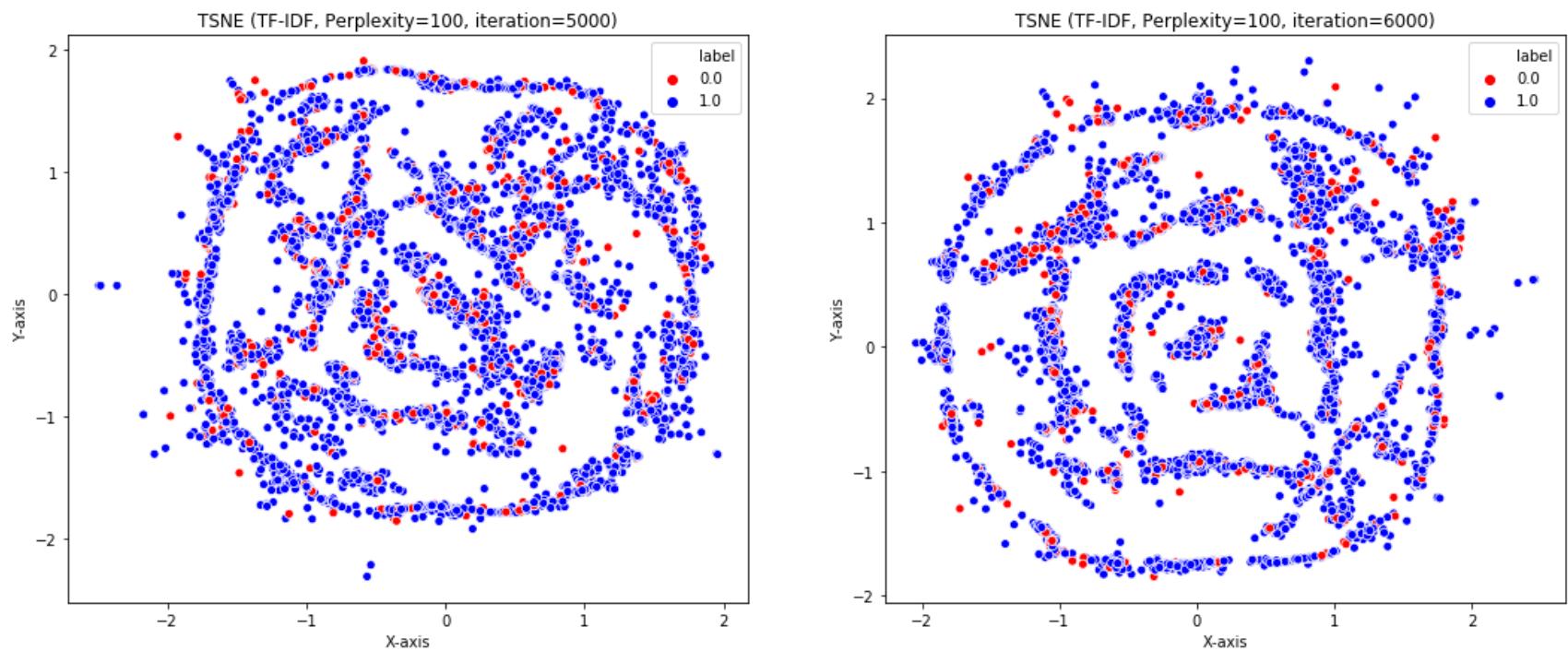


CPU times: user 3h 19min 52s, sys: 1h 32min 10s, total: 4h 52min 3s  
Wall time: 18min 51s

```
In [12]: 1 print('2nd run :\n')
2 %time tsne_plots(data, y, params, vect[1])
```

2nd run :





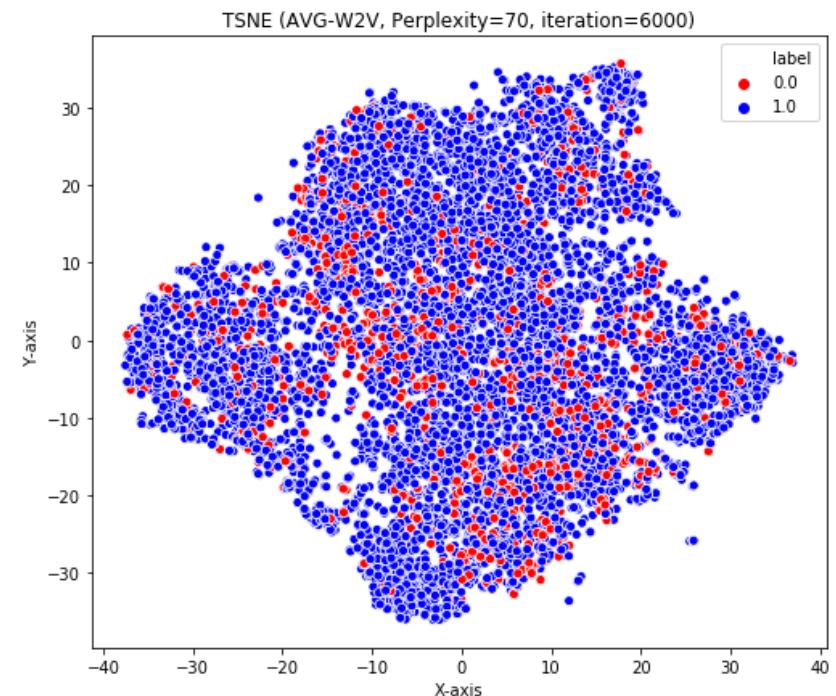
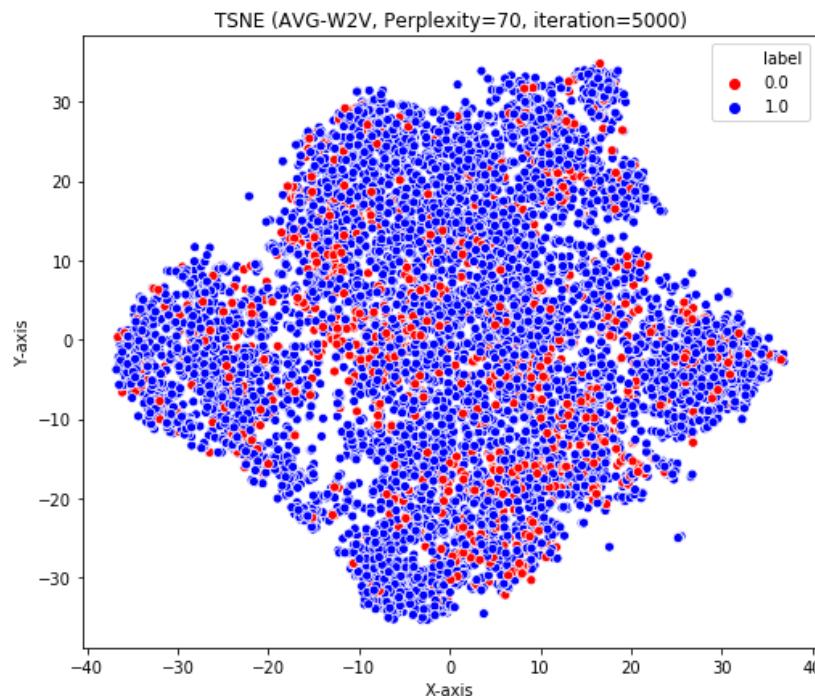
CPU times: user 3h 20min 10s, sys: 1h 32min 30s, total: 4h 52min 40s  
Wall time: 18min 52s

## AVG-W2V

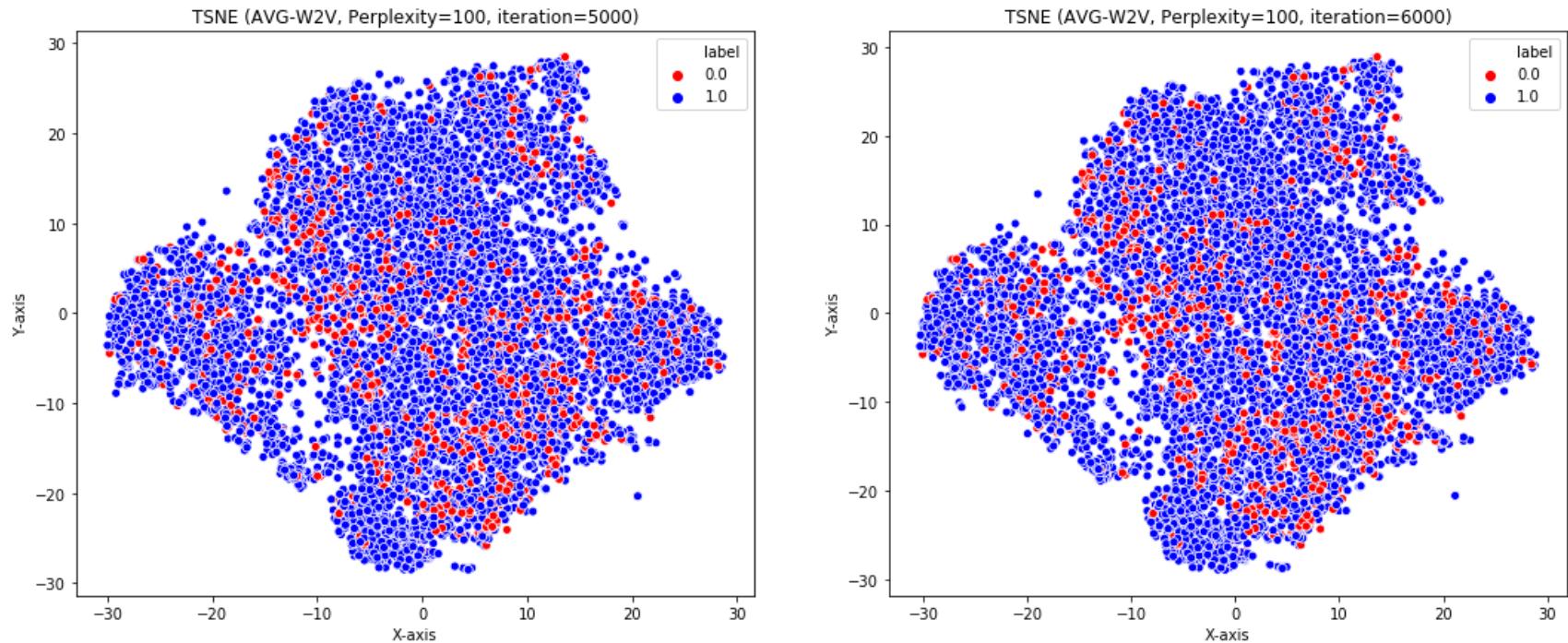
In [13]:

```
1 print('1st run :\n')
2 data=std_data(avg_sent_vectors,mean=True)
3 %time tsne_plots(data, y, params, vect[2])
```

1st run :



## TSNE\_AFFR\_8K\_POINTS



CPU times: user 1h 28min 29s, sys: 1h 16min 14s, total: 2h 44min 44s  
 Wall time: 10min 23s

### Observation:

1. from the above TSNE plots we observe that for :
  - [a.] perplexity=70, the plot almost stabilized when the iterations increased from 5000 to 6000.
2. from the above TSNE plots we observe that for :
  - [a.] perplexity=100, the plot is still changing its shape(not stable), when the iterations increased from 5000 to 6000.
3. so we choose optimal perplexity=70

### Note:

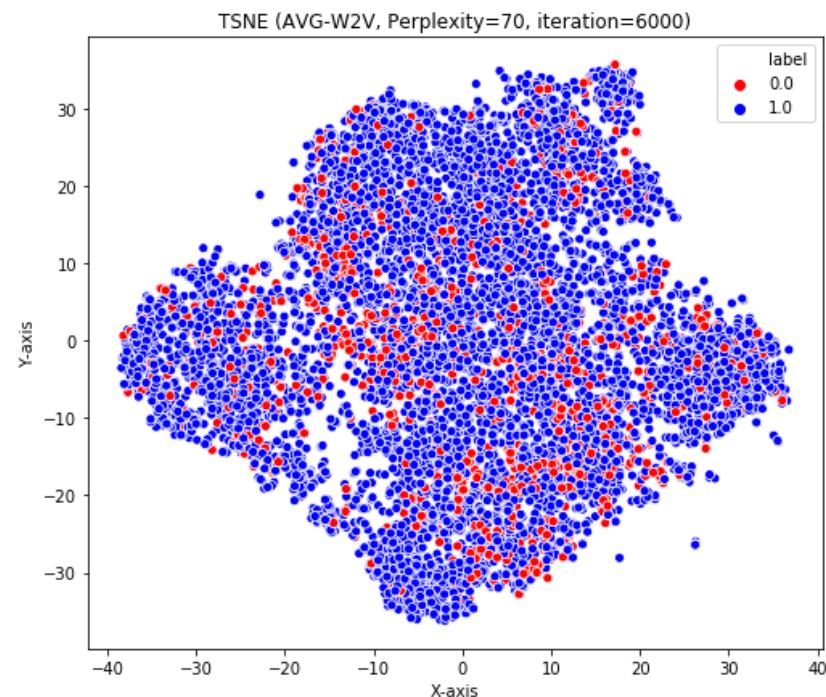
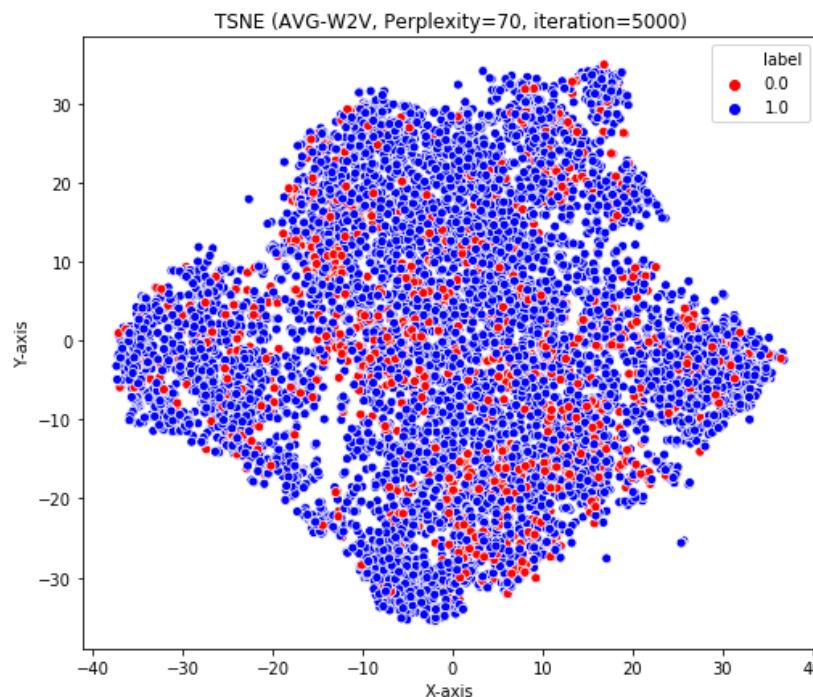
1. To check whether our result is correct or not we rerun TSNE.

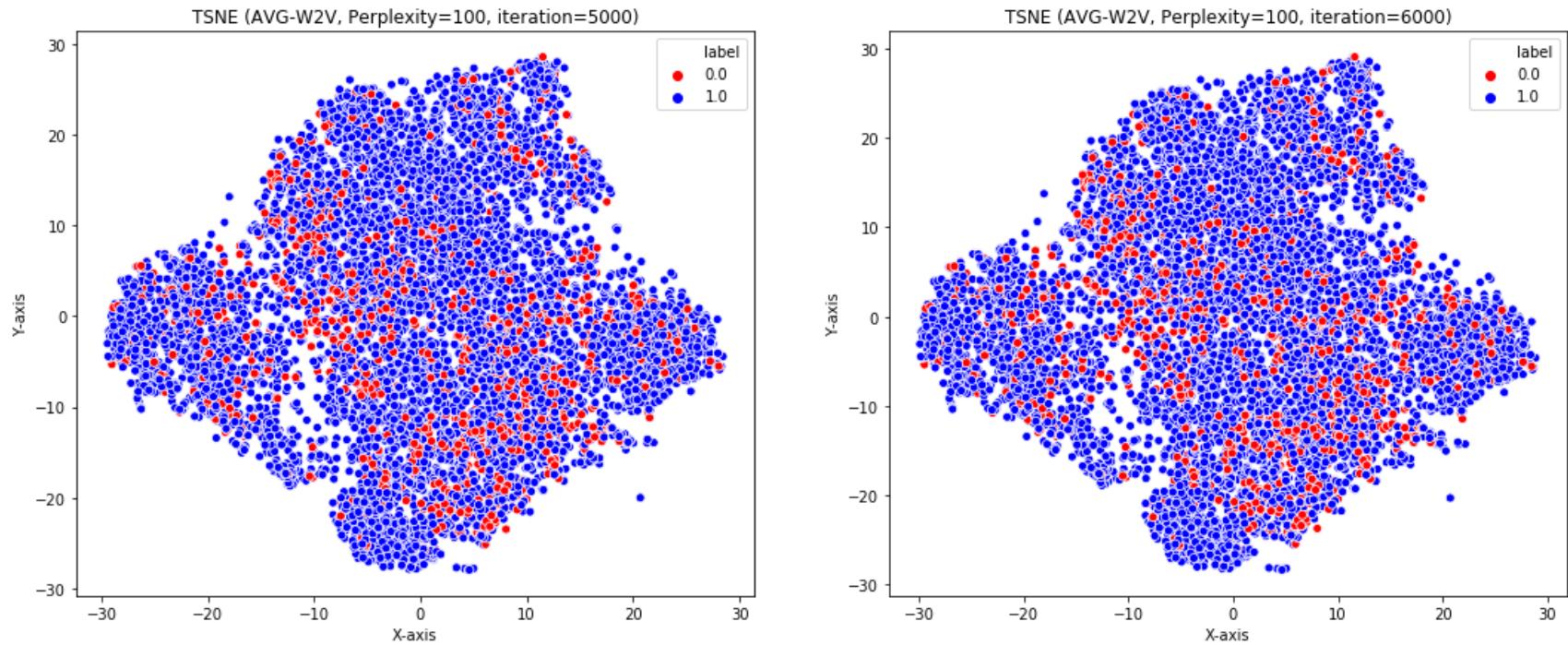


In [14]:

```
1 print('2nd run :\n')
2 %time tsne_plots(data, y, params, vect[2])
```

2nd run :





CPU times: user 1h 28min 21s, sys: 1h 16min, total: 2h 44min 22s  
 Wall time: 10min 22s

### Observation:

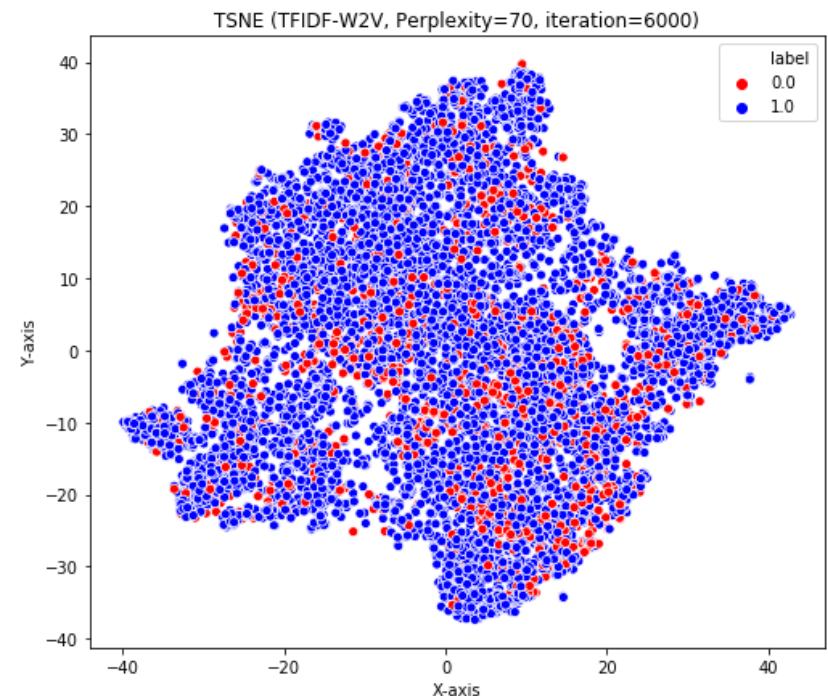
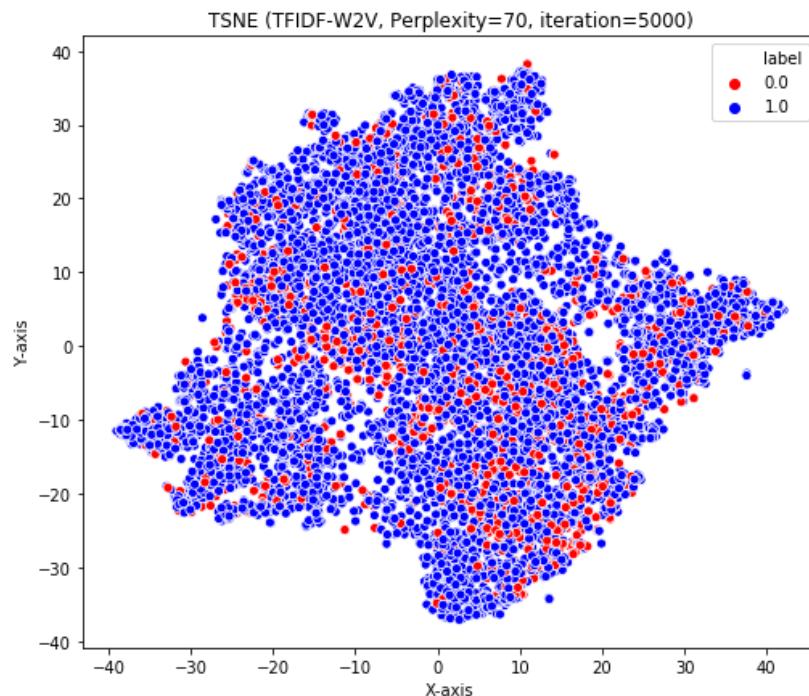
1. from the above TSNE plots we observe that for multiple run of TSNE we get same plot so finally we choose optimal perplexity as:  
 [a.] perplexity=70

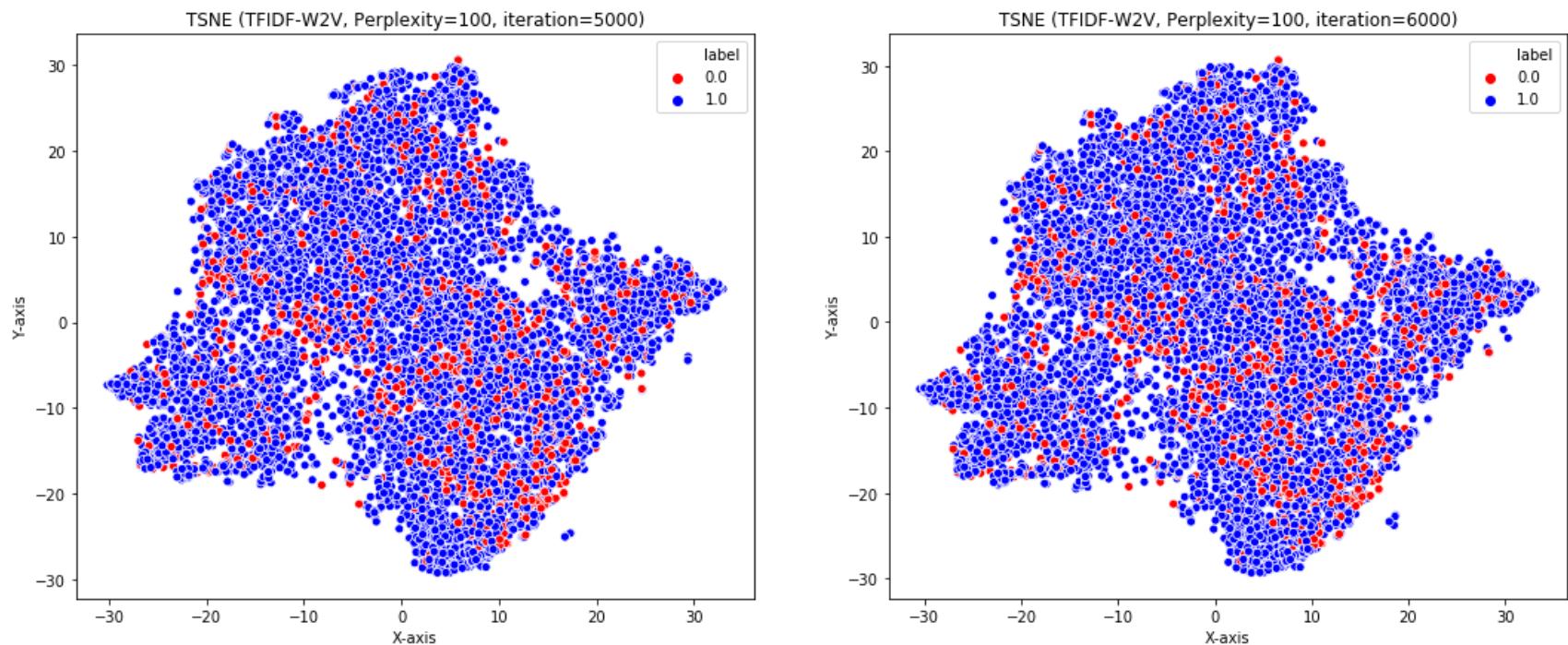
## TFIDF-W2V

In [15]:

```
1 print('1st run :\n')
2 data=std_data(tfidf_sent_vectors,mean=True)
3 %time tsne_plots(data, y, params, vect[3])
```

1st run :





CPU times: user 1h 29min 8s, sys: 1h 16min 20s, total: 2h 45min 29s  
 Wall time: 10min 26s

## Observation:

1. from the above TSNE plots we observe that for :
  - [a.] perplexity=70, the plot almost stabilized when the iterations increased from 5000 to 6000.
2. from the above TSNE plots we observe that for :
  - [a.] perplexity=100, the plot is still changing its shape(not stable), when the iterations increased from 5000 to 6000.
3. so we choose optimal perplexity=70

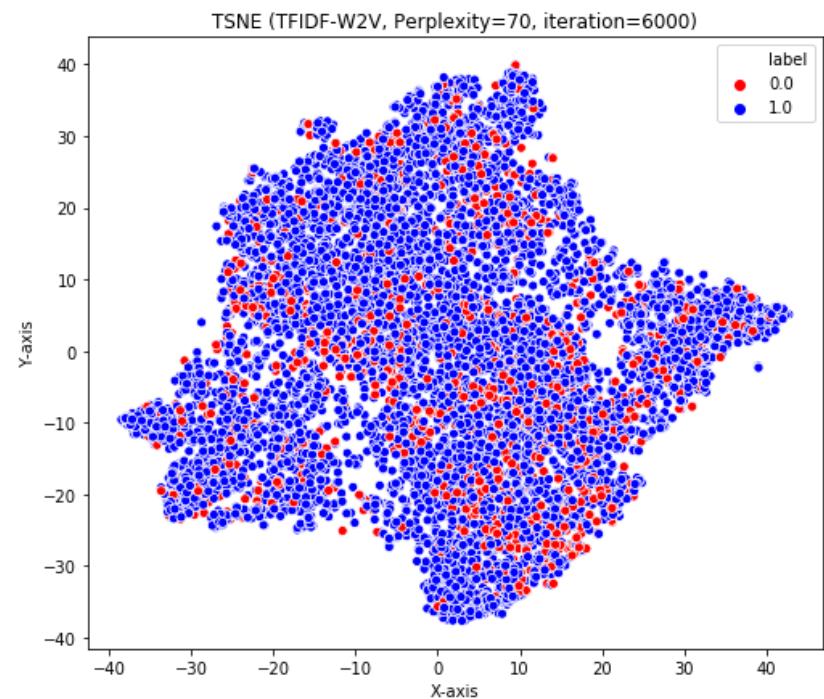
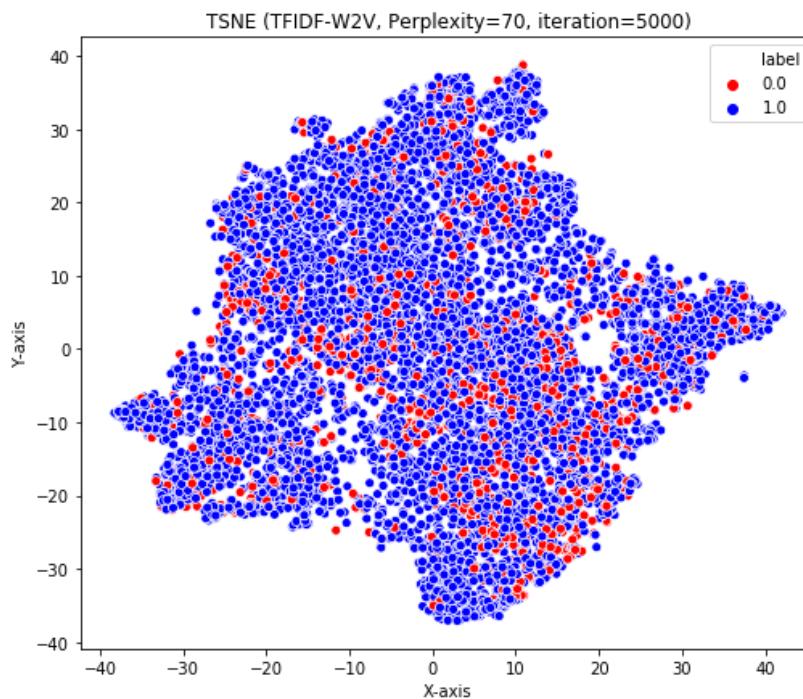
## Note:

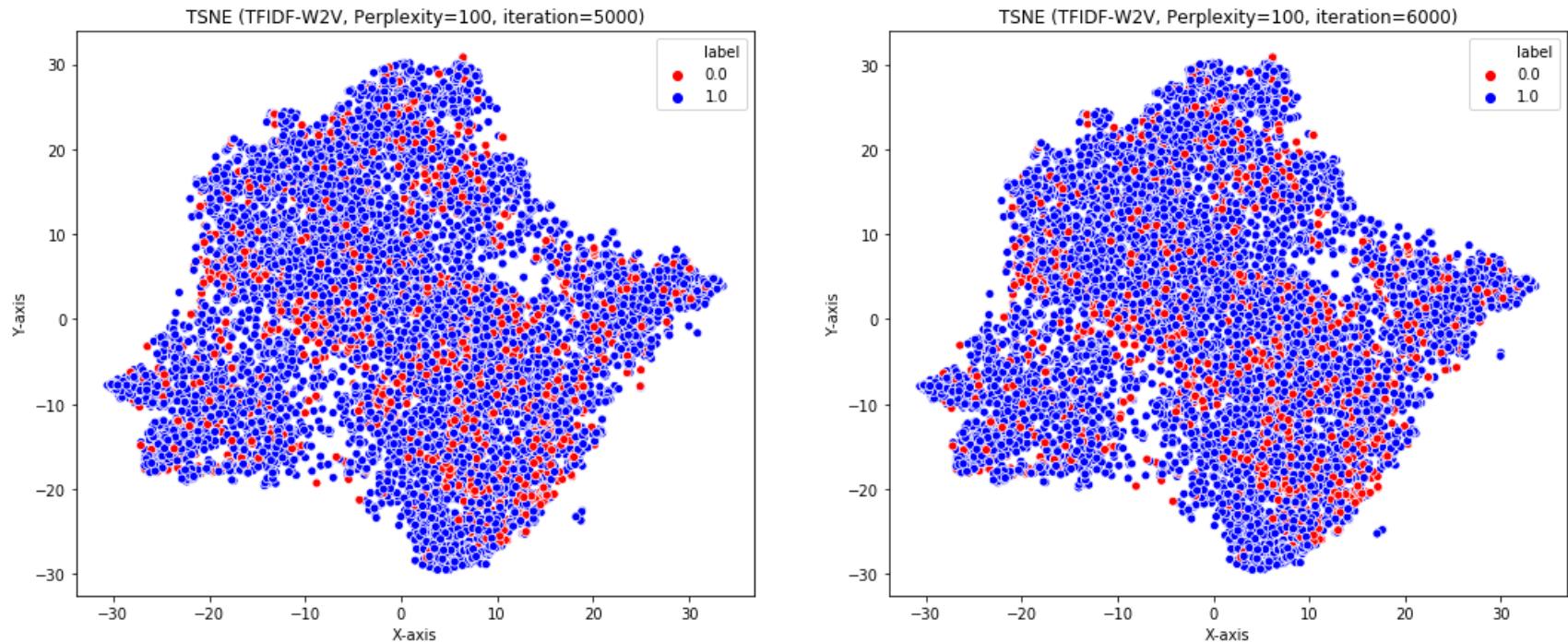
1. To check whether our result is correct or not we rerun TSNE.



```
In [16]: 1 print('2nd run :\n')
2 %time tsne_plots(data, y, params, vect[3])
```

2nd run :





CPU times: user 1h 29min 7s, sys: 1h 16min 15s, total: 2h 45min 22s  
 Wall time: 10min 25s

## Observation:

1. from the above TSNE plots we observe that for multiple run of TSNE we get same plot so finally we choose optimal perplexity as:  
 [a.] perplexity=70

## NOTE:

1. For clear explanation of perplexity tuning saw AVGW2V and TFIDFW2V section where i clearly explained perplexity tuning.
2. As the TSNE is computationally expensive algorithm tuning perplexity for BOW and TFIDF is not easy, for that reason i am just showing approximate value of perplexity .

## Reference Links:

1. <https://distill.pub/2016/misread-tsne/> (<https://distill.pub/2016/misread-tsne/>)
2. <http://colah.github.io/posts/2014-10-Visualizing-MNIST/> (<http://colah.github.io/posts/2014-10-Visualizing-MNIST/>)
3. <https://www.appliedaicourse.com/> (<https://www.appliedaicourse.com/>)