

CSCI 5922 Fall 2022–Lab 1

Abulitibu Tuguluke

September 11, 2022

The Data I used

Through the semester, I would probably used data generated from work, for this lab particularly, I used job posted from dice.com categorized tech job (e.g.: software engineer, data scientist) and non-tech job (e.g.: nurse, writer). It was carefully hand labeled by SMEs. I want to use Feed-forward network to predict whether a job post a tech job or not.

Table 1: Data Description

	text	accept	answer
count	968	968	968
unique	968	2	1
		Tech	accept
count	1	876	968
		Non-Tech	accept
count	1	92	968

Which stated there are 876 jobs labeled as **Tech**, and 92 as **NonTech**, which makes the total number of sample 968. Since I hand-picked the **accepted** answer, we only need two columns of information from this data-frame: **text** and **answer**. Here are the 2 examples:

- Non Tech (Job): Deliver with Uber Eats Deliver with Uber. Earn on your schedule. Work on your schedule. Deliver for a few hours in the mornings, every night, or just on weekends-it's up to you. Earn Good Money. You'll earn by bringing people the food they love from local restaurants. Choose your wheels. Use your car, scooter, or bike to make deliveries.* Delivery requirements Car delivery: Be at least 19 years old Have a 2-door or 4-door car made after 2000 Have a valid driver's license, vehicle registration and vehicle insurance Have at least one year of driving experience in the U.S. Bike delivery: Be at least 18 years old Have a state-issued ID or Driver's License When signing up be sure to choose 'Biking' under transportation method. Scooter delivery: Be at least 19 years old Have a valid driver's license and vehicle insurance Have a 2-wheel scooter made after 2000 Ready to get started? Sign up today and start earning. No experience necessary. But, if you have previous employment experience in delivery (such as a delivery driver, food service, food delivery, delivery runner, or courier) you may enjoy delivering with Uber Eats! *Vehicles allowed for delivery vary by city.
- Tech (Job): Hadoop Consultant In depth and strong knowledge and hands on experience Hadoop, Horton Works, Scoop, Oozie workflowsShould have done complex data ingestions in HadoopKnowledge in Data modeling and DWH naming patternData profiling knowledgeKnowledge in DWH and ETL concepts

Word embedding (Pre-processing)

Since the data is in text form, we need to convert it into numeric, **Tech** label is converted as 1, **Non-Tech** as 0. For description, we first used TF-IDF method, yet the reslut is not idea, as shown in Figure 1,

There are no major improvement in the accuracy, my guess is that the TF-IDF vector is sparse, as we seen below:

[0. , 0.07048313, 0. , ..., 0. , 0. ,0.],

There is not much to learn, so we used another method similar to Word2Vec, called Doc2Vec, which embeds text into dense vector, e.g.:

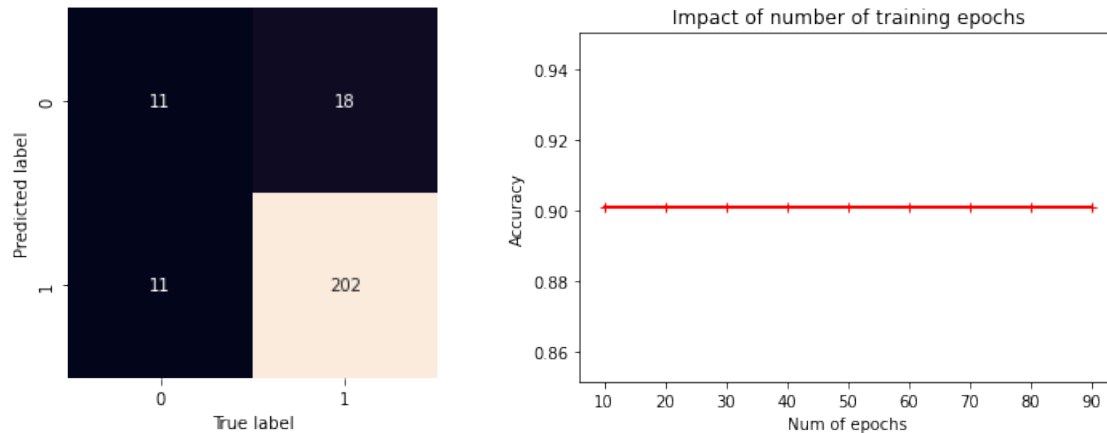


Figure 1: Tf-Idf emvedding for text.

```
array([ 0.69282454, -0.767151 , 1.1698718 , 0.0125812 , 0.04911822,
-0.6022604 , -0.7765408 , 0.14662024, -0.5277134 , -1.1415157 ,
0.9658156 , 0.12996148, -0.7017618 , 0.25222412, 0.2311553 ,
-0.4849749 , 0.460774 , -0.6721754 , 1.1252933 , 0.5912803 ,
1.066204 , 0.91162 , -0.13604069, -0.46164566, -0.22352642,
-0.02144009, -0.55692935, 0.12265404, -0.6255455 , 0.02748431,
-0.05019109, 0.65613216, 1.3773457 , 0.55956566, 0.5098652 ,
1.1547431 , -0.18142179, 0.1515532 , -0.4199732 , 0.51574713,
0.3102894 , 1.2909027 , 0.5681546 , 0.90260255, -0.82004976,
-0.76896405, -0.3467222 , 0.12360989, -0.8401419 , 0.09505583,
1.5170844 , 0.887022 , 0.88092905, -0.45412895, -0.07270601,
0.00651041, 0.0348754 , 0.5411873 , 0.07916429, -0.17342937,
0.39807764, -0.91081697, -0.55455333, -0.6905511 ], dtype=float32)
```

with which we will be using in this assignment.

1 Neural Network Hyper-parameters

We ran 9 experiments: 3 activation functions of **ReLU**, **Tanh**, **Logistic(Sigmoid)**, with 3 sets of hidden layers of 10×10 , 20×20 , 30×30 . We set rest of the Hyper-parameters as default in `MLPClassifier`, such as: batch size=200, learning rate=0.001, max iteration=200, optimization method=adam. The result is shown in Figure 2:

According to the accuracy data, all 3 layers with **ReLU** performs the best with 0.92 accuracy even with less neuron layers. **Tanh** is the least accurate with 0.90 after all three layers. But what I found more interesting is the True-Positive for **Non-Tech(0)** for **Sigmoid**, since the first 2 set of layers of 10×10 , 20×20 did not able to spot a single **Non-Tech(0)** correctly until 30×30 , I suspect is it due to the high cost of non-linearity of the logistic functions which demands more computing to get to a better result. So my bets would be on **Tanh** performs better on this particular data set due to it's ability to label **Non-Tech(0)** correctly than **Sigmoid**, which confirms the textbook's claim that **ReLU** is the most popular activation function, for now.

Also, if we ignore the activation function, and solely focused on Layer architecture, then the bigger layer the better results, 30×30 performs the best, which again confirms the book's claim that more computing yields better result. The Evaluation Metrics in Table 2 tells a similar story, but if we focus more on **Non-Tech(0)**'s precision, then clearly **ReLU** with 20 layers performs better than **ReLU** with 30 layers, I am still not sure if this is due to 'stretch pants' approach, which mean 30 layers is a bit over-fitting, hence the accuracy is dropped by 0.01 point, we can not confirm this until we try with a bigger size layer.

Judging by Figure 3 and Table 3, that may be the case, as we can see the accuracy dropped by 0.02 while increasing the layer size for all three activation functions, along with the precision for **Non-Tech(0)**. And increasing the size does not matter that much after certain size.

	Non-Tech(0)	Tech(1)	accuracy
Activaton Function: relu, Hidden layer: 10	{'precision': 0.5454545454545454, 'recall': 0....	{'precision': 0.9405204460966543, 'recall': 0....	0.910653
Activaton Function: logistic, Hidden layer: 10	{'precision': 0.0, 'recall': 0.0, 'f1-score': ...	{'precision': 0.9037800687285223, 'recall': 1....	0.90378
Activaton Function: tanh, Hidden layer: 10	{'precision': 0.6, 'recall': 0.642857142857142...	{'precision': 0.9616858237547893, 'recall': 0....	0.924399
Activaton Function: relu, Hidden layer: 20	{'precision': 0.7619047619047619, 'recall': 0....	{'precision': 0.9555555555555556, 'recall': 0....	0.941581
Activaton Function: logistic, Hidden layer: 20	{'precision': 0.0, 'recall': 0.0, 'f1-score': ...	{'precision': 0.9037800687285223, 'recall': 1....	0.90378
Activaton Function: tanh, Hidden layer: 20	{'precision': 0.48148148148148145, 'recall': 0....	{'precision': 0.9431818181818182, 'recall': 0....	0.900344
Activaton Function: relu, Hidden layer: 30	{'precision': 0.6818181818181818, 'recall': 0....	{'precision': 0.9516728624535316, 'recall': 0....	0.931271
Activaton Function: logistic, Hidden layer: 30	{'precision': 0.6363636363636364, 'recall': 0....	{'precision': 0.9479553903345725, 'recall': 0....	0.924399
Activaton Function: tanh, Hidden layer: 30	{'precision': 0.5, 'recall': 0.464285714285714...	{'precision': 0.9433962264150944, 'recall': 0....	0.90378

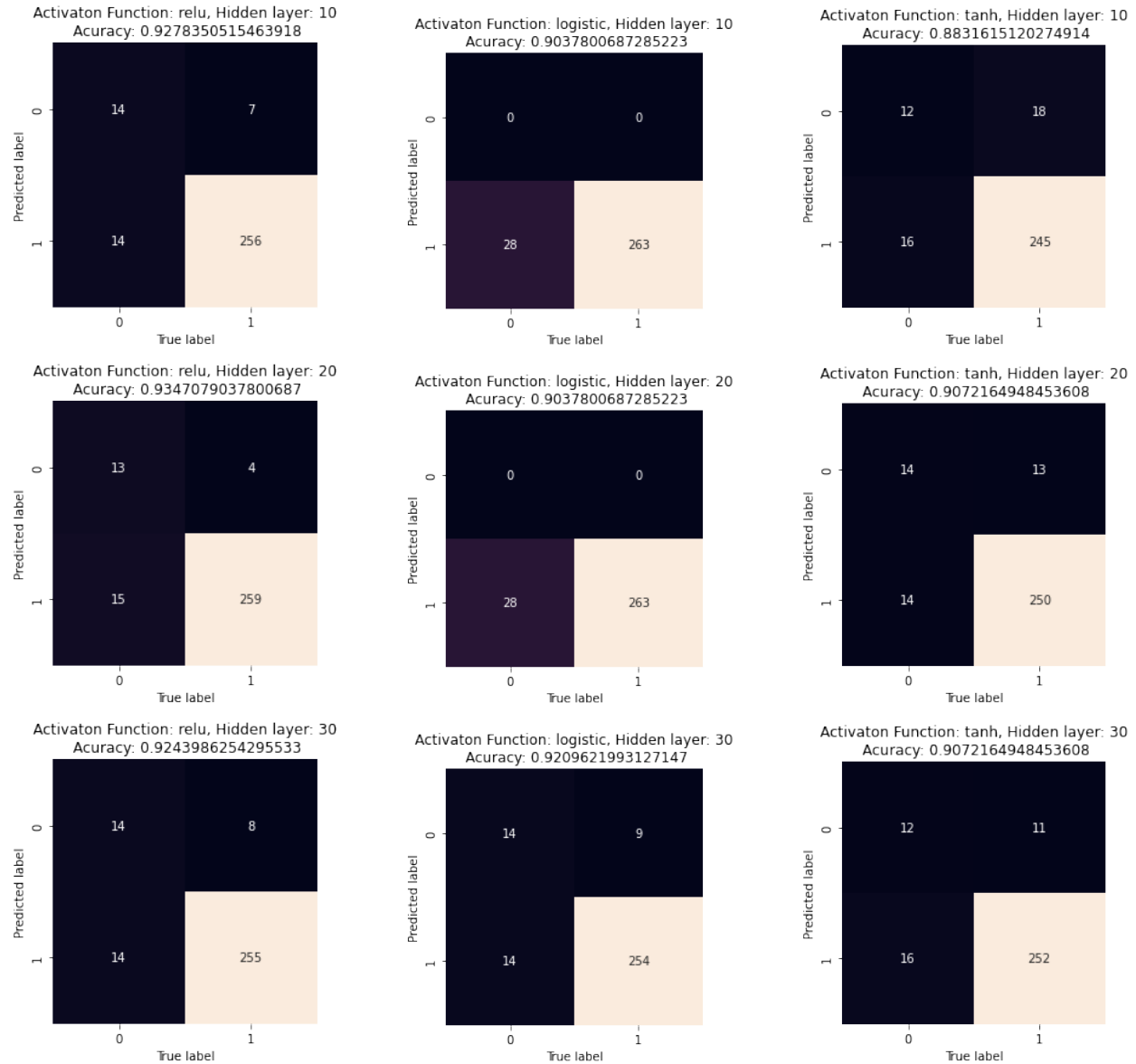
Table 2: Classification report for hidden layers of 10×10 , 20×20 , 30×30 

Figure 2: Part 1 Neural Network Hyper-parameters: 3 activation functions of **ReLU**, **Tanh**, **Logistic(Sigmoid)**, with 3 sets of hidden layers of 10×10 , 20×20 , 30×30 , Hyper-parameters as default in MLPClassifier, such as: batch size=200, learning rate=0.001, max iteration=200, optimization method=adam.

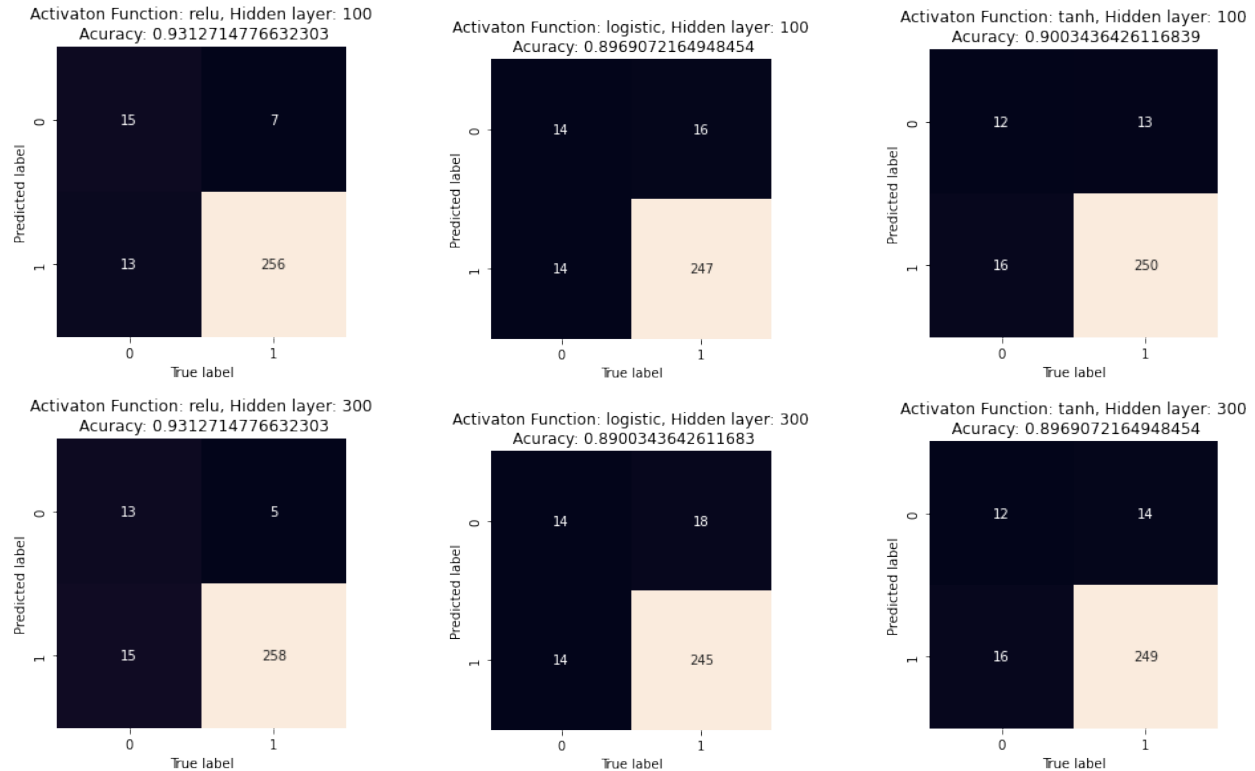


Figure 3: Part 1 Neural Network Hyper-parameters: 3 activation functions of **ReLU**, **Tanh**, **Logistic(Sigmoid)**, with 2 sets of hidden layers of 100×100 , 300×300 , Hyper-parameters as default in MLP-Classifer, such as: batch size=200, learning rate=0.001, max iteration=200, optimization method=adam.

	0	1	accuracy
Activaton Function: relu, Hidden layer: 100	{'precision': 0.6818181818181818, 'recall': 0....	{'precision': 0.9516728624535316, 'recall': 0....	0.931271
Activaton Function: logistic, Hidden layer: 100	{'precision': 0.4666666666666667, 'recall': 0....	{'precision': 0.946360153256705, 'recall': 0.9...	0.896907
Activaton Function: tanh, Hidden layer: 100	{'precision': 0.48, 'recall': 0.42857142857142...	{'precision': 0.9398496240601504, 'recall': 0....	0.900344
Activaton Function: relu, Hidden layer: 300	{'precision': 0.7222222222222222, 'recall': 0....	{'precision': 0.945054945054945, 'recall': 0.9...	0.931271
Activaton Function: logistic, Hidden layer: 300	{'precision': 0.4375, 'recall': 0.5, 'f1-score...	{'precision': 0.9459459459459459, 'recall': 0....	0.890034
Activaton Function: tanh, Hidden layer: 300	{'precision': 0.46153846153846156, 'recall': 0...	{'precision': 0.939622641509434, 'recall': 0.9...	0.896907

Table 3: Classification report for hidden layers of 100×100 , 300×300

2 Impact of Training Duration and Training Data

We use the same data set as part 1 with the same random 70/30 splits and Doc2Vec embedding, while setting rest of the parameters as constants: hidden layer =100, batch size=200, learning rate=0.001, max iteration=200, optimization method=adam, and activation function with **ReLU** and epoch from 10 to 200 with 10 increment each iteration. The learning curves are shown in Figure 4. Before I ran this experiment, I assume that train with 20%, 40%, 60%, 80%, and 100% of the training data will gradually be increased by the accuracy rate, yet the results they are pretty much converge to 0.94 to 0.95. What is interesting is the accuracy during the iterations: with 20% started with the lowest accuracy around 0.60 to full training set starting with 0.92. My guess is that with lesser the training, it take the model longer time to learning from, but eventually improving the accuracy with each epoch, this again confirms the textbook's claim that the more data the better, along with bigger iteration.

We can clearly see the influence of bigger the data percentage selected for the training the better accuracy showing in the early stage. The question of influence of the training duration not only depend on epoch size but also the maximum iteration setting. As shown in Figure 5, after 200 epoch, the accuracy stay constant with a high starting point as we set the starting epoch as 50 instead of 10.

Additional notes

The **Doc2Vec** is already a matured supervised learning algorithm, so it helped this lab with majority of the effort since the embedding of the text info already been 'learned', by resending dense matrix. The data is skewed toward **Tech(1)**, so once we collect enough **Non-Tech(0)** into our training, or even tune down the volume of the Tech set, the precision for **Non-Tech(0)** substantially, in my opinion.

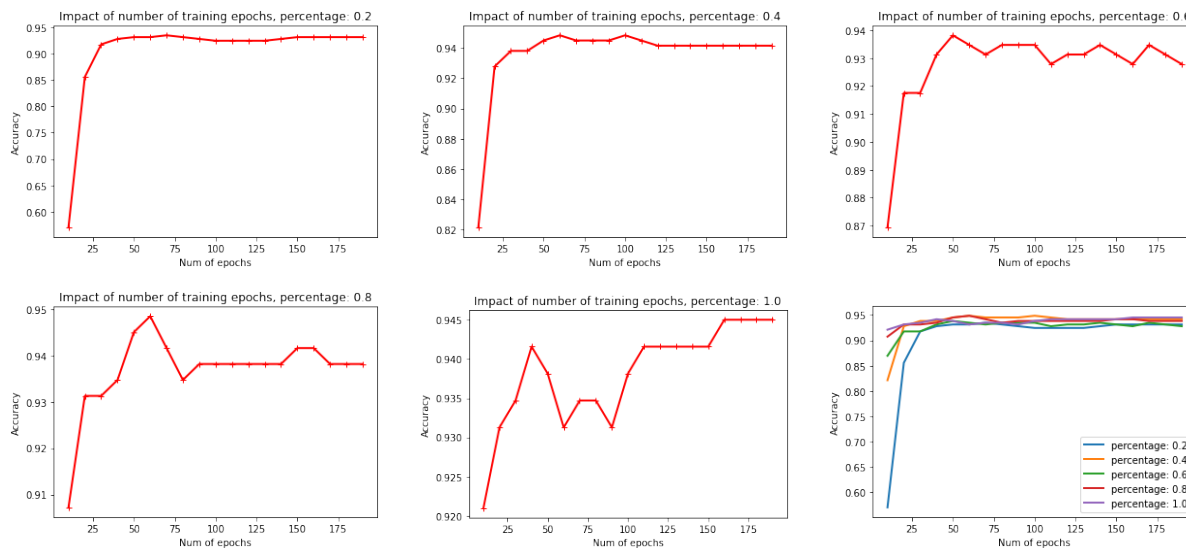


Figure 4: Five learning curves

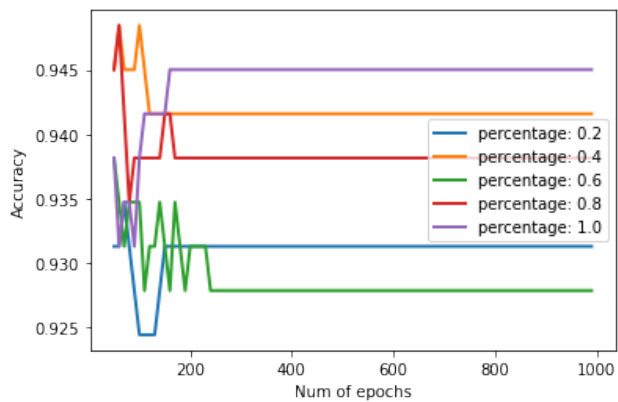


Figure 5: 50 to 500 epoch

Code

```

1 import matplotlib.pyplot as plt
2 import numpy as np
3 import pandas as pd
4 import seaborn as sns
5 from gensim.models.doc2vec import Doc2Vec, TaggedDocument
6 from sklearn import metrics
7 from sklearn.model_selection import train_test_split
8 from sklearn.neural_network import MLPClassifier
9 from sklearn.preprocessing import StandardScaler
10
11
12 def create_categories(x):
13     '''
14     change tech column to 1,
15     non-tech to 0
16     '''
17     if x == "Tech":
18         return 1
19     else:
20         return 0
21
22
23 def add_embedding(data):
24     '''
25     adding doc2vec column based on text
26     '''
27     # tokenize and tag the card text
28     card_docs = [TaggedDocument(doc.split(' '), [i])
29                  for i, doc in enumerate(data.text)]
30
31     # instantiate model
32     model = Doc2Vec(vector_size=64, window=2,
33                    min_count=1, workers=8, epochs=40)
34     # build vocab
35     model.build_vocab(card_docs)
36     # train model
37     model.train(card_docs, total_examples=model.corpus_count,
38               epochs=model.epochs)
39     card2vec = [model.infer_vector((data['text'][i].split(' ')))
40                for i in range(0, len(data['text']))]
41     # Create a list of lists
42     dtv = card2vec
43     # set list to dataframe column
44     data['doc2vec'] = dtv
45     return data
46
47
48 def split_data(test_ratio=None):
49     '''
50     split the data with
51     test_ratio input: .3
52     '''
53     return train_test_split(X, y, test_size=test_ratio, random_state=1)
54
55
56 def train_with_activation_funciton_nnodes(activation_function=None, n_hidden_nodes=None):
57     # Standardize data
58     stdsc = StandardScaler()
59     stdsc.fit(X_train)
60     X_train_std = stdsc.transform(X_train)
61     X_test_std = stdsc.transform(X_test)
62     # activation{'identity', 'logistic', 'tanh', 'relu'}, default='relu'
63     mlp = MLPClassifier(activation=activation_function, hidden_layer_sizes=[
64                        n_hidden_nodes, n_hidden_nodes])
65
66     mlp.fit(X_train_std, y_train)

```

```

67 y_predicted = mlp.predict(X_test_std)
68 # other metric
69 report = metrics.classification_report(
70     y_test, y_predicted, output_dict=True)
71 # plot
72 mat = metrics.confusion_matrix(y_test, y_predicted)
73 sns.heatmap(mat.T, square=True, annot=True, fmt='d', cbar=False)
74 plt.xlabel("True label")
75 plt.ylabel("Predicted label")
76 plt.title("Activaton Function: {}, Hidden layer: {}\n Acuracy: {}".format(
77     activation_function, n_hidden_nodes, mlp.score(X_test_std, y_test)))
78 plt.show()
79
80 return report
81
82
83 def training_by_percentage(percentage=None):
84     start_num_epochs = 10
85     finish_num_epochs = 200
86     inc_amt = 10
87
88     pred_scores = []
89     num_epochs = []
90     # standard then slice
91     X_train_slice = X_train[:round(percentge*len(X_train))]
92     X_train_std_slice = stdsc.transform(X_train_slice)
93     X_test_std = stdsc.transform(X_test)
94     y_train_slice = y_train[:round(percentge*len(y_train))]
95
96     for epoch_count in range(start_num_epochs, finish_num_epochs, inc_amt):
97         my_classifier = MLPClassifier(
98             activation='relu', random_state=20, max_iter=epoch_count)
99         my_classifier.fit(X_train_std_slice, y_train_slice)
100         score = my_classifier.score(X_test_std, y_test)
101         pred_scores.append(score)
102         num_epochs.append(epoch_count)
103
104     # plt.figure()
105     plt.plot(num_epochs, pred_scores, 'r--', linewidth=2)
106     plt.xlabel("Num of epochs")
107     plt.ylabel("Accuracy")
108     plt.title(
109         "Impact of number of training epochs, percentage: {}".format(percentge))
110     plt.show()
111
112
113 if __name__ == "__main__":
114
115     tech_data = pd.read_csv('tech_nontech_classification_annotations.csv')
116     tech_data = add_embedding(tech_data)
117     tech_data["binary"] = tech_data['accept'].apply(create_categories)
118     X = (np.vstack(np.array(tech_data['doc2vec']))))
119     y = np.array(tech_data['binary'])
120     X_train, X_test, y_train, y_test = split_data(test_ratio=.3)
121     # part1:
122     reports = {}
123     for hidden in [10, 20, 30]:
124         for act_func in ['relu', 'logistic', 'tanh']:
125             report = train_with_activation_funciton_nnodes(
126                 activation_function=act_func, n_hidden_nodes=hidden)
127             reports["Activaton Function: {}, Hidden layer: {}".
128                 .format(act_func, hidden)] = report
129     # print(pd.DataFrame(reports).T[['0', '1', 'accuracy']].to_latex())
130     # part2:
131     percent_slice = [.2, .4, .6, .8, 1.0]
132     for percentge in percent_slice:
133         training_by_percentage(percentge=percentge)

```