

IE406 Machine Learning

Lab Assignment - 7

Group 28

201801015: Shantanu Tyagi
201801076: Shivani Nandani
201801407: Pratvi Shah
201801408: Arkaprabha Baerjee

Question 1

Cars24 is the most popular website of used vehicles for sale, yet it's very difficult to collect all of them in the same place. Among all cities, data from 5 major cities which include Hyderabad, New Delhi, Mumbai, Bangalore, and Chennai is collected. Develop an algorithm for predict price of car.

Answer

For this dataset we performed preliminary analysis on the data and cleaned the data for making it a good fit for model creation. After that we standardized the data and implemented multiple models to find an algorithm which predicts price of car accurately.

Code

```
1 # Import libraries
2 import numpy as np
3 import pandas as pd
4 from sklearn.model_selection import train_test_split
5 from sklearn.linear_model import LogisticRegression
6 from sklearn.metrics import mean_squared_error
7 from sklearn.decomposition import PCA
8 import matplotlib.pyplot as plt
9 from sklearn.linear_model import SGDRegressor
10 from sklearn.pipeline import make_pipeline
11 from sklearn.preprocessing import StandardScaler
12 import tensorflow as tf
13 from tensorflow import keras
14 from tensorflow.keras.layers import Conv2D, MaxPooling2D, Dropout, Flatten, Dense
15 #Read data
16 df = pd.read_csv('Cars24.csv', index_col=0)
17 df.head()
18 df.dropna(inplace=True)
19
20 print(len(np.unique(df['Car Brand'])))
21 print(len(np.unique(df['Model'].astype(str))))
22 print(len(np.unique(df['Model Year'])))
23 print(len(np.unique(df['Location'])))
24 print(len(np.unique(df['Fuel'])))
25 print(len(np.unique(df['Gear'].astype(str))))
26
27 # Categorical to numeric
28 def getVecForm(vocab, df):
29     for i in range(len(vocab)):
30         df.replace(vocab[i], i, inplace=True)
31     return df
32
33 df['Fuel'] = getVecForm(np.unique(df['Fuel']), df['Fuel'])
34 df['Car Brand'] = getVecForm(np.unique(df['Car Brand']), df['Car Brand'])
35 df['Model'] = getVecForm(np.unique(df['Model'].astype(str)), df['Model'])
36 df['Location'] = getVecForm(np.unique(df['Location']), df['Location'])
37 df['Gear'] = getVecForm(np.unique(df['Gear'].astype(str)), df['Gear'])
38 print(df)
```

```

39
40 #Final data
41 X = df.drop(['Price'], axis=1).to_numpy()
42 y = df['Price'].to_numpy()
43
44 #Standardize
45 scaler = StandardScaler()
46 X = scaler.fit_transform(X)
47
48 # Logistic Regression
49 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
50 reg = LogisticRegression(random_state=0, solver='liblinear', max_iter=1000).fit(X_train,
    y_train)
51 y_pred = reg.predict(X_test)
52 reg.score(X_test, y_test)
53 print(mean_squared_error(y_test, y_pred, squared=False))
54 plt.figure(figsize=[10,8])
55 plt.plot(y_test, 'b', alpha=0.5, label='Actual')
56 plt.plot(y_pred, 'k', alpha=0.5, label='Predicted')
57 plt.legend()
58 plt.show()
59 plt.figure(figsize=[8,8])
60 plt.scatter(y_test, y_pred, color='b')
61 plt.xlabel('True Values [Price]')
62 plt.ylabel('Predictions [Price]')
63 plt.show()
64 error = abs(y_test-y_pred)
65 plt.figure(figsize=[8,8])
66 plt.hist(error, bins = 10, color='b', alpha=0.5)
67 plt.xlabel("Prediction Error [Price]")
68 plt.ylabel("Count")
69 plt.show()
70
71 # PCA followed by Logistic Regression
72 pca = PCA(n_components=X.shape[1])
73 pca.fit_transform(X)
74 ex_var_ratio = pca.explained_variance_ratio_
75
76 plt.figure(figsize=[8,8])
77 plt.plot(ex_var_ratio, color='b', linestyle='--', marker='x')
78 plt.show()
79 pca = PCA(n_components=3)
80 X_new = pca.fit_transform(X)
81 print(np.sum(pca.explained_variance_ratio_))
82 X_train, X_test, y_train, y_test = train_test_split(X_new, y, test_size=0.2, random_state=42)
83
84 reg = LogisticRegression(max_iter=1000).fit(X_train, y_train)
85 y_pred = reg.predict(X_test)
86 reg.score(X_test, y_test)
87 print(mean_squared_error(y_test, y_pred, squared=False))
88 plt.figure(figsize=[8,6])
89 plt.plot(y_test, 'b', alpha=0.5, label='Actual')
90 plt.plot(y_pred, 'k', alpha=0.5, label='Predicted')
91 plt.legend()
92 plt.show()
93 plt.figure(figsize=[8,8])
94 plt.scatter(y_test, y_pred, color='b')
95 plt.xlabel('True Values [Price]')
96 plt.ylabel('Predictions [Price]')
97 plt.show()
98 error = abs(y_test-y_pred)
99 plt.figure(figsize=[8,8])
100 plt.hist(error, bins = 10, color='b', alpha=0.5)
101 plt.xlabel("Prediction Error [Price]")
102 plt.ylabel("Count")
103 plt.show()
104
105 # Stochastic Gradient Descent
106 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
107 n_samples, n_features = X_train.shape[0], X_train.shape[1]
108 reg = make_pipeline(StandardScaler(),
109                     SGDRegressor(max_iter=1000, tol=1e-3))
110 reg.fit(X_train, y_train)
111 y_pred = reg.predict(X_test)
112 reg.score(X_test, y_test)

```

```

113 plt.figure(figsize=[10,8])
114 plt.plot(y_test, 'b', alpha=0.5, label='Actual')
115 plt.plot(y_pred, 'k', alpha=0.5, label='Predicted')
116 plt.legend()
117 plt.show()
118
119 plt.figure(figsize=[8,8])
120 plt.scatter(y_test, y_pred, color='b')
121 plt.xlabel('True Values [Price]')
122 plt.ylabel('Predictions [Price]')
123 plt.show()
124
125 error = abs(y_test-y_pred)
126 plt.figure(figsize=[8,8])
127 plt.hist(error, bins = 10, color='b', alpha=0.5)
128 plt.xlabel("Prediction Error [Price]")
129 plt.ylabel("Count")
130 plt.show()
131
132 # Neural Network
133 def build_model():
134     model = tf.keras.Sequential([
135         tf.keras.layers.Dense(256, activation='relu', input_shape=[X_train.shape[1]]),
136         tf.keras.layers.Dense(256, activation='relu'),
137         tf.keras.layers.Dense(256, activation='relu'),
138         tf.keras.layers.Dense(128, activation='relu'),
139         tf.keras.layers.Dense(128, activation='relu'),
140         tf.keras.layers.Dense(128, activation='relu'),
141         tf.keras.layers.Dense(64, activation='relu'),
142         tf.keras.layers.Dense(64, activation='relu'),
143         tf.keras.layers.Dense(1)
144     ])
145
146     optimizer = tf.keras.optimizers.Adam()
147
148     model.compile(loss='mse',
149                 optimizer=optimizer,
150                 metrics=['mae', 'mse'])
151     return model
152
153
154 class PrintDot(keras.callbacks.Callback):
155     def on_epoch_end(self, epoch, logs):
156         if epoch % 100 == 0: print('')
157         print('.', end='')
158
159 model = build_model()
160 history = model.fit(X_train,
161                    y_train,
162                    epochs=1000, validation_split = 0.2, verbose=0,
163                    callbacks=[PrintDot()])
164
165 hist = pd.DataFrame(history.history)
166 hist['epoch'] = history.epoch
167
168 def plot_history():
169     plt.figure(figsize=[10,8])
170     plt.xlabel('Epoch')
171     plt.ylabel('Mean Square Error')
172     plt.plot(hist.epoch[1:], hist.loss[1:], label='Train Error')
173     plt.plot(hist.epoch[1:], hist.val_loss[1:], label = 'Val Error')
174     plt.legend()
175
176 plot_history()
177 mse, _, _ = model.evaluate(X_test, y_test)
178 rmse = np.sqrt(mse)
179 print('Root Mean Square Error on test set: {}'.format(round(rmse, 3)))
180 y_pred = model.predict(X_test).flatten()
181 print(mean_squared_error(y_test, y_pred, squared=False))
182 plt.figure(figsize=[10,8])
183 plt.plot(y_test, 'b', alpha=0.5, label='Actual')
184 plt.plot(y_pred, 'k', alpha=0.5, label='Predicted')
185 plt.legend()
186 plt.show()
187

```

```

188 plt.figure(figsize=[8,8])
189 plt.scatter(y_test, y_pred, color='b')
190 plt.xlabel('True Values [Price]')
191 plt.ylabel('Predictions [Price]')
192 plt.show()
193
194 error = abs(y_test - y_pred)
195 plt.figure(figsize=[8,8])
196 plt.hist(error, bins = 10, color='b', alpha=0.5)
197 plt.xlabel("Prediction Error [Price]")
198 plt.ylabel("Count")
199 plt.show()

```

Listing 1: Question 1

Result

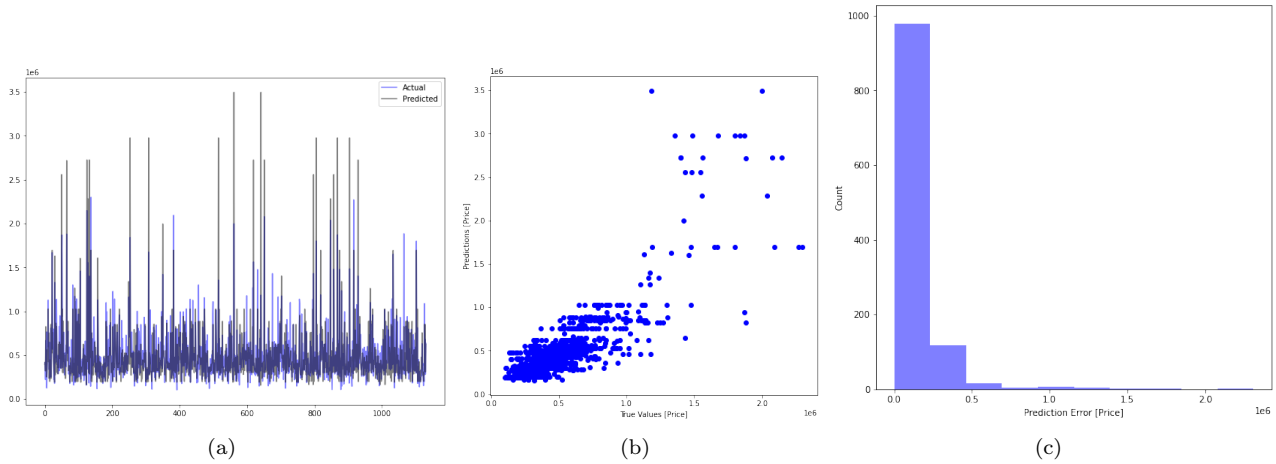


Figure 1: Logistic Regression

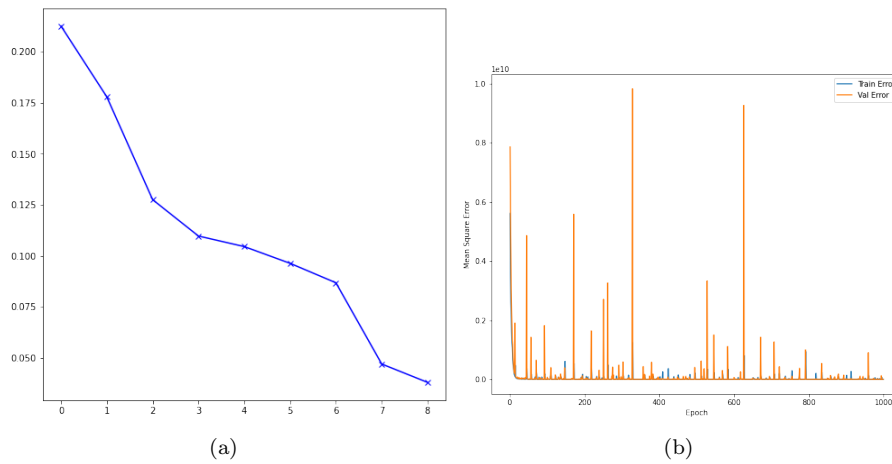


Figure 2: (a) Elbow curve for PCA (b) MSE v/s Epoch (Model loss plot)

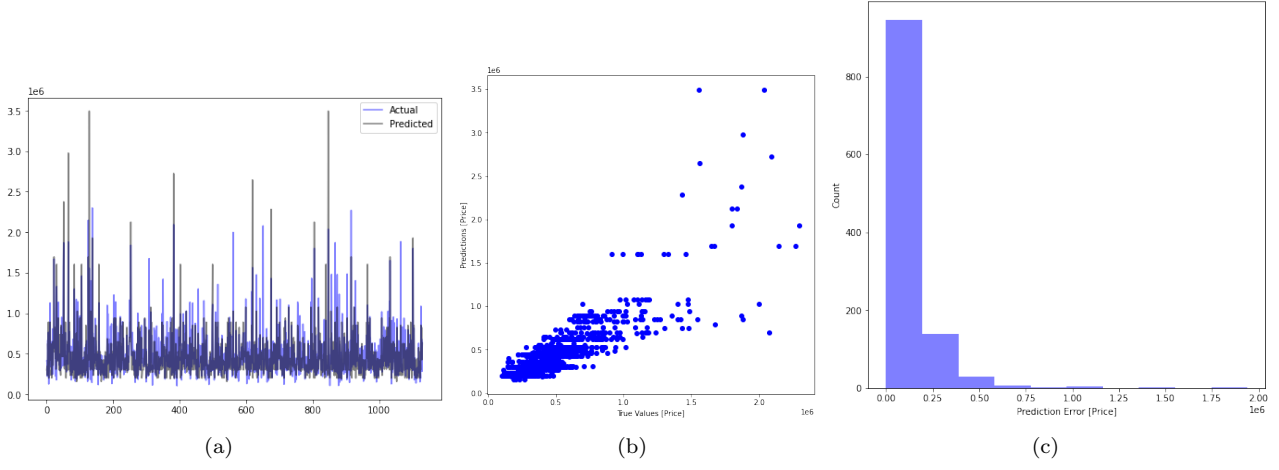


Figure 3: Logistic Regression post PCA

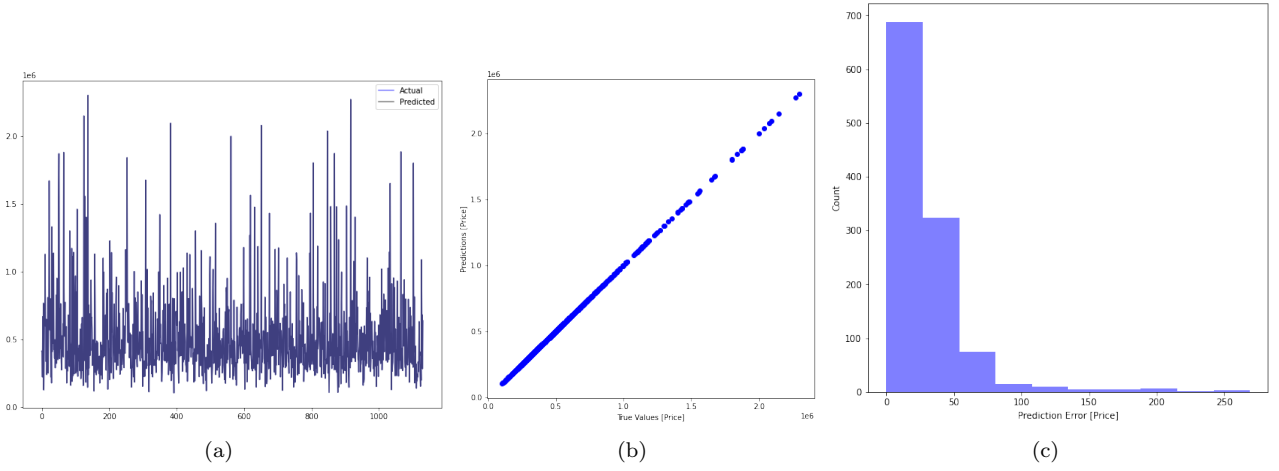


Figure 4: Stochastic Gradient Descent

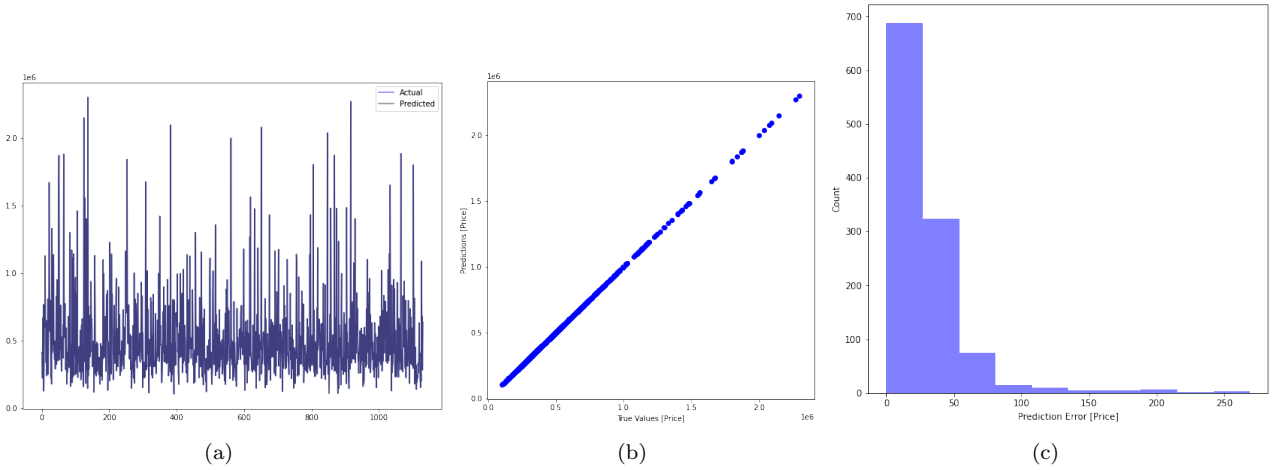


Figure 5: Neural Network

Observation/ Justification

We can make following observations from above:

1. From Figure.2(a), we can see that we do not get a clear cutoff for number of components that best describes the price and this can be due to the inherent categorical nature of majority of the features.
2. Figure.2(b) shows the Training and Validation set error for each epoch. It can be observed that, on an average, the error for Validation set for each epoch has a decreasing trend. This shows that the model is not over-fitting the data.

3. We could clearly observe from the plot (c) in the above figures that the error in prediction of car price is highest for Logistic regression model with and without PCA.
4. Due to the transformation from categorical to numeric data mean squared error after performing PCA is(= 189427.88) lesser than the original Logistic Regression model(= 223357.67).
5. From plot (b) in the above figures we can see that Stochastic GD and Neural Network models perform much better than original Logistic Regression model evident from the linear line whose slope is almost 1.

Question 2

Fashion-MNIST is a dataset of Zalando's article images—consisting of a training set of 60,000 examples and a test set of 10,000 examples. Each example is a 28x28 grayscale image, associated with a label from 10 classes. That means one image contains 784 pixel and pixel-value is an integer between 0 and 255. Make a classification model to classify the product.

Answer

The following classification models have been used :

1. Logistic Regression
2. Linear Discriminant Analysis
3. SVM
 - (a) Linear Kernel
 - (b) Polynomial Kernel
4. Neural Network architecture with 512 neurons in dense layer and 10 nodes in output layer. The Dense layer uses relu activation function and the final layer employs the softmax function. An adam optimizer has been used along categorical cross-entropy loss function.

Code

```
1 import numpy as np
2 import pandas as pd
3 import time
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 import plotly.graph_objects as go
7 from sklearn.model_selection import train_test_split
8 from sklearn.metrics import accuracy_score, confusion_matrix
9 from sklearn.decomposition import PCA
10 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
11 from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay, classification_report
12 from sklearn.preprocessing import StandardScaler
13 from sklearn.linear_model import LogisticRegression
14 from sklearn.naive_bayes import GaussianNB
15 from sklearn.ensemble import RandomForestClassifier
16 from sklearn.metrics import accuracy_score, recall_score, precision_score, f1_score
17 from sklearn import svm
18 import tensorflow as tf
19 import keras as keras
20
21 # Load data
22 training_images = pd.read_csv('fashion-mnist_train.csv')
23 training_labels = training_images['label']
24 training_images = training_images[training_images.columns[1:]]
25 test_images = pd.read_csv('fashion-mnist_test.csv')
26 test_labels = test_images['label']
27 test_images = test_images[test_images.columns[1:]]
28 training_images = StandardScaler().fit_transform(training_images)
29 test_images = StandardScaler().fit_transform(test_images)
30
31 # NN architecture with 512 neurons and 10 output nodes.
32 # Dense layer uses relu activation function and the final layer employs the softmax function.
33 # An adam optimizer has been used.
34 model = tf.keras.models.Sequential([tf.keras.layers.Flatten(), tf.keras.layers.Dense(512,
35                                     activation=tf.nn.relu), tf.keras.layers.Dense(10, activation=tf.nn.softmax)])
36 model.compile(optimizer = tf.optimizers.Adam(),
37               loss = 'sparse_categorical_crossentropy',
38               metrics=['accuracy'])
39 model.fit(training_images, training_labels, epochs=10)
40
41 model.evaluate(test_images, test_labels)
42 print(model.metrics_names)
43
44 # SVM
45 def plot_confusion_matrix(title, y_test, y_pred):
```

```

44
45     confusionMatrix = confusion_matrix(y_test,y_pred)
46     sns.heatmap(confusionMatrix,annot=True, fmt='.4g')
47     plt.ylabel('True',fontsize=15)
48     plt.xlabel('Predicted',fontsize=15)
49     plt.title(title,fontsize=15)
50     plt.show()
51     print(classification_report(y_test,y_pred))
52
53
54 def svm_func(x_train, x_test, y_train, y_test):
55     k = ['linear', 'poly']
56
57     for i in range(len(k)):
58         clf = svm.SVC(kernel=k[i])
59         clf.fit(x_train, y_train)
60         y_pred = clf.predict(x_test)
61         plot_confusion_matrix(k[i]+' kernel', y_test, y_pred)
62
63 svm_func(training_images, test_images, training_labels, test_labels)
64
65 # Logistic regression
66 clf = LogisticRegression(max_iter=10000).fit(training_images, training_labels)
67 accuracy_inbuilt = accuracy_score(test_labels, clf.predict(test_images))
68
69 cm = confusion_matrix(test_labels, pred_data)
70 disp = ConfusionMatrixDisplay(confusion_matrix=cm,display_labels=clf.classes_)
71 fig, ax = plt.subplots(figsize=(10,10))
72 disp.plot(ax=ax)
73 plt.show()
74 print(classification_report(test_labels,pred_data))
75
76 # LDA
77 clf = LinearDiscriminantAnalysis()
78 clf.fit(training_images, training_labels)
79 pred_data = clf.predict(test_images)
80
81 cm = confusion_matrix(test_labels, pred_data)
82 disp = ConfusionMatrixDisplay(confusion_matrix=cm,display_labels=clf.classes_)
83 fig, ax = plt.subplots(figsize=(10,10))
84 disp.plot(ax=ax)
85 plt.show()
86 print(classification_report(test_labels,pred_data))

```

Listing 2: Question 2

Result

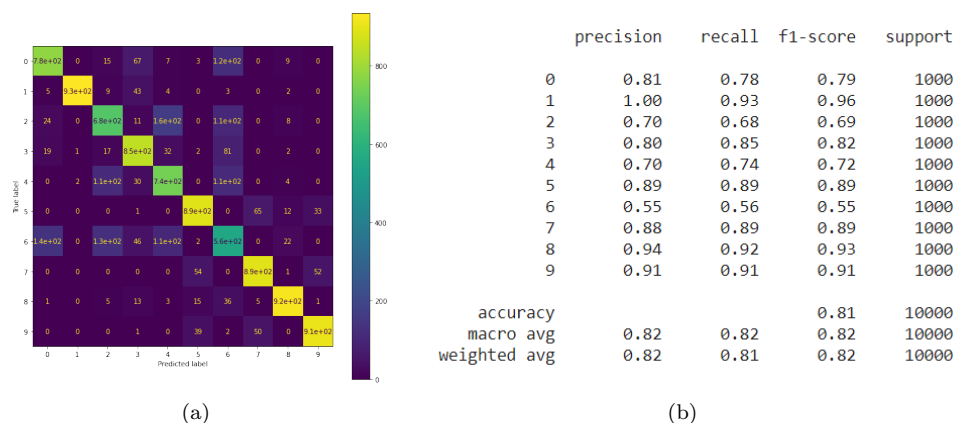
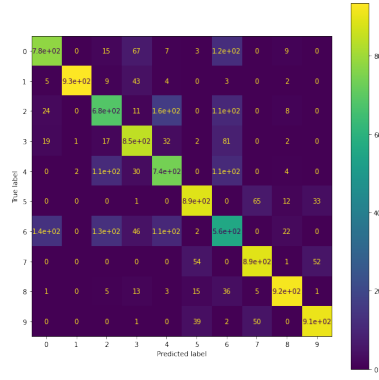


Figure 6: Logistic Regression

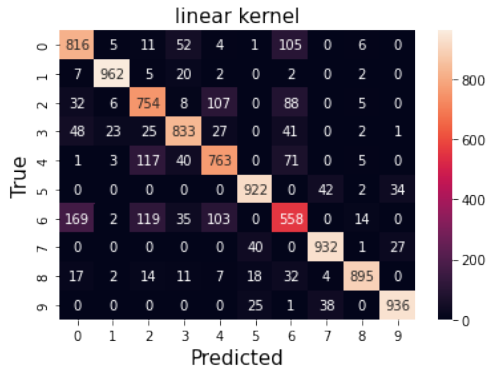


(a)

	precision	recall	f1-score	support
0	0.81	0.78	0.79	1000
1	1.00	0.93	0.96	1000
2	0.70	0.68	0.69	1000
3	0.80	0.85	0.82	1000
4	0.70	0.74	0.72	1000
5	0.89	0.89	0.89	1000
6	0.55	0.56	0.55	1000
7	0.88	0.89	0.89	1000
8	0.94	0.92	0.93	1000
9	0.91	0.91	0.91	1000
accuracy			0.81	10000
macro avg	0.82	0.82	0.82	10000
weighted avg	0.82	0.81	0.82	10000

(b)

Figure 7: Linear Discriminant Analysis

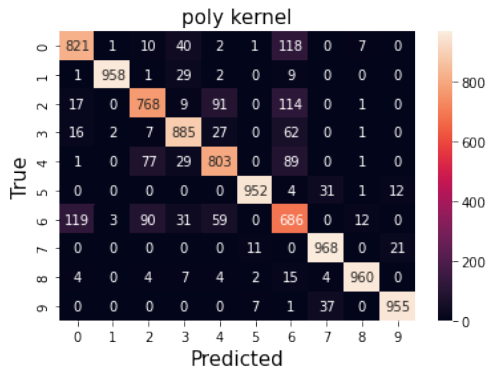


(a)

	precision	recall	f1-score	support
0	0.75	0.82	0.78	1000
1	0.96	0.96	0.96	1000
2	0.72	0.75	0.74	1000
3	0.83	0.83	0.83	1000
4	0.75	0.76	0.76	1000
5	0.92	0.92	0.92	1000
6	0.62	0.56	0.59	1000
7	0.92	0.93	0.92	1000
8	0.96	0.90	0.93	1000
9	0.94	0.94	0.94	1000
accuracy			0.84	10000
macro avg	0.84	0.84	0.84	10000
weighted avg	0.84	0.84	0.84	10000

(b)

Figure 8: SVM Linear kernel



(a)

	precision	recall	f1-score	support
0	0.84	0.82	0.83	1000
1	0.99	0.96	0.98	1000
2	0.80	0.77	0.78	1000
3	0.86	0.89	0.87	1000
4	0.81	0.80	0.81	1000
5	0.98	0.95	0.97	1000
6	0.62	0.69	0.65	1000
7	0.93	0.97	0.95	1000
8	0.98	0.96	0.97	1000
9	0.97	0.95	0.96	1000
accuracy			0.88	10000
macro avg	0.88	0.88	0.88	10000
weighted avg	0.88	0.88	0.88	10000

(b)

Figure 9: SVM Polynomial kernel

Model	Accuracy (Testing Dataset)
Logistic Regression	0.81
LDA	0.81
Neural Network	0.8721
SVM(linear kernel)	0.84
SVM(poly kernel)	0.88

Table 1: accuracy

Observation/ Justification

We could draw the following observations:

- From the above results we can see that polynomial kernel for support vector machine performs the best out of all other methods.
- The Neural network achieves a saturation point for its accuracy within 10 epochs.
- Logistic Regression and LDA perform the worst among all models chosen.
- Shirt (Number 6) has often been confused with T-shirt/top (Number 0), Pullover (Number 2) and Coat (Number 4).
- Coat (Number 4) has also often been confused with a pullover (Number 2).