

Name:

Jeffrey Hui

Netid:

jhui8

CS 441 - HW2: PCA and Linear Models

Complete the sections below. You do not need to fill out the checklist.

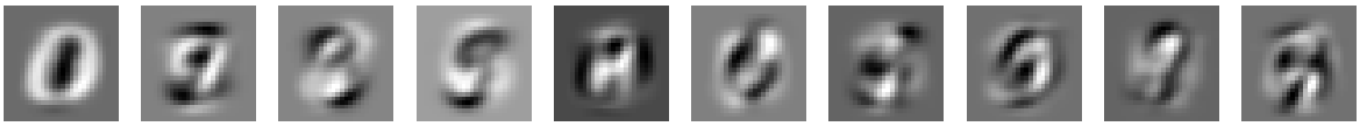
Total Points Available

[] / 160

1. PCA on MNIST
 - a. Display 10 principal component vectors [] / 5
 - b. Display scatterplot [] / 5
 - c. Plot cumulative explained variance [] / 5
 - d. Compression and 1-NN experiment [] / 15
2. MNIST Classification with Linear Models
 - a. LLR / SVM error vs training size [] / 20
 - b. Error visualization [] / 10
 - c. Parameter selection experiments [] / 15
3. Temperature Regression
 - a. Linear regression test [] / 10
 - b. Feature selection results [] / 15
4. Stretch Goals
 - a. PR and ROC curves [] / 10
 - b. Visualize weights [] / 10
 - c. Other embeddings [] / 15
 - d. One city is all you need [] / 15
 - e. SVM with RBF kernel [] / 10

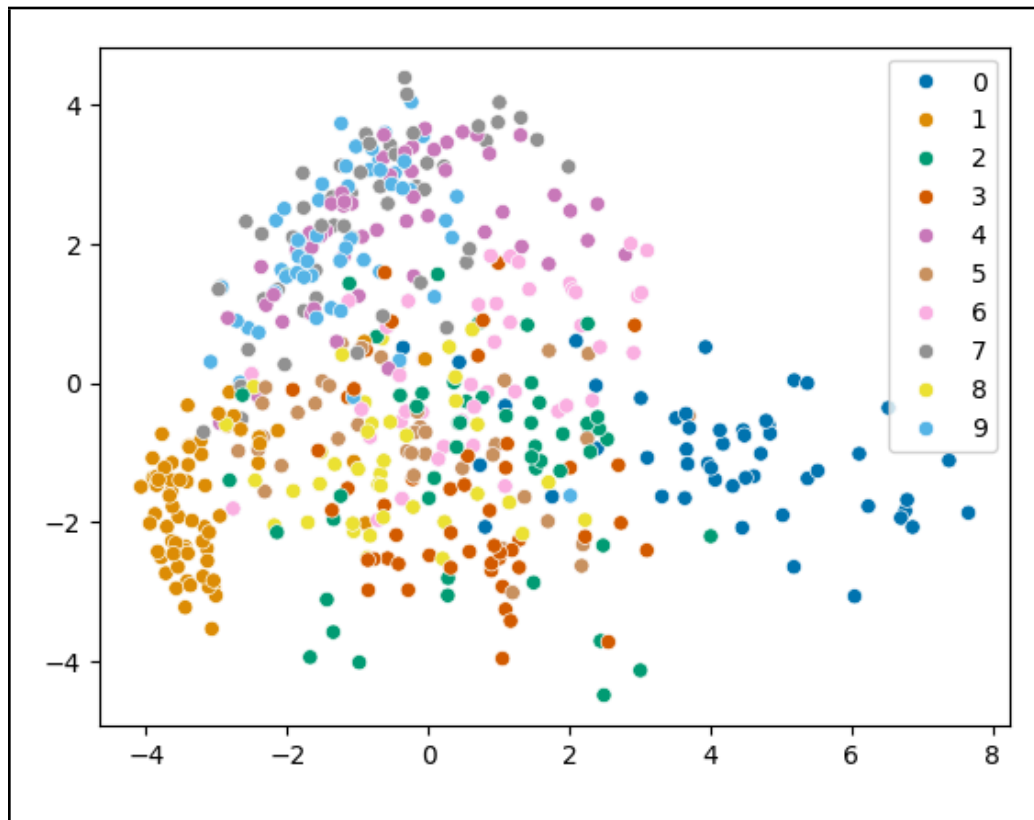
1. PCA on MNIST

a. Display 10 principal component vectors

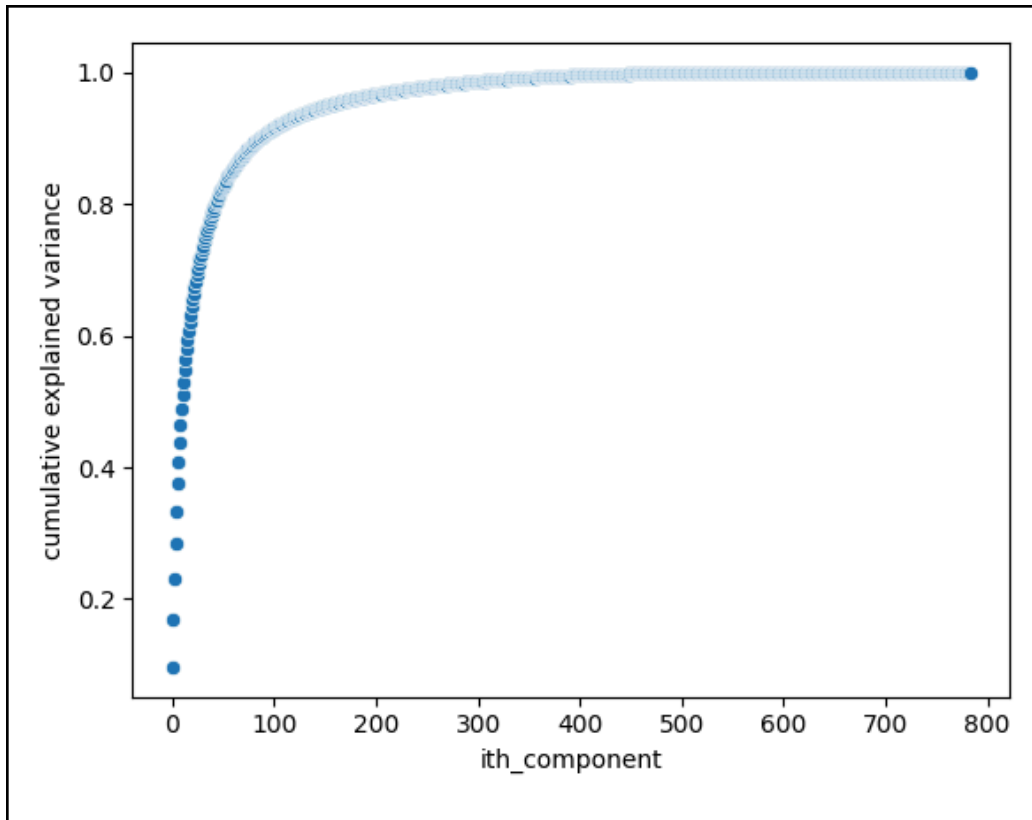


b. Display scatterplot

Scatterplot `x_train[:500]` for the first two PCA dimensions. Show a different color for each label.



c. Plot cumulative explained variance



d. Compression and 1-NN experiment

Number of components selected

	Total Time (s)	Test Error (%)	Dimensions
Brute Force (PCA)	1.28	2.68	87
Brute Force	5.49	3.09	784

2. MNIST Classification with Linear Models

a. LLR / SVM error vs training size

Test error (%)

# training samples	LLR	SVM
100	32.5%	32.4%
1,000	13.64	16.11

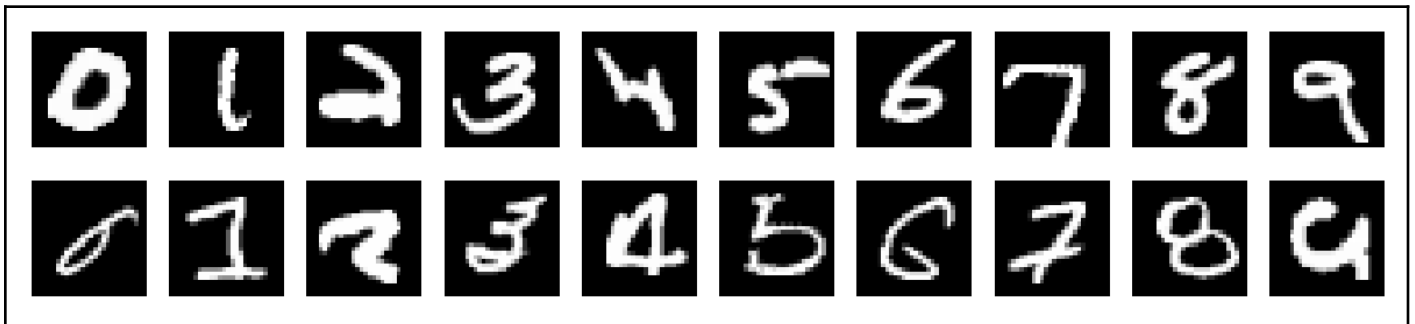
10,000	9.5	11.12
60,000	7.39	8.17

b. Error visualization

LLR



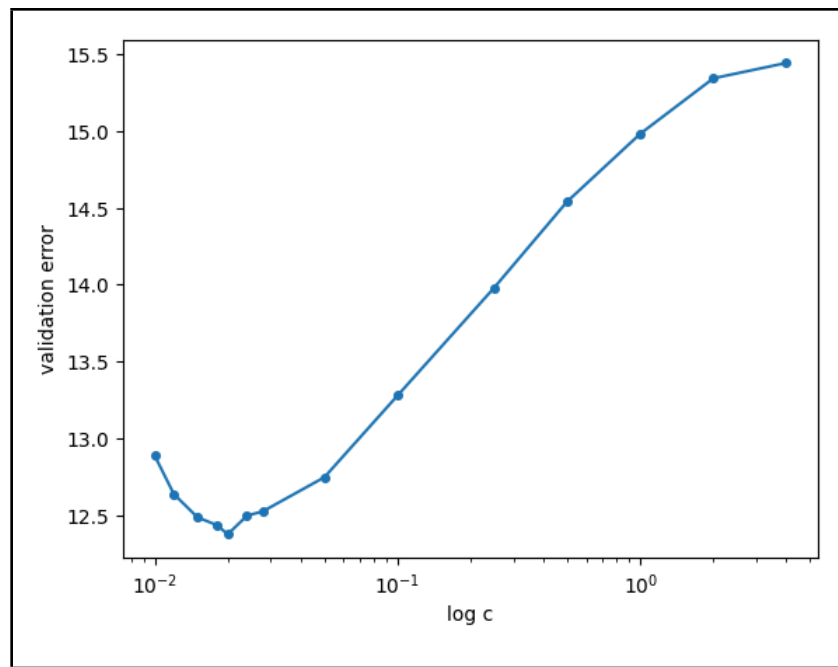
SVM



c. Parameter selection experiments

	Logistic Regression
Best C value	0.02
Validation error (%)	12.38
Test error (%)	13.57

Plot C value vs validation error for values tested



3. Temperature Regression

a. Linear regression test

Test RMSE

	Linear regression
Original features	2.161
Normalized features	2.163

Why might normalizing features in this way not be as helpful as it is for KNN?

Linear regression is not sensitive to the scale of the features as much as KNN. Additionally, if the relationship between the features and the target variable is non-linear, normalizing the data might not necessarily improve performance and could even lead to an increase in RMSE.

b. Feature selection results

Feature Rank	Feature number	City	Day
1	334	Chicago	-1
2	347	Minneapolis	-1
3	405	Grand Rapids	-1
4	336	Kansas City	-1
5	361	Cleveland	-1
6	307	Omaha	-2
7	367	Indianapolis	-1
8	264	Minneapolis	-2
9	9	Boston	-5
10	236	Springfield	-3

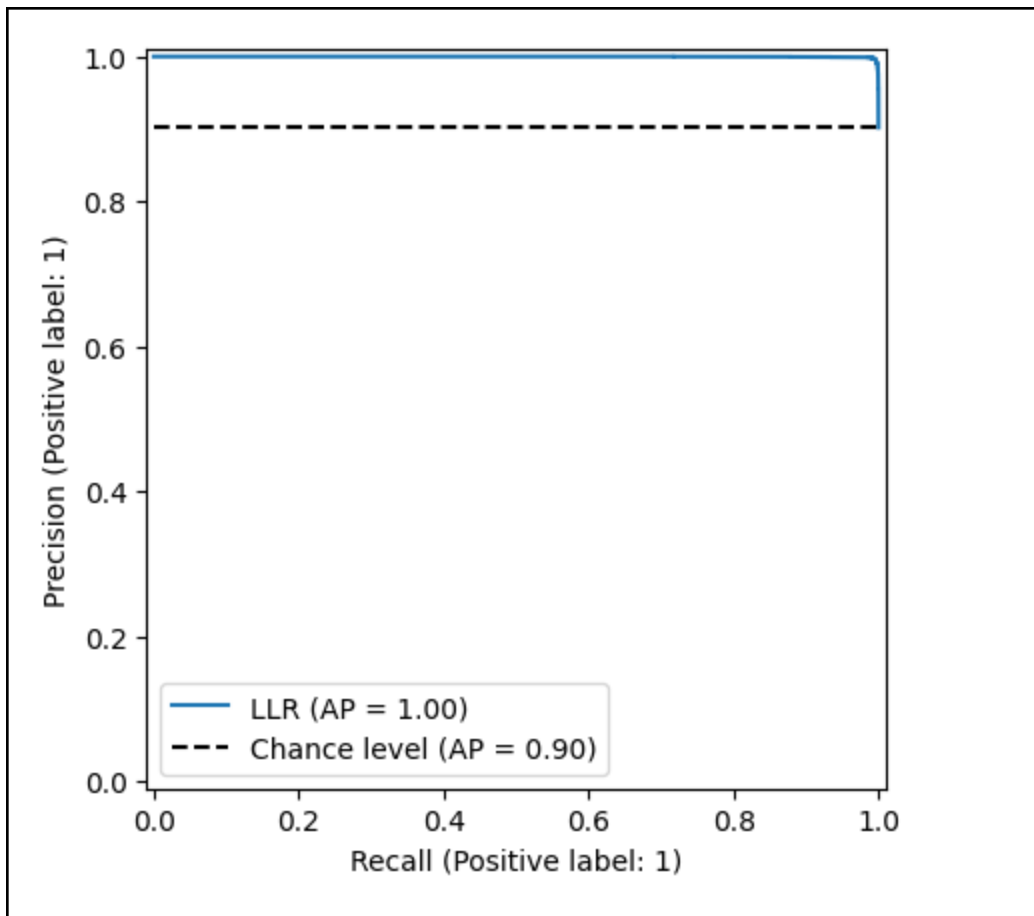
Test error using only the 10 most important features for regression

	Linear Regression
RMS Error	2.058

4. Stretch Goals

a. PR and ROC curves

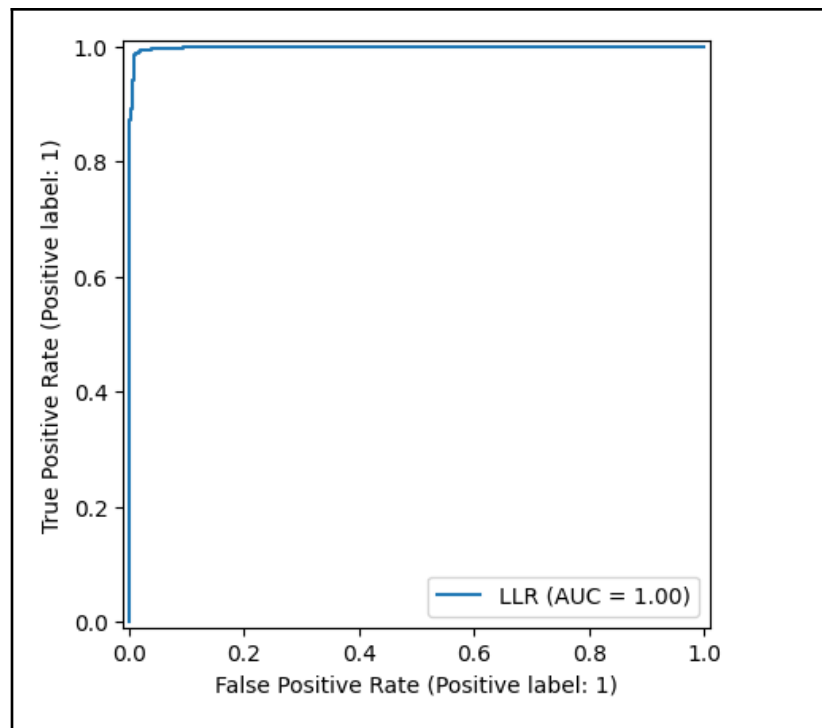
PR plot



Average Precision

1.0

ROC plot

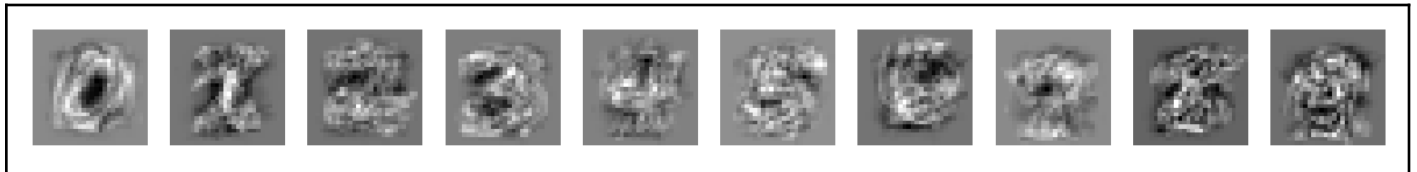


Area under the curve (AUC)

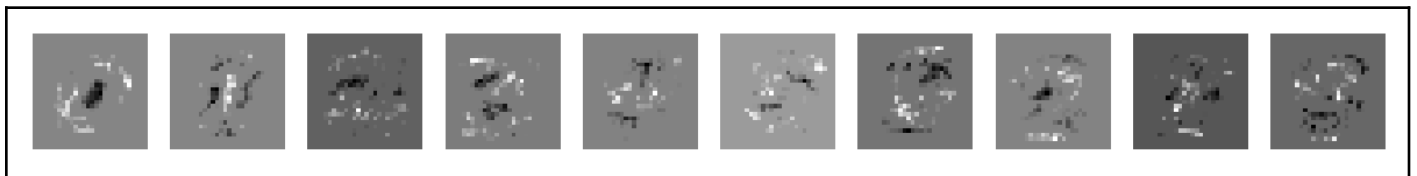
1.0

b. Visualize weights

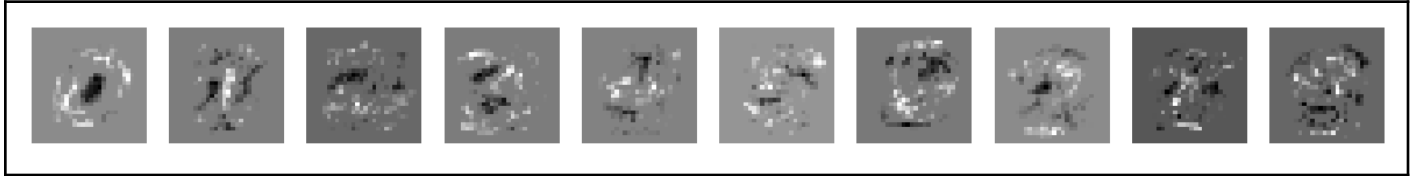
LLR - L2



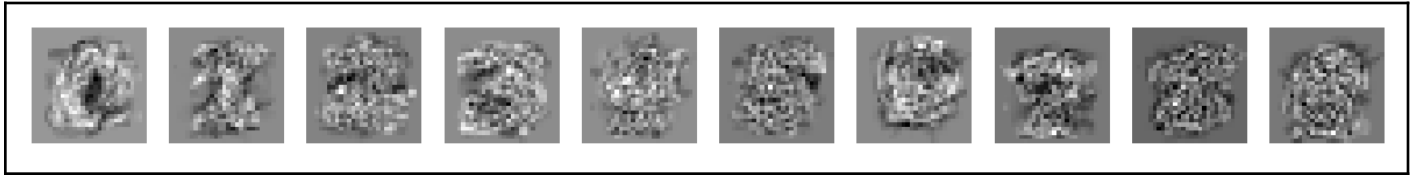
LLR - L1



LLR - elastic



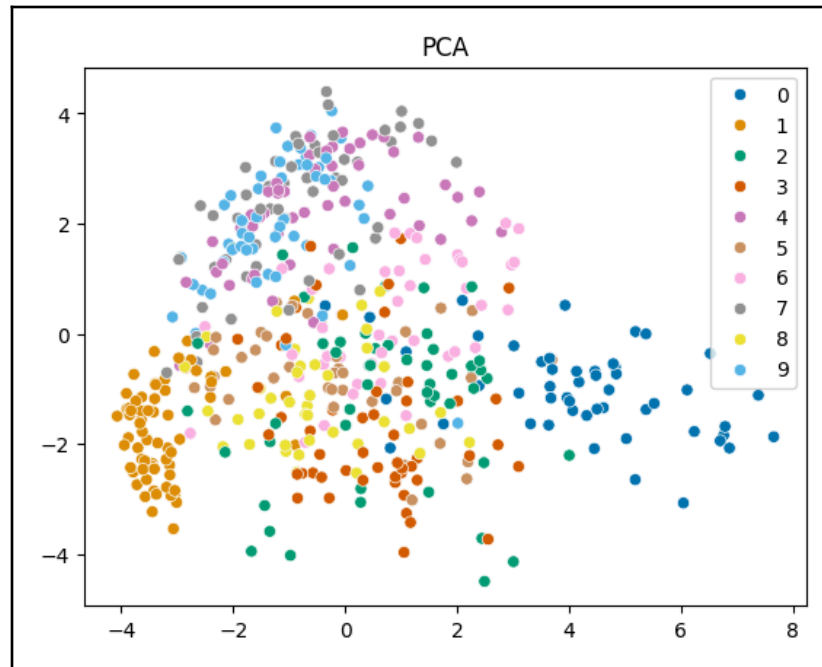
SVM



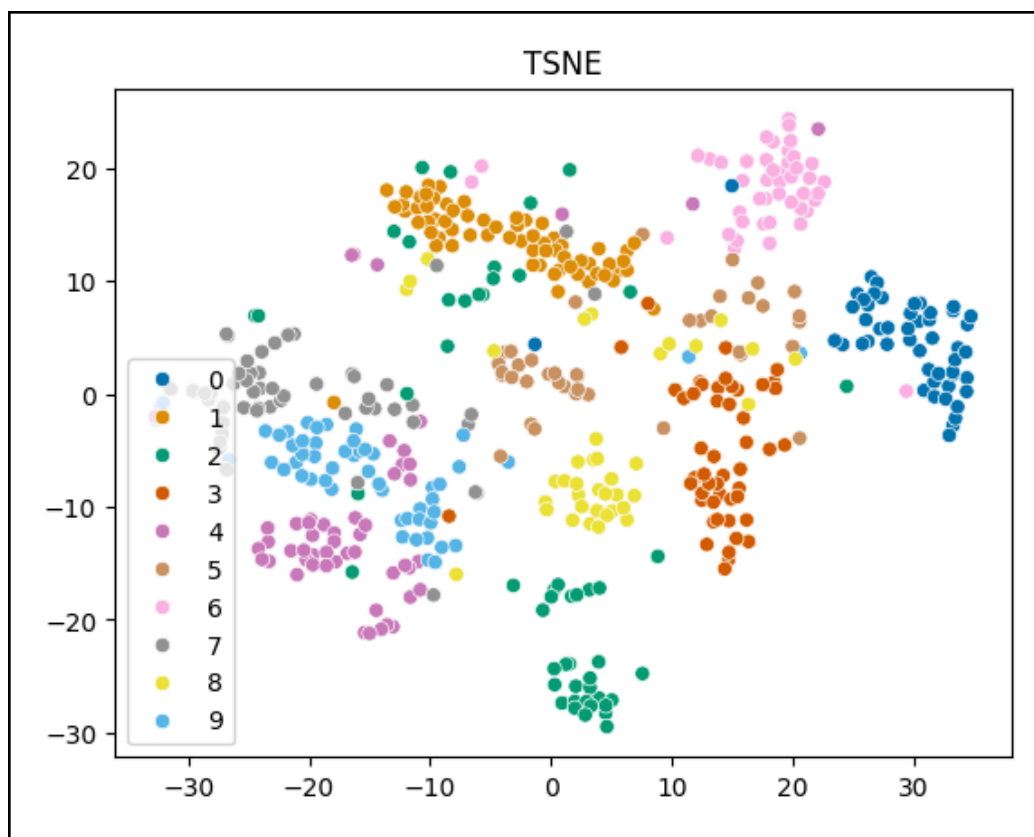
c. Other embeddings

Display 2+ plots for TSNE, MDA, and/or LDA, and copy PCA plot from 1b here.

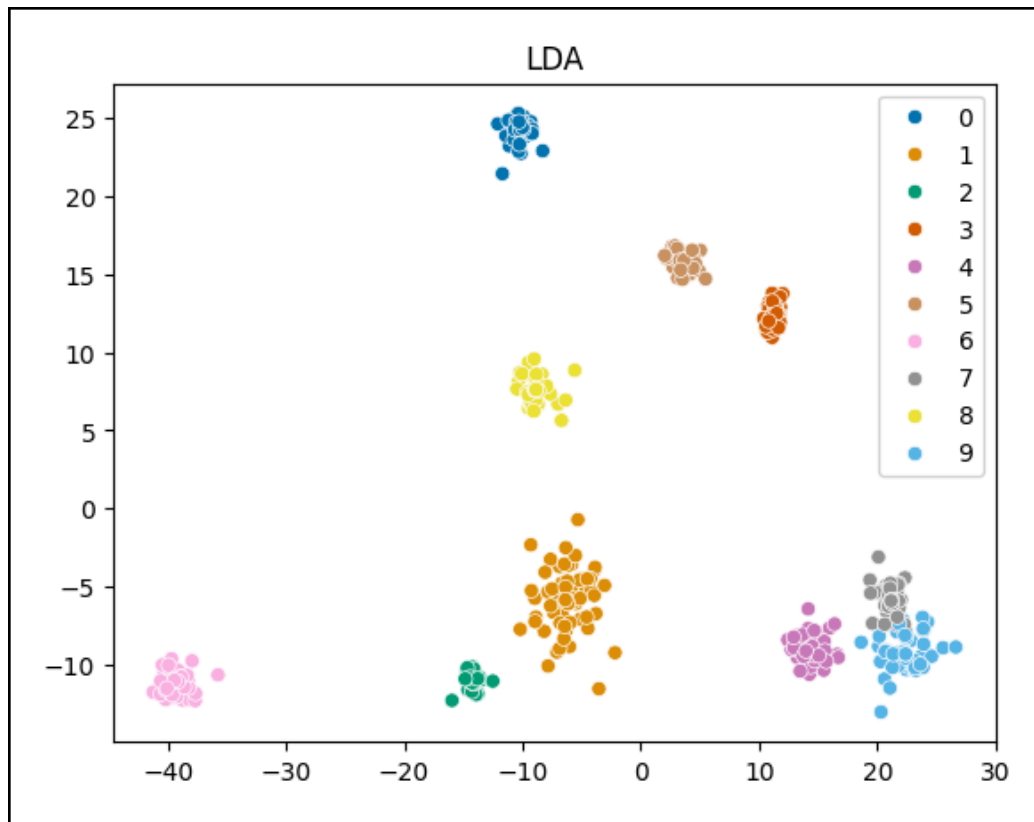
PCA



[TSNE]



[LDA]



d. One city is all you need

City

Test error using features only from that city

Explain your process (in words):

e. Compare linear SVM and SVM with RBF kernel

Test accuracy (%)

# training samples	SVM-Linear	SVM-RBF
100	67.5%	65.59%
1,000	86.36%	90.83%
10,000	90.5%	95.94%
60,000	92.61%	97.92%

Acknowledgments / Attribution

List any outside sources for code or ideas or “None”.

None