

Homework Assignment 6 [30 pts]

STAT437 Unsupervised Learning – Fall 2023

Due: Friday, October 6 on Canvas at 11:59pm CST.

Problems	Points
1.1	0.25
1.2	0.25
1.3	1
1.4	1
2.1	1.5
2.2	2
2.3	2
2.4	1.5
3.1	1.5
3.2	1.5
3.3	1.5
3.4	0.5
3.5	0.75
3.6.1.	1
3.6.2	0.5
3.6.3	1
3.6.4	0.5
4.1.1	0.75
4.1.2	0.75
4.2.1	0.75
4.2.2	0.75
4.3.1	0.75
4.3.2	1
4.4	1
5	3
6	3

Questions 1-4: See Jupyter notebook

Question 5: The dataset below is comprised of 10 professors in the UIUC Statistics department. Suppose we'd like to cluster this dataset using k-modes with k=2 clusters. The *current* cluster modes are shown below. What would be the *new cluster modes* found in the next iteration of the k-modes clustering algorithm. Show your work.

- Hint: If there's a tie in the cluster assignment step, assign the person to the mode with the lowest index.*
- Hint: If there's a tie in the centroid/mode update step, select the attribute value with the highest alphabetical order (ie. A>B).*

	PhD	Sex	Generation
Tori Ellison	Operations Research 1 1	Female 1 0	millennial 0 0
Karle Flanagan	Statistics Education 1 0	Female 1 0	millennial 0 0
Kelly Findley	Statistics Education 1 0	Male 0 1	millennial 0 0
Julie Deeke	Statistics 0 1	Female 1 0	millennial 0 0
Chris Kinson	Statistics 0 1	Male 0 1	millennial 0 0
Jeff Douglas	Statistics 0 1	Male 0 1	Gen X 1 1
Bo Li	Statistics 0 1	Female 1 0	Gen X 1 1
Steve Culpepper	Educational Psychology 1 1	Male 0 1	Gen X 1 1
Dave Zhao	Statistics 0 1	Male 0 1	millennial 0 0
Vimal Rao	Educational Psychology 1 1	Male 0 1	millennial 0 0

Hamming Distance	Mode
Mode 1	Mode
2	Mode 2
2	Mode 2
1	Mode 1
1	Mode 1
0	Mode 1
1	Mode 1
2	Mode 1
2	Mode 1
0	Mode 1
1	Mode 1
Mode 2	Mode
1	Mode 2
0	Mode 2
1 (TIE)	Mode 1
1 (TIE)	Mode 1
2	Mode 1
3	Mode 1
2 (TIE)	Mode 1
3	Mode 1
2	Mode 1
2	Mode 1

	PhD	Sex	Generation
Mode 1	Statistics	Male	millennial
Mode 2	Statistics Education	Female	millennial

New Mode

Mode 1 statistics Male millennial

Mode 2 operational Research Female millennial

Note: For mode 2,
TIE in PhD.

Question 6: In the page below 4 categorical datasets are listed. Each of these datasets was put into a Hamming distance matrix. Each hamming distance matrix was then used as input into the t-SNE algorithm producing the following 4 sets of t-SNE plots shown in the pages below.

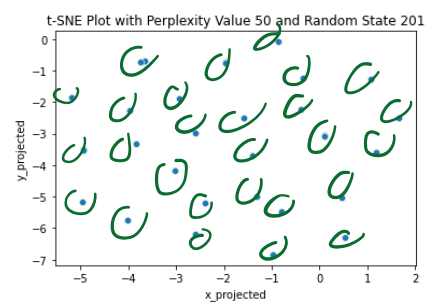
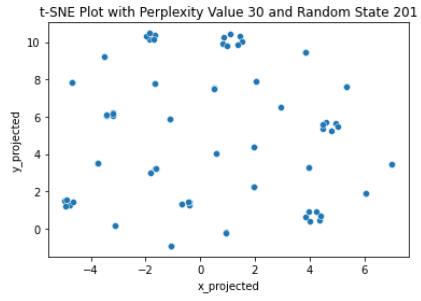
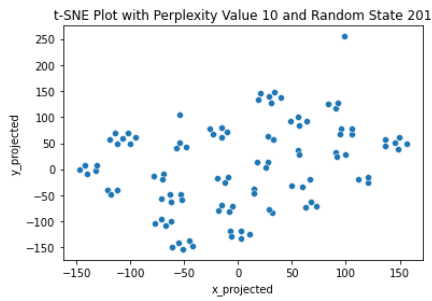
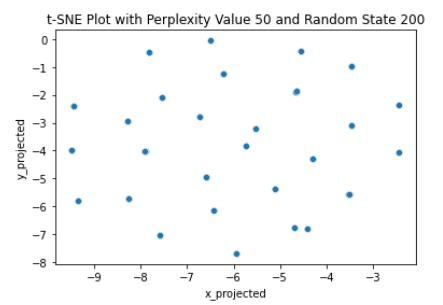
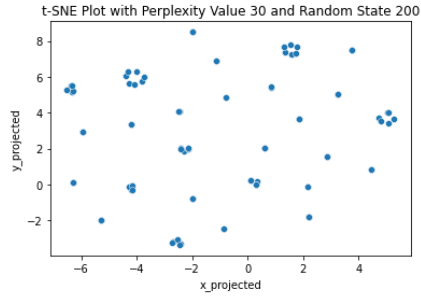
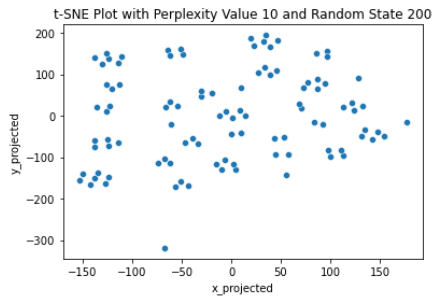
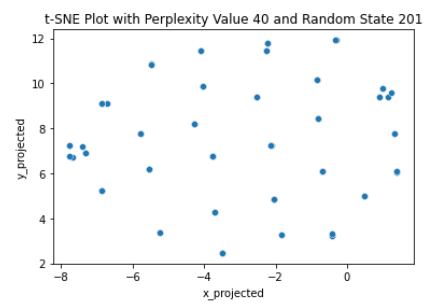
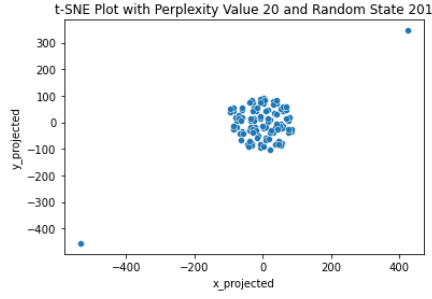
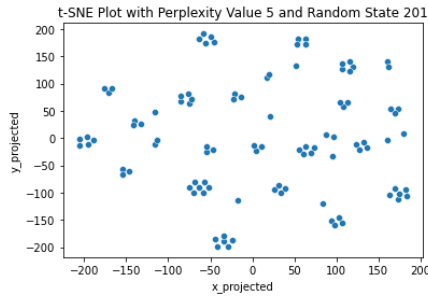
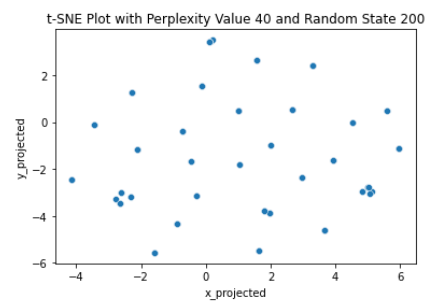
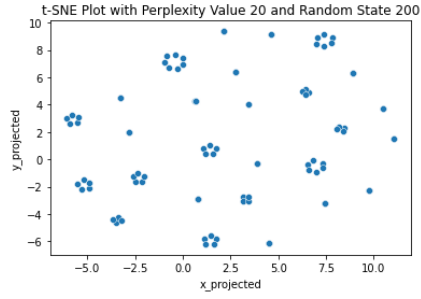
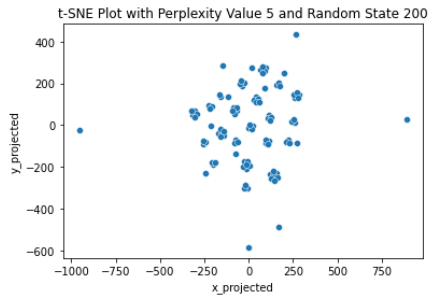
Match each of the 4 datasets 1-4 to the corresponding t-SNE plot sets A-D.

Hint: Some points may be completely overlapping in some of the t-SNE plots below.

Explanations not required, but may help for partial credit if you get it wrong.

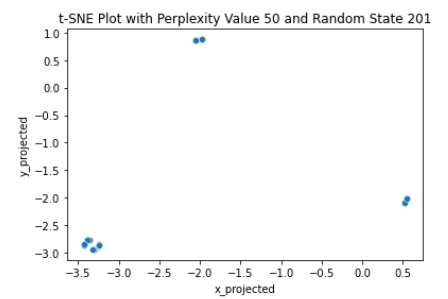
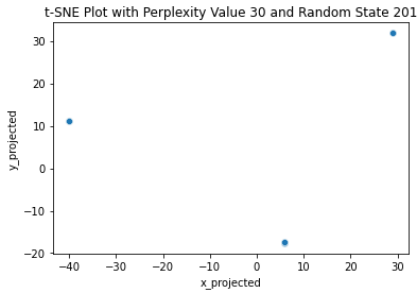
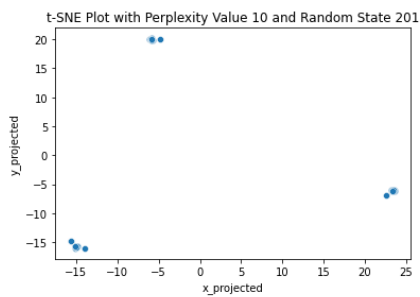
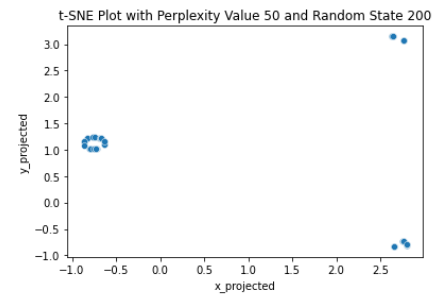
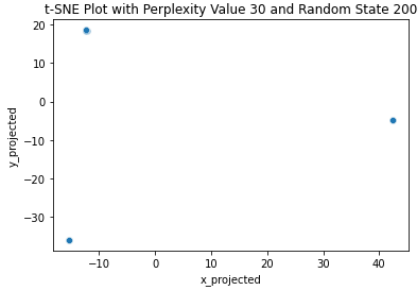
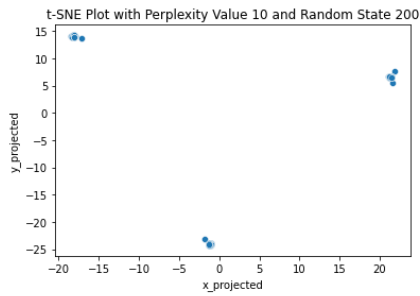
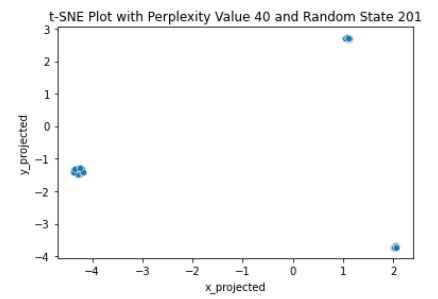
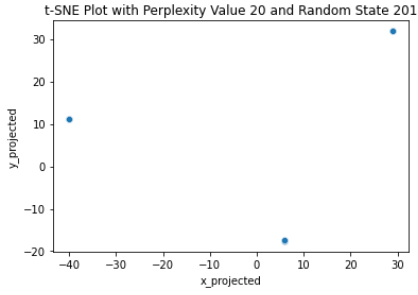
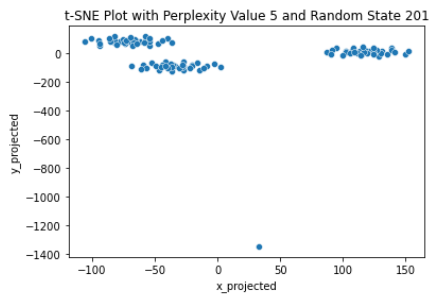
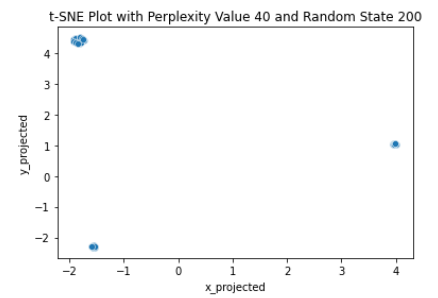
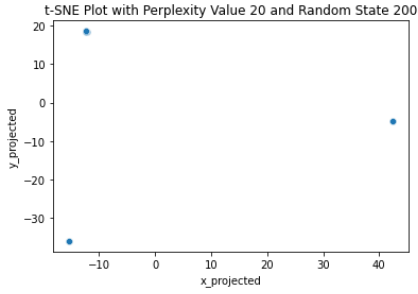
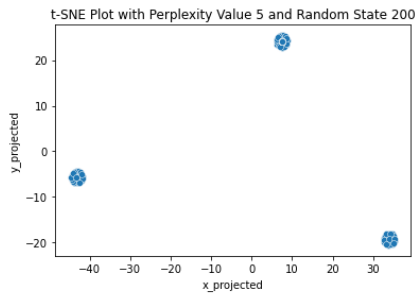
Plot Set A

DATASET 3

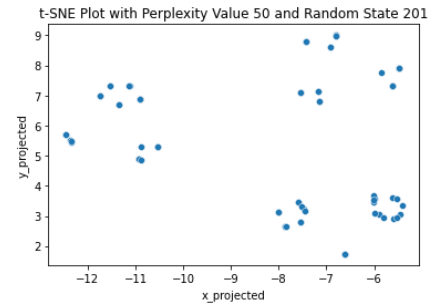
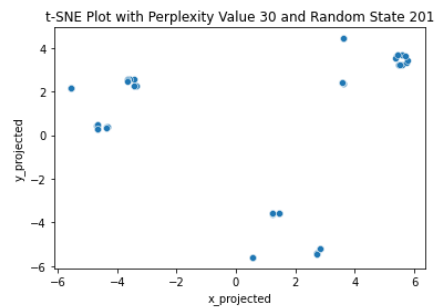
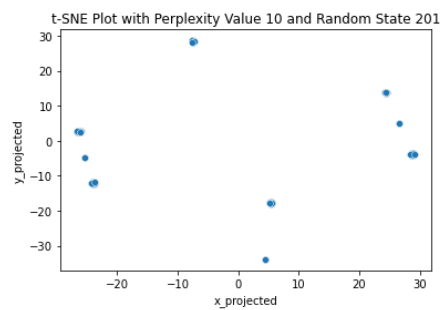
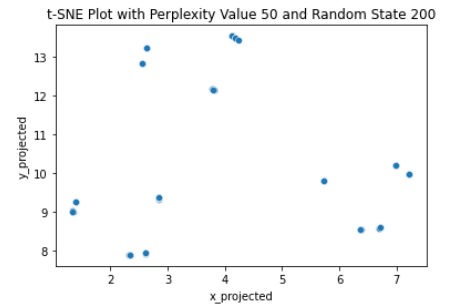
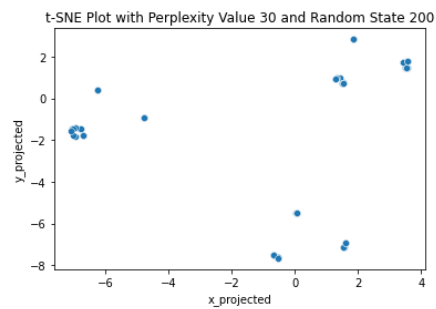
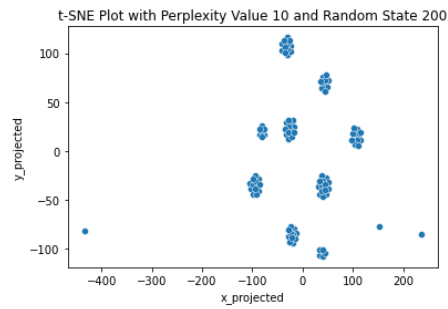
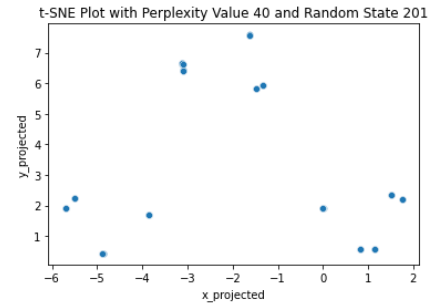
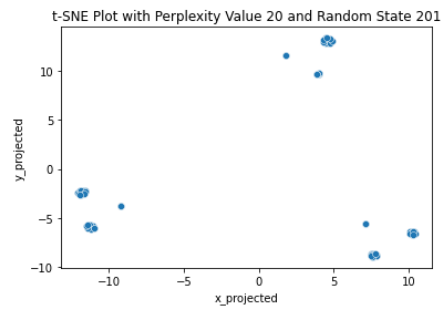
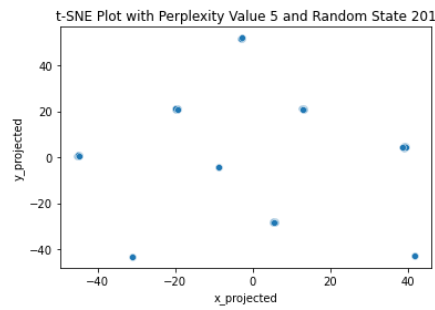
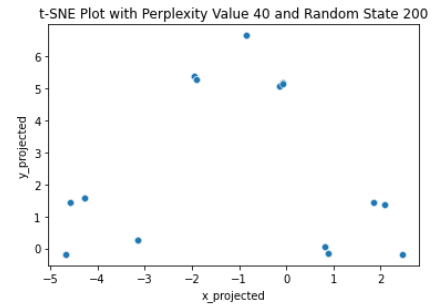
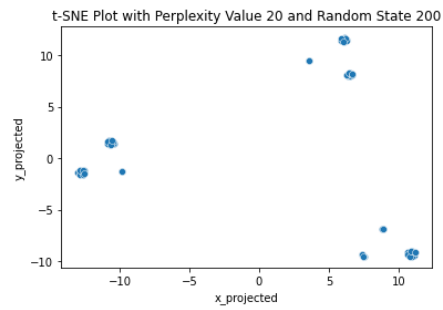
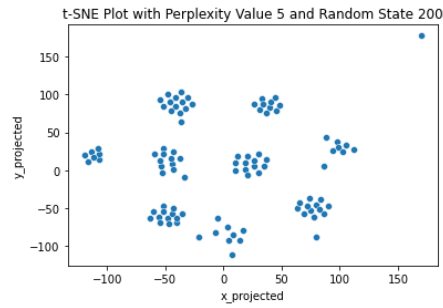


Plot Set B

DATASET 1.



Plot Set C



Plot Set D

DATASET 4

