

Name :- Prateek Mahajan

EE511:- Homework 5

Problem 1:-

I have read and understood the general instructions at the top of HWS and I formally declare that all work I turn in for everything in this course will not contain or involve any cheating at all.

Problem 2:-

- a) Note that all printed evaluation metrics are w/o matching
(i) Please refer to the end of this pdf for this gr.
(ii) Please refer to the end of this pdf for this gr.

- b) For the MNIST dataset with 10 classes & 10 clusters, there are $(10 \times 9 \times 8 \times 7 \times 6 \times 5 \times 4 \times 3 \times 2 \times 1) = 10!$
 $= 3628800$ matchings possible.
(The ~~10~~ above can be deduced while matching by matching clusters to a class 1 by 1.)
First, match cluster 1 to a class (there are 10 possible classes), then cluster 2 which would have 9 possible ~~and~~ till cluster 10, which would have only 1 class remaining.)

There are 2 parts of the Hungarian algorithm that are responsible for its worst case complexity :- equality graph construction & augmented path search.

The equality graph construction assigns labels to rows & columns of the cost matrix, and has a worst case complexity of $O(n^2)$ per iteration. The same also applies to the ~~any~~ path search in the EGr, using BFS/DFS.

- ∴ the net worst case complexity would be $O(n^2 \times n)$, where n is the number of tasks/workers in the Hungarian algo..
⇒ net worst case time complexity = $O(n^3)$.

(ii) Please refer to the end of this ^{pdf} for this gn.

c) (i) Please refer to the end of this pdf for this gn.

(ii) Yes, spectral clustering improves the accuracy of K-means from $\approx 50\%$. $\rightarrow \approx 60\%$.

This is mostly because spectral clustering is good at representing complex relationships b/w data & transforming data from a high dimensional manifold to a low dimensional one with these relationships.

K-means on the other hand is biased towards spherical relationships, and so does not handle datasets like MNIST as well.

d) (i) Please refer to the end of this pdf for this gn.

(ii) The random sampled ~~$x_{dataset}$~~ performs the worst, which is expected as doing any form of clustering would give you better reference points for K-NN than random sampling.

However, spectral KNN show a strange pattern:-

- $K=1 \rightarrow$ KNN performs better
- $K=3 \rightarrow$ Similar perf.
- $K=5 \rightarrow$ Spectral performs better

My ~~intuition~~ intuition says this is because spectral takes a more nuanced approach to representing relationships amongst various points in the dataset, that make it better for higher K s, as it's more thorough in its "clustering analysis". However, for lower K s (like $K=1$) this kind of thorough analysis in the data generalises it too much, which is why K-NN performs better

Problem 3:-

- a) Please ~~update this~~ refer to the end of this pdf for this qn.
- b) Please refer to the end of this pdf for this qn.
- c) Please refer to the end of this pdf for this qn.
- d) Please refer to the end of this pdf for this qn.
- e) (next page)

Validation Accuracy:-

logistic Regression	Random Forest	K-NN	Multi-layer Perceptrons
0.90264	0.04565	0.9428	0.9664

Overall, it seems like ~~K-NN > MLP >~~ logistic > RF.

Unfortunately, in RFs, I suspect it's an implementation issue causing a poor accuracy. For the others, I like to view MLP as a complex version of logistic regression in a way. Since the given dataset is also complex & likely has ~~over~~ a ~~very~~ complex decision boundary, MLP performs better than the logistic regression. K-NN also has a tendency to simplify decision boundaries and so, MLP outperforms K-NN in this dataset.

Further, while K-NN may not generate complex decision boundaries, as MLP, logistic regression tends to them closer to a linear form while K-NN is still good at representing non-linear boundaries.
∴ K-NN outperforms the logistic regression in this dataset, which clearly has complex non-linear decision boundaries.

- f) As MLP performed best in my previous run, I chose MLP as the classifier for this go. I used sklearn to implement it & tuned its hyperparameters in 3d) (which I reused in 3f) using GridSearchCV. Please find the .bat file in the submission & the code in the attached .zip file.