



PRML Mid Sem Exam

27.10.2020

Vipul Yuvraj Nikam

180010041

3rd year, CSE

Question 1:

a. Let

$$f_1(x_1, x_2) = x_1^2 + x_2^2 + 2x_1 + 2x_2, -10 \leq x_1 \leq 10, -10 \leq x_2 \leq 10$$

$$f_2(x_1, x_2) = x_1 \sin(x_1) + x_2 \sin(x_2), -10 \leq x_1 \leq 10, -10 \leq x_2 \leq 10$$

For both the functions:

- Draw the surface using appropriate libraries and comment on the type of surface (convex/ non-convex).
- Find the minimum value of x_1 and x_2 analytically and verify the same from the plot.
- Use gradient descent and normal equation method, and obtain the minimum value of x_1 and x_2 . From the observation clearly illustrate the pros. and cons. of both techniques.
- Using various learning rates and initialization tabulate the results and justify your observation using appropriate theoretical illustration.

b. What is the difference between regression and k-means clustering tasks.

Based on the cost function of both the techniques, comment on the convergence (i.e will they reach the global minima). Justify your solution analytically (using appropriate mathematical equations or intuition).

- Take an appropriate example and illustrate the same (coding is expected).
- If you think that any of the methods will not reach the global minima, propose a technique that may help (to some extent) to reach the global minima. Justify your proposed method by extending your earlier illustration.?

Answers:

I. Code:

[Click here](#) :

https://colab.research.google.com/drive/13aTif8VL6M_8b3_vkcSC57GP8hDd1GqD?usp=sharing

II. Graphs & Plots:

All graphs are attached above in the code file.

III. Explanation:

1) In f_1 we can see the global minimum or its presence. So f_1 is convex & f_2 is non-convex function.
 as for f_1 we have $f_1 = x_1^2 + x_2^2 + 2x_1 + 2x_2$
 by using partial derivatives,

$$\frac{\partial f_1}{\partial x_1} = 2x_1 + 1 \quad \& \quad \frac{\partial f_1}{\partial x_2} = 2x_2 + 2.$$

 by $\frac{\partial f_1}{\partial x_1} = 0 \quad \& \quad \frac{\partial f_1}{\partial x_2} = 0$

$$x_1 = -1 \quad \& \quad x_2 = -1$$

 so for f_1 , we get minimum value of -2 at $x_1 = -1$ & $x_2 = -1$

2) $x_1 = -\tan x_2$ & $x_2 = -\tan x_1$
 from graph of $x = -\tan x_2$ we get if we differentiate it partially.
 for f_1 , we know that surface is convex, so it will converge to minimum value in gradient descent. If learning rate is greater than or equal to 1, it diverges & will not lead to minimum value.
 for f_2 , we know that surface is non-convex so according to initial value of x_1 & x_2 .

→ 1b) Regression: it predicts the continuous value & their output. regression analysis is the statistics model that is used to predict the numeric data instead of lab it can also identify the distribution trends based on the available data.

k-means clustering: it is grouping method of data according to the similarity of data points & data pattern. the aim of this is to separate similar categories of data.

working:-

Basic goal to reduce cost function in k means clustering cost function is same of euclidean distance from point to their mean. by cluster

$$J = \sum_{i=1}^n \sum_{k=1}^K w_{ik} \|x_i - c_k\|^2$$

$w_{ik} = 0$ data pt belong to the cluster.

$w_{ik} = 1$ if data pt belong to cluster

i) minimize J by the c_k fixed & w_{ik} variable:

$$\frac{\partial J}{\partial w_{ik}} = \sum_{i=1}^n \sum_{k=1}^K \|x_i - c_k\|^2$$

reducing this J will get:

$$\sum_{k=1}^K 1 \text{ if } k = \arg \min_i \|x_i - c_k\|$$

otherwise 0.

2) minimize by treating w_{ik} as constant & c_k variable.

$$\frac{\partial I}{\partial c_k} = 2 \frac{k}{2} w_{ik} \|x_i - c_k\|$$

we will get the following by reducing it

$$c_k = \frac{\sum_{i=1}^m w_{ik} x_i}{\sum_{i=1}^m w_{ik}}$$

The above two steps will be carried out iteratively until we get the optimal segmentation of data point.

Question 2:

Consider the credit card fraud detection task (downloaded the data-set from: kaggle- link).

- Follow these steps to pre-process the data and obtain its subset:
 1. Export the data in .csv file using any library of your choice.
 2. As the dataset is imbalanced (i.e less no. of fraud examples compared to genuine), choose randomly the same number of genuine examples as fraud examples and obtain a balanced dataset .
 3. Remove the outliers from the selected examples (use threshold of 2.5 IQR).
 4. Normalize the data (features should have zero mean and unit standard deviation).
 5. Use this data for further processing.
- Use t-SNE plot (refer: t-SNE link) to visualize the subset of data in 2D. Intuitively comment which clustering algorithm (k-means, GMM, agglomerative hierarchical clustering, DBSCAN) will work better and provide justification for the same using your theoretical knowledge.
- Propose and clearly explain an evaluation measure to evaluate the clustering algorithms. Using a piece of code, prove your intuition of clustering algorithm selection using the proposed evaluation measure (use the same subset of data obtained in the first part).
- Take the whole (imbalanced) data, do the data pre-processing task (out-lier removal, data normalization), perform clustering using all the above mentioned techniques. If degradation in performance is observed, comment on the possible causes for it. Propose a possible solution to overcome this issue (without reducing the size of the dataset). Justify the same through a code.
- Split the whole data in a training and testing set (90%, 10%, perform class specific division as the data is imbalanced). Use any one of the above mentioned clustering techniques and perform classification tasks (using a piece of python code). In the report, clearly mention the steps involved in training and testing, and the performance measure with appropriate mathematical equations.?

Answers:

I. Code:

[Click here](#) :

<https://colab.research.google.com/drive/1HxdPtYtILWWgbZlfoahkPzkO-LV98s4f?usp=sharing>

II. Graphs & Plots:

All graphs are attached above in the code file.

III. Explanation:

→ 2) After seeing the t-SNE plot, intuitively we come to conclusion that the best algorithm would be k-means. But we can't use the fact that agglomerative Hierarchical gives us Holistic visualization of the dendrogram whilst considering all the data points in the clusters.

But since k-means iterates through all the items k-clusters each item to measure clusters using similarity measures we chose k-means to better one.

or But DBSCAN & GMM are not viable here because GMM is very high computation. & DBSCAN is more favorable.

$$a(i) = \frac{1}{k_i - 1} \cdot \sum_{c \in C_i} d(i, c) \quad b_i = \min_j \left(\frac{1}{k_j - 1} \sum_{c \in C_j} d(i, c) \right)$$

$$S(i) = \frac{k(i) - a(i)}{\max(a(i), b(i))} \quad \text{where } C(i) \text{ is cluster index } i \text{th data.}$$

so, from the plot we can see that DBSCAN & GMM have performed poorly than K-means & agglomerative ~~the~~ comparison but DBSCAN was better than GMM when whole data set was considered for clustering, the fraud clean wasn't visible. So we built all the algo to get idea of clusters formed.

To overcome this problem, we can increase our dataset size also we can generate synthetic example. In we saw unbalanced dataset module in python helps us to implement SMOTE which helps us to generate synthetic samples. To increase our performance of clusters we used above method.