

HW4 HADI SHARIFI

Q1-a

Check the code bundled with the report

Q1-b

Check the code bundled with the report

Q1-c

Check the code bundled with the report

Q1-d

Check the code bundled with the report.

I used the KMean function from ScikitLearn package. The best accuracy I got was 80%.

	Accuracy	Misclassification for 2	Misclassification for 6
KMean	94%	0.05	0.07
GMM-EM	81%	0.19	0.15

Q2 - a

To derive the gradient of cost function. The cost function is the log likelihood

$$\ell(\theta) = \sum_{i=1}^n \{-\log(1 + e^{-\theta x_i}) + (y_i - 1)\theta x_i\}$$

$$\frac{\partial \ell(\theta)}{\partial \theta} = \sum_{i=1}^n (y_i x_i - x_i - (\frac{e^{-\theta x_i}}{1 + e^{-\theta x_i}}) x_i)$$

$$= \sum_{i=1}^n (y_i - \frac{1}{1 + e^{-\theta x_i}}) x_i = \text{Fdn} \quad \text{this is our gradient descent steps.}$$

for gradient descent, we need define the steps which is the formula above
And we need the learning rate λ .

the pseudo code will be:

Procedure GD (Fdn, $\theta(0)$)

$\theta \leftarrow \theta(0)$

while not converged do

$\theta \leftarrow \theta + \lambda \text{Fdn}(\theta)$

return θ

Fdn is a function that returns a matrix of form

$$\text{Fdn}(\theta) = \begin{bmatrix} \frac{\partial \ell(\theta)}{\partial \theta_1} \\ \frac{\partial \ell(\theta)}{\partial \theta_2} \\ \vdots \\ \frac{\partial \ell(\theta)}{\partial \theta_n} \end{bmatrix}$$

$\theta \in \mathbb{R}$

$\theta(0)$ is a Random initialization

the returned θ is a point in which the log likelihood is the optimum. We run the Algorithm for each and points.

Q2-b

Stochastic gradient decent is little different.
gradient decent iterates over all points and for large data that is very time consuming.

In SGD, we shuffle the data, Pick Random K ^{Points} data of them and run the gradient steps. In SGD, we also have the learning rate λ .
The gradient step is similar to Q2-a ($F_{\text{data}}(\theta)$).

Procedure SGD ($F_{\text{data}}, \theta(0)$)

$\theta \leftarrow \theta(0)$

while not converged do

for shuffle data $\{1, 2, \dots, N\}$ pick K data points

for $k \in \{1, 2, \dots, K\}$ do

$\theta_k \leftarrow \theta_k + \lambda F_{\text{data}}(\theta_k)$

return θ_k

Q2-C

In order to identify a training problem as concave or convex, we need to derive the Hessian matrix of $\ell(\theta)$.

For data points of type \mathbb{R}^n the Hessian matrix will be:

$$H = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \dots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

for problem Q2, the data point are in \mathbb{R} therefore

Taking second derivative would do the work.

The Hessian matrix has only one item.

From Q1-a

$$\frac{\partial \ell(\theta)}{\partial \theta} = \sum_{i=1}^n \left(y_i - \frac{1}{1 + e^{-\theta x_i}} \right) x_i$$

$$\frac{d}{dx} \frac{1}{1 + e^{ax}} = \frac{-ae^{ax}}{(1 + e^{ax})^2}$$

$$\frac{\partial^2 \ell(\theta)}{\partial \theta^2} = \sum_{i=1}^n \frac{x_i^T e^{\theta x_i}}{(1 + e^{\theta x_i})^2} \times x_i = \sum_{i=1}^n \frac{e^{\theta x_i}}{(1 + e^{\theta x_i})^2}$$

$$= \sum_{i=1}^n \frac{1}{1 + e^{\theta x_i}} \cdot \left(\frac{e^{\theta x_i}}{1 + e^{\theta x_i}} \right) = \sum_{i=1}^n \frac{1}{1 + e^{\theta x_i}} \left(\frac{1}{1 + e^{\theta x_i}} - 1 \right)$$

⊕

⑧ $\Rightarrow \frac{1}{1+e^{\theta x_i}}$ is always $0 < \square < 1$ which makes

$\sum_{i=1}^n \left(\frac{1}{1+e^{\theta x_i}} \times \left(\frac{1}{1+e^{\theta x_i}} - 1 \right) \right) < 0 \Rightarrow$ hence the
Problem is Concave
and NOT Convex