

Name - Md Farhan Istham
180041120

Name : Md Farhan Istham

ID : 180041120

Semester : Sixth

Date : 07 - Feb - 22

Course Code : CSE 4621

Course Name : Machine Learning

Exam : Mid Semester

Ans. to Qno. 3(a)

$$\frac{d\lambda}{dz^{[2]}} = \frac{d\lambda}{da^{[2]}} \times \frac{da^{[2]}}{dz^{[2]}} \times \frac{dz^{[2]}}{da^{[1]}} \times \frac{da^{[1]}}{dz^{[1]}} \quad [\text{Chain - Rule}]$$

$$\text{Now, } \lambda = -y \log \hat{y} - (1-y) \log (1-\hat{y}) \quad [\hat{y} = a^{[2]}]$$

$$\begin{aligned} \frac{d\lambda}{da^{[2]}} &= \frac{d}{da^{[2]}} [-y \log a^{[2]} - (1-y) \log (1-a^{[2]})] \\ &= -\frac{y}{a^{[2]}} + \frac{1-y}{1-a^{[2]}} \end{aligned} \quad (i)$$

$$\begin{aligned} \frac{da^{[2]}}{dz^{[2]}} &= \frac{d}{dz^{[2]}} (\sigma(z^{[2]})) \\ &= \sigma(z^{[2]}) \{1 - \sigma(z^{[2]})\} \\ &= a^{[2]} (1-a^{[2]}) \end{aligned} \quad (ii)$$

$$\begin{aligned} \frac{dz^{[2]}}{da^{[1]}} &= \frac{d}{da^{[1]}} (W^{[2]} a^{[1]} + b^{[2]}) \\ &= W^{[2]} \end{aligned} \quad (iii)$$

$$\begin{aligned} \frac{da^{[1]}}{dz^{[1]}} &= \frac{d}{dz^{[1]}} \sigma(z^{[1]}) \\ &= a^{[1]} (1-a^{[1]}) \end{aligned} \quad (iv)$$

Multiplying (i), (ii), (iii) and (iv),

$$\begin{aligned}\frac{\partial \delta}{\partial z^{[1]}} &= \left(-\frac{y}{a^{[1]}} + \frac{1-y}{1-a^{[1]}}\right) \{a^{[1]}(1-a^{[1]})\} \cdot w^{[2]}.a^{[1]}(1-a^{[1]}) \\ &= (-y + a^{[2]}) \cdot w^{[2]}.a^{[1]}(1-a^{[1]}) \\ &= w^{[2]}.a^{[1]}(1-a^{[1]}) (a^{[2]} - y) \quad (\text{Ans.})\end{aligned}$$

This can be rewritten as

$$d_z^{[1]} = d_z^{[2]}. \sigma'(a^{[1]}). w^{[2]}$$

The update equation for $w^{[1]}$ is

$$\frac{\partial \delta}{\partial w^{[1]}} = \frac{\partial \delta}{\partial z^{[1]}} \times \frac{\partial z^{[1]}}{\partial w^{[1]}}$$

$$\begin{aligned}\text{Now, } \frac{\partial z^{[1]}}{\partial w^{[1]}} &= \frac{\partial}{\partial w^{[1]}} \sigma(w^{[1]}x + b^{[1]}) \\ &= x\end{aligned}$$

$$\begin{aligned}\therefore \frac{\partial \delta}{\partial w^{[1]}} &= w^{[2]}.a^{[1]}(1-a^{[1]})(a^{[2]} - y) \cdot x \quad (\text{Ans.}) \\ &= d_z^{[1]}.x\end{aligned}$$

Update equation.

$$w^{[1]} := w^{[1]} - \alpha d_w^{[1]}$$

$$\text{or, } w^{[1]} := w^{[1]} - \alpha(d_z^{[1]}.x)$$

$$\text{or, } w^{[1]} := w^{[1]} - \alpha \{w^{[2]}a^{[1]}x(1-a^{[1]})(a^{[2]} - y)\}$$

(Ans.)

Anuto Q.no. 3(b)

$$\text{Given, } x = (1, 1) = \begin{bmatrix} 1 \\ 1 \end{bmatrix}; \quad y = 0.$$

$$(i) \quad a_1^{[1]} = \sigma(w_1^{[1]}x + b_1)$$

$$= \sigma(0.8 - 0.5 \times 1 + 0.4 \times 1)$$

$$= \sigma(0.7)$$

$$= 0.332$$

$$a_2^{[1]} = \sigma(w_2^{[1]}x + b_2)$$

$$= \sigma(0.9 \times 1 + 1.0 \times 1 - 0.1)$$

$$= \sigma(1.8)$$

$$= 0.858$$

$$\hat{y} = a^{[2]} = \sigma(w_0^{[2]}x + b^{[2]})$$

$$= \sigma(w^{[2]}a^{[1]} + b^{[2]})$$

$$= \sigma(0.3 + 0.332 \times (-1.2) + 1.1 \times 0.858)$$

$$= \sigma(0.8454)$$

$$= 0.7$$

$$\therefore \hat{y} = 0.7.$$

Ashwani

$$(ii) \quad dW^{[2]} = \frac{d\alpha}{d\alpha^{[1]}} \times \frac{da^{[1]}}{dx} \times \frac{dz^{[2]}}{dw^{[2]}}$$

$$= (a^{[2]} - y) \times a^{[1]} \quad \text{--- (1)}$$

$$db^{[2]} = a^{[2]} - y \quad \text{--- (2)}$$

$$W^{[2]} = W^{[1]} - \alpha \times dW^{[2]}$$

$$b^{[2]} = b^{[2]} - \alpha \times db^{[2]} \quad [\text{Assume, } \alpha = 1]$$

$$W^{[2]} = \begin{bmatrix} -1.2 \\ 1.1 \end{bmatrix} - (0.7 - 0) \times \begin{bmatrix} 0.332 \\ 0.858 \end{bmatrix}$$

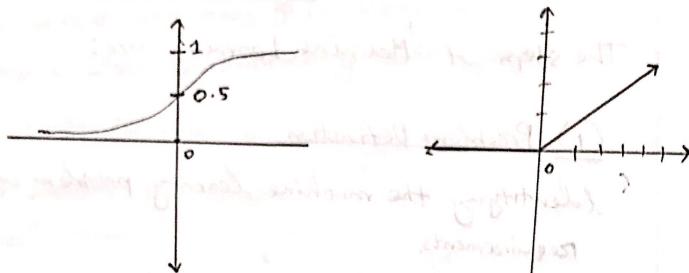
$$= \begin{bmatrix} -1.2 \\ 1.1 \end{bmatrix} - \begin{bmatrix} 0.2324 \\ 0.6 \end{bmatrix}$$

$$= \begin{bmatrix} -1.432 \\ 0.5 \end{bmatrix} \quad (\text{Ans})$$

$$(iii) \quad db^{[2]} = b^{[2]} - \alpha \times db^{[2]}$$

$$= 0.3 - 1 \times (0.7 - 0)$$

$$= 0.4 \quad (\text{Ans.})$$

Ans to Q.no. 3(c)

Sigmoid

$$\sigma(z) = \frac{1}{1+e^{-z}} \in (0, 1)$$

ReLU

$$\text{ReLU}(z) = \max(0, z)$$

Benefits of ReLU over sigmoid:

- (i) Sigmoid punishes negative ~~cos~~ values by taking it close to zero. It serves the purpose while a ReLU assigns exactly 0 to negative costs, which speeds up training.
- (ii) ReLU is simpler and faster. Easier to train as it is less computationally expensive.
- (iii) ReLU is better for Image Processing using CNN as the negative cost of sigmoid makes it harder to train such models. They can generalize easier using ReLU.
- (iv) Sigmoid is more prone to vanishing and exploding gradient problem.

Ans to Ques. 1(a)

The steps of Machine Learning are:

(i) Problem Definition

Identifying the machine learning problem with its requirements.

(ii) Data Collection

Gathering the data required for the machine learning task through third-party sources, survey, databases and so on.

(iii) Data pre-processing

Cleaning the data, scaling and normalizing it to make it suitable for training. Data splitting is performed.

(iv) Model Selection

Selecting an appropriate model, and learning algorithm for the task. This includes - linear regression, neural network, support vector machine etc.

(v) Training the Model

Using whole data or portion of it, the model is trained by minimizing the cost through iterations.

(vi) Hyperparameter Tuning

The learning rate, regularization parameter, network architecture are changed and fine-tuned to improve the results in cross-validation dataset.

(vii) Evaluating the Model

The model accuracy is measured on Test set. This step is done concurrently with hyperparameter tuning.

(viii) Model deployment

The model is deployed and can be used for practical applications.

Major types of learning —

(i) Supervised — Training is performed with labeled dataset. Ex - linear, logistic regression, neural networks.

(ii) Unsupervised — Trained using unlabeled dataset.

Ex - Principal Component Analysis (PCA)
Anomaly Detection.

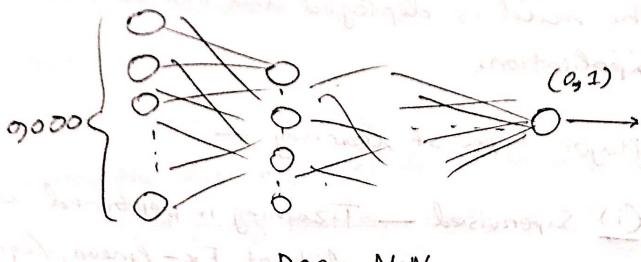
(iii) Semi-supervised — Dataset contains both labeled and unlabeled data. Trained using labeled data, fine-tuned / refined using unlabeled data.

(iv) Reinforcement Learning — Reward / punishment based learning. If ~~the~~ model does wrong it gets punished and tries to do something new.

Ex - AlphaZero (deep N.N. to play chess).

Ans to Ques 1(b)

As the number of input features is huge, I believe using a deep neural network will be the best choice.



This model is chosen because ~~huge~~ inputs is a huge feature set and at least 9,000 will be used. Deep Neural networks work better on datasets with large input dimensions. A 10-12 layered neural network with a sufficient number of neurons in each layer should give good accuracy.

Regularization -

I will use Dropout Regularization, as there are 10,000 genes, some genes like 9,000 will be used. So, rest 1,000 can be removed using dropout regularization.

Ans. to Q. no. 1(c)

Log-loss cost

$$\lambda(\hat{y}, y) = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)}) \right]$$

$$\nabla_{\theta} \lambda = -\frac{1}{m}$$

The vectorized form is

$$\lambda = -\frac{1}{m} (y \log \hat{y} + (1-y) \log (1-\hat{y}))$$

$$= -\frac{1}{m} \sum_j \log \sigma(\theta x) + (1-y) \log (1-\sigma(\theta x))$$

$$\nabla_{\theta} \lambda = -\frac{1}{m} \frac{\partial \lambda}{\partial \theta}$$

$$\nabla_{\theta} \lambda = \frac{\partial \lambda}{\partial \theta} = \frac{\partial \lambda}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial \theta z} \times \frac{\partial z}{\partial \theta} \quad [\because \sigma'(z) = \sigma(z)(1-\sigma(z))]$$

$$= \left(-\frac{y}{\hat{y}} + \frac{1-y}{1-\hat{y}} \right) \times \hat{y}(1-\hat{y}) \times \dots \times$$

$$= (\hat{y} - y) \cdot x$$

$$\nabla_{\theta}^2 \lambda = \frac{\partial^2 \lambda}{\partial \theta^2} = \left(\frac{\partial}{\partial \theta} \hat{y} \right) \cdot x$$

$$= \hat{y}(1-\hat{y}) \cdot x^T x$$

This term is always positive as $\hat{y} \geq 0$ for $\hat{y} = \sigma(\theta x)$ ∴ $\nabla_{\theta}^2 \lambda \geq 0$

$$\hat{y} = \sigma(\theta x)$$

This is the second derivative test for convexity.

∴ Log loss function is convex (proved.)

Ans. to Q. no. 2(a)

Suppose, we have $x_1 \in [0, 10000]$

$x_2 \in [0, 5]$.

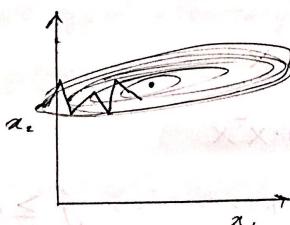
When using linear or logistic regression, the model will be more dependent on x_1 as it has higher values. For similar weights model will prefer x_1 and will try to optimize it first. This prevents the model from generalizing and makes it harder to train.

Instead if $x_1 \in [0, 1]$
 $x_2 \in [0, 1]$

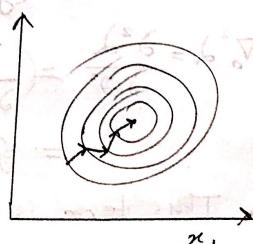
each feature will be given similar priority.

Using feature scaling, $x_i = \frac{x_i - x_{\min}}{x_{\max} - x_{\min}}$

we can keep the features in a small range $[0, 1]$.



Unscaled Dataset
Features.



Scaled Features.

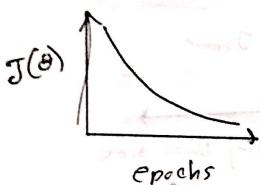
Scaled features make it easier to reach the minimum.

(Training) comes a lot faster and good fit.

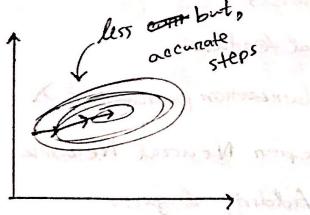
Ans to Q no 2(b)Batch

(i) Whole training data is used in a single epoch.

(ii) Cost decreases in each step.



(iii) Each step takes longer to train but accurate.

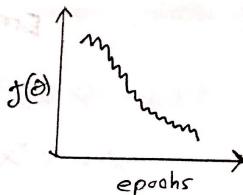


(iv) Usually takes longer time to train.

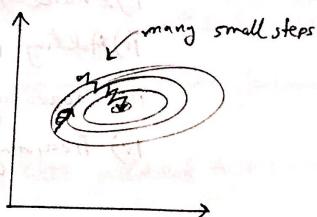
Stochastic

(i) A single training data sample is used.

(ii) Cost might increase, but in multiple steps overall cost decreases.



(iii) Each step takes shorter time to train but can be inaccurate.



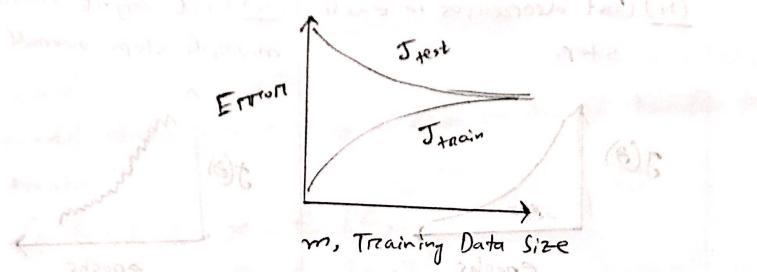
(iv) Usually faster to train.

Ques. 2(e)

Ans. to Q. no. 2(e)

If the training and test error — both are high, then the model is suffering from high bias problem.

No, the increasing the ~~training~~ training data size won't help.



As, we can see increasing data size won't help reduce the ~~error~~ error.

Solution →

- Increasing features
- Adding polynomial features
- Decreasing regularization parameter, λ
- Training a deeper Neural Network

→ Adding layers

→ Increases no. of neurons / layer.

Ans. to Q.no. 2(d)

$$\text{Given, } J(\theta) = \frac{1}{2} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)})^2 + \lambda \sum_{j=1}^n |\theta_j|$$

$$\nabla_{\theta} J = \frac{\partial J}{\partial \theta} = \frac{1}{2} \cdot 2$$

The vector form of the given equation is

$$J(\theta) = \frac{1}{2} (\cancel{x}^T \theta - y)^2 + \lambda \sum |\theta_j|$$

$$\nabla_{\theta} J = \frac{\partial J}{\partial \theta} = \frac{1}{2} \times 2 (\cancel{x}^T \theta - y) \cancel{x} + \lambda$$

$$\Rightarrow \theta = \cancel{x} \times {}^T (\cancel{x} \theta - y) + \lambda \quad [\because \nabla_{\theta} J = 0]$$

$$\Rightarrow \theta = x^T x \theta - x^T y + \lambda$$

$$\Rightarrow x^T y - \lambda = \cancel{x}^T \cancel{x} \times {}^T x \theta$$

$$\Rightarrow (x^T x)^{-1} (x^T y - \lambda) = \theta$$

$$\Rightarrow \hat{\theta} = (x^T x)^{-1} (x^T y - [0, 1, \dots, 1] \lambda) \quad \underline{\text{(Ans.)}}$$

Hence, the $\begin{bmatrix} 0 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$ term is given because

regularization is NOT added to the bias.

We shall avoid normal equation if the number of training data is high in a dataset or the number of features is too high. Because, normal ~~data~~ equation is ~~computationally~~ more expensive and has complexity $O(n^3)$.