

CS-577 ASSIGNMENT - 3

THEORETICAL QUESTIONS:

LOSS

1. Equations for L_1 , L_2 and hubber loss, cosh loss

$$L_1: L_1(\theta) = \sum_{j=1}^k |y_j^{(i)} - \hat{y}_j^{(i)}|$$

$$L_2: L_2(\theta) = \sum_{j=1}^k (\hat{y}_j^{(i)} - y_j^{(i)})^2$$

Hubber loss: $g(d) = \begin{cases} \frac{1}{2}d^2 & \text{if } |d| \leq \delta \\ \delta(d - \frac{1}{2}\delta) & \text{otherwise} \end{cases}$

$$L_3(\theta) = \sum_{j=1}^k g(\hat{y}_j^{(i)} - y_j^{(i)})$$

Log cosh loss: $L_4(\theta) = \sum_{j=1}^k \log(\cosh(\hat{y}_j^{(i)} - y_j^{(i)}))$

$$\log(\cosh(d)) \approx \begin{cases} d^2 & \text{if } d \text{ is small} \\ |d| - \log 2 & \text{otherwise} \end{cases}$$

L_1 loss is generally used in regression problems where we are certain that there aren't many outliers. This loss function is not sensitive to outliers. Median is the optimal predicted value.

L_2 loss provides the sum of squared errors and is affected by large error values. This loss is comparatively more sensitive to outliers than the L_1 loss and the mean is the optimal prediction in this loss.

Hubber loss is used as a loss function for robust regression. This loss handles outliers the best. The curve is more flatter compared to L_2 loss.

(2)

log(cosh) loss: A smoother loss function compared to ℓ_2 used for regression problems. The function behaves differently for small ~~dependencies~~ & large d , representing inliers and outliers respectively.

2. Maximizing the log likelihood function:

$$L(\theta) = \prod_{i=1}^m \prod_{j=1}^k p(y=j|x^{(i)})^{y_j^{(i)}}$$

Minimizing the log likelihood f :

$$\begin{aligned} l(\theta) &= -\log L(\theta) = -\sum_{i=1}^m \sum_{j=1}^k \log(p(y=j|x^{(i)})) \\ &= -\sum_{i=1}^m \sum_{j=1}^k y_i^{(i)} \log(\hat{y}_j^{(i)}) \end{aligned}$$

For a random assignment of k class

$p(y=j|x^{(i)}) = \frac{1}{k} = \log(k)$. If the cross entropy loss becomes greater than $\log(k)$ this loss is considered to be the worst loss.

3. Softmax loss $l(\theta) = -\sum_{i=1}^m \sum_{j=1}^k y_i^{(i)} \log(\hat{y}_j^{(i)})$.

When multiclass classification needs to be performed we use softmax loss to perform the updatums of weights. Gradients of the loss are calculated in order to manipulate the initial weights to achieve multiclass classification.

(3)

4. Kullback - Leibler loss :

$$L(\theta) = - \sum_{i=1}^m \sum_{j=1}^k y_j^{(i)} \log \left(\frac{y_j^{(i)}}{\hat{y}_j^{(i)}} \right)$$

$$\text{K.L. loss } L(\theta) = \underbrace{\sum_{j=1}^k (y_j^{(i)} \log y_j^{(i)})}_{\text{- entropy}} - \underbrace{\sum_{j=1}^k (y_j^{(i)} \log \hat{y}_j^{(i)})}_{\text{cross entropy}}$$

Given identically independent samples (ID) $\{x_i\}_{i=1}^m$ and distributions $P(n)$ and $Q(n)$, $P(n)$ is the true distribution of labels and $Q(n)$ being the predicted distribution, Kullback - Leibler loss measures the similarity of these distributions.

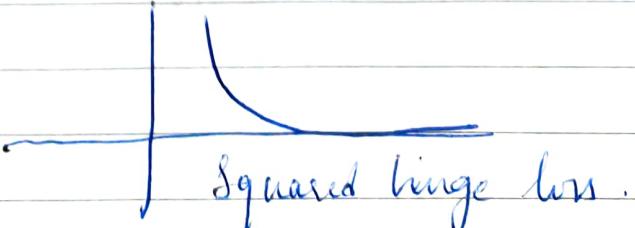
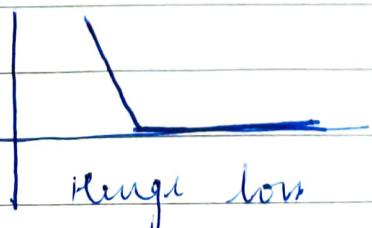
Log likelihood for the entire distribution being

$$\sum_{i=1}^m \log \left(\frac{p(x_i)}{q(x_i)} \right) \begin{cases} > 0 & p(x) \text{ is better method} \\ < 0 & q(x) \text{ is better method} \end{cases}$$

5. Binge loss : $L_i(\theta) = \sum_{j \neq t} \max(0, y_j^{(i)} - \hat{y}_t^{(i)} + 1)$

sum of incorrect class labels.

Squared binge loss : $L_i(\theta) = \frac{1}{2} \sum_{j \neq t} \max(0, y_j^{(i)} - \hat{y}_t^{(i)} + 1)^2$



(4)

Both of the above terms add a penalty to the terms that are misclassified. Here the main idea is to add the terms that are incorrectly classified. The term $\hat{y}_j^{(i)} - \hat{y}_t^{(i)} + 1$ will be positive if there is correct classification and negative when there is correct classification.

	$x^{(1)}$	$x^{(2)}$	$x^{(3)}$
$\hat{y}^{(1)}$	0.5	0.4	0.3
$\hat{y}^{(2)}$	1.3	0.8	-0.6
$\hat{y}^{(3)}$	1.4	-0.4	2.7

$$L_1(\theta) = \sum_{j \in t} \max(0, \hat{y}_j^{(1)} - \hat{y}_t^{(1)} + 1)$$

$$L_1 = \max(0, (1.3 - 0.5) + 1) + \max(0, (0.4 - 0.5) + 1)$$

$$L_1 = 1.8 + 1.9 = 3.7$$

$$L_2 = \max(0, (0.8 - 0.4) + 1) + \max(0, (0.9 - 0.8) + 1)$$

$$L_2 = 0.6 + 0 = 0.6$$

$$L_3 = \max(0, (0.3 - 2.7) + 1) + \max(0, (-0.6 - 2.7) + 1)$$

$$L_3 = 0 + 0 = 0$$

from L_1 and L_2 the predictions are better than L_3 since the values are much higher than zero.

(5)

7. Regularization is performed on a model to yield a simpler model than the original model, this is done because simpler explanations are better the coefficients are made smaller to obtain a model that generalizes better.

$$L(\theta) = \underbrace{\frac{1}{m} \sum_{i=1}^m L_i(f(x_i, w), y_i)}_{\text{loss}} + \lambda R(\theta)$$

Regularization

L_1 regularization $R(\theta) = \sum_{i,j} |\theta_{ij}|$

L_2 regularization $R(\theta) = \sum_{i,j} \theta_{ij}^2$

Adding a regularization term to loss results in weight L_1 regularization makes the weights sparse, i.e., concentrate the weight, L_2 regularization the weights smaller while spreading.

λ should be chosen within a range that its not too high or low, which might lead to the model that is too general for use.

8. we have the gradient terms

$$\frac{\partial L}{\partial \theta} = \frac{1}{m} \sum_{i=1}^m L_i(f(x^{(i)}, \theta), y^{(i)}) + \lambda \frac{\partial}{\partial \theta} R(\theta)$$

L_2 Regularization: $R(\theta) = \sum \theta_i^2 \quad \theta^T \theta = \sum \theta_i^2$
 during gradient descent subtract $\lambda \theta / \partial \theta$ (weight decay)

L1 regularization: $R(\theta) = \sum |θ_i| \rightarrow \frac{\partial}{\partial θ} R(\theta) = \text{sign}(θ)$
 During gradient descent subtract $\frac{\partial}{\partial θ} R(\theta)$.
 (add or subtract depending on the sign).

9. We have the predicted output equation
 $\hat{y} = σ(wx + b)$

Kernal regularization involves regularizing the weights wherein the weights are shrank.
 Bias regularization involves in regularizing the bias term where the bias term is shrank.
 Activity regularization involves in regularizing ~~the~~ both weights and the bias term which intern reduces \hat{y} as \hat{y} is a combination of both w and b .

Optimization

1. Numerical computation of gradients is ~~not~~ not a feasible option when we deal with non linear systems. During which, we use the backpropagation method to find the gradients using chain rule. Using back propagation we can multiply a chain of partial derivatives to obtain the gradients in a faster way compare to numerical computation of gradients.
 Numerical method can be effective while calculating the gradients of a linear system and an activation function is complex that backpropagation is difficult.

2. Gradient descent involves in calculating the optimal weights for the model designed. the update information of the weights is given by the equation

$$\theta^{(i+1)} \leftarrow \theta^{(i)} - \eta \nabla L(\theta^{(i)})$$

here the gradient of the loss is calculated over the entire dataset. We stop updating the weights when the difference in old and new loss is lower than a threshold.

Stochastic gradient involves in calculating the weights in a similar way as we do in the gradient decent method only difference being the loss captured will be calculated using a randomly selected data point. The update equation is as follows:

$$\theta^{(i+1)} \leftarrow \theta^{(i)} - \eta (x^{(i)}, \theta - y^{(i)}) x^{(i)}$$

Stochastic gradient decent converges quickly as compared with gradient decent method because in SGD loss is calculated over one point and weights are upgraded instead of computing it over all datapoints.

3. while performing SGD, for updation of weights if all the data points are considered then

The resulting update of weight would be more accurate when done with a single data point. But performing gradient descent on a single data point is faster computationally when compared to performing it over the entire dataset, which leads to the conclusion, larger the batch size time taken to compute increases with increase in accuracy of the updated weight.

4 main problems with SGD are:

- ① what should be the learning rate?
- ② what happens if loss is more sensitive to me param?
- ③ How to avoid getting stuck at a local minima?
- ④ Minibatch gradient estimates are noisy.

4. Momentum fixes poor conditioning by averaging the gradients over previous gradients:

As far as local minima is concerned (saddle points), momentum reduces the gradient decent with a velocity so that there is no stop

in calculating the gradients.

SGD with momentum smoothens changes to the gradients so that the resulting output is not noisy.

5. In momentum the gradient is computed at the old location whereas in NAG the gradient are computed by predicting the values of the gradient.

while computing the value for $\theta^{(i+1)}$
we get

$$\theta^{(i+1)} \leftarrow \theta^{(i)} + v^{(i+1)} + \underbrace{\gamma(v^{(i+1)} - v^{(i)})}_{\text{Acceleration}}$$

There is an acceleration term added to the momentum equation which leads to the form Nesterov accelerated Gradient (NAG).

- NAG converges faster than SGD + momentum and with better results for parameter update.

6. Strategies for learning rate decay.

(i) Step decay:

every k iterations reduce the learning rate by 50% $n \rightarrow n/2$ over k iterations.

(ii) Exponential decay

decrease the learning rate exponentially over time $\eta = \eta_0 e^{-kt}$.

(iii) fractional decay:

Reduce the ~~average~~ learning rate using the following equation $\eta = \eta_0 / (1 + kt)$.

7. Newton's method is as follows:

Find x such that $f(x) = 0$ by taking an initial guess of x as x_0 .

Find the update such that $f'(x_0 + \Delta x) = 0$.

Expanding the above equation using Taylor

(10)

series expansion

$$f(x_0 + \Delta x) = f(x_0) + \Delta x f'(x_0) + \dots = 0$$

$$\Rightarrow f(x_0) + \Delta x f'(x_0) = 0 \Rightarrow \Delta x = -\frac{f'(x_0)}{f''(x_0)} \text{ and}$$

continue until $f(x_0 + \Delta x) = 0$.we have the objective to compute $\nabla J(\theta) = 0$
where $\nabla J(\theta)$ is the loss

$$\Rightarrow f(x_0 + \Delta x) = f(x_0) + \Delta x^T \nabla f(x_0) + \dots = 0$$

$$\Rightarrow f(x_0) + \Delta x^T \nabla f(x_0) = 0$$

$$\Rightarrow \nabla J(\theta) + \nabla(\nabla J(\theta_0)) \Delta \theta = 0$$

The term $\nabla(\nabla J(\theta_0))$ is known as the hessian matrix denoted by H .

$$H = \begin{bmatrix} \frac{\partial^2}{\partial \theta_1 \partial \theta_1}, & \dots & \frac{\partial^2}{\partial \theta_1 \partial \theta_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2}{\partial \theta_n \partial \theta_1}, & \dots & \frac{\partial^2}{\partial \theta_n \partial \theta_n} \end{bmatrix}$$

$$\Rightarrow \nabla J(\theta_0) + H \Delta \theta = 0$$

$$\Rightarrow \Delta \theta = -H^{-1} \nabla J(\theta_0),$$

As per newton's method

$$\theta^{(i+1)} \leftarrow \theta^{(i)} - H^{-1} \nabla J(\theta^{(i)})$$

here we use H^{-1} as the learning rate.Hessian matrix is a matrix comprising of all the gradients. We compute H^{-1} to find the learning rate η .

8. Condition numbers are generally used to describe the sensitivity of the parameters using the singular values of the hessian matrix. Condition number is the ratio of the largest singular values in the hessian matrix to the smallest singular value.

$$\text{condition number} = \frac{\sigma_{\text{largest}}}{\sigma_{\text{smallest}}}$$

In case of poor conditioning the condition number increases

9. In Adagrad we replace H with different preconditioners. This is done as follows

$$B^{(i)} = \text{diag} \left(\sum_{j=1}^i \nabla J(\theta^{(i)}) \nabla J(\theta^{(i)})^\top \right)^{1/2}$$

$$B^{(i)} = \begin{bmatrix} \sqrt{\sum (\frac{\partial J}{\partial \theta_1})^2} \\ \vdots \\ \sqrt{\sum (\frac{\partial J}{\partial \theta_n})^2} \end{bmatrix}$$

$$\theta^{(i+1)} \leftarrow \theta^{(i)} - \eta B^{(i)}^{-1} \nabla J(\theta^{(i)})$$

so we approximate the hessian matrix as follows

$$\frac{\partial^2 J}{\partial \theta^2} = \sum \left(\frac{\partial J}{\partial \theta_i} \right)^2$$

15. In original we use elementwise scaling of gradients based on history of gradients due to having a large squared sum of past gradients that can be handled in small ~~steps~~ steps.

RUS prop works this by using a array just (eq. 0.9) while addressing new gradient to the gradient sum.

The $\mathcal{E}^{(t+1)}$ feature looks as follows

$$\mathcal{E}^{(t+1)} = \left(f \sum_{j=1}^K (g_j^{(t)})^2 + (1-f)(g_k^{(t)})^2 \right)^{\frac{1}{2}}$$

$\underbrace{\phantom{\sum_{j=1}^K}}_{k^{\text{th}} \text{ component of } \nabla \mathcal{L}(\theta^{(t)})}$

16. To use sometimes RUS prop with momentum in the following way.

Let assume the first momentum velocity with momentum:

$$m_1^{(t+1)} = \beta_1 \cdot m_1^{(t)} + (1-\beta_1) \nabla L(\theta^{(t)})$$

$$m_2^{(t+1)} = \beta_2 \cdot m_2^{(t)} + (1-\beta_2) \nabla L(\theta^{(t)}) \odot \nabla L(\theta^{(t)})$$

$$\theta^{(t+1)} = \theta^{(t)} - \eta \frac{m_1^{(t+1)}, 1}{\|m_2^{(t+1)}\| + \epsilon}$$

Because the momentum are initialized to zero while dividing by the second momentum we get a large step.

$$\hat{m}_1^{(t+1)} = \frac{m_1^{(t)}}{1-\beta_1^t}$$

$$\hat{m}_2^{(t+1)} = \frac{m_2^{(t)}}{1-\beta_2^t}$$

$$\text{so set } \theta^{(t+1)} = \theta^{(t)} - \eta \hat{m}_1^{(t+1)}, \Gamma$$

12. Gradient descent with line search involves in finding the best step size instead of fixing the first step size.

For a given direction $u = \nabla f(x)$ best step size is given by $\eta^* = \underset{\eta}{\operatorname{argmin}} f(x + \eta u)$ and perform gradient descent $\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta^{(t)} \nabla f(\theta^{(t)})$. Then value for η^* to find explicit computation or use gradient descent or use simple line search.

Bracketing is performed by grouping points into a bracket like $[a, b, c] \quad x = (b+c)/2$
 if $f(x) \leq f(b) \Rightarrow [b, x, c]$
 else if $f(x) > f(b) \Rightarrow [a, b, x]$.

Continue this process to produce smaller and smaller brackets until the brackets are small enough.

An alternate to bracketing is to perform a successive line search: Start with θ_0 and direction set $\{u^{(1)}\}$. Compute $\eta^{(1)} = \underset{\eta}{\operatorname{argmin}} f(\theta^{(1)} + \eta u^{(1)})$
 $\theta^{(1+\dagger)} \leftarrow \theta^{(1)} + \eta^{(1)} u^{(1)}$.

The above equation is not as efficient as gradient descent because the step calculations are not as accurate but QD with line search is less heavy computationally and gives faster results.

13. Quasi Newton methods approximate the Hessian matrix inverse using gradient evaluations. Although this method is efficient this method requires large number of datapoints.

Compute quasi newton direction

$$\Delta \theta = (-H^{(i-1)})^{-1} \nabla J(\theta^{(i-1)})$$

Determine step size η^* with bracketing line search and then compute parameter update with $\theta^{(i)} = \theta^{(i-1)} + \eta^* \Delta \theta$

Finally compute updated hessian approximate $H^{(k)}$.

For BFGS update $\Delta \theta = \theta^{(i)} - \theta^{(i-1)}$

$$\nabla J = \nabla J(\theta^{(i)}) - \nabla J(\theta^{(i-1)})$$

inverse update

$$(H^{(i)})^{-1} = \left(I - \frac{\Delta \theta \Delta J^T}{\Delta J^T \Delta \theta} \right) (H^{(i-1)})^{-1} \left(I - \frac{\Delta J \Delta \theta^T}{\Delta J^T \Delta \theta} \right) + \frac{\Delta \theta}{\Delta J^T \Delta \theta}$$

BFGS is comparatively faster than the newton algorithm but less accurate.

BFGS is comparatively slower than ADAM but produces more accurate results.

Regularization:

1. For weight decay we multiply each coefficient by $\rho \in [0, 1]$. As iterations progress weights that are not reinforced decay to 0 which is equivalent to adding regularization term to the loss function.

2. Early stopping is done by stopping when validation error increases instead of training error stops decreasing. This can be interpreted as ~~is~~ ~~less~~ regularization on weights in loss.

Strategies :

① Run validation data - I

Retrain on all data using the number of iterations determined through validation.

We need more data so that we can confidently say that the number of iterations computed is the correct number of iterations.

② Run validation data II

Retrain on all data from previous weights until validation loss is bigger than training loss. The idea here is to check if both training and validation loss converges and hence checking if the loss obtained is correct.

3. To prevent overfitting of data, synthetic data is added to increase variability in training leading to better generalization.

The following tasks can be performed:

① Augment by interpolating b/w examples or by adding noise.

② Augment in feature or data domain.

③ Augment by ~~interpolating~~ transforming data using crop / rotate / scale / change intensity of images

④ Popular in image classification to introduce scale / illumination / rotation invariance.

4. At each training stage dropout layers are added to fully connected layers with a probability of $(1-p)$ where ' p ' is the hyperparameter.

The removed nodes are reinitialized with the original weights in the subsequent stage.

Advantages: Reduce node interactions, node ~~fat~~ overfitting, reduce dependency on single node, distribute features across multiple nodes and ~~increase~~ training speed.

Disadvantages: longer training due to dropout

5. Multiplying the output of each node by p is equivalent to computing expected value for 2^n dropped out networks

$$\hat{y} = E_D [f(x, D)] = \int p(D) f(x, D) dD$$

During tuning all nodes are utilized and hence leads to a higher values. To approximate this, the o/p is multiplied by θ or $(1-p)$.

- 6 Batch normalization normalizes the output of the previous activation layer by subtracting the batch mean and dividing by the batch standard deviation.

During training of deep neural networks one or many weights value can become high which might affect the flow in the network this is where Batch normalization comes into picture. Here when each batch is normalized at each layer this imbalance in weights are normalized to obtain a stable network.

Normalizing data induces randomness to the model.

Advantages :

- ① Activations are not saturated.
- ② normalizes values compact for each batch
- ③ Suitable for back propagation

Dis The model generated can loose information of some key info, the model can get too general.

7. Saturations is required during learning for which we scale and shift after normalization.

$$\{z^{(0)}\}_{j=1}^q \rightarrow \{\tilde{z}^{(1)}\}_{j=1}^q \rightarrow \{z^{(2)}\}_{j=1}^q$$

$$\tilde{z}_j^{(1)} = \sigma_j z_j + \beta_j$$

σ_j and β_j are first learnt then the network can learn to cancel BN if there is no need for it.

During training we do Batch normalization because batches are random, the added randomness reduces overfitting.

8. Ensemble classifiers track multiple independent models and uses majority vote or average during testing.

Different models generated from ensemble overfit in different ways and hence cancel out each other's method of overfitting when a collective decision is made.

Different models can be obtained by

- changing data
- changing parameters.
- Record multiple snapshots of the model during training.