

CS 577 - F22 - Assignment 3 Writing Portion

Due date: 10/25/2022

Anas Puthawala - A20416308
Professor Gady Agam

Loss



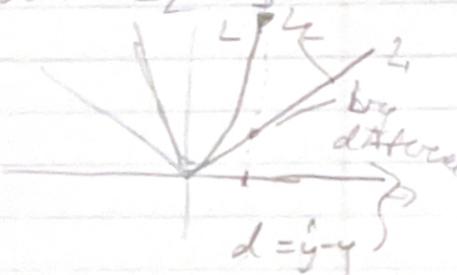
L₁ loss - regression problems, minimize distance between known & pred. value

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

L₂ Loss - same scenario as above, except it's squared distance

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

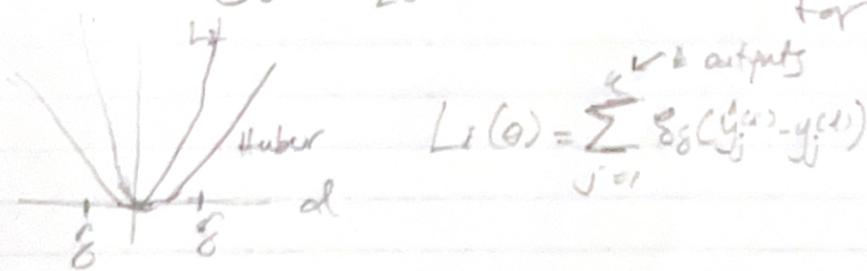
difference between L₁ & L₂ is that L₂ is more sensitive to outliers



Huber loss - between L₁ & L₂ loss, caps the losses of outliers

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

& quadratic for small d, linear for large d.



Log-Cosh loss - similar to huber loss

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

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

reduce sensitivity to outliers

- 1) The main purpose of these loss functions as mentioned above are for regression tasks. They all differ in terms of sensitivity to outliers like for example L1 loss is a lot less sensitive to outliers as opposed to L2 loss. Huber, log-cosh loss are derivatives of the same type of loss function in terms of differing sensitivities to outliers. Each of them can be configured to the task at hand and the data at hand (i.e. choose a delta to properly penalize the outliers / data distribution) to improve upon learning.

2)

2) cross entropy loss.

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

you want to maximize the likelihood argmax(\mathcal{L})

but you want to use a monotonically inc. or dec function so you use $\log(\dots)$

$$\therefore \log(\mathcal{L}) = \text{log-likelihood}$$

now you want to minimize the negative log likelihood *

$$= -\sum_{i=1}^m \sum_{j=1}^k y_j^{(i)} \log(p(y=j|x^{(i)}))$$

The worst cross entropy value is ∞ b/c you will have $\log(\dots)$ very small # and so the actual loss will be very high

Additionally, when the classifier is really bad then it'll just output the same value for each of the units and in that case the output will be $\log(k)$.

3)

3) softmax loss

$$L_s(\theta) = -\sum_{j=1}^k y_j^{(i)} \log(p(y=j|x^{(i)}))$$

You would use softmax loss when dealing with multi-class classification problems. It turns the output vector into a probability that it belongs to a class in the class list of $j=1$ to K

4) The KL loss measures the difference between two probability distributions. The KL divergence is basically the expected value of the log of the likelihood ratio of the two distributions p and q. If p and q are the same we have the minimum value, 0 (no error). The KL divergence is the same as cross-entropy when the class labels do not change. In other words, you can expand the KL loss to have two main terms; negative entropy and cross entropy and if you assume that the negative entropy is always the same since the dataset doesn't change then the loss is effectively identical to the cross entropy

$$4) \text{KL loss eqn } L(\theta) = -\sum_{i=1}^m \sum_{j=1}^k y_i^{(i)} \log \left(\frac{y_i^{(i)}}{q(x_i)} \right)$$

$$\frac{P(x_1, \dots, x_m)}{q(x_1, \dots, x_m)} = \prod_{i=1}^m \frac{P(x_i)}{q(x_i)}$$

≥ 1

p(x) better
q(x) better

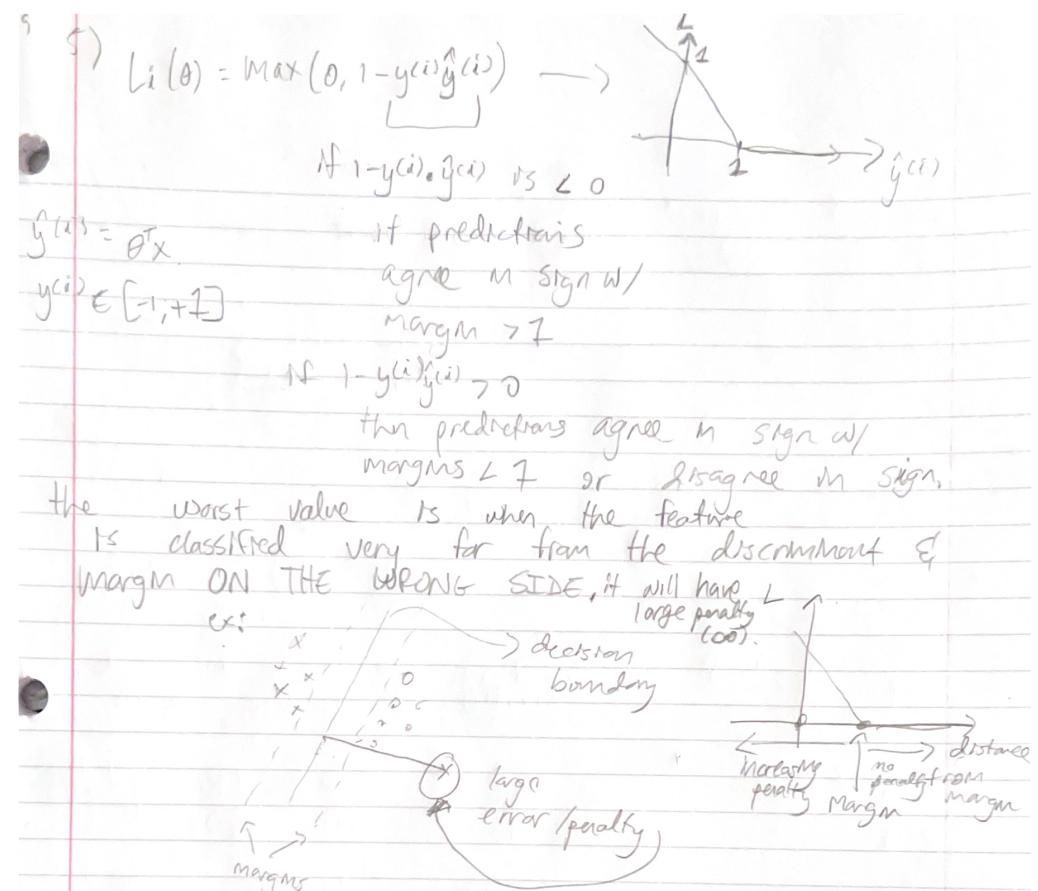
$$\log \left(\prod_{i=1}^m \frac{P(x_i)}{q(x_i)} \right)$$

$$\sum_{i=1}^m \log \left(\frac{P(x_i)}{q(x_i)} \right)$$

$$\mathbb{E} \left[\frac{P(x_i)}{q(x_i)} \right] = \int p(x) \log \left(\frac{P(x)}{q(x)} \right) dx$$

$$\sum_{i=1}^m P(x_i) \log \left(\frac{P(x_i)}{q(x_i)} \right) = \text{KL}(P || Q) = \sum_{i=1}^m P(x_i) \log(P(x_i)) - \sum_{i=1}^m P(x_i) \log(q(x_i))$$

5) The motivation behind hinge loss is to have the neurons delineate a margin between the two classes. In other words, you want the discriminant to be of better certainty due to margins on either side. The y_i hat is the distance from the decision boundary (no longer a probability as in the previous losses)



Ultimately the idea is to introduce a margin in order to train the discriminant to separate the classes as best as possible and that the classes are separated as far as possible.

Multi-class hinge loss:

$$\text{Multi-class } L_i(\alpha) = \sum_{j \neq t} \max(0, \hat{y}_j^{(i)} - y_t^{(i)} + 1)$$

$$y_j^{(i)} = \alpha_j^T x \quad y^{(i)} \in [1, k]$$

Squared hinge loss & worst case expected value:

Squared hinge loss b'may:

$$L_i(\alpha) = \max(0, 1 - y^{(i)} \hat{y}^{(i)})^2$$

multi-class:

$$L_i(\alpha) = \frac{1}{2} \sum_{j=1}^k \max(0, \hat{y}_j^{(i)} - y_t^{(i)} + 1)^2$$

worst case value

before learning

is $k-1$, where k is # classes.

$$L_i(\alpha) = (k-1) * \max(0, 0 - 0 + 1)$$

$$= k-1 * \max(0, 1)$$

$$= \boxed{k-1}$$

6)

$$\begin{array}{c} \hat{y}^{(1)} \\ \hat{y}^{(2)} \\ \hat{y}^{(3)} \end{array} \quad \begin{array}{c} x^{(1)}, \\ 0.5 \\ 0.4 \\ 0.3 \end{array} \quad \begin{array}{c} x^{(2)}, \\ 1.3 \\ 0.8 \\ -0.6 \end{array} \quad \begin{array}{c} x^{(3)}, \\ 1.4 \\ -0.4 \\ 2.7 \end{array} \quad L_i = \sum_{j \neq t} \max(0, \hat{y}_j^{(i)} - y_t^{(i)} + 1)$$

y 1 2 3

$$L_1 = \max(0, 0.4 - 0.5 + 1) + \max(0, 0.3 - 0.5 + 1) = \max(0, 0.9) + \max(0, 0.2) = 1.7 \#$$

$$L_2 = \max(0, 1.3 - 0.8 + 1) + \max(0, -0.6 - 0.8 + 1) = 1.5 + 0 = 1.5 \#$$

$$L_3 = \max(0, 1.4 - 2.7 + 1) + \max(0, -0.4 - 2.7 + 1) = 0 + 0 = 0 \#$$

7) Adding a regularization term can tackle overfitting. The purpose of the regularization term is to introduce an additional penalty to the loss term. By doing this it increases generalizability (more bias). The primary difference between the two are that L1 regularization makes the weights more sparse whereas L2 regularization spreads the weights. L1 reg. is just the absolute value of the weights matrix whereas L2 reg. is the squared value. The way to choose a regularization coefficient is simply by trial and error (i.e. hyper parameter tuning). You can just use a scheme such as grid-search or something to find the optimal regularization term by trial and error (pick a term, train the model, evaluate the performance, change the term and repeat and choose the one that results in the best performance).

8) The L1 and L2 loss have additional terms when dealing with gradient descent. For L1 it's $\lambda \text{sign}(\theta)$ and L2 it's $2\lambda\theta$. The regularization terms change the penalty based on the weights. The way they actually change the loss is mentioned up above in my answer to question 7, and here is the derivation:

$$L_2 = R(\theta) : \sum_i \theta_i^2 = \theta^T \theta \quad \frac{\partial R}{\partial \theta} = 2\theta$$

$$L_1 = R(\theta) \sum_i |\theta_i| = \theta \quad \frac{\partial R}{\partial \theta} = \text{sign}(\theta)$$

You add λ , regularization coefficient on both of the terms:

$$L_2 = 2\lambda\theta, L_1 = \lambda \text{sign}(\theta)$$

weight decay

9) Kernel Regularization:

Regularization of the parameters that reduce the weight and introduce a weight decay on the weights 'w'.

Bias Regularization:

Regularization of the bias terms, aims to reduce the bias and expects smaller output values.

Activity Regularization:

Regularization of the parameter y_{hat} by introducing a decay on the output y_{hat} and reducing the output value, bringing it closer to lower values.

Optimization

1. The advantage of using back-prop to compute gradients numerically is that they're easier to compute and use symbolic derivatives. The numerical approach may also in some cases be slower. One possible use of the direct computation of gradients is perhaps for verification purposes.
2. Stochastic gradient descent is when you update the weights with batches from the training dataset. Whereas in normal gradient descent you update the weights after every epoch (i.e. seeing the entire training dataset). SGD is faster because you're updating the weights sooner as you progressively see the batches.
3. Larger batch sizes serve to decrease the variance for the SGD updates but may end up deterring the model performance the same way a model can over-fit to the data. Smaller batch sizes may take longer to converge and will simply approach the GD algorithm but would also improve the bias of the model. The four main problems with SGD are as follows:

- (1) What should be the optimal LR?
 - (2) What happens if the loss is too sensitive to our parameters, i.e. updates too fast in one direction or too slow in another.
 - (3) How can you avoid getting stuck in a local minima / saddle point
 - (4) Mini-batch gradient estimates may be noisy
4. SGD w/ momentum takes care of poor conditioning by smoothing it out by averaging the gradient previous gradients.
 It addresses saddle points because there is a velocity vector that is also calculated at each step, so if it reaches a local minimum (where the gradient may be 0) the idea is that the velocity vector will still be an element that will remain and will help to continue the iteration and move us out of the local minimum.
 And lastly, it addresses noisy gradients by smoothing them out via moving average.
5. NAG is more accurate than simple momentum in the sense that it provides a more precise update for the gradient because it's calculated at the future time-step. It's called accelerated because of the following term in the update equation:

$$\beta(v^{(t+1)} - v^{(t)})$$

Difference of two velocities --> acceleration.

6. There is a step decay where by a certain step you simply may choose to half the learning rate (for example).
 There is exponential decay, where the learning rate starts off high and exponentially decays. You may tune the decay rate and the iteration index of when the decay will occur.
 And lastly there is fractional decay whereby the learning rate decays by some given fraction. Similar to the exponential decay case, you may tune the decay rate and iteration index of when the decay will occur.

1) step decay:
 every k iterations $\eta \leftarrow \eta/2$



2) exponential decay
 $\eta = \eta_0 e^{-k/t}$
 decay rate
 iteration index

3) fractional decay
 $\eta = \eta_0 / (1 + k^t)$

7. The learning rate in the newton's method is computed by taking the inverse of the hessian matrix which is the second order partial derivative matrix with respect to components in the theta vector. The basic meaning behind the hessian matrix is that it is a good estimate of the sensitivity of the LR with respect to each of the different parameters. It helps to resolve issue number 2 (Loss too sensitive to parameters) from problem (3).

8. The condition number helps to define the sensitivity using the singular value of the Hessian matrix:
 condition number = largest_singular value / smallest_singular value

The problem gets more difficult when the condition number is high because that means that the matrix is more sensitive to error in the input which ultimately may lead to improper optimization. This is often the case when there is poor conditioning.

9. Adagrad approximates the inverse of the hessian matrix by computing the matrix $B(i)^{-1}$:

Adagrad Hessian approx

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

$$B^{(i)} = \begin{bmatrix} \sqrt{\sum \frac{(\partial J)^2}{\partial \theta_1}} \\ \vdots \\ \sqrt{\sum \frac{(\partial J)^2}{\partial \theta_n}} \end{bmatrix} \quad B^{(i)^{-1}} = \begin{bmatrix} \frac{1}{\sqrt{\sum \frac{(\partial J)^2}{\partial \theta_1}}} \\ \vdots \\ \frac{1}{\sqrt{\sum \frac{(\partial J)^2}{\partial \theta_n}}} \end{bmatrix}$$

approx. for hessian

update rule:

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

10. In Adagrad since we normalize by element-wise sum of the square gradients, the step size will become smaller as the iterations progress. To combat this, RMSprop introduces a decay factor when adding new gradients to the gradient sum.

RMSProp

- In Adagrad:

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

- In RMSprop add decay:

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

- Th. k-th feature

$$\tilde{B}_{kk}^{(i+1)} = \left(S \sum_{j=1}^i (\tilde{g}_k^{(j)})^2 + (1-\gamma) (\tilde{g}_k^{(i)})^2 \right)^{1/2}$$

k-th component
of $\nabla J(\theta^{(i)})$

RMSProp

$$\begin{cases} S_j^{(i+1)} = \gamma S_j^{(i)} + (1-\gamma) \|\nabla_j L(\theta^{(i)})\|^2 \\ \theta_j^{(i+1)} \leftarrow \theta_j^{(i)} - \eta \nabla_j L(\theta^{(i)}) \cdot \frac{1}{S_j^{(i+1)} + \epsilon} \end{cases}$$

11. You use a bias correction term because in the beginning when $m_1 = m_2 = 0$ will result in too large of a step from the second moment in the update equation. So to prevent the m_2 from being 0 we add a bias correction term that divides the number based on the iteration number so as to increase the initial moment and therefore prevent the update from being too large since the moment we are interested in is in the denominator and a larger m_2 will serve to decrease the overall update:

Also, Adam combines RMSprop and momentum by incorporating velocity w/ momentum (first moment) and element-wise step scaling (second moment) to the update equation:

Adam

$$\begin{aligned} M_1^{(i+1)} &= \beta_1 \cdot m_1^{(i)} + (1-\beta_1) \nabla L(\theta^{(i)}) && \text{first moment:} \\ M_2^{(i+1)} &= \beta_2 \cdot m_2^{(i)} + (1-\beta_2) (\nabla L(\theta^{(i)}) \odot \nabla L(\theta^{(i)})) && \text{velocity w/ momentum,} \\ \theta^{(i+1)} &\leftarrow \theta^{(i)} - \eta \frac{m_1^{(i+1)}}{\sqrt{M_2^{(i+1)} + \epsilon}} && \text{second moment,} \\ &&& \text{elementwise step scale} \\ m_1^{(i+1)} &= \frac{m_1^{(i)}}{1-\beta_1} && \\ m_2^{(i+1)} &= \frac{m_2^{(i)}}{1-\beta_2} && \end{aligned}$$

12. The idea behind line search is that instead of a fixed step size try and find the 'optimal' step size in a given direction.

given direction: $\alpha = \nabla f(x)$

best step size: $\alpha^T \nabla f(x + \alpha t)$

algorithms

1. expect step decay
2. broadest descent
3. simple line search

G.D.: $\theta^{(i+1)} \leftarrow \theta^{(i)} - \eta \nabla f(\theta^{(i)})$

Bracketing is a process for searching the minimum point

Given bracket $[a, b, c]$

$$x = \frac{b+c}{2}$$

If $f(x) \leq f(b) \Rightarrow [b, x, c]$

If $f(x) > f(b) \Rightarrow [a, b, x]$

and continue w/ smaller & smaller brackets until the bracket is small enough

An alternative to bracketing can be coordinate descent. The advantage is that it's more accurate but the disadvantage is that it's computationally more expensive.

13. Quasi-newton methods basically approximate the inverse of the hessian matrix via gradient evaluation. The advantage of BFGS algorithm over Newtons is the time complexity. In BFGS it's $O(N^2)$ whereas in Newtons its $O(N^3)$. The advantage of BFGS when comparing with Adam is that it's more accurate and but the disadvantage is that it's more computationally expensive than Adam.

Regularization

1. You multiply each coefficient by p , where $p \in [0,1]$ and as the iterations progress, weights that are not multiplied by the coefficient and therefore reinforced, decay to 0. This process is effectively equivalent to adding a regularization term to the loss.
2. Early stopping basically halts the training process once the model begins to overfit, i.e. when the models validation loss begins increasing and the validation accuracy begin decreasing but the training loss continues to decrease and the training accuracy continue to increase.

Strategy to reuse validation data:

1. Once early stopping kicks in and you know the number of epochs where the model begins to overfit, you retrain the model and combine the validation and training set and train only up to the number of epochs before it overfits.
2. You can continue to train from the previous weight using the entire data when the validation loss is bigger than the training loss.
3. Data augmentation is creating data with transformation etc. applied from the existing feature data set and introducing this newly created data into the feature set for model training. The goal is to improve variability, generalizability, and overall performance in the model by increasing the data at hand for the model to learn with.
4. Dropout is performed by dropping (or a better term is deactivating) random nodes (with probability of $1-p$) in fully connected layers where p is a hyper-parameter. The main idea behind this is that it addresses overfitting by improving the generalizability of the model by forcing other neurons (nodes) to learn features when the various other neurons are dropped (turned off / deactivated by dropout). This helps reinforce the neurons' feature detection capability and just overall helps the model to overfit less.

The disadvantage is that it may take longer for the model to properly train since not all units are being used in the training stage (since some are being dropped due to dropout).

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

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

*↑
mask for all n nodes*

6. In training batch normalization is performed by taking the mean and standard deviation of the mini batch at the respective layer / neuron and then using the normalization equation of subtracting by the mean and dividing by the standard deviation. So that'll leave us with values created (mean and standard deviation) for all the respective times we normalize a batch, and when we get to inference we just average the values out (for the mean and standard deviation for the batches) that were created and use that mean and standard deviation.

The randomness for batch normalization gets introduced by the input batch.

7. So some saturation is actually needed for the training to stop and the purpose of batch norm is to reduce saturation but not so much so that the model never trains. The scale and shift parameters serve to introduce

some saturation that is needed for the model to converge. So when γ_j and β_j (the parameters used to scale and shift) are equal to standard deviation and the mean (respectively) the normalization will be cancelled. The values can be learned simply by calculating gradients with respect to the parameters (γ_j and β_j) and updating the parameters with optimization process. And a good initial value is to set the mean to 0 and the standard deviation to 1, so γ_j and $\beta_j = 1, 0$ respectively.

8. When we train multiple models that are independent and use the average to make a prediction we reduce overfitting because the prediction is now an average of multiple models which is much less sensitive to the bias and variance of a single model. There are some strategies to obtain multiple models:

1. change the data
2. change the parameters
3. you can record multiple snapshots of the model during the training phase (and vary learning rate)