CS2109S Tutorial 9 Back Propagation (Last F2F Tutorial)

(AY 24/25 Semester 2)

April 11, 2025

(Prepared by Benson)

Contents

Back Propagation

Recap Q1. Back Propagation Q2. Back Propagation for a Deep(er) Network

Training Neural Networks

- Q3. Training Deep Neural Networks
- Q4. Dying ReLU Problem

Bonus. Behind Pytorch Autograd (Practical)

Recap. Back Propagation

v Warm-Up Exercise: Let $\hat{\boldsymbol{y}} = \boldsymbol{W}^{\top} \boldsymbol{x}$.





Recap. Back Propagation

Forward propagation helps us to evaluate $\hat{\mathbf{y}}$ in a neural network given \mathbf{x} . To perform gradient descent, we also need to evaluate $\frac{\partial \mathcal{E}}{\partial \mathbf{W}^{[i]}}$ for all layers (to update the weights).



Recap. Back Propagation

Suppose a neural network takes time T for inference. What's the closest estimate of the training time using the same data?

A. *T*

B. 2*T*

C. 3*T*

D. 4*T*

E. None of the other options

Solution: Forward propagation takes time T. Backward propagation requires time 2T (2 matrix multiplications are involved for input gradient & weight gradient). Hence the total time is 3T.

$$\mathcal{E} = -\frac{1}{n} \sum_{i=0}^{n-1} \left\{ [Y_{0i} \log(\hat{Y}_{0i})] + [(1 - Y_{0i}) \log(1 - \hat{Y}_{0i})] \right\}$$

$$\mathcal{E} = -\left\{ [Y_{00} \log(\hat{Y}_{00})] + [(1 - Y_{00}) \log(1 - \hat{Y}_{00})] \right\} \blacktriangleleft n = 1$$

(i) When $n = 1$, show that $\frac{\partial \mathcal{E}}{\partial \hat{Y}} = \left[-\frac{Y_{00}}{\hat{Y}_{00}} + \frac{1 - Y_{00}}{1 - \hat{Y}_{00}} \right].$
 $\frac{\partial \mathcal{E}}{\partial \hat{Y}} = \left[\frac{\partial \mathcal{E}}{\partial \hat{Y}_{00}} \right] = \left[-\left(\frac{Y_{00}}{\hat{Y}_{00}} + \frac{1 - Y_{00}}{1 - \hat{Y}_{00}} \times (-1) \right) \right]$
 $= \left[-\frac{Y_{00}}{\hat{Y}_{00}} + \frac{1 - Y_{00}}{1 - \hat{Y}_{00}} \right]$



(ii) When
$$n = 1$$
, show that $\frac{\partial \mathcal{E}}{\partial f^{[1]}} = \hat{\mathbf{Y}} - \mathbf{Y}$.
Chain Rule
 $\hat{\partial} \mathcal{E}$
 $\frac{\partial \mathcal{E}}{\partial f^{[1]}} = \begin{bmatrix} \frac{\partial \mathcal{E}}{\partial f^{[1]}_{00}} \end{bmatrix} = \begin{bmatrix} \partial \hat{\mathcal{E}} & \partial \hat{Y}_{00} \\ \partial \hat{Y}_{00} & \partial \hat{f}^{[1]}_{00} \end{bmatrix}$
from (i)
from (i)
 $= \begin{bmatrix} -\hat{Y}_{00} + 1 - \hat{Y}_{00} \\ \partial \hat{Y}_{00} + 1 - \hat{Y}_{00} \end{pmatrix} + \hat{Y}_{00}(1 - \hat{Y}_{00}) \end{bmatrix}$
 $= \begin{bmatrix} -Y_{00}(1 - \hat{Y}_{00}) + (1 - Y_{00}) \hat{Y}_{00} \end{bmatrix}$
 $= \begin{bmatrix} -Y_{00} + \hat{Y}_{00} \end{bmatrix} = \hat{\mathbf{Y}} - \mathbf{Y}$

(iii) When
$$n = 1$$
, show that $\frac{\partial \mathcal{E}}{\partial W_{20}^{[1]}} = \left(\frac{\partial \mathcal{E}}{\partial f^{[1]}}\right)_{00} X_{20}$.

$$\frac{\partial \mathcal{E}}{\partial W_{20}^{[1]}} = \left|\frac{\partial \mathcal{E}}{\partial f_{00}^{[1]}}\right| = \left|\frac{\partial f_{00}^{[1]}}{\partial W_{20}^{[1]}}\right| = \left|\left(\frac{\partial \mathcal{E}}{\partial f^{[1]}}\right)_{00}\right| X_{20}$$

$$\frac{\partial \mathcal{E}}{\partial W_{20}^{[1]}} = \left|\frac{\partial \mathcal{E}}{\partial f_{00}^{[1]}}\right| = \left|\frac{\partial \mathcal{E}}{\partial f^{[1]}_{00}}\right| = \left|\frac{\partial \mathcal{E}}{\partial$$





$$\begin{array}{c} \mathbf{x}_{0i} = 1 \\ \mathbf{x}_{0i} = 1 \\ \mathbf{x}_{1i} \\ \mathbf{w}_{10}^{[1]} \\ \mathbf{w}_{20}^{[1]} \\ \mathbf{w}_{20}^{[1]$$

w

$$\mathcal{E} = -rac{1}{n} \sum_{i=0}^{n-1} \left\{ [Y_{0i} \log(\hat{Y}_{0i})] + [(1 - Y_{0i}) \log(1 - \hat{Y}_{0i})]
ight\}$$

(c) Consider a general case where $n \in \mathbb{N}$. Find $\frac{\partial \mathcal{E}}{\partial f^{[1]}}$. $\frac{\partial \mathcal{E}}{\partial \hat{\mathbf{Y}}} = \left[\frac{\partial \mathcal{E}}{\partial \hat{\mathbf{Y}}_{0i}}\right]_{i=0}^{n-1} = \left[\frac{1}{n}\left(-\frac{Y_{0i}}{\hat{\mathbf{Y}}} + \frac{1-Y_{0i}}{1-\hat{\mathbf{Y}}}\right)\right]^{n-1}$ Chain Rule $\frac{\partial \mathcal{E}}{\partial \boldsymbol{f}^{[1]}} = \left[\frac{\partial \mathcal{E}}{\partial f_{0i}^{[1]}}\right]_{i=0}^{n-1} = \left[\frac{\partial \mathcal{E}}{\partial \hat{Y}_{0i}}\right] \left[\frac{\partial \hat{Y}_{0i}}{\partial f_{0i}^{[1]}}\right]_{i=0}^{n-1} \hat{Y}_{0i} = \sigma(f_{0i}^{[1]})$ $= \begin{bmatrix} \frac{1}{n} \left(-\frac{Y_{0i}}{\hat{Y}_{0i}} + \frac{1 - Y_{0i}}{1 - \hat{Y}_{0i}} \right) & \left| \hat{Y}_{0i}(1 - \hat{Y}_{0i}) \right|_{i=0}^{n-1} \\ = \begin{bmatrix} \frac{1}{n} (-Y_{0i} + \hat{Y}_{0i}) \end{bmatrix}_{i=0}^{n-1} & \sigma'(x) = \sigma(x)(1 - \sigma(x)) \\ = \frac{1}{n} (\hat{Y} - Y) \end{bmatrix}$



(d) Let's say that Grace has 100 samples from cultivar A and 1000 samples from cultivar B that make up the 1100 total samples. To deal with the imbalanced data set, she decided to introduce two hyper-parameters α and β in the loss function, as shown below:

$$\mathcal{E} = -rac{1}{n}\sum_{i=0}^{n-1} \left\{ lpha [Y_{0i} \cdot \log(\hat{Y}_{0i})] + eta [(1 - Y_{0i})\log(1 - \hat{Y}_{0i})]
ight\}$$

Why do you think that she introduced the hyper-parameters α and β ? How should she set their values?

- We would like the weight of cultivar A \approx the weight of cultivar B.
- $\alpha = 10\beta$ (since the dataset has 10 times more samples from cultivar B than those from cultivar A).

Q2. Back Propagation for a Deep(er) Network





$$f^{[1]} = (W^{[1]})^{\top} X$$

$$a^{[1]} = g^{[1]}(f^{[1]}) = \text{ReLU}(f^{[1]})$$

$$f^{[2]} = (W^{[2]})^{\top} a^{[1]}$$

$$\hat{Y} = g^{[2]}(f^{[2]}) = \sigma(f^{[2]})$$

Suppose we use σ (the sigmoid function) as our activation function in a neural network with 50 hidden layers, as per the code in the accompanying Python notebook.

(a) Play around with the code. Notice that when performing back propagation, the gradient magnitudes of the first few layers are extremely small. What do you think causes this problem?



Q3. Training Deep Neural Networks



Suppose we use σ (the sigmoid function) as our activation function in a neural network with 50 hidden layers, as per the code in the accompanying Python notebook.

(b) Based on what we have learnt thus far, how can we **mitigate** this problem? Test out your solution by modifying the code and checking the gradient magnitudes.

Sigmoid Function:



ReLU Function:





ReLU Function:

Q4. Dying ReLU Problem

This problem occurs when majority of the activations are 0 (meaning the underlying pre-activations are mostly non-positive), resulting in the network dying midway. The gradients passed back are also 0 which leads to poor gradient descent performance and hence poor learning. Refer to the figure below on the ReLU and Leaky ReLU activation functions.



How does Leaky ReLU fix this? What happens if we set a = 1 in Leaky ReLU?

Q4. Dying ReLU Problem ReLU Function:



Leaky ReLU Function:

Beyond CS2109S

Introductory

CS2109S Intro to AI and ML

Cool 5k mods

CS5339 Theory and Algo for ML

Algos, e.g. Perceptron, SVM, Kernels ML Theory, e.g. Conc. Measures, VC-dim **Prerequisite**: CS3264

CS5275 Algo Designer's Toolkit

Math Tools for Algo/ML: Randomized, Optimization, Info Theory, and more **Prerequisite**: CS3230 Theoretical Foundations

CS3263 Foundations of AI

↔ CS4246 AI Decision Making

CS3264 Foundations of ML

Sem 1 (Harold): Math-intensive version of CS2109S (linear algebra?) Sem 2 (Bryan): Focuses on proving / more "old-school" techniques

Applications

CS4243 Computer Vision

CS4248 Natural Language P.



Student Feedback Exercise: Your Voice Matters!

Provide your feedback now >>

https://blue.nus.edu.sg/blue/



Bonus. Behind Pytorch Autograd (Practical)

Copy the template code from the website.

Complete the function visualize_autograd_graph(loss) to print out the Pytorch autograd graph. Get started by reading the provided starter code and the comments.

```
Sample neural network with 2 linear layers
  class TwoLayerNet(nn.Module):
      def __init__(self, D_in, H, D_out):
          super(TwoLayerNet, self).__init__()
4
          self.linear1 = nn.Linear(D_in, H)
          self.relu = nn.ReLU()
          self.linear2 = nn.Linear(H, D_out)
      def forward(self, x):
0
          x = self.linear1(x)
          x = self.relu(x)
          x = self.linear2(x)
          return x
14
  model = TwoLayerNet(1000, 100, 10)
  loss = nn.MSELoss()
```



Extra Slide

Bonus. Behind Pytorch Autograd (Practical)

Solution.¹

```
def add_nodes(var):
      if var in visited:
2
          return
      visited.add(var)
4
5
      if hasattr(var, "variable"):
6
7
           create_parameter_node(var, var.variable.size())
      else:
8
9
           create_func_node(var)
11
      if hasattr(var, "next_functions"):
           for u in var.next_functions:
               if u[0] is not None:
                   create_edge(u[0], var)
14
                   add_nodes(u[0])
```

Extra Slide

¹Adapted from https://www.kaggle.com/code/ehsanmjz/torchviz-example.