

Homework 1

Guilherme Albertini

September 22, 2022

Theory

Let $Linear_1 \rightarrow f \rightarrow Linear_2 \rightarrow g$ be a two-layer neural net architecture whereby $Linear_i(x) = \mathbf{W}^{(i)}\mathbf{x} + \mathbf{b}^{(i)}$ is the i^{th} affine transformation and f, g are element-wise nonlinear activation functions (else must use transposes and/or Hadamard products). When an input $\mathbf{x} \in \mathbb{R}^n$ is fed into the network, $\hat{\mathbf{y}} \in \mathbb{R}^K$ is obtained as output.

Problem 1: Regression Task

We would like to perform regression task. We choose $f(\cdot) = 5(\cdot)^+ = 5ReLU(\cdot)$ and g to be the identity function. To train, we choose MSE loss function, $\ell_{MSE}(\hat{\mathbf{y}}, \mathbf{y}) = ||(\hat{\mathbf{y}} - \mathbf{y})||^2$.

1. Name and mathematically describe the 5 programming steps you would take to train this model with PyTorch using SGD on a single batch of data.

- (a) First, we initialize the parameters to random values (and can optionally save the original parameters in another variable), and tell PyTorch that we want to track their gradients.
- (b) Second, we calculate the predicted values and visualise the actual and predicted values.
- (c) Third, we calculate the average loss of the model using back-propagation, seeking to reduce MSE of our existing model. To do that, we compute the current gradients to approximate how to change update the existing parameters.
- (d) Fourth, we update the parameters via the chosen learning rate.
- (e) Finally, iterate until early stopping or another criterion (i.e. certain number of epochs, target MSE threshold, etc.) is applied. It is best to monitor the training and validation losses and chosen metrics when deciding to stop.

2. For a single data point (x, y) , write down all inputs and outputs for forward pass of each layer. You can only use variables and mechanics specified prior in your answer.

— $Linear_1$ —
 Input: \mathbf{x}
 Output (z_1): $\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$
 — f —
 Input: $Linear_1(\mathbf{x}) = \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$
 Output (z_2): $f(Linear_1(\mathbf{x})) = 5ReLU(Linear_1(\mathbf{x}))$
 $= 5ReLU(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$
 $= 5\max(0, \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$
 — $Linear_2$ —
 Input: $f(Linear_1(\mathbf{x}))$
 Output (z_3): $Linear_2(f(Linear_1(\mathbf{x}))) = \mathbf{W}^{(2)}f(Linear_1(\mathbf{x})) + \mathbf{b}^{(2)}$
 $= 5\mathbf{W}^{(2)}(\max(0, \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})) + \mathbf{b}^{(2)}$
 — g —
 Input: $Linear_2(f(Linear_1(\mathbf{x})))$
 Output: $g(Linear_2(f(Linear_1(\mathbf{x})))) =$
 $5(\mathbf{W}^{(2)}(\max(0, \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})) + \mathbf{b}^{(2)})\mathbf{I} =$
 $5(\mathbf{W}^{(2)}(\max(0, \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})) + \mathbf{b}^{(2)}) = \hat{\mathbf{y}}$

3. Write down the gradients calculated from the backward pass. You can only use the following variables: $\mathbf{x}, \mathbf{y}, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{W}^{(2)}, \mathbf{b}^{(2)}, \frac{\partial \ell}{\partial \mathbf{y}}, \frac{\partial z_2}{\partial \mathbf{y}}, \frac{\partial \hat{\mathbf{y}}}{\partial z_3},$

where z_1, z_2, z_3, \hat{y} are outputs of $Linear_1, f, Linear_2, g$, respectively.

$$\begin{aligned}
\frac{\partial z_3}{\partial W^{(2)}} &= \frac{\partial(5W^{(2)} \max(0, W^{(1)}x + b^{(1)}) + b^{(2)})}{\partial W^{(2)}} \\
&= 5 \max(0, W^{(1)}x + b^{(1)}) \\
\frac{\partial z_3}{\partial b^{(2)}} &= \frac{\partial(5W^{(2)} \max(0, W^{(1)}x + b^{(1)}) + b^{(2)})}{\partial b^{(2)}} \\
&= 1 \\
\frac{\partial \ell}{\partial W^{(2)}} &= \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_3} \frac{\partial z_3}{\partial W^{(2)}} \\
&= 5 \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_3} \max(0, W^{(1)}x + b^{(1)}) \\
\frac{\partial \ell}{\partial b^{(2)}} &= \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_3} \frac{\partial z_3}{\partial b^{(2)}} \\
&= \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_3}
\end{aligned}$$

$$\begin{aligned}
\frac{\partial z_3}{\partial z_2} &= \frac{\partial(W^{(2)} 5 \max(0, W^{(1)}x + b^{(1)}))}{\partial 5 \max(0, W^{(1)}x + b^{(1)})} = W^{(2)} \\
\frac{\partial z_1}{\partial W^{(1)}} &= \frac{\partial(W^{(1)}x + b^{(1)})}{\partial W^{(1)}} = x \\
\frac{\partial z_1}{\partial b^{(1)}} &= \frac{\partial(W^{(1)}x + b^{(1)})}{\partial b^{(1)}} = 1 \\
\frac{\partial \ell}{\partial W^{(1)}} &= \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_3} \frac{\partial z_3}{\partial z_2} \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial W^{(1)}} \\
&= \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_3} W^{(2)} \frac{\partial z_2}{\partial z_1} x \\
\frac{\partial \ell}{\partial b^{(1)}} &= \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_3} \frac{\partial z_3}{\partial z_2} \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial b^{(1)}} \\
&= \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_3} W^{(2)} \frac{\partial z_2}{\partial z_1}
\end{aligned}$$

4. Show the elements of $\frac{\partial z_2}{\partial z_1}, \frac{\partial \hat{y}}{\partial z_3}, \frac{\partial \ell}{\partial \hat{y}}$. Be careful about dimensionality.

$$\frac{\partial z_2}{\partial z_1} = \frac{\partial(5 \max(0, \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}))}{\partial(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})}$$

$$= \begin{cases} 0, & z_1 < 0 \\ 5, & z_1 > 0 \\ \text{undefined (or assigned a value in code)}, & z_1 = 0 \end{cases}$$

Note: $z_1 = \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$

$$\frac{\partial \hat{\mathbf{y}}}{\partial z_3} = 1 \text{ (Due to identity transformation)}$$

$$\frac{\partial \ell}{\partial \hat{\mathbf{y}}} = \frac{\partial(\|\hat{\mathbf{y}} - \mathbf{y}\|^2)}{\partial \hat{\mathbf{y}}} = 2(\hat{\mathbf{y}} - \mathbf{y})^T$$

Problem 2: Classification Task

We would like to perform multi-class classification task, so we set $f = \tanh$ and $g = \sigma$, the logistic sigmoid function, $\sigma(z) = \frac{1}{1+\exp(-z)}$.

1. If you want to train this network, what do you need to change in the equations of (b), (c) and (d), assuming we are using the same MSE loss function.

ya

Problem 3

Proof.

□

Section 2.2

Problem 6

Blah

Problem 7

Blah

Problem 10

Blah