

# Homework1 Backprop

Rakesh Kumar

December 21, 2022

---

## 1 Theory

### 1.1 Two-Layer Neural Nets

$Linear_1 \rightarrow f \rightarrow Linear_2 \rightarrow g$  where  $Linear_i(x) = \mathbf{W}^{(i)}\mathbf{x} + \mathbf{b}^{(i)}$  is the  $i$ -th transformation, and  $f, g$  are elementwise nonlinear activation functions. where input  $\mathbf{x} \in \mathbb{R}^n$  and  $\hat{y} \in \mathbb{R}^n$ .

### 1.2 Regression Task

$f(.) = (. )^+ = ReLU(.)$  and  $g$  is to be identity function.

MSE loss function  $l_{MSE}(y, \hat{y}) = ||\hat{y} - y||^2$  where  $y$  is the target output.

1. Name and mathematically define the 5 programming steps to train the above model architecture with PyTorch using SGD on a single batch of data.

Solution:

step 1. Feed Forward to get the logits.

$$y_{pred} = model(X) \quad (1)$$

$$\mathbf{x} \in \mathbb{R}^n \rightarrow \mathbf{h} \in \mathbb{R}^d \rightarrow \hat{\mathbf{y}} \in \mathbb{R}^k$$

Where:

$$\mathbf{h} = f(\mathbf{W}_h\mathbf{x} + \mathbf{b}_h) \quad (2)$$

$$\hat{\mathbf{y}} = g(\mathbf{W}_y\mathbf{h} + \mathbf{b}_y) \quad (3)$$

$$\mathbf{W}_h \in \mathbb{R}^{d \times n}, \mathbf{b}_h \in \mathbb{R}^d, \mathbf{W}_y \in \mathbb{R}^{K \times d}, \mathbf{b}_y \in \mathbb{R}^K$$

step 2. Compute the loss using different criterion.

$$loss = criterion(y_{pred}, y) \quad (4)$$

Example of loss function:  $l_{MSE}(y, \hat{y}) = \frac{1}{n} \sum ||y_{pred} - y||^2$

**Cross entropy or neagtive likelihood:**

$$L(\hat{Y}, c) = \frac{1}{m} \sum_{i=1}^m (l(y^{\hat{i}}, c_i)) \quad (5)$$

Where:  $l(y^{\hat{i}}, c_i) = -\log(\hat{y}[c])$

step 3. Zeroing the gradients before going through the backward pass.  
`optimizer.zero_grad()`

step 4. Fourth step after zeroing the gradients, is Bakward pass to compute the gradient of loss w.r.t our learnable parameters.

$$\Theta = (\mathbf{W}_h, \mathbf{b}_h, \mathbf{W}_y, \mathbf{b}_y) \quad (6)$$

$$J(\Theta) = L(\hat{Y}(\Theta), c) \in \mathbb{R}^+ \quad (7)$$

$$\frac{\partial J(\Theta)}{\partial \mathbf{W}_y} = \frac{\partial J(\Theta)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{W}_y} \quad (8)$$

$$\frac{\partial J(\Theta)}{\partial \mathbf{W}_h} = \frac{\partial J(\Theta)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{W}_h} \quad (9)$$

step 5. Updating the parameters.

Based on the SGD optimization using backprop we update the learnable parameters.

2. Write doen the forward pass of each layer:

Layer	Input	Output
<i>Linear</i> <sub>1</sub>	$\mathbf{x}, \mathbf{W}^1, \mathbf{b}^1$	$\mathbf{z}_1 = \mathbf{W}^1 \mathbf{x} + \mathbf{b}^1$
f	$\mathbf{W}^1 \mathbf{x} + \mathbf{b}^1$	$\mathbf{z}_2 = f(\mathbf{W}^1 \mathbf{x} + \mathbf{b}^1)$
<i>Linear</i> <sub>2</sub>	$\mathbf{z}_2, \mathbf{W}^2, \mathbf{b}^2$	$\mathbf{z}_3 = \mathbf{W}^2 \mathbf{z}_2 + \mathbf{b}^2$
g	$\mathbf{W}^2 \mathbf{h} + \mathbf{b}_y$	$\hat{\mathbf{Y}} = g(\mathbf{W}^2 \mathbf{h} + \mathbf{b}^2)$
Loss	$\hat{\mathbf{Y}}, \mathbf{c}$	$L(\hat{Y}(\Theta), c)$

3. Gradient calulated from the backward pass:

Parameters	Gradient
$\mathbf{W}^1$	
$\mathbf{b}^1$	
$\mathbf{W}^2$	
$\mathbf{b}^2$	