

Problem 41.  $x \begin{matrix} \nearrow y \\ \rightarrow z \end{matrix}$ ,  $\bar{x} = \bar{y} \cdot \frac{\partial y}{\partial x} + \bar{z} \cdot \frac{\partial z}{\partial x}$

$$\bar{1} = \frac{\partial 1}{\partial 1} = \underline{1}, \quad \bar{y}_1 = \bar{1} \cdot \frac{\partial 1}{\partial y_1}, \quad \bar{y}_2 = \bar{1} \cdot \frac{\partial 1}{\partial y_2}.$$

$$\bar{z}_1 = \bar{y}_1 \cdot \frac{\partial y_1}{\partial z_1} + \bar{y}_2 \cdot \frac{\partial y_2}{\partial z_1}, \quad \bar{z}_2 = \bar{y}_1 \cdot \frac{\partial y_1}{\partial z_2} + \bar{y}_2 \cdot \frac{\partial y_2}{\partial z_2}$$

$$\bar{w}_{11} = \bar{z}_1 \cdot \frac{\partial z_1}{\partial w_{11}}, \quad \bar{w}_{12} = \bar{z}_1 \cdot \frac{\partial z_1}{\partial w_{12}}, \quad \bar{b}_1 = \bar{z}_1 \cdot \frac{\partial z_1}{\partial b_1}$$

$\downarrow$   $\bar{w}_{21}$        $\downarrow$   $\bar{w}_{22}$        $\downarrow$   $b_2$

Problem 26. backprop provides gradients, that are used for GD.  
 $\alpha$  too high = too much change, unable to converge.


$$w - \alpha \cdot \nabla = w - \alpha \cdot \frac{\partial y_1}{\partial z_1}$$



$$w_{11} - \alpha \cdot \nabla = w_{11} - \alpha \cdot \frac{\partial z_1}{\partial w_{11}}$$

$\alpha$  too low = convergence takes a long time.

Problem 27:

$$x_i = \text{~~at~~ } 1 \times d$$

input layer neurons =  $d$ , output layer after 1<sup>st</sup> layer: 

2<sup>nd</sup> layer neurons = , output layer after 2<sup>nd</sup> layer = 

3<sup>rd</sup> layer neurons = ,

$\vdots$

$N$  layers

Output layer:  $\text{input } X$ , output layer's final output = 10,

dimensionality from  $N-1$ 's output

or the # of classes.

the explanation on

Problem 28. In essence, GD is used to train the MLP.

So  $\alpha$ 's being high/low doesn't change if/when the architecture is MLP/CNN/others. ~~xxxx~~ ↓

Check Problem 26.

Problem 29. backprop v. SGD.

SGD = using 1 sample at a time for training

GD = using entire training set at a time for training.

Problem 30. Introduces non-linearity to the otherwise linear layers in MLP.

Problem 31.  $b^{(\text{filter})} = [1, 1, 2]$ ,  $a^{(\text{vector/input})} = [2, -1, 1]$ .

$$b * a = 1 \times [2, -1, 1, 0, 0]$$

$$+ 1 \times [0, 2, -1, 1, 0]$$

$$+ 2 \times [0, 0, 2, -1, 1]$$

$$= [-, -, -, -, -] \quad \& \quad d = 5.$$

Problem 32 ①  $x = [-1, 0, 1, 2, 3]$  ②  $x = [1, 2, 3, 4, 5]$ .

$$\text{ReLU}(\textcircled{1}) = [0, 0, 1, 2, 3] \quad \text{ReLU}(\textcircled{2}) = [1, 2, 3, 4, 5].$$

Problem 33: kernels = applying transformations onto input, so that we get different feature maps.

Multiple kernels = multiple representations.

Problem 34. Dropout: = randomly deactivate certain neurons  
Normalization / LayerNorm.



Problem 35. CNN lecture 21-23.

max pooling: taking the maximum of a given set of activations

avg pooling: " " average of a given region ↗

Problem 36: self attention:  $Q \cdot K^T$

Q: a version (a.k.a. linearly processed "Hi How Are You")

K: a version (a.k.a. linearly processed "Same")

different.

same ↗

Problem 37. Multi-Head Attention = multiple scaled dot product attention.  
Attention slides 15.

Problem 38. ① the self-attention needs to be masked in the decoders.  
② the decoder has two distinct "blocks" of multi-head attention.

Problem 39. Image Captioning = look at (pay attention to) a subset region of the input picture to produce one word/phrase at a time.

problem 40. ① Get rid of the decoder

② add linear layers to the encoder

optional ↙ ③ do a softmax, assuming the dimensionality is correct  
⇒ if not, address it with linear layer(s).

④. get the output