

1. (15 points) **Backpropagation for autoencoders.** In an autoencoder, we seek to reconstruct the original data after some operation that reduces the data's dimensionality. We may be interested in reducing the data's dimensionality to gain a more compact representation of the data.

For example, consider $\mathbf{x} \in \mathbb{R}^n$. Further, consider $\mathbf{W} \in \mathbb{R}^{m \times n}$ where $m < n$. Then $\mathbf{W}\mathbf{x} \in \mathbb{R}^{(m \times n)(n \times 1)} = (m \times 1)$ is of lower dimensionality than \mathbf{x} . One way to design \mathbf{W} so that $\mathbf{W}\mathbf{x}$ still contains key features of \mathbf{x} is to minimize the following expression

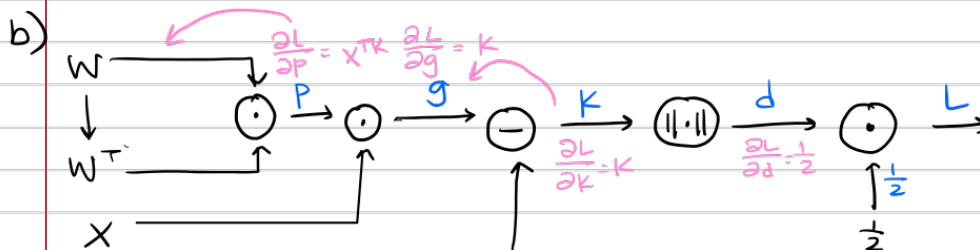
$$\mathcal{L} = \frac{1}{2} \|\mathbf{W}^T \mathbf{W} \mathbf{x} - \mathbf{x}\|^2$$

with respect to \mathbf{W} . (To be complete, autoencoders also have a nonlinearity in each layer, i.e., the loss is $\frac{1}{2} \|f(\mathbf{W}^T f(\mathbf{W}\mathbf{x})) - \mathbf{x}\|^2$. However, we'll work with the linear example.)

- (3 points) In words, describe why this minimization finds a \mathbf{W} that ought to preserve information about \mathbf{x} .
- (3 points) Draw the computational graph for \mathcal{L} . **Hint:** You can set up the computational graph to this problem in a way that will allow you to solve for part (d) without taking 4D tensor derivative.
- (3 points) In the computational graph, there should be two paths to \mathbf{W} . How do we account for these two paths when calculating $\nabla_{\mathbf{W}} \mathcal{L}$? Your answer should include a mathematical argument.
- (6 points) Calculate the gradient: $\nabla_{\mathbf{W}} \mathcal{L}$.

a) In $\mathcal{L} = \frac{1}{2} \|\underbrace{\mathbf{W}^T \mathbf{W} \mathbf{x}}_{\text{reduction into lower dimensionality}} - \mathbf{x}\|_{\text{squared loss}}^2$

The objective of the minimization function is to find a matrix \mathbf{W} that when reconstructed via $\mathbf{W}^T \mathbf{W} \mathbf{x}$, preserves the most information of \mathbf{x} through minimizing the loss \mathcal{L} .



- c) To account for the 2 paths leading to \mathbf{W} , we can use the law of total derivatives and sum the gradients of the 2 paths.

$$\nabla_{\mathbf{W}} \mathcal{L} = \frac{\partial \mathcal{L}}{\partial \mathbf{W}} + \left(\frac{\partial \mathcal{L}}{\partial \mathbf{W}^T} \right)^T$$

d) $\mathcal{L} = \frac{1}{2} \|\mathbf{K}\|^2 \rightarrow \frac{\partial \mathcal{L}}{\partial \mathbf{K}} = \mathbf{K}$

$$\mathbf{K} = \mathbf{g} - \mathbf{x} \rightarrow \frac{\partial \mathcal{L}}{\partial \mathbf{g}} = \frac{\partial \mathcal{L}}{\partial \mathbf{g}} \cdot \frac{\partial \mathcal{L}}{\partial \mathbf{K}} = 1 \cdot \mathbf{K}$$

$$\mathbf{g} = \mathbf{P} \mathbf{x} \rightarrow \frac{\partial \mathcal{L}}{\partial \mathbf{P}} = \frac{\partial \mathcal{L}}{\partial \mathbf{g}} \cdot \frac{\partial \mathcal{L}}{\partial \mathbf{g}} = \mathbf{x}^T \cdot \mathbf{K}$$

$$\mathbf{P} = \mathbf{W}^T \mathbf{W} \rightarrow \frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \frac{\partial \mathcal{L}}{\partial \mathbf{P}} \cdot \frac{\partial \mathcal{L}}{\partial \mathbf{P}} = \mathbf{W}^T \mathbf{x}^T \mathbf{K} = \mathbf{W} \mathbf{x} \mathbf{K}^T$$

$$\rightarrow \frac{\partial \mathcal{L}}{\partial \mathbf{W}^T} = \frac{\partial \mathcal{L}}{\partial \mathbf{W}^T} \cdot \frac{\partial \mathcal{L}}{\partial \mathbf{P}} = (\mathbf{W}^T)^T \mathbf{x}^T \mathbf{K} = \mathbf{W} \mathbf{x}^T \mathbf{K} = \mathbf{W} \mathbf{K} \mathbf{x}^T$$

In the architecture shown, D represents the number of neurons in input layer, H represents the number of neurons in hidden layer, and C represents the number of neurons in the output layer (in our design $C=7$). The output is then passed through a softmax classifier. Although we learned about the ReLU activation in class, we decided to use the Swish activation function (introduced by Google Brain) for the hidden layer. The Swish activation function for any scalar input k is defined as,

$$\text{swish}(k) = \frac{k}{1 + e^{-k}} = k\sigma(k),$$

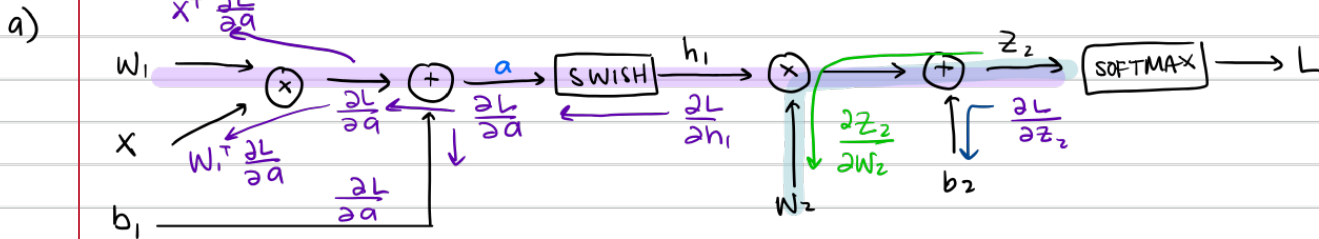
$D = \# \text{ input neurons}$
 $H = \# \text{ hidden neurons}$
 $C = \# \text{ output neurons} = 7$

where $\sigma(k)$ is the sigmoid activation function you have seen in lecture.

You will train the 2-layer FC net using gradient descent and for that you will need to compute the gradients. For the gradient computations, you are allowed to keep your final answer in terms of $\frac{\partial L}{\partial z_2}$.

- (3 points) Draw the computational graph for the 2-layer FC net.
- (5 points) Compute $\nabla_{W_2} L, \nabla_{b_2} L$.
- (7 points) Compute $\nabla_{W_1} L, \nabla_{b_1} L$.

3



b) $z_2 = W_2 h_1 + b_2$ $h_1 = a \cdot \sigma(a) = \frac{a}{1+e^{-a}}$ $a = W_1 X + b_1$

upstream grad = $\frac{\partial L}{\partial z_2}$

$\nabla_{W_2} L = \frac{\partial z_2}{\partial W_2} \cdot \frac{\partial L}{\partial z_2} = \boxed{h_1 \cdot \frac{\partial L}{\partial z_2}}$ (switcharoo?)

$\hookrightarrow \frac{\partial z_2}{\partial W_2} (W_1 h_1 + b_2) = h_1$

$\frac{\partial L}{\partial a} = \frac{\partial h_1}{\partial a} \cdot \frac{\partial L}{\partial h_1}$

$\nabla_{b_2} L = \frac{\partial z_2}{\partial b_2} \cdot \frac{\partial L}{\partial z_2} = \boxed{\frac{\partial L}{\partial z_2}}$
 $\hookrightarrow = 1$

$\frac{\partial h_1}{\partial a} \left(\frac{a}{1+e^{-a}} \right)$

c) $\nabla_{W_1} L = \frac{\partial L}{\partial a} \cdot X^T = \boxed{\sigma(a) + a [\sigma(a)(1-\sigma(a))]} X^T$

$\hookrightarrow \frac{\partial L}{\partial a} = \frac{\partial h_1}{\partial a} \cdot \frac{\partial L}{\partial h_1}$

$\hookrightarrow h_1 = \frac{a}{1+e^{-a}} = a \cdot \sigma(a)$

$\frac{\partial h_1}{\partial a} = \sigma(a) + a \sigma'(a)$
 $\sigma'(a) = \frac{\partial}{\partial a} \frac{a}{1+e^{-a}} = \sigma(a)(1-\sigma(a))$
 $= \sigma(a) + a [\sigma(a)(1-\sigma(a))]$

$\nabla_{b_2} L = \frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial a} = \nearrow$

2. I am a C147 student is