

Brair Tiloon Elberier

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 $\mathbf{W}\mathbf{x}$ 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

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

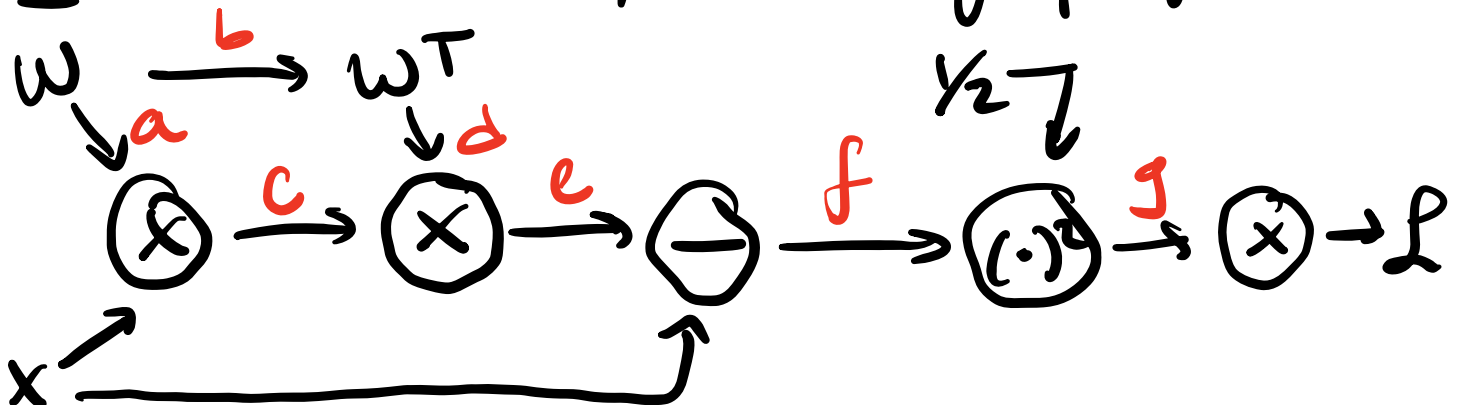
with respect to \mathbf{W} . (To be complete, autoencoders also have a nonlinearity in each layer, i.e., the loss is $\frac{1}{2} \|f(\mathbf{W}^T f(\mathbf{W}\mathbf{x})) - \mathbf{x}\|^2$. However, we'll work with the linear example.)

- (a) (3 points) In words, describe why this minimization finds a \mathbf{W} that ought to preserve information about \mathbf{x} .
- (b) (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.
- (c) (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.
- (d) (6 points) Calculate the gradient: $\nabla_{\mathbf{W}} \mathcal{L}$.

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

The loss function shows us that if the difference between $\mathbf{W}^T \mathbf{W} \mathbf{x}$ and \mathbf{x} is minimized $\Rightarrow f(\mathbf{W}^T f(\mathbf{W} \mathbf{x})) \rightarrow \mathbf{x}$ then $\mathbf{W} \mathbf{x}$ contains information about \mathbf{x} .

b. Draw the computational graph for \mathcal{L} .



C. In the computational graph, there should be 2 paths to W . When calculating $\nabla_W L$, how do we account for these two paths?

We must sum the partial derivatives from the two paths. This will factor both paths into the formulation of $\nabla_W L$.
For example in this problem there are 2 paths

$$p_1: a \rightarrow c \rightarrow e \rightarrow f \rightarrow g \rightarrow L$$

$$p_2: b \rightarrow d \rightarrow e \rightarrow f \rightarrow g \rightarrow L$$

given this example, then in the calculation of $\nabla_W L$

$$\frac{dL}{dp_1} + \frac{dL}{dp_2} \Rightarrow \text{where each term represents the partial derivative product of each subpaths}$$

d. Calculate the gradient $\nabla_w \mathcal{L}$:

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

back-propagate to w through $p1$:

$$\frac{d\mathcal{L}}{dy} = \frac{1}{2} \quad , \quad \frac{d\mathcal{L}}{df} = \frac{1}{2} \cdot 2f = f$$

$$\frac{d\mathcal{L}}{de} = \frac{d}{de}(e^{-x}) \cdot f = f$$

$$\frac{d\mathcal{L}}{dc} = \frac{d(W^T W x)}{d(Wx)} \cdot f = W f$$

$$\frac{d\mathcal{L}}{dW} = \frac{d(W^T W x)}{d(W^T)} \cdot f = (Wx)^T \cdot f$$

$$\frac{d\mathcal{L}}{db} = W f \cdot \frac{d(Wx)}{d(W)} = W f x^T$$

$$\frac{d\mathcal{L}}{da} = W x f^T$$

$$\nabla_w \mathcal{L} = W (W^T W x - x) x^T + W x (W^T W x - x)^T$$

2. (20 points) **Backpropagation for Gaussian-process latent variable model.** (Optional for students in C147: Please write 'I am a C147 student' in the solution and you will get full credit for this problem). An important component of unsupervised learning is visualizing high-dimensional data in low-dimensional spaces. One such nonlinear algorithm to do so is from Lawrence, NIPS 2004, called GP-LVM. GP-LVM optimizes the maximum-likelihood of a probabilistic model. We won't get into the details here, but rather to the bottom line: in this paper, a log-likelihood has to be differentiated with respect to a matrix to derive the optimal parameters.

To do so, we will apply the chain rule for multivariate derivatives via backpropagation. The log-likelihood is:

$$\mathcal{L} = -c - \frac{D}{2} \log |\mathbf{K}| - \frac{1}{2} \text{tr}(\mathbf{K}^{-1} \mathbf{Y} \mathbf{Y}^T)$$

where $\mathbf{K} = \alpha \mathbf{X} \mathbf{X}^T + \beta^{-1} \mathbf{I}$ and c is a constant. The $|\cdot|$ symbol in this context refers to the determinant of a matrix. To solve this, we'll take the derivatives with respect to the two terms with dependencies on \mathbf{X} :

$$\begin{aligned} \mathcal{L}_1 &= -\frac{D}{2} \log |\alpha \mathbf{X} \mathbf{X}^T + \beta^{-1} \mathbf{I}| \\ \mathcal{L}_2 &= -\frac{1}{2} \text{tr}((\alpha \mathbf{X} \mathbf{X}^T + \beta^{-1} \mathbf{I})^{-1} \mathbf{Y} \mathbf{Y}^T) \end{aligned}$$

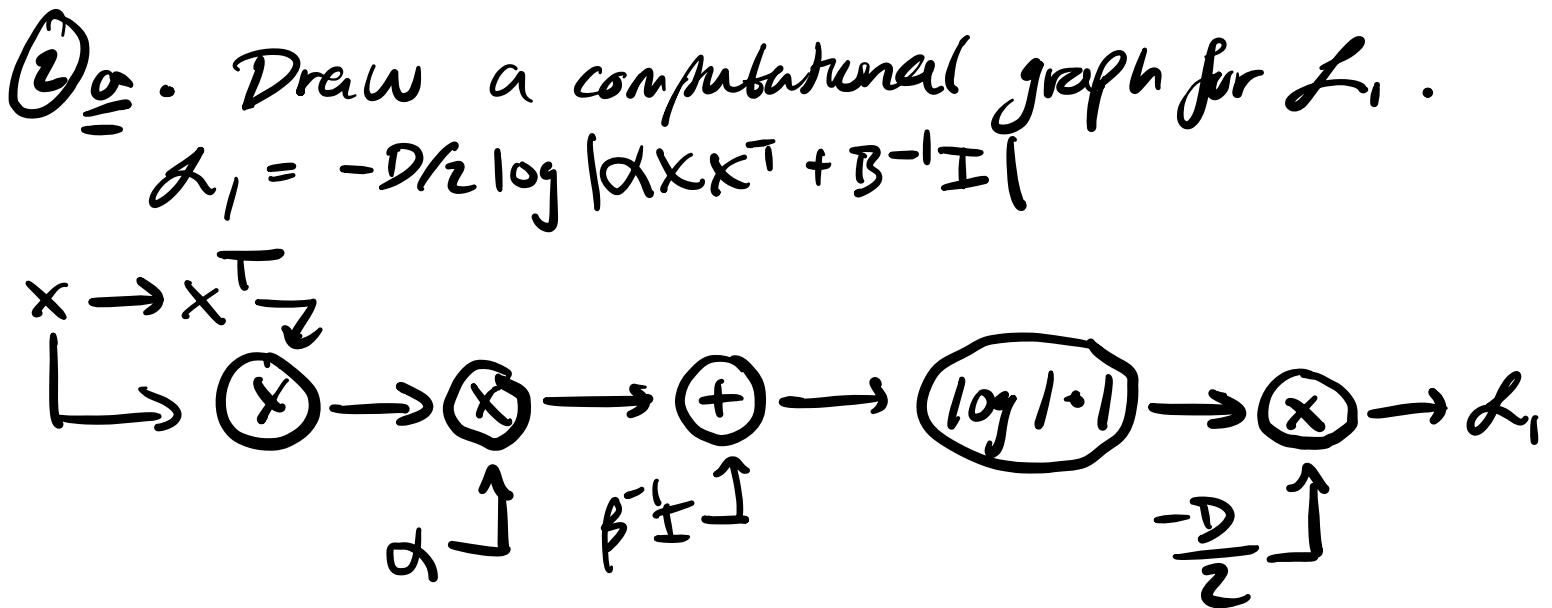
Hint: To receive full credit, you will be required to show all work. You may use the following matrix derivative without proof:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{K}} = -\mathbf{K}^{-T} \frac{\partial \mathcal{L}}{\partial \mathbf{K}^{-1}} \mathbf{K}^{-T}.$$

Also, consider the matrix operation, $\mathbf{Z} = \mathbf{X} \mathbf{Y}$. If we have an upstream derivative, $\partial \mathcal{L} / \partial \mathbf{Z}$, then backpropagate the derivatives in the following way:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{X}} &= \frac{\partial \mathcal{L}}{\partial \mathbf{Z}} \mathbf{Y}^T \\ \frac{\partial \mathcal{L}}{\partial \mathbf{Y}} &= \mathbf{X}^T \frac{\partial \mathcal{L}}{\partial \mathbf{Z}} \end{aligned}$$

- (3 points) Draw a computational graph for \mathcal{L}_1 .
- (6 points) Compute $\frac{\partial \mathcal{L}_1}{\partial \mathbf{X}}$.
- (3 points) Draw a computational graph for \mathcal{L}_2 .
- (6 points) Compute $\frac{\partial \mathcal{L}_2}{\partial \mathbf{X}}$.
- (2 points) Compute $\frac{\partial \mathcal{L}}{\partial \mathbf{X}}$.



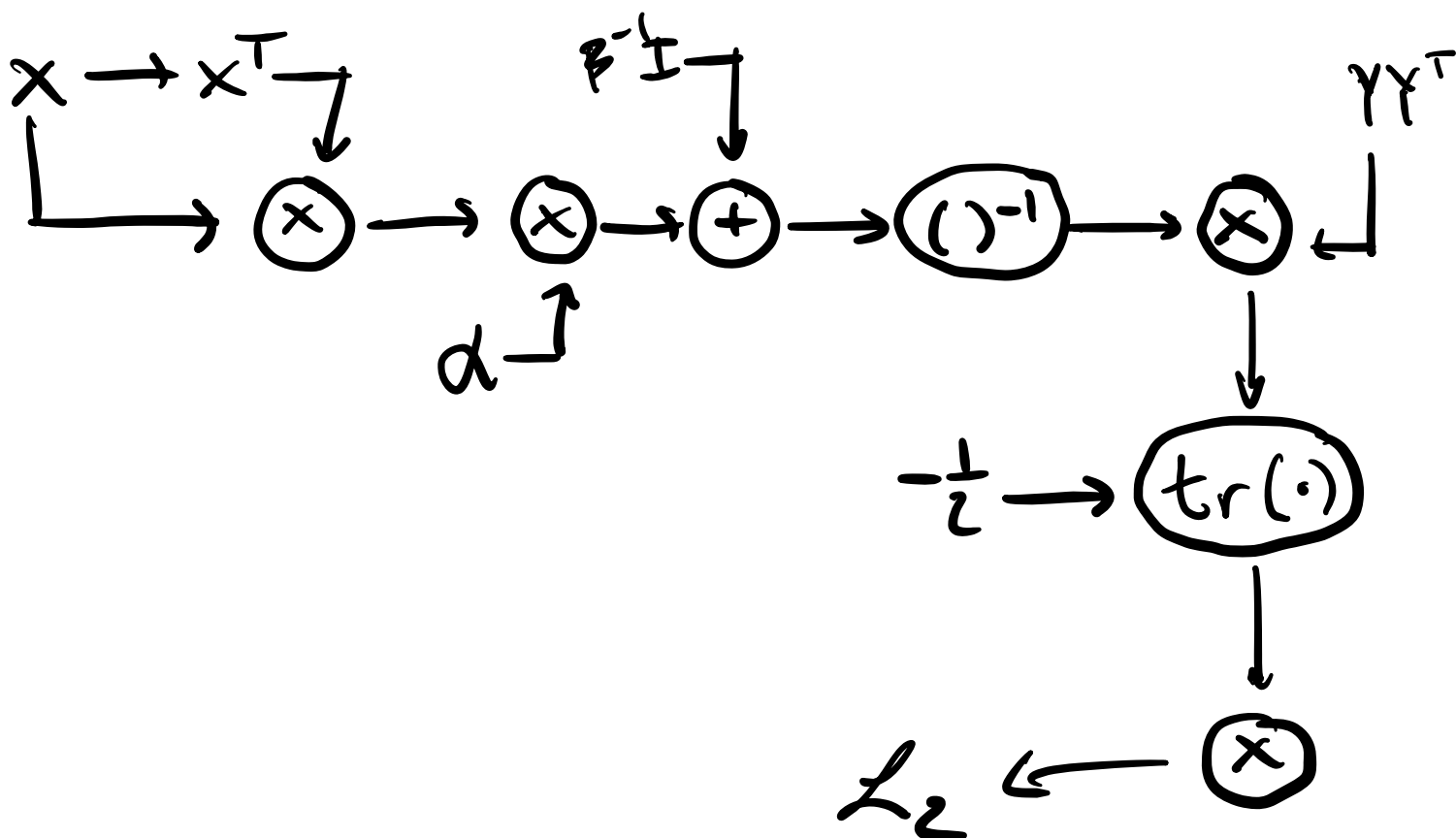
b. compute $\frac{d\mathcal{L}_1}{dX}$:

$$\left[\begin{aligned} \frac{d\mathcal{L}_1}{dK} &= \frac{d\mathcal{L}_1}{d(\log K)} \cdot \frac{d(\log(K))}{d(K)} \cdot f_1 = -\frac{D}{2} \log(K) \\ &= -\frac{D}{2} \cdot (K^{-1})^T \end{aligned} \right.$$

$$\frac{d\mathcal{L}_1}{dX} = \frac{d\mathcal{L}_1}{dK} \cdot \frac{dK}{dX} \quad \downarrow$$

$$\boxed{\frac{d\mathcal{L}_1}{dX} = -\alpha D (K^{-1})^T X}$$

c. $\mathcal{L}_2 = -\frac{1}{2} \text{tr}[(\alpha X X^T + \beta^{-1} I)^{-1} Y Y^T]$



d. compute $\frac{d\mathcal{L}_2}{dx}$:

$$\frac{d\mathcal{L}_2}{d(k^{-1}YY^T)} = -\frac{1}{2}I \quad \rightarrow YY^T$$

$$\frac{d(k^{-1}YY^T)}{dk} = -k^{-T} \frac{d(k^{-1}YY^T)}{dk^{-1}} \cdot k^{-T}$$

$$\rightarrow -k^{-T}YY^Tk^{-T}$$

$$\left[\frac{dk}{d(XX^T)} \cdot \frac{d(XX^T)}{dx} \right] = d-2x$$

$$\boxed{\frac{d\mathcal{L}_2}{dx} = d k^{-T}YY^Tk^{-T}x}$$

e. compute $\frac{d\mathcal{L}}{dx}$:

to compute $\frac{d\mathcal{L}}{dx}$, we must sum the expressions we have computed previously

$$\frac{d\mathcal{L}_1}{dx} \text{ and } \frac{d\mathcal{L}_2}{dx} \dots$$

$$\boxed{\frac{d\mathcal{L}}{dx} = \frac{d\mathcal{L}_1}{dx} + \frac{d\mathcal{L}_2}{dx} = \left[-dD(k^{-1})^T x + d k^{-T}YY^Tk^{-T}x \right]}$$

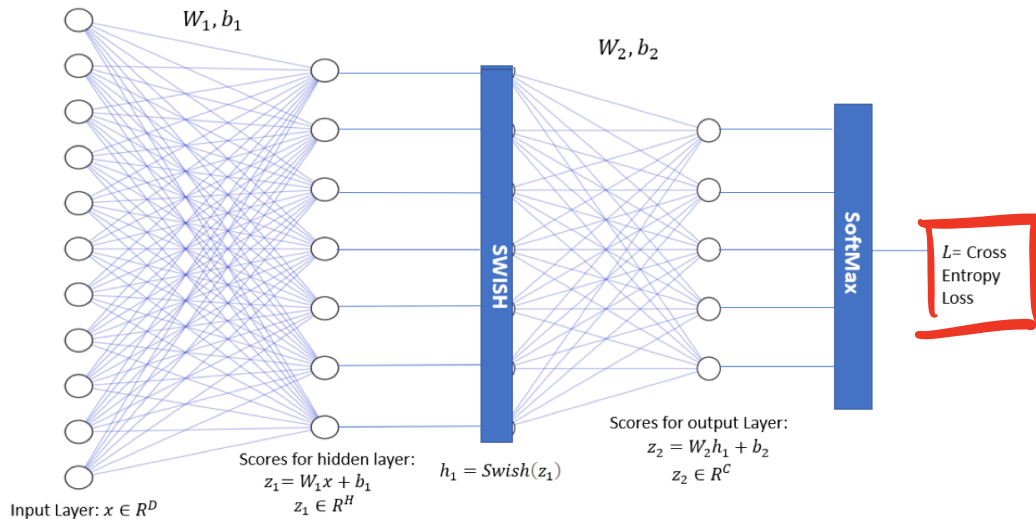
3. (15 points) **NNDL to the rescue!!**

It looks like a calm Monday morning and you are almost done with NNDL HW for the week (sigh)! But then suddenly (tring tring ...) your phone starts buzzing, you pick up the call, and the person from the other end sounds tense. The person exclaims ... there is a national emergency!!

7 different Pandora creature species (from Avatar) have been spotted in 1000's of numbers across various places in the country. They are having a hard time adjusting to the earth's climate and are causing chaos. As a result there has been a power outage in many cities. Luckily LA is an exception. UCLA's engineering division is helping out with this emergency, and you have been summoned to help.

You quickly take a bird to the secret facility and meet with director in charge of this operation. The director gives you a dataset consisting of images of these creatures along with their species type and instructs you to design a machine learning model to classify the images into species type. The only design constraint that the director has imposed is that the model should not have a very large number of parameters because some of UCLA's compute facilities are overloaded due to the power outages.

You just learned about fully connected neural networks (FC net) in class and decide to use it for accomplishing the task. To satisfy the design constraint, you decide to build a 2-layer FC net and train it using the provided dataset. The trained model will not only enable you to classify the images into species type but the hidden representations (outputs of intermediate layers) can be used to analyze the various properties of the species. A pictorial representation of the 2-layer FC net is shown below:



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,

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

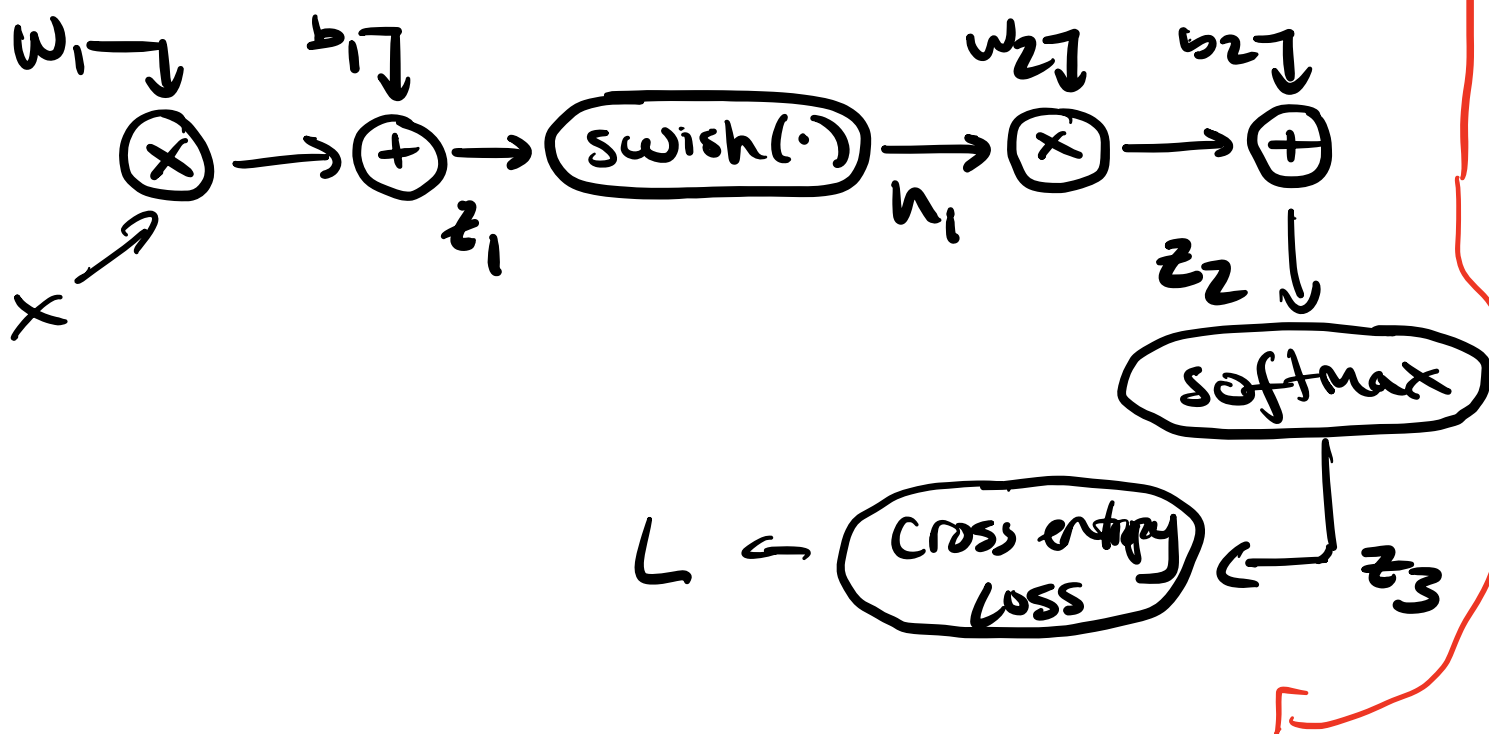
where $\sigma(k)$ is the sigmoid activation function you have seen in lecture.

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 $\nabla_{W_1} L$, $\nabla_{b_1} L$.

⑧. $\text{swish}(k) = \frac{k}{1 + e^{-k}} = k \sigma(k)$

a. Draw the computational graph for the 2-layer FC net...



b. Compute $\nabla_{w_2} L, \nabla_{b_2} L \dots$

$$\boxed{\frac{dL}{dw_2} = \frac{dL}{dz_2} \cdot h_1^T} \Rightarrow \frac{dz_2}{dw_2} \cdot \frac{dL}{dz_2}$$

$$\boxed{\frac{dL}{db_2} = \frac{dL}{dz_2} \cdot \frac{dz_2}{db_2} = \frac{dL}{dz_2}}$$

C. Compute $\nabla w, L, \nabla b, L \dots$

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

$$\text{swish}(z_1) = z_1 \sigma(z_1)$$

from lecture ... $\sigma(k)^2 = 1/(1 + e^{-k})$

$$\rightarrow [\sigma(z_1) + z_1 \sigma(z_1)(1 - \sigma(z_1))] = \underline{\underline{B}}$$

$$\frac{dL}{dw_1} = \frac{dz_1}{dw_1} \cdot \frac{dh_1}{dz_1} \cdot \frac{dz}{dh_1} \cdot \frac{dL}{dz_2}$$

$$= \frac{dz_1}{dw_1} \cdot B^T \cdot w_2^T \cdot \frac{dL}{dz_2}$$

$$\frac{dL}{db_1} = \frac{dz_1}{db_1} \cdot \frac{dh_1}{dz_1} \cdot \frac{dL}{dz_2}$$

$$= \frac{dz_1}{db_1} \cdot (B)^T w_2^T \cdot \frac{dL}{dz_2}$$