# HW5 Recitation

Wednesday, March 1, 2023          9:24 AM

## RECITATION 5
## NEURAL NETWORKS

10-301/10-601: INTRODUCTION TO MACHINE LEARNING

2022-10-14

## 1 Matrix Calculus

Consider $\mathbf{x} \in \mathbb{R}^p$, $\mathbf{y} \in \mathbb{R}^r$, $\mathbf{z} \in \mathbb{R}^n$ where $\mathbf{z} = g(\mathbf{y})$, and $\mathbf{y} = f(\mathbf{x})$. We want to derive $d\mathbf{z}/d\mathbf{x}$ (a vector form of the scalar chain rule).

1. If $\mathbf{x}$, $\mathbf{y}$, and $\mathbf{z}$ were all scalars, what would $dz/dx$ be?

$$\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx} = \frac{dy}{dx} \cdot \frac{dz}{dy}$$

$$x \xrightarrow{f} y \xrightarrow{g} z$$

**Shape matching:**

2. Fill in the following shapes:

$x \in \mathbb{R}^p$

$y \in \mathbb{R}^r$

$z \in \mathbb{R}^n$

$\left( \frac{d\mathbf{y}}{d\mathbf{x}} : p \times n \right.$

$\frac{d\mathbf{z}}{d\mathbf{y}} : r \times n$

$\frac{d\mathbf{z}}{d\mathbf{x}} : p \times n$

3. Therefore, the correct derivative is

$$\frac{dz}{dy} @ \frac{dy}{dx}$$

$$(r \times n) \quad (p \times r)$$

$$\frac{d\mathbf{z}}{d\mathbf{x}} = \frac{dy}{dx} \quad \frac{dz}{dy}$$

$$(p \times r) \cdot (r \times n)$$

**Generalizing a single element:** In order to ensure your derivatives are correct, we recommend you use matrix calculus rules whenever possible. When you're not sure how to apply a rule or if one applies, use the method of generalizing a single element.

4. Example: Suppose $\mathbf{A} \in \mathbb{R}^{n \times n}, \mathbf{x} \in \mathbb{R}^n$, and we want to compute $\frac{d\mathbf{Ax}}{d\mathbf{x}}$. Note that our numerator and denominator are both length-$n$ vectors, so our derivative has shape $\mathbb{R}^{n \times n}$.

5. Is $\left(\frac{d\mathbf{Ax}}{d\mathbf{x}}\right)_{ij}$ equal to $\frac{d\mathbf{Ax}_i}{d\mathbf{x}_j}$ or $\frac{d\mathbf{Ax}_j}{d\mathbf{x}_i}$? Why?

$$\frac{d\,Ax}{dx} \leftarrow n$$
$$\uparrow n$$

$$\left(\frac{d\,Ax}{dx}\right)_{ij} = \frac{d(Ax)_i}{dx_i} \quad or \quad \boxed{\frac{d(Ax)_j}{dx_i}}$$

6. Compute $\left(\frac{d\mathbf{Ax}}{d\mathbf{x}}\right)_{ij}$:

$$\left(\frac{d\,Ax}{dx}\right)_{ij} = \frac{d\,(Ax)_j}{dx_i}$$

$$= \frac{d\,A_j \cdot x}{dx_i} \quad \frac{d\,A_j^T x}{dx_i}$$

$$= \frac{d}{dx_i} \sum_{k=1}^{n} A_{jk} x_k$$

$$= A_{ji}$$

$$A \qquad \boxed{A^T} \qquad \frac{d\,Ax}{dx}$$

**Applying Matrix Calculus** For example, suppose we are finding the closed-form solution to linear regression: given $\mathbf{X} \in \mathbb{R}^{n \times d}, \mathbf{y} \in \mathbb{R}^n$, we wish to find $\boldsymbol{\theta} \in \mathbb{R}^d$ that minimizes $\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2.$ ←

7. What is the shape of $\mathbf{X}^\top \mathbf{y}$? Is this our solution?   $X^T y$ is not the right soln

   ← $d \times 1$

8. What is the closed-form solution? Use matrix calculus to derive the solution.

$(X^T X)^{-1} X^T y$

$\nearrow$

$$\frac{d}{d\theta} \|y - X\theta\|_2^2 = \frac{d}{d\theta} (y-X\theta)^T (y-X\theta)$$

$$= \frac{d}{d\theta} (\cancel{y^T y} - y^T(X\theta) - (X\theta)^T y + (X\theta)^T(X\theta))$$

$$= \frac{d}{d\theta} -2y^T X\theta + (X\theta)^T(X\theta)$$

$-2\frac{d}{d\theta} y^T(X\theta)$          $= -2X^T y + 2X^T X\theta$

$-2\frac{d}{d\theta} (y^T X)\theta$

$-2\frac{d}{d\theta} (X^T y)^T \theta$      $0 = -2X^T y + 2X^T X\theta$

$\underset{n \times 1}{\uparrow} \quad \underset{n \times 1}{\uparrow} \quad \underset{d \times 1}{\uparrow}$      $X^T y = X^T X\theta$

$(X^T X)^{-1} X^T y = \theta$

$-2 X^T y$

$\frac{d}{d\theta} (X\theta)^T(X\theta)$

$g'(x) \cdot f'(g(x))$

$\left(\frac{d}{d\theta} \overset{\downarrow}{X\theta}\right)(2X\theta)$

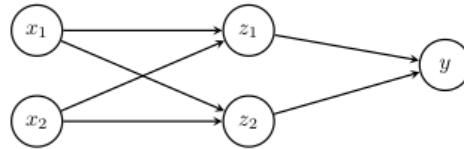$\uparrow$

$X^T 2X\theta$

# 2 Forward Propagation



Figure 1: Neural Network For Example Questions

**Forward Propagation** is the process of calculating the value of your loss function, given data, weights and activation functions. Given the input data $\mathbf{x}$, we can transform it by the given weights, $\boldsymbol{\alpha}$, then apply the corresponding activation function to it and finally pass the result to the next layer. Forward propagation does not involve taking derivatives and proceeds from the input layer to the output layer.

**Network Overview** Consider the neural network with one hidden layer shown in Figure 2. The input layer consists of 2 features $\mathbf{x} = [x_1, x_2]^T$, the hidden layer has 2 nodes with output $\mathbf{z} = [z_1, z_2]^T$, and the output layer is a scalar $\hat{y}$. We also add an intercept to the input, $x_0 = 1$ and the output of the hidden layer $z_0 = 1$, both of which are fixed to 1.

$\boldsymbol{\alpha}$ is the matrix of weights from the inputs to the hidden layer and $\boldsymbol{\beta}$ is the matrix of weights from the hidden layer to the output layer. $\alpha_{j,i}$ represents the weight going *to* the node $z_j$ in the hidden layer *from* the node $x_i$ in the input layer (e.g. $\alpha_{1,2}$ is the weight from $x_2$ to $z_1$), and $\boldsymbol{\beta}$ is defined similarly. We will use a **ReLU** activation function for the hidden layer and no activation for the output layer.

**Network Details** Equivalently, we define each of the following.

The input:

$$\mathbf{x} = [x_0, x_1, x_2]^T \tag{1}$$

Linear combination at the first (hidden) layer:

$$a_j = \sum_{i=0}^{2} \alpha_{j,i} \cdot x_i, \ \forall j \in \{1, \ldots, 2\} \tag{2}$$

Activation at the first (hidden) layer:

$$z_j = \mathrm{ReLU}(a_j) = \max(0, a_j), \forall j \in \{1, \ldots, 2\} \tag{3}$$

$$\mathbf{z} = [z_0, z_1, z_2]^T \tag{4}$$

Linear combination at the second (output) layer:

$$\hat{y} = \sum_{i=0}^{2} \beta_j \cdot z_j, \tag{5}$$

Here we fold in the intercept term $\alpha_{j,0}$ by thinking of $x_0 = 1$, and fold in $\beta_0$ by thinking of $z_0 = 1$.
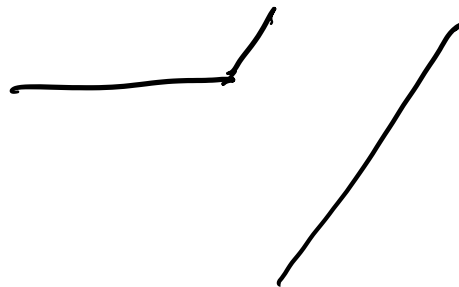
**Loss**   We will use Squared error loss, $\ell(\hat{y}, y)$:

$$\ell(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2 \tag{6}$$

1. Why and how do we include an intercept term in the input and in the hidden-layer?

   To encode shifts from the origin to help better fit data

2. Why do we need to use nonlinear activation functions in our neural net?

   to prevent it from being linear regression

$$a = ( W_1 x )$$

$$b = W_2 (a) = W_2 f(W_1 (x))$$

$$= (W_2 W_1) x$$

We initialize the network weights as:

$$\alpha = \begin{bmatrix} 2 & 1 & 0 \\ 0 & 2 & 0 \end{bmatrix}$$

→α=

$$\beta = \begin{bmatrix} 0 & 1 & 2 \end{bmatrix}$$

For the following questions, we use $\gamma = 3$.

1. **Scalar Form:**

   • Given $x_1 = 1$, $x_2 = 2$, What are the values of $a$?

$$a_2 = \sum_{i=0}^{2} \alpha_{2,i} x_i = \quad 1 \cdot 0 + 1 \cdot 2 + 2 \cdot 0 = 2$$

$$x = [1, 1, 2]$$

   • Given $z_1 = 0$, $z_2 = 1$ calculate $\hat{y}, l$

$$z = [1, 0, 1]$$

$$\hat{y} = \sum_{i=0}^{2} \beta_i \cdot z_i = \quad 1 \cdot 0 + 0 \cdot 1 + 1 \cdot 2 = 2$$

$$\tfrac{1}{2}(2-3)^2 = \tfrac{1}{2}$$

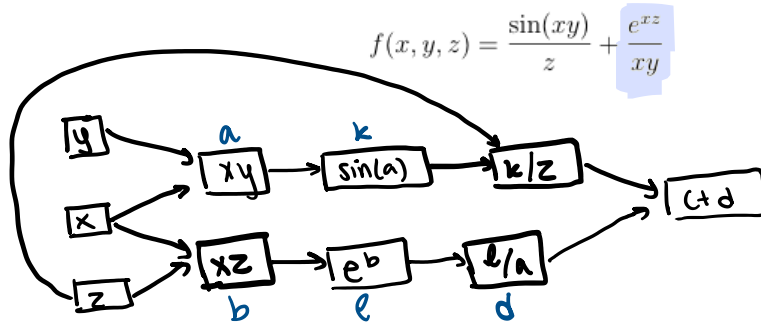2. **Vector Form:** Find the vector form of forward computation, given **x** is a column vector.

$$\alpha \in \mathbb{R}^{2 \times 3} \qquad x \in \mathbb{R}^3 \qquad a \in \mathbb{R}^2$$

$$\begin{bmatrix} a_1 \\ a_2 \end{bmatrix} \rightarrow a = \alpha x$$

$$\begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \rightarrow z = ReLU(a)$$

$$\rightarrow \begin{bmatrix} z_0 = 1 \\ z \end{bmatrix} \qquad \hat{y} = \beta z$$
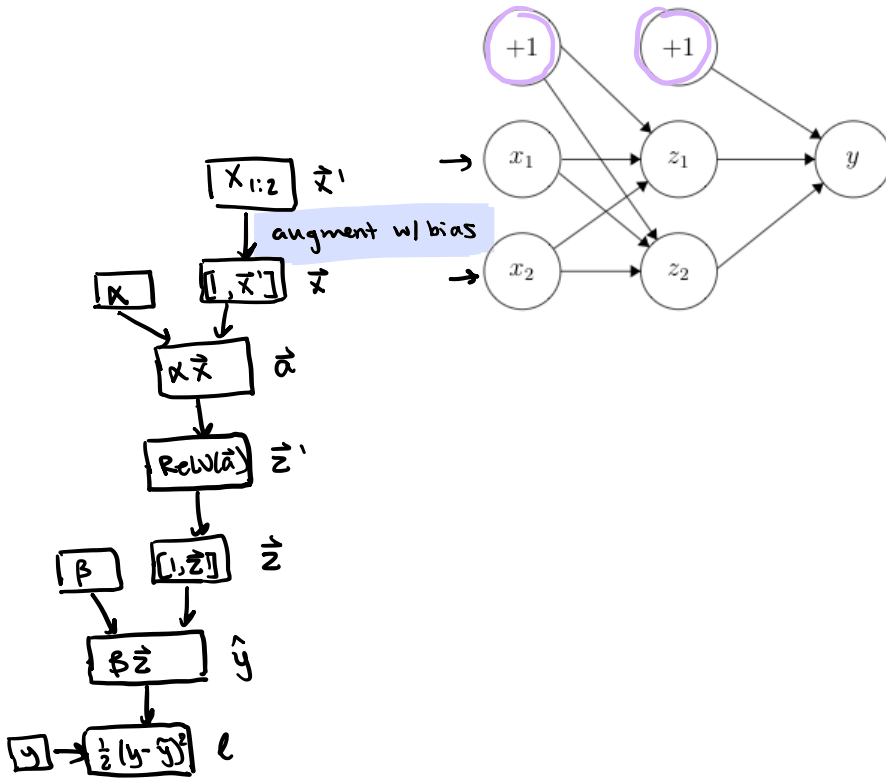
$$\left\lfloor \begin{matrix} 7 \\ z \end{matrix} \right\rfloor$$

# 3 Computation Diagrams

1. For the following function $f$, create the computation graph using the conventions defined in lecture.

$$f(x, y, z) = \frac{\sin(xy)}{z} + \frac{e^{xz}}{xy}$$



2. For the following neural network, draw the corresponding computation graph. Assume that all hidden units use the ReLU function as the activation function and that the loss is mean squared error. Provide the shape of all parameters defined in the computation graph. Assume the weights for the first layer and second layers are respectively the matrices $\alpha$ and $\beta$.
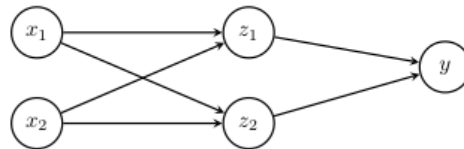
# 4   Backward Propagation



Figure 2: Neural Network For Example Questions

Given a Neural Network and a corresponding loss function $J(\theta)$, backpropagation gives us the gradient of the loss function with respect to the weights of the neural network. The method is called *backward* propagation because we calculate the gradients of the final layer of weights first, then proceed backward to the first layer. In a simple neural network with one hidden layer, the partial derivatives that we need for learning are $\frac{\partial \ell}{\partial \alpha_{ij}}$ and $\frac{\partial \ell}{\partial \beta_{kj}}$, and we need to apply chain rule recursively to obtain these. Note that in implementation, it is easier to use matrix/vector forms to conduct computations.

1. Many gradients are calculated in back propagation. Which of these gradients are used to update the weights? Do not include intermediate value(s) used to calculate these gradient(s).

the gradients used to update the weights are the gradients of the loss func w/ respect to $\alpha$ and $\beta$    ($\frac{\partial \ell}{\partial \alpha}$, $\frac{\partial \ell}{\partial \beta}$)

2. **Scalar Form:** Given

- $x_1 = 1$, $x_2 = 2$
- $a_1 = 3$, $a_2 = 2$
- $z_1 = 3$, $z_2 = 2$
- $\alpha = \begin{bmatrix} 2 & 1 & 0 \\ 0 & 2 & 0 \end{bmatrix}$
- $\beta = \begin{bmatrix} 0 & 1 & 2 \end{bmatrix}$
- $y = 3$,

what are the values of $\frac{\partial \ell}{\partial \beta_1}$, $\frac{\partial \ell}{\partial \alpha_{1,1}}$?

**Hint:** Derive expressions for $\frac{\partial \ell}{\partial \beta_i}$ and $\frac{\partial \ell}{\partial \alpha_{i,j}}$ first, then substitute in values.

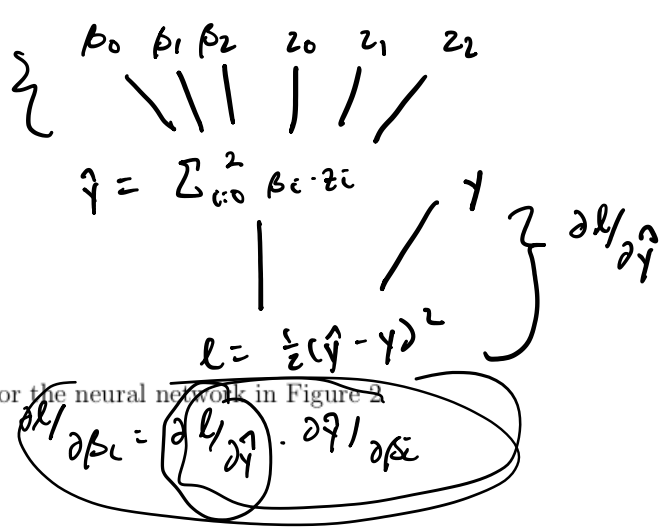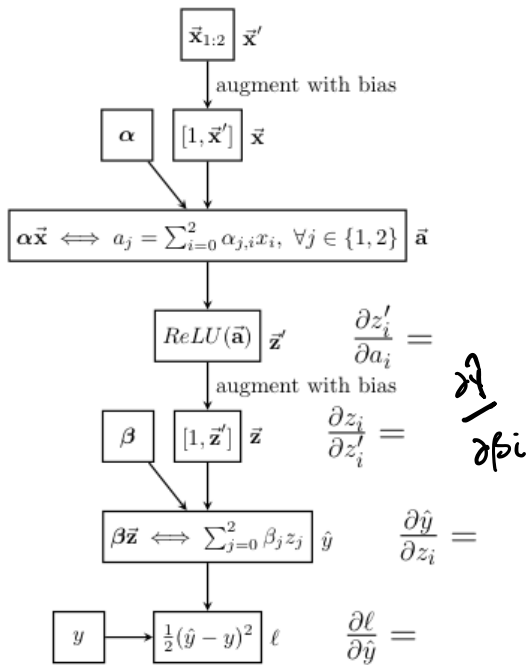For convenience, the computation graph for the neural network is displayed below:



Figure 3: Computation graph for the neural network in Figure 2

$$\frac{\partial z_i'}{\partial a_i} =$$

$$\frac{\partial z_i}{\partial z_i'} =$$

$$\frac{\partial \hat{y}}{\partial z_i} =$$

$$\frac{\partial \ell}{\partial \hat{y}} =$$

Handwritten notes:

$\frac{\partial \hat{y}}{\partial \beta_i} \Big\{ \quad \beta_0 \ \beta_1 \ \beta_2 \quad z_0 \ z_1 \ z_2$

$\hat{y} = \sum_{i=0}^{2} \beta_i \cdot z_i$

$\Big\} \frac{\partial \ell}{\partial \hat{y}}$

$\ell = \frac{1}{2}(\hat{y} - y)^2$

$\frac{\partial \ell}{\partial \beta_i} = \left(\frac{\partial \ell}{\partial \hat{y}}\right) \cdot \frac{\partial \hat{y}}{\partial \beta_i}$

**Hint:** $\frac{\partial ReLU(x)}{\partial x} = 1$ if $x > 0$, 0 if $x <= 0$

$$\frac{\partial \ell}{\partial \beta_i} = \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \beta_i}$$

$$\frac{\partial \ell}{\partial \beta_1} =$$

$\frac{\partial \ell}{\partial \hat{y}} = \frac{\partial}{\partial \hat{y}}\left(\frac{1}{2}(\hat{y}-y)^2\right)$

$$\partial \ell /_{\partial \beta i} = \begin{pmatrix} \hat{y} - y \end{pmatrix}$$

$$\partial \ell /_{\partial \beta i} = \partial \ell /_{\partial \hat{y}} \cdot \partial \hat{y} /_{\partial \beta i}$$

$$= \partial \ell /_{\partial \hat{y}} \cdot \partial /_{\partial \beta i} \left( \sum_{j=0}^{2} \beta_j \cdot z_j \right)$$

$$= \partial \ell /_{\partial \hat{y}} \cdot z_i$$

$$\frac{\partial \ell}{\partial \alpha_{i,j}} = \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_i} \frac{\partial z_i}{\partial a_i} \frac{\partial a_i}{\partial \alpha_{i,j}} \qquad = (\hat{y} - y) \cdot z_i \qquad \frac{\partial \ell}{\partial \alpha_{1,1}} =$$

$$\frac{\partial \ell}{\partial \hat{y}} = \partial /_{\partial \hat{y}} \left( \tfrac{1}{2}(\hat{y} - y)^2 \right) \qquad\qqu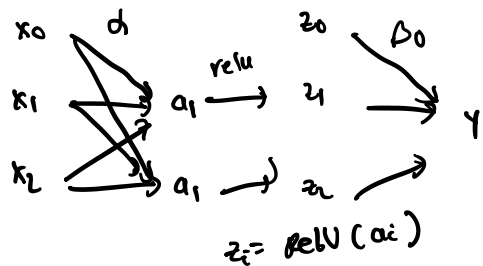ad \frac{\partial \ell}{\partial \alpha_{ij}} = \partial \ell /_{\partial a_i} \cdot \partial a_i /_{\partial \alpha_{ij}}$$

$$= \hat{y} - y \qquad\qquad\qquad = \frac{\partial \ell}{\partial a_i} \cdot \partial /_{\partial \alpha_{ij}} \left( \sum_{k=0}^{2} \alpha_{i,k} \cdot x_k \right)$$

$$\frac{\partial \ell}{\partial z_i} = \partial \ell /_{\partial \hat{y}} \cdot \partial /_{\partial z_i} \left( \sum_{j=0}^{2} \beta_j \cdot z_j \right) \qquad = \frac{\partial \ell}{\partial a_i} \cdot x_j$$

$$= \partial \ell /_{\partial \hat{y}} \cdot \beta_i \qquad\qquad \frac{\partial \ell}{\partial \alpha_{ij}} = (\hat{y} - y)(\beta_i)\left( \partial /_{\partial a_i} ReLU(a_i) \right)$$
$$\cdot (x_j)$$

$$\partial \ell /_{\partial a_i} = \partial \ell /_{\partial z_i} \cdot \partial z_i /_{\partial a_i}$$

$$= \partial \ell /_{\partial z_i} \cdot \partial /_{\partial a_i} \left( ReLU(a_i) \right)$$

3. **Vector Form:** What are the values of $\frac{\partial \ell}{\partial \beta}, \frac{\partial \ell}{\partial \alpha}$?

$$\partial \ell /_{\partial \alpha_{ij}} = \frac{\partial \ell}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_i} \cdot \frac{\partial z_i}{\partial a_i} \cdot \frac{\partial a_i}{\partial \alpha_{ij}}$$

$$\partial \ell /_{\partial \beta i} = (\hat{y} - y) \cdot z_i \quad_R$$

$$\partial \ell /_{\partial \vec{\beta}} = (\hat{y} - y) \begin{pmatrix} z_1 \\ z_2 \end{pmatrix}$$

$$= (\hat{y} - y) \vec{z}$$

$z_i = ReLU(a_i)$

$\hat{\beta} = \beta$ w/o the bias term

$\partial \ell /_{\partial \hat{y}} = (\hat{y} - y) \leftarrow$ scalar

$\partial z_i /_{\partial a_j} = $ if $i \ne j$, $0$

$\partial \hat{y} /_{\partial z_i} = \partial /_{\partial z_i} \left( \sum_{j=0}^{2} \beta_i \cdot z_i \right)$    else, $\partial /_{\partial a_i} ReLU(a_i)$

$$= \beta_i \qquad\qquad = \begin{pmatrix} \partial /_{\partial a_1} ReLU(a_1) & 0 \\ 0 & \partial /_{\partial a_2} ReLU(a_2) \end{pmatrix}$$

$$\partial \hat{y} /_{\partial \vec{z}} = \hat{\beta} \leftarrow 2\times 1$$

$$\partial \ell /_{\partial \vec{a}} = (\hat{y} - y) \cdot \hat{\beta} \cdot \begin{pmatrix} \partial /_{\partial a_1} ReLU(a_1) & 0 \\ 0 & \partial /_{\partial a_2} ReLU(a_2) \end{pmatrix}$$

$2\times1$

$2\times1$     $1\times2$

$$= (\hat{p} - y) \begin{pmatrix} \partial/_{\partial a_1} \text{ReLU}(a_1) & 0 \\ 0 & \partial/_{\partial_2}\text{ReLU}(a_2) \end{pmatrix} \hat{p}$$

$1\times1$     $2\times2$     $2\times1$

$\underbrace{\phantom{xxxxxxxxxxxxx}}$
$1\times1$

$$\frac{\partial \ell}{\partial d_{ij}} = \frac{\partial \ell}{\partial a_i} \cdot \frac{\partial a_i}{\partial d_{ij}}$$

$$= \frac{\partial \ell}{\partial a_i} \cdot \partial/_{\partial_1 a_{ij}} \left( \sum_{k=0}^{2} d_{ik} \cdot x_k \right)$$

$(k\times1$

$$= \frac{\partial \ell}{\partial a_i} \cdot x_j$$

$\frac{\partial \ell}{\partial d}$ =
$2\times 3$

$$\begin{pmatrix} \frac{\partial \ell}{\partial a_1} \cdot x_0 & \frac{\partial \ell}{\partial a_1} \cdot x_1 & \frac{\partial \ell}{\partial a_1} \cdot x_2 \\ \frac{\partial \ell}{\partial a_2} \cdot x_0 & \frac{\partial \ell}{\partial a_2} \cdot x_1 & \frac{\partial \ell}{\partial_2} \cdot x_2 \end{pmatrix}$$

$$= \begin{pmatrix} \frac{\partial \ell}{\partial a_1} \\ \frac{\partial \ell}{\partial a_2} \end{pmatrix} \begin{pmatrix} x_0 & x_1 & x_2 \end{pmatrix}$$

$$= \frac{\partial \ell}{\partial a} \vec{x}^T$$