

# Gradient Descent Method

## Backpropagation

Derivative of Loss function w.r.t. Vectors and Matrices

Sunmook Choi

Dept. of Mathematics  
Korea University

# Settings

- For practical reasons, we express each data  $\mathbf{x}$  into a **row vector**.
- There are different layouts to express derivatives with respect to vectors or matrices.
- In this presentation, we use the ‘denominator layout’.  
For more information, you can see more details in [this Wikipedia page](#).

# Linear Regression

# Gradients

For **real-valued** function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ , the derivative of  $f$  is the **gradient** of  $f$

$$\nabla f = \left[ \frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \quad \cdots \quad \frac{\partial f}{\partial x_n} \right]$$

where  $\mathbf{x} = \begin{bmatrix} x_1 & x_2 & \cdots & x_n \end{bmatrix} \in \mathbb{R}^n$  ( $\mathbf{x}$  is a **row vector**).

## Examples

- Let  $\mathbf{a} \in \mathbb{R}^n$  and  $f(\mathbf{x}) = \mathbf{x}\mathbf{a}^\top = \sum_{i=1}^n a_i x_i$ , then the gradient is  $\nabla_{\mathbf{x}} f = \mathbf{a}$ .
- Let  $A = [a_{ij}] \in \mathbb{R}^{n \times n}$  and  $f(\mathbf{x}) = \mathbf{x}A\mathbf{x}^\top$ , then the gradient is  $\nabla_{\mathbf{x}} f = \mathbf{x}(A + A^\top)$ .
  - \* Why?  $\rightarrow$  Homework
  - \* Hint)  $f(\mathbf{x}) = \mathbf{x}A\mathbf{x}^\top = \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j$ .

# Linear Regression

Consider a dataset  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  where  $\mathbf{x}_i = [x_{i1} \ x_{i2} \ \dots \ x_{id}] \in \mathbb{R}^d$  and  $y_i \in \mathbb{R}$ .

We assume that the relation is linear, that is, there exists  $\mathbf{w} = [w_1 \ \dots \ w_d] \in \mathbb{R}^d$  such that

$$y_i = \mathbf{x}_i \mathbf{w}^\top = \mathbf{w} \mathbf{x}_i^\top = \sum_{j=1}^d w_j x_{ij} \text{ where } x_{i1} = 1.$$

- $X = [\mathbf{x}_1^\top \ \mathbf{x}_2^\top \ \dots \ \mathbf{x}_n^\top]^\top \in \mathbb{R}^{n \times d}$  : *the design matrix* of  $\mathcal{D}$  (each row is an input vector).
- $\mathbf{y} = [y_1 \ y_2 \ \dots \ y_n] \in \mathbb{R}^n$  : the row vector containing the labels.

Then we want to find  $\mathbf{w} \in \mathbb{R}^d$  that satisfies

$$\mathbf{y} = \mathbf{w} X^\top,$$

so we will try to minimize the mean squared loss function: (*Ordinary Least Squares*)

$$L(\mathbf{w}) = \frac{1}{2n} \sum_{i=1}^n (y_i - \mathbf{w} \mathbf{x}_i^\top)^2 = \frac{1}{2n} (\mathbf{y} - \mathbf{w} X^\top)(\mathbf{y} - \mathbf{w} X^\top)^\top.$$

# Analytic Solution for Linear Regression

MSE loss of Linear Regression:

$$L(\mathbf{w}) = \frac{1}{2n} \sum_{i=1}^n (y_i - \mathbf{w} \mathbf{x}_i^\top)^2 = \frac{1}{2n} (\mathbf{y} - \mathbf{w} \mathbf{X}^\top) (\mathbf{y} - \mathbf{w} \mathbf{X}^\top)^\top.$$

Since  $L(\mathbf{w})$  is convex w.r.t.  $\mathbf{w}$ , the minimum is the point where  $\nabla_{\mathbf{w}} L(\mathbf{w}) = 0$  (gradient!).

$$\begin{aligned} \nabla_{\mathbf{w}} L(\mathbf{w}) &= \frac{1}{2n} \nabla_{\mathbf{w}} (\mathbf{y} \mathbf{y}^\top - \mathbf{y} \mathbf{X} \mathbf{w}^\top - \mathbf{w} \mathbf{X}^\top \mathbf{y}^\top + \mathbf{w} \mathbf{X}^\top \mathbf{X} \mathbf{w}^\top) \\ &= \frac{1}{2n} \nabla_{\mathbf{w}} (\mathbf{y} \mathbf{y}^\top - 2 \mathbf{w} \mathbf{X}^\top \mathbf{y}^\top + \mathbf{w} \mathbf{X}^\top \mathbf{X} \mathbf{w}^\top) \\ &= \frac{1}{2n} \left( -2 \mathbf{y} \mathbf{X} + \mathbf{w} (\mathbf{X}^\top \mathbf{X} + (\mathbf{X}^\top \mathbf{X})^\top) \right) \\ &= \frac{1}{n} \left( -\mathbf{y} \mathbf{X} + \mathbf{w} (\mathbf{X}^\top \mathbf{X}) \right) \end{aligned}$$

# Analytic Solution for Linear Regression

MSE loss of Linear Regression:

$$L(\mathbf{w}) = \frac{1}{2n} \sum_{i=1}^n (y_i - \mathbf{w} \mathbf{x}_i^\top)^2 = \frac{1}{2n} (\mathbf{y} - \mathbf{w} X^\top) (\mathbf{y} - \mathbf{w} X^\top)^\top,$$

and

$$\nabla_{\mathbf{w}} L(\mathbf{w}) = \frac{1}{n} \left( -\mathbf{y} X + \mathbf{w} (X^\top X) \right) = 0$$

Hence, if  $X^\top X$  is invertible, then we have an optimal solution:  $\hat{\mathbf{w}} = \mathbf{y} X (X^\top X)^{-1}$ .

- The rank of  $X^\top X$  is equal to the rank of  $X$ .
- Hence,  $X^\top X \in \mathbb{R}^{d \times d}$  is invertible if and only if  $X \in \mathbb{R}^{n \times d}$  has a full rank of  $d$ , that is, when each column vector (feature) of  $X$  is linearly independent of each other ( $\because d \ll n$ ).

# Numerical Solution for Linear Regression

MSE loss of Linear Regression:

$$L(\mathbf{w}) = \frac{1}{2n} \sum_{i=1}^n (y_i - \mathbf{w} \mathbf{x}_i^\top)^2 = \frac{1}{2n} (\mathbf{y} - \mathbf{w} \mathbf{X}^\top) (\mathbf{y} - \mathbf{w} \mathbf{X}^\top)^\top.$$

and

$$\nabla_{\mathbf{w}} L(\mathbf{w}) = \frac{1}{n} (-\mathbf{y} \mathbf{X} + \mathbf{w} \mathbf{X}^\top \mathbf{X}) = \frac{1}{n} (\mathbf{w} \mathbf{X}^\top - \mathbf{y}) \mathbf{X}$$

so that

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta \nabla_{\mathbf{w}} L(\mathbf{w}^t) = \mathbf{w}^t - \frac{\eta}{n} (\mathbf{w} \mathbf{X}^\top - \mathbf{y}) \mathbf{X}$$

or equivalently, for  $j = 1, \dots, d$ ,

$$w_j^{t+1} = w_j^t - \frac{\eta}{n} \sum_{i=1}^n (\mathbf{w} \mathbf{x}_i^\top - y_i) x_{ij}.$$



# MLP Classification

# Jacobian Matrix

For **vector-valued** function  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ , the derivative of  $\mathbf{y} = f(\mathbf{x})$  is the **Jacobian matrix**

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \left[ \frac{\partial y_i}{\partial x_j} \right]_{ij} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix} \in \mathbb{R}^{m \times n}$$

## Example

- Let  $A = [a_{ij}] \in \mathbb{R}^{n \times m}$ , and  $\mathbf{y} = \mathbf{x}A$ . Then the Jacobian matrix is

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \left[ \frac{\partial y_i}{\partial x_j} \right]_{ij} = A^\top.$$

## Derivative w.r.t. matrices

Let  $X = [x_{ij}] \in \mathbb{R}^{p \times q}$  and  $y$  be a **real-valued** function of  $X$  of independent variables.

- Here, independent variables of  $X$  means the matrix  $X$  has no special structure, e.g., not symmetric nor positive definite, etc.
- Then, the derivative of  $y$  with respect to  $X$  is given by

$$\frac{\partial y}{\partial X} = \left[ \frac{\partial y}{\partial x_{ij}} \right]_{ij} \in \mathbb{R}^{p \times q}.$$

- **Example:** Define  $f(X) = \mathbf{a}X\mathbf{b}^\top$  for  $\mathbf{a} \in \mathbb{R}^{1 \times p}$  and  $\mathbf{b} \in \mathbb{R}^{1 \times q}$ , we have

$$\frac{\partial f}{\partial X} = \left[ \frac{\partial f}{\partial x_{ij}} \right]_{ij} = \mathbf{a}^\top \mathbf{b} \in \mathbb{R}^{p \times q}.$$

## Derivative w.r.t. matrices

Let  $\mathbf{y} = \mathbf{x}W + \mathbf{b}$  where  $W \in \mathbb{R}^{n \times m}$ ,  $\mathbf{x} \in \mathbb{R}^n$ , and  $\mathbf{b} \in \mathbb{R}^m$ .

- Then, the derivative  $\frac{\partial \mathbf{y}}{\partial W}$  of  $\mathbf{y}$  with respect to  $W$  should be three dimensional!
- If  $L$  is a real-valued function of  $\mathbf{y}$ , then the derivative of  $L$  w.r.t.  $W$  is

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial W} = \frac{\partial L}{\partial \mathbf{y}} \frac{\partial(\mathbf{x}W + \mathbf{b})}{\partial W} = \mathbf{x}^\top \frac{\partial L}{\partial \mathbf{y}}$$

by the chain rule.

$$* \mathbf{x}^\top \in \mathbb{R}^{n \times 1}, \frac{\partial L}{\partial \mathbf{y}} \in \mathbb{R}^{1 \times m} \quad \text{so that} \quad \frac{\partial L}{\partial W} \in \mathbb{R}^{n \times m}.$$

## MLP Classification Forward Pass

Let  $\mathcal{D} = \{(\mathbf{x}^n, y^n)\}_{n=1}^N$  be a dataset of  $K$  classes. Consider 2-layer Multi-layer Perceptron.

- Input dim:  $d$ , Hidden dim:  $h$ , Output dim:  $K$  (# of classes)

$$\mathbf{x}^n \mapsto \mathbf{z}^{n,1} = \mathbf{x}^n W^1 + \mathbf{b}^1 \mapsto \mathbf{a}^n = \sigma(\mathbf{z}^{n,1}) \mapsto \mathbf{z}^{n,2} = \mathbf{a}^n W^2 + \mathbf{b}^2 \mapsto \hat{\mathbf{y}}^n = \text{Softmax}(\mathbf{z}^{n,2})$$

- Here, the parameters are  $W^1 \in \mathbb{R}^{d \times h}$ ,  $\mathbf{b}^1 \in \mathbb{R}^h$ ,  $W^2 \in \mathbb{R}^{h \times K}$ ,  $\mathbf{b}^2 \in \mathbb{R}^K$ .
- The activation function  $\sigma(\cdot)$  is the sigmoid function, that is,  $\sigma(z) = \frac{1}{1 + e^{-z}}$ .
- Notice that  $\frac{d}{dz}\sigma(z) = \sigma(z)(1 - \sigma(z))$ .

## MLP Classification Loss

$$\mathbf{x}^n \mapsto \mathbf{z}^{n,1} = \mathbf{x}^n W^1 + \mathbf{b}^1 \mapsto \mathbf{a}^n = \sigma(\mathbf{z}^{n,1}) \mapsto \mathbf{z}^{n,2} = \mathbf{a}^n W^2 + \mathbf{b}^2 \mapsto \hat{\mathbf{y}}^n = \text{Softmax}(\mathbf{z}^{n,2})$$

- Categorical Cross Entropy loss with respect to  $\theta = (W^1, \mathbf{b}^1, W^2, \mathbf{b}^2)$

$$L(\theta) = \frac{1}{N} \sum_{n=1}^N \sum_{j=1}^K [-y_j^n \log(\hat{y}_j^n)]$$

- We use one-hot encoding so that if the data  $\mathbf{x}^n$  belongs to  $j_*$ th class, then

$$\mathbf{y}^n = [y_1^n \ y_2^n \ \cdots \ y_K^n] \quad \text{where} \quad y_j^n = \begin{cases} 1 & \text{if } j = j_* \\ 0 & \text{if } j \neq j_* \end{cases}.$$

- Notice that the summation over  $j$  is actually not a summation due to one-hot encoding.

## Then gradients are...

$$\mathbf{x}^n \mapsto \mathbf{z}^{n,1} = \mathbf{x}^n W^1 + \mathbf{b}^1 \mapsto \mathbf{a}^n = \sigma(\mathbf{z}^{n,1}) \mapsto \mathbf{z}^{n,2} = \mathbf{a}^n W^2 + \mathbf{b}^2 \mapsto \hat{\mathbf{y}}^n = \text{Softmax}(\mathbf{z}^{n,2}) \mapsto L$$

We want to find the following gradients:  $\frac{\partial L}{\partial W^1}$ ,  $\frac{\partial L}{\partial \mathbf{b}^1}$ ,  $\frac{\partial L}{\partial W^2}$ ,  $\frac{\partial L}{\partial \mathbf{b}^2}$ .

Using the chain rule, we first find that

$$\frac{\partial L}{\partial W^2} = \sum_n \frac{\partial L}{\partial \hat{\mathbf{y}}^n} \frac{\partial \hat{\mathbf{y}}^n}{\partial \mathbf{z}^{n,2}} \frac{\partial \mathbf{z}^{n,2}}{\partial W^2} \quad \text{and} \quad \frac{\partial L}{\partial \mathbf{b}^2} = \sum_n \frac{\partial L}{\partial \hat{\mathbf{y}}^n} \frac{\partial \hat{\mathbf{y}}^n}{\partial \mathbf{z}^{n,2}} \frac{\partial \mathbf{z}^{n,2}}{\partial \mathbf{b}^2}.$$

$$\mathbf{x}^n \mapsto \mathbf{z}^{n,1} = \mathbf{x}^n W^1 + \mathbf{b}^1 \mapsto \mathbf{a}^n = \sigma(\mathbf{z}^{n,1}) \mapsto \mathbf{z}^{n,2} = \mathbf{a}^n W^2 + \mathbf{b}^2 \mapsto \hat{\mathbf{y}}^n = \text{Softmax}(\mathbf{z}^{n,2}) \mapsto L$$

For  $L(\theta) = \frac{1}{N} \sum_{n=1}^N \sum_{j=1}^K [-y_j^n \log(\hat{y}_j^n)]$ , we have  $\frac{\partial L}{\partial \hat{\mathbf{y}}^n} = \frac{1}{N} \begin{bmatrix} -y_1^n/\hat{y}_1^n & -y_2^n/\hat{y}_2^n & \cdots & -y_K^n/\hat{y}_K^n \end{bmatrix}$  and

$$\frac{\partial \hat{\mathbf{y}}^n}{\partial \mathbf{z}^{n,2}} = \left[ \frac{\partial \hat{y}_i^n}{\partial z_j^{n,2}} \right]_{ij} \text{ where } \frac{\partial \hat{y}_i^n}{\partial z_j^{n,2}} = \begin{cases} \hat{y}_i^n(1 - \hat{y}_i^n) & \text{if } i = j \\ -\hat{y}_i^n \hat{y}_j^n & \text{if } i \neq j \end{cases} \text{ because } \hat{y}_i^n = \frac{e^{z_i^{n,2}}}{\sum_{k=1}^K e^{z_k^{n,2}}}.$$

Using the fact that  $\mathbf{y}^n$  is one-hot encoded, we have  $\frac{\partial L}{\partial \mathbf{z}^{n,2}} = \frac{\partial L}{\partial \hat{\mathbf{y}}^n} \frac{\partial \hat{\mathbf{y}}^n}{\partial \mathbf{z}^{n,2}} = \frac{1}{N} (\hat{\mathbf{y}}^n - \mathbf{y}^n) \in \mathbb{R}^{1 \times K}$ .

$$\therefore \frac{\partial L}{\partial W^2} = \sum_n \frac{\partial L}{\partial \mathbf{z}^{n,2}} \frac{\partial \mathbf{z}^{n,2}}{\partial W^2} = \frac{1}{N} \sum_n (\hat{\mathbf{y}}^n - \mathbf{y}^n) \frac{\partial (\mathbf{a}^n W^2 + \mathbf{b}^2)}{\partial W^2} = \frac{1}{N} \sum_n (\mathbf{a}^n)^\top (\hat{\mathbf{y}}^n - \mathbf{y}^n)$$

$$\frac{\partial L}{\partial \mathbf{b}^2} = \sum_n \frac{\partial L}{\partial \mathbf{z}^{n,2}} \frac{\partial \mathbf{z}^{n,2}}{\partial \mathbf{b}^2} = \frac{1}{N} \sum_n (\hat{\mathbf{y}}^n - \mathbf{y}^n) I_K = \frac{1}{N} \sum_n (\hat{\mathbf{y}}^n - \mathbf{y}^n)$$



$$\mathbf{x}^n \mapsto \mathbf{z}^{n,1} = W^1 \mathbf{x}^n + \mathbf{b}^1 \mapsto \mathbf{a}^n = \sigma(\mathbf{z}^{n,1}) \mapsto \mathbf{z}^{n,2} = W^2 \mathbf{a}^n + \mathbf{b}^2 \mapsto \hat{\mathbf{y}}^n = \text{Softmax}(\mathbf{z}^{n,2}) \mapsto L$$

We now compute  $\frac{\partial L}{\partial W^1} = \sum_n \frac{\partial L}{\partial \mathbf{z}^{n,1}} \frac{\partial \mathbf{z}^{n,1}}{\partial W^1}$ ,  $\frac{\partial L}{\partial \mathbf{b}^1} = \sum_n \frac{\partial L}{\partial \mathbf{z}^{n,1}} \frac{\partial \mathbf{z}^{n,1}}{\partial \mathbf{b}^1}$ .

We have  $\frac{\partial L}{\partial \mathbf{a}^n} = \frac{\partial L}{\partial \mathbf{z}^{n,2}} \frac{\partial \mathbf{z}^{n,2}}{\partial \mathbf{a}^n} = \frac{\partial L}{\partial \mathbf{z}^{n,2}} (W^2)^\top$ ,  $\frac{\partial L}{\partial \mathbf{z}^{n,1}} = \frac{\partial L}{\partial \mathbf{a}^n} \frac{\partial \mathbf{a}^n}{\partial \mathbf{z}^{n,1}}$ , and  $\frac{\partial \mathbf{a}^n}{\partial \mathbf{z}^{n,1}} = \text{diag}[a_i^n(1 - a_i^n)]$ .

Then  $\frac{\partial L}{\partial \mathbf{z}^{n,1}} = \frac{\partial L}{\partial \mathbf{z}^{n,2}} (W^2)^\top \text{diag}[a_i^n(1 - a_i^n)] \in \mathbb{R}^{1 \times h}$ .

$$\begin{aligned} \therefore \frac{\partial L}{\partial W^1} &= \sum_n \frac{\partial L}{\partial \mathbf{z}^{n,1}} \frac{\partial \mathbf{z}^{n,1}}{\partial W^1} = \sum_n (\mathbf{x}^n)^\top \frac{\partial L}{\partial \mathbf{z}^{n,1}} = \sum_n (\mathbf{x}^n)^\top \frac{\partial L}{\partial \mathbf{z}^{n,2}} (W^2)^\top \text{diag}[a_i^n(1 - a_i^n)] \in \mathbb{R}^{d \times h} \\ \frac{\partial L}{\partial \mathbf{b}^1} &= \sum_n \frac{\partial L}{\partial \mathbf{z}^{n,1}} \frac{\partial \mathbf{z}^{n,1}}{\partial \mathbf{b}^1} = \sum_n \frac{\partial L}{\partial \mathbf{z}^{n,1}} = \sum_n \frac{\partial L}{\partial \mathbf{z}^{n,2}} (W^2)^\top \text{diag}[a_i^n(1 - a_i^n)]. \end{aligned}$$

## Gradient Descent Methods

Finally, we apply the gradient descent method to update the parameters.

- Recall that  $W^1 \in \mathbb{R}^{d \times h}$ ,  $\mathbf{b}^1 \in \mathbb{R}^{1 \times h}$ ,  $W^2 \in \mathbb{R}^{h \times K}$ ,  $\mathbf{b}^2 \in \mathbb{R}^{1 \times K}$  and

$$\frac{\partial L}{\partial W^1} \in \mathbb{R}^{d \times h}, \quad \frac{\partial L}{\partial \mathbf{b}^1} \in \mathbb{R}^{1 \times h}, \quad \frac{\partial L}{\partial W^2} \in \mathbb{R}^{h \times K}, \quad \frac{\partial L}{\partial \mathbf{b}^2} \in \mathbb{R}^{1 \times K}.$$

- Hence, the update equations should be

$$W^1 \leftarrow W^1 - \eta \frac{\partial L}{\partial W^1}, \quad \mathbf{b}^1 \leftarrow \mathbf{b}^1 - \eta \frac{\partial L}{\partial \mathbf{b}^1}$$

$$W^2 \leftarrow W^2 - \eta \frac{\partial L}{\partial W^2}, \quad \mathbf{b}^2 \leftarrow \mathbf{b}^2 - \eta \frac{\partial L}{\partial \mathbf{b}^2}$$