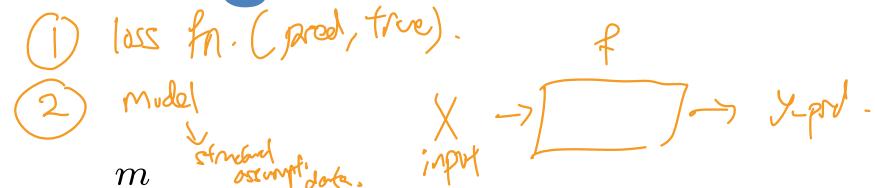


Building Blocks of ML algorithms

- From previous slide



How to find?

* Gradient Descent gives improved estimates of w

w a.k.a. $C(w)$ should decrease

- To make it even more general, we can write

→ good "w"
 ~ small $C(w)$

→ the best "w"
 ~ w that minimizes $C(w)$

$$\underset{w}{\operatorname{argmin}} C(w) = \underset{w}{\operatorname{argmin}} \sum_{i=1}^m L(f(\mathbf{x}_i, \mathbf{w}), y_i) + \lambda R(\mathbf{w})$$

no. of samples: m

Error

pred

true output

inputs

parameters

e.g. in ridge (in) reg $w^T w$.

- Learning model f reflects our belief about the relationship between the features \mathbf{x}_i & target y_i
- Loss function L is the penalty for predicting $f(\mathbf{x}_i, \mathbf{w})$ when the true value is y_i
- Regularization R encourages less complex models
- Cost function C is the final optimization criterion we want to minimize
- Optimization routine to find solution to cost function

Derivative and Gradient

The gradient of a function is a vector of **partial derivatives**

Differentiation of a scalar function w.r.t. a vector

If $f(\mathbf{x})$ is a **scalar function** of d variables, \mathbf{x} is a $d \times 1$ vector.

Then differentiation of $f(\mathbf{x})$ w.r.t. \mathbf{x} results in a $d \times 1$ vector

$$f: \mathbb{R}^d \rightarrow \mathbb{R} . \quad \frac{df(\mathbf{x})}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix}$$

$$\frac{df}{d\mathbf{x}} : \mathbb{R}^d \rightarrow \mathbb{R}^d \\ \mathbf{x}_0 \mapsto \begin{pmatrix} \frac{\partial f}{\partial x_1}(\mathbf{x}_0) \\ \vdots \\ \frac{\partial f}{\partial x_d}(\mathbf{x}_0) \end{pmatrix}$$

This is referred to as the **gradient** of $f(\mathbf{x})$ and often written as $\nabla_{\mathbf{x}} f$.

$$\text{E.g. } f(\mathbf{x}) = ax_1 + bx_2 \quad \nabla_{\mathbf{x}} f = \begin{bmatrix} a \\ b \end{bmatrix}$$

Ref: Duda, Hart, and Stork, "Pattern Classification", 2001 (Appendix)

Derivative and Gradient

Partial Derivatives

Differentiation of a vector function w.r.t. a vector

If $\mathbf{f}(\mathbf{x})$ is a **vector function** of size $h \times 1$ and \mathbf{x} is a $d \times 1$ vector.

Then differentiation of $\mathbf{f}(\mathbf{x})$ results in a $h \times d$ matrix

$$f: \mathbb{R}^d \rightarrow \mathbb{R}^h$$

$$\frac{df}{dx}: \mathbb{R}^d \rightarrow \mathbb{R}^{h \times d}$$

$$\frac{df}{dx} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_d} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_h}{\partial x_1} & \dots & \frac{\partial f_h}{\partial x_d} \end{bmatrix}$$

$h=1$

$$f: \mathbb{R}^d \rightarrow \mathbb{R}^1$$

$$\frac{df}{dx}: \mathbb{R}^d \rightarrow \mathbb{R}^{1 \times d}$$

$$x_0 \mapsto \left[\frac{df}{dx}(x_0) \right]$$

The matrix is referred to as the **Jacobian of $f(\mathbf{x})$**
 (total) derivative.

Ref: Duda, Hart, and Stork, "Pattern Classification", 2001 (Appendix)

Linear Alg. recap

- ① $L = h \begin{pmatrix} L_1 & L_2 & \dots & L_d \end{pmatrix}$: $h \times d$ matrix \sim linear map
 $L(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^h$
 $v \mapsto L(v)$
- ② The function $x^T(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$
- Fix $x = \begin{pmatrix} x^1 \\ \vdots \\ x^d \end{pmatrix}$ is just a linear map
 $w \mapsto x^T w$
- $\begin{pmatrix} w^1 \\ \vdots \\ w^d \end{pmatrix} \mapsto (x^1 \dots x^d) \begin{pmatrix} w^1 \\ \vdots \\ w^d \end{pmatrix}$
- $= x^1 w^1 + \dots + x^d w^d$
- $= \begin{pmatrix} x^1 \\ \vdots \\ x^d \end{pmatrix} \bullet \begin{pmatrix} w^1 \\ \vdots \\ w^d \end{pmatrix}$
- ↑
not product
- sending
 $\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \mapsto L_1$
 $\begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \mapsto L_2$
 \vdots
 $\begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \mapsto L_d$

Terminology

Fix a function $f: \mathbb{R}^d \rightarrow \mathbb{R}^h$

$\xleftarrow[\text{dummy var.}]{\text{position}} x_0 \mapsto f(x_0) = \begin{pmatrix} f_1(x_0) \\ f_2(x_0) \\ \vdots \\ f_h(x_0) \end{pmatrix}$

The Jacobian or (total) derivative is a function

$$\frac{df}{dx}: \mathbb{R}^d \rightarrow \mathbb{R}^{h \times d}$$

$$x_0 \mapsto h \left(\frac{df}{dx}(x_0) \right)$$

dummy variable for position

this is called the (position) derivative at x_0 . It is a linear map

$$= \begin{bmatrix} \frac{\partial f_1}{\partial x_1}(x_0) & \cdots & \frac{\partial f_1}{\partial x_d}(x_0) \\ \vdots & \ddots & \vdots \\ \frac{\partial f_h}{\partial x_1}(x_0) & \cdots & \frac{\partial f_h}{\partial x_d}(x_0) \end{bmatrix}$$

key idea:

$$f(x_0 + \epsilon) - f(x_0) \approx h \left(\frac{df}{dx}(x_0) \right) (\epsilon)$$

output in \mathbb{R}^h

must be small direction change in \mathbb{R}^d

h x d matrix, aka linear map

direction vector; \mathbb{R}^d

Special case : $h=1$

$$f : \mathbb{R}^d \rightarrow \mathbb{R}$$

$$x_0 \mapsto f(x_0)$$

The gradient is a function

$$\nabla_x f : \mathbb{R}^d \rightarrow \mathbb{R}^d$$

$$x_0 \mapsto \nabla_x f(x_0) = \begin{pmatrix} \frac{\partial f}{\partial x_1}(x_0) \\ \vdots \\ \frac{\partial f}{\partial x_d}(x_0) \end{pmatrix}$$

this is called the

gradient at x_0

The Jacobian or (total) derivative is a function

$$\frac{df}{dx} : \mathbb{R}^d \rightarrow \mathbb{R}^{1 \times d}$$

$$x_0 \mapsto \begin{pmatrix} \frac{df}{dx}(x_0) \end{pmatrix} = \left(\frac{\partial f}{\partial x_1}(x_0), \frac{\partial f}{\partial x_2}(x_0), \dots, \frac{\partial f}{\partial x_d}(x_0) \right)$$

$$= \nabla_x f(x_0)^T$$

this is called the derivative at x_0 .

If it is a

linear map

key idea :

$$f(x_0 + \varepsilon) - f(x_0) \approx \left[\frac{df}{dx}(x_0) \right] (\varepsilon) = \nabla_x f(x_0)^T (\varepsilon) = \nabla_x f(x_0) \cdot \bar{\varepsilon}$$

where $\varepsilon \in \mathbb{R}^d$

$f(x_0 + \varepsilon)$ position in \mathbb{R}^d

$f(x_0)$ position in \mathbb{R}^d

ε small step size in \mathbb{R}^d

$\frac{df}{dx}(x_0)$ linear map

ε small step in \mathbb{R}^d

$\nabla_x f(x_0)^T (\varepsilon)$ dot product

Relationship between gradient & derivative ($h=1$)

$$\frac{df}{dx}(x_0) (-) = \left(\frac{\partial f}{\partial x_1}(x_0), \frac{\partial f}{\partial x_2}(x_0), \dots, \frac{\partial f}{\partial x_d}(x_0) \right) (-)$$



linear map

$\mathbb{R}^d \rightarrow \mathbb{R}$.

$$= \begin{pmatrix} \frac{\partial f}{\partial x_1}(x_0) \\ \vdots \\ \frac{\partial f}{\partial x_d}(x_0) \end{pmatrix}^T (-)$$



linear map

$$= \nabla_x f(x_0)^T (-) : \mathbb{R}^d \rightarrow \mathbb{R}^d$$

i.e.: the gradient is the transpose of the $(1 \times d)$ matrix

$$\stackrel{1}{\left(\frac{df}{dx} \right)}$$

Example 1: Fix $x = \begin{pmatrix} x^1 \\ x^2 \\ \vdots \\ x^d \end{pmatrix}$. Consider the function

$$x^T(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$$

$$w = \begin{pmatrix} w_1 \\ \vdots \\ w_d \end{pmatrix} \mapsto x^T w = (x^1 \ x^2 \ \dots \ x^d) \begin{pmatrix} w_1 \\ \vdots \\ w_d \end{pmatrix}$$

$$= x^1 w_1 + x^2 w_2 + \dots + x^d w_d$$

The derivative is

$$\frac{dx^T(\cdot)}{dw} : \mathbb{R}^d \rightarrow \mathbb{R}^{1 \times d}$$

$$w_0 \mapsto \left(\frac{\partial x^T(\cdot)}{\partial w_1} \Big|_{w=w_0} \quad \frac{\partial x^T(\cdot)}{\partial w_2} \Big|_{w=w_0} \quad \dots \quad \frac{\partial x^T(\cdot)}{\partial w_d} \Big|_{w=w_0} \right)$$

$$= (x^1 \ x^2 \ \dots \ x^d) = x^T$$

$$= x \bullet w$$

dot product.

The gradient is $\nabla_w x^T : \mathbb{R}^d \rightarrow \mathbb{R}^d$

$$w_0 \mapsto x$$

Example 2 : $\sigma: \mathbb{R} \rightarrow \mathbb{R}$

$$a_0 \mapsto \frac{1}{1 + \exp(-\beta a_0)}$$

where $\beta \in \mathbb{R}^+$
is some positive real number

$$\begin{aligned}\frac{\partial \sigma(a)}{\partial a} &= \frac{\partial}{\partial a} \left(\frac{1}{1 + \exp(-\beta a)} \right) \\ &= -\frac{1}{(1 + e^{-\beta a})^2} \frac{\partial(1 + e^{-\beta a})}{\partial a} \\ &= \frac{\beta}{(1 + e^{-\beta a})^2} e^{-\beta a} \\ &= \frac{\beta}{(1 + e^{-\beta a})^2} (1 + e^{-\beta a} - 1) \\ &= \beta \left(\frac{1}{1 + e^{-\beta a}} - \frac{1}{(1 + e^{-\beta a})^2} \right) \\ &= \beta (\sigma(a) - \sigma^2(a)) \\ &= \beta \sigma(a)(1 - \sigma(a))\end{aligned}$$

$$\begin{aligned}\nabla_a \sigma: \mathbb{R} &\rightarrow \mathbb{R}^1 \\ a_0 &\mapsto \left[\frac{\partial}{\partial a} \frac{1}{1 + \exp(-\beta a)} \right] \Big|_{a=a_0} \\ &= \beta \sigma(a_0)(1 - \sigma(a_0)) \\ &= \beta \frac{1}{1 + \exp(-\beta a_0)} \cdot \left(1 - \frac{1}{1 + \exp(-\beta a_0)} \right)\end{aligned}$$

$$\frac{d\sigma}{da}: \mathbb{R} \rightarrow \mathbb{R}^{1 \times 1}$$

$$a_0 \mapsto [\nabla_a \sigma(a_0)]^\top \cdot (-)$$

1x1 matrix.

Example 3 : $f : \mathbb{R}^2 \rightarrow \mathbb{R}$

$$\begin{pmatrix} w_1 \\ w_2 \end{pmatrix} \mapsto w_1^2 + w_2^2$$

$\nabla_w f : \mathbb{R}^2 \rightarrow \mathbb{R}^2$

$$\begin{pmatrix} w_1 \\ w_2 \end{pmatrix} \mapsto \begin{pmatrix} 2w_1 \\ 2w_2 \end{pmatrix}$$

$\frac{df}{dw} : \mathbb{R}^2 \rightarrow \mathbb{R}^{1 \times 2}$

$$\begin{pmatrix} w_1^\circ \\ w_2^\circ \end{pmatrix} \mapsto \left(\frac{\partial f}{\partial w_1} \left(\begin{pmatrix} w_1^\circ \\ w_2^\circ \end{pmatrix} \right), \frac{\partial f}{\partial w_2} \left(\begin{pmatrix} w_1^\circ \\ w_2^\circ \end{pmatrix} \right) \right)$$

$$= \left(2w_1 \begin{pmatrix} w_1^\circ \\ w_2^\circ \end{pmatrix}, 2w_2 \begin{pmatrix} w_1^\circ \\ w_2^\circ \end{pmatrix} \right)$$

$$= \begin{pmatrix} 2w_1^\circ & 2w_2^\circ \end{pmatrix}$$

$$\begin{pmatrix} 2 \\ 2 \end{pmatrix} \mapsto \begin{pmatrix} 4 & 6 \end{pmatrix}$$

Gradient Descent Algorithm

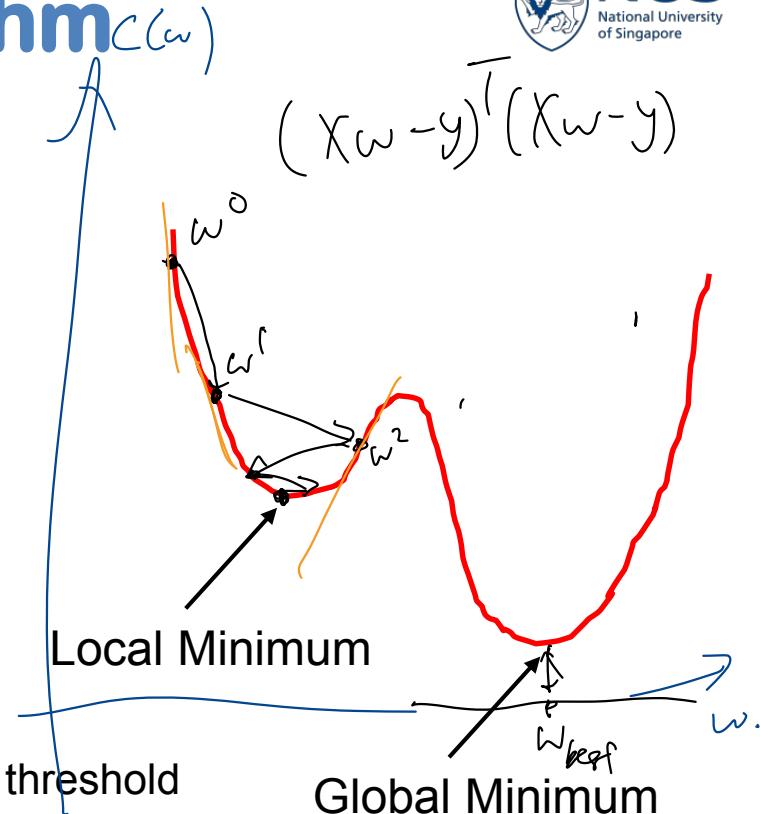
- Gradient Descent:

```

Initialize  $w_0$  and learning rate  $\eta$ ;
while true do
    Compute  $w_{k+1} \leftarrow w_k - \eta \nabla_w C(w_k)$ 
    if converge then
        return  $w_{k+1}$ 
    end
end
  
```

\downarrow
 opposite direction
 gradient
 adjust for
 suitable step size

- Possible convergence criteria
 - Set maximum iteration k
 - Check percentage or absolute change in C below a threshold
 - Check percentage or absolute change in w below a threshold
- Gradient descent can only find local minimum
 - Because gradient = 0 at local minimum, so w won't change after that
- Many variations of gradient descent, e.g., change how gradient is computed or learning rate η decreases with increasing k



Gradient Descent Algorithm

- Suppose we want to minimize $C(\mathbf{w})$ with respect to $\mathbf{w} = [w_1, \dots, w_d]^T$

- Gradient $\nabla_{\mathbf{w}} C(\mathbf{w}) = \begin{pmatrix} \frac{\partial C}{\partial w_1} \\ \frac{\partial C}{\partial w_2} \\ \vdots \\ \frac{\partial C}{\partial w_d} \end{pmatrix}$

- $\nabla_{\mathbf{w}} C(\mathbf{w})$ is vector & function of \mathbf{w}
 - $\nabla_{\mathbf{w}} C(\mathbf{w})$ is direction at \mathbf{w} where C is increasing most rapidly, so $-\nabla_{\mathbf{w}} C(\mathbf{w})$ is direction at \mathbf{w} where C is decreasing most rapidly

- Gradient Descent:

Initialize \mathbf{w}_0 and learning rate η ;

while true **do**

Compute $\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k - \eta \nabla_{\mathbf{w}} C(\mathbf{w}_k)$

if converge **then**

return \mathbf{w}_{k+1}

end

end

According to multi-variable calculus, if eta is not too big, then $C(\mathbf{w}_{k+1}) < C(\mathbf{w}_k) \Rightarrow$ we get better \mathbf{w} after each iteration

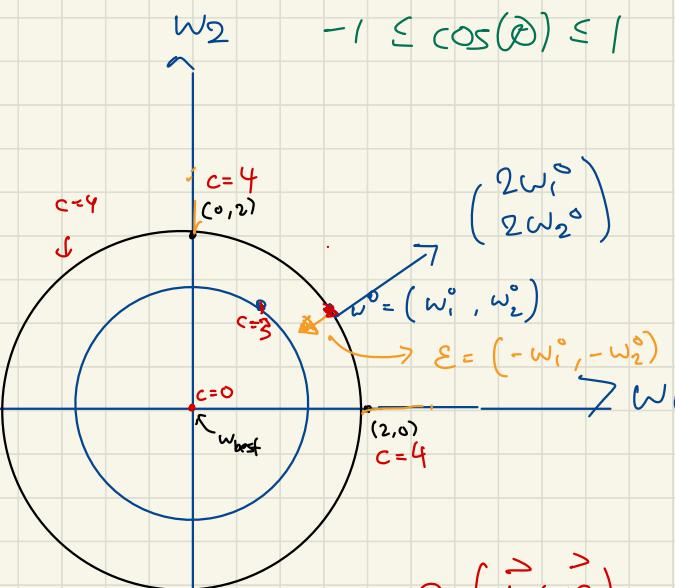
gradient descent (example)

$$C(\vec{w}) = w_1^2 + w_2^2$$

$$C: \mathbb{R}^2 \rightarrow \mathbb{R}$$

punkt

$$\begin{pmatrix} w_1 \\ w_2 \end{pmatrix} \mapsto w_1^2 + w_2^2$$



$$C(\vec{w}_0 + \vec{\epsilon}) - C(\vec{w}_0)$$

$$\approx \nabla_w C(w_0) \cdot \vec{\epsilon}$$

if unit direction vector

$$< \|\nabla_w C(w_0)\| \|\vec{\epsilon}\| \cos \theta.$$

$$\nabla_w C: \mathbb{R}^2 \rightarrow \mathbb{R}^2.$$

$$\begin{pmatrix} w_1 \\ w_2 \end{pmatrix} \rightarrow \begin{pmatrix} \frac{\partial C}{\partial w_1} \left(\begin{pmatrix} w_1 \\ w_2 \end{pmatrix} \right) \\ \frac{\partial C}{\partial w_2} \left(\begin{pmatrix} w_1 \\ w_2 \end{pmatrix} \right) \end{pmatrix} = \begin{pmatrix} 2w_1 \\ 2w_2 \end{pmatrix}$$

punkt

$\vec{a} \cdot \vec{b} = \|\vec{a}\| \|\vec{b}\| \cos \theta$

 $\cos(0) = 1$

 $\cos(180) = -1$

 $\cos(90) = 0$

$$x_0 = 2 ; \eta = 0.1$$

EE2211 Tutorial 8

$$\hat{x}^{(1)} = x^{(0)} - \eta \nabla_x f(x^{(0)})$$

$$= 2 - (0.1) \cdot 4(2)^3$$

$$\nabla_x f(x_0) = 4x_0^3$$

Question 1

Suppose we are minimizing $f(x) = x^4$ with respect to x . We initialize x to be 2. We perform gradient descent with learning rate 0.1. What is the value of x after the first iteration?

Question 2

Please consider the csv file (government-expenditure-on-education.csv), which depicts the government's educational expenditure over the years. We would like to predict expenditure as a function of year. To do this, fit an exponential model $f(\mathbf{x}, \mathbf{w}) = \exp(-\mathbf{x}^T \mathbf{w})$ with squared error loss to estimate \mathbf{w} based on the csv file and gradient descent. In other words, $C(\mathbf{w}) = \sum_{i=1}^m (f(\mathbf{x}_i, \mathbf{w}) - y_i)^2$.

Note that even though year is one dimensional, we should add the bias term, so $\mathbf{x} = [1 \text{ year}]^T$. Furthermore, optimizing the exponential function is tricky (because a small change in \mathbf{w} can lead to large change in f). Therefore for the purpose of optimization, divide the "year" variable by the largest year (2018) and divide the "expenditure" by the largest expenditure, so that the resulting normalized year and normalized expenditure variables have maximum values of 1. Use a learning rate of 0.03 and run gradient descent for 2000000 iterations.

- (a) Plot the cost function $C(\mathbf{w})$ as a function of the number of iterations.
- (b) Use the fitted parameters to plot the predicted educational expenditure from year 1981 to year 2023.
- (c) Repeat (a) using a learning rate of 0.1 and learning rate of 0.001. What do you observe relative to (a)?

The goal of this question is for you to code up gradient descent, so I will provide you with the gradient derivation. First, please note that in general, $\nabla_{\mathbf{w}}(\mathbf{x}^T \mathbf{w}) = \mathbf{x}$. To see this:

$$\nabla_{\mathbf{w}}(\mathbf{x}^T \mathbf{w}) = \begin{bmatrix} \frac{\partial(\mathbf{x}^T \mathbf{w})}{\partial w_1} \\ \frac{\partial(\mathbf{x}^T \mathbf{w})}{\partial w_2} \\ \vdots \\ \frac{\partial(\mathbf{x}^T \mathbf{w})}{\partial w_d} \end{bmatrix} = \begin{bmatrix} \frac{\partial(w_1 x_1 + w_2 x_2 + \dots + w_d x_d)}{\partial w_1} \\ \frac{\partial(w_1 x_1 + w_2 x_2 + \dots + w_d x_d)}{\partial w_2} \\ \vdots \\ \frac{\partial(w_1 x_1 + w_2 x_2 + \dots + w_d x_d)}{\partial w_d} \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} = \mathbf{x} \quad (1)$$

The above equality will be very useful for the other questions as well. Now, going back to our question,

$$\nabla_{\mathbf{w}} C(\mathbf{w}) = \nabla_{\mathbf{w}} \sum_{i=1}^m (f(\mathbf{x}_i, \mathbf{w}) - y_i)^2 \quad (2)$$

$$= \sum_{i=1}^m \nabla_{\mathbf{w}} (f(\mathbf{x}_i, \mathbf{w}) - y_i)^2 \quad (3)$$

$$= \sum_{i=1}^m 2(f(\mathbf{x}_i, \mathbf{w}) - y_i) \nabla_{\mathbf{w}} f(\mathbf{x}_i, \mathbf{w}) \quad \text{chain rule} \quad (4)$$

$$= \sum_{i=1}^m 2(f(\mathbf{x}_i, \mathbf{w}) - y_i) \nabla_{\mathbf{w}} \exp(-\mathbf{x}_i^T \mathbf{w}) \quad (5)$$

$$= - \sum_{i=1}^m 2(f(\mathbf{x}_i, \mathbf{w}) - y_i) \exp(-\mathbf{x}_i^T \mathbf{w}) \nabla_{\mathbf{w}} (\mathbf{x}_i^T \mathbf{w}) \quad \text{chain rule} \quad (6)$$

$$= - \sum_{i=1}^m 2(f(\mathbf{x}_i, \mathbf{w}) - y_i) \exp(-\mathbf{x}_i^T \mathbf{w}) \mathbf{x}_i \quad (7)$$

$$= - \sum_{i=1}^m 2(f(\mathbf{x}_i, \mathbf{w}) - y_i) f(\mathbf{x}_i, \mathbf{w}) \mathbf{x}_i \quad (8)$$

+

Question 3

Given the linear learning model $f(\mathbf{x}, \mathbf{w}) = \mathbf{x}^T \mathbf{w}$, where $\mathbf{x} \in \mathbb{R}^d$. Consider the loss function $L(f(\mathbf{x}_i, \mathbf{w}), y_i) = (f(\mathbf{x}_i, \mathbf{w}) - y_i)^2$, where i indexes the i -th training sample. The final cost function is $C(\mathbf{w}) = \sum_{i=1}^m L(f(\mathbf{x}_i, \mathbf{w}), y_i)$, where m is the total number of training samples. Derive the gradient of the cost function with respect to \mathbf{w} .

Question 4

Repeat Question 3 using $f(\mathbf{x}, \mathbf{w}) = \sigma(\mathbf{x}^T \mathbf{w})$, where $\sigma(a) = \frac{1}{1 + \exp(-\beta a)}$

Question 5

Repeat Question 3 using $f(\mathbf{x}, \mathbf{w}) = \sigma(\mathbf{x}^T \mathbf{w})$, where $\sigma(a) = \max(0, a)$

$$\begin{array}{c}
 \vec{\omega} \rightarrow \boxed{f(-, \vec{x}_i)} \rightarrow f(\vec{x}_i, \omega) \rightarrow \boxed{L(-, y_i)} \rightarrow L(f(\vec{x}_i, \omega), y_i) \\
 \downarrow \\
 f(-, \vec{x}_i) : \mathbb{R}^d \rightarrow \mathbb{R} \\
 L(-, y_i) : \mathbb{R} \rightarrow \mathbb{R} \\
 C : \mathbb{R}^d \rightarrow \mathbb{R} \\
 \end{array}$$

* differentiation is linear
 $\nabla \sum \text{stuff} = \sum \nabla \text{stuff}$

Question 3

Given the linear learning model $f(\mathbf{x}, \mathbf{w}) = \mathbf{x}^T \mathbf{w}$, where $\mathbf{x} \in \mathbb{R}^d$. Consider the loss function $L(f(\mathbf{x}_i, \mathbf{w}), y_i) = (f(\mathbf{x}_i, \mathbf{w}) - y_i)^4$, where i indexes the i -th training sample. The final cost function is $C(\mathbf{w}) = \sum_{i=1}^m L(f(\mathbf{x}_i, \mathbf{w}), y_i)$, where m is the total number of training samples. Derive the gradient of the cost function with respect to \mathbf{w} .

For each

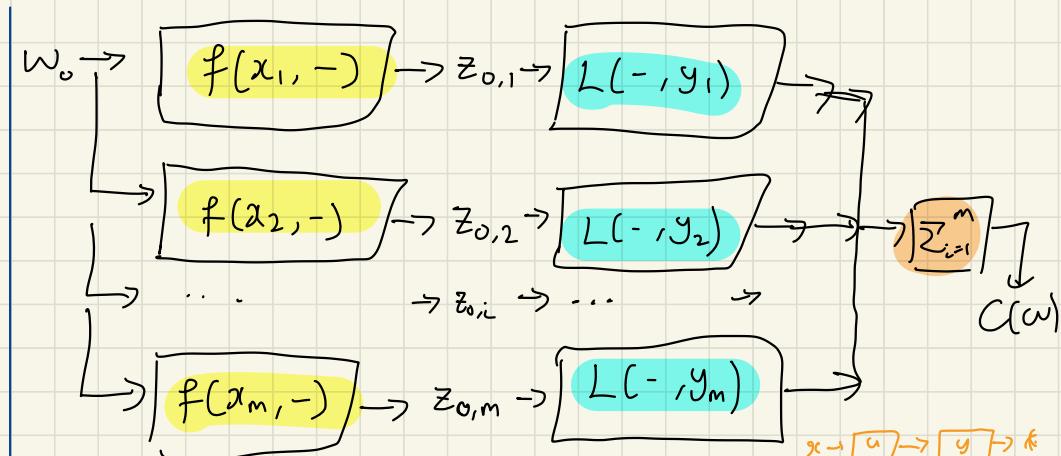
$$\mathbf{x}_i = \begin{pmatrix} x_i^1 \\ x_i^2 \\ \vdots \\ x_i^d \end{pmatrix} \in \mathbb{R}^d, \quad y_i \in \mathbb{R}$$

HAVE

$$f(x_i, -) : \mathbb{R}^d \rightarrow \mathbb{R} \quad w \mapsto x_i^T w$$

$$L(-, y_i) : \mathbb{R} \rightarrow \mathbb{R} \quad z \mapsto (z - y_i)^4$$

$$\sum_{i=1}^m : \mathbb{R}^m \rightarrow \mathbb{R} \quad \begin{pmatrix} c_1 \\ c_2 \\ \vdots \\ c_m \end{pmatrix} \mapsto c_1 + c_2 + \dots + c_m$$



How to differentiate?

Slogan : $\frac{dy}{dz} = \left[\frac{dy}{du} \right] \left[\frac{du}{dx} \right]$

matrix = matrix \cdot matrix
 linear map = linear map (compos) linear map

Intuition : $y \circ u(x_0 + \varepsilon) - y \circ u(x_0) \approx \underbrace{\left[\frac{dy}{du}(u(x_0)) \right]}_h \underbrace{\left[\frac{du}{dx}(x_0) \right]}_d (\varepsilon)$

$x_0 \mapsto u(x_0)$
 $u: \mathbb{R}^d \rightarrow \mathbb{R}^h$
 $\frac{du}{dz}: \mathbb{R}^d \rightarrow \mathbb{R}^{h \times d}$
 $x_0 \mapsto \left[\frac{\partial u}{\partial x} \right]$

$y: \mathbb{R}^h \rightarrow \mathbb{R}^k$
 $y(u_0) \mapsto \left[\frac{\partial y}{\partial u} \right]$
 $\frac{dy}{du}: \mathbb{R}^h \rightarrow \mathbb{R}^{k \times h}$

$\mathbb{R}^d \rightarrow \mathbb{R}^h$

Let's calculate the derivative of w_0 $\frac{dc}{dw}(w_0) \in \mathbb{R}^{1 \times d}$

$$\left[\frac{df(x_i, -)}{dw}(w_0) \right] = \begin{pmatrix} \frac{\partial x_i^T(-)}{\partial w_1} & \frac{\partial x_i^T(-)}{\partial w_2} & \dots & \frac{\partial x_i^T(-)}{\partial w_d} \end{pmatrix}$$

$$= 1 \begin{pmatrix} x_i^1 & x_i^2 & \dots & x_i^d \end{pmatrix} = x_i^T \in \mathbb{R}^{1 \times d}$$

$$\left[\frac{dL(-, y_i)}{dz}(z_0) \right] = \left[\frac{d(z-y)^4}{dz} \Big|_{z=z_0} \right] = \left[4(z_0-y)^3 \right] \in \mathbb{R}^{1 \times 1}$$

What is z_0 in chain rule?

$$= \left[4(x_i^T w_0 - y)^3 \right] \in \mathbb{R}^{1 \times 1}$$

must be $z_0 = f(x_i, w_0) = x_i^T w_0$

$$\left[\frac{d \sum_{i=1}^m}{dc}(c_0) \right] = 1 \left[1 \quad 1 \quad \dots \quad 1 \right] = \text{"sum over all entries"}$$

OBV! Because differentiation is linear!

$$\frac{d}{dw} \sum_{i=1}^m F_i(w_0) = \sum_{i=1}^m \frac{dF_i}{dw}(w_0)$$

$$\frac{df(x_i, -)}{dw}: \mathbb{R}^d \rightarrow \mathbb{R}^{1 \times d}$$

$$\frac{dL(-, y_i)}{dz}: \mathbb{R}^1 \rightarrow \mathbb{R}^{1 \times 1}$$

Everything TGR:

$$\left[\frac{dC}{dw}(w_0) \right] = \sum_{i=1}^M \left[d \left(\frac{L(-, y_i) \circ f(x_i, -)}{dw} \right) (w_0) \right]$$

differentiation
is linear

$$\begin{aligned}
 &= \sum_{i=1}^M \mathbf{1} \left[\frac{dL(-, y_i)}{dz} \left(z_0 = f(x_i, w_0) \right) \right] \mathbf{1} \left[\frac{df(x_i, -)}{dw} (w_0) \right] \\
 &= \sum_{i=1}^M \mathbf{1} \left[4(z_0 - y)^3 \Big|_{z_0 = x_i^T w_0} \right] \mathbf{1} \left[x_i^T \right] \\
 &= \mathbf{1} \left[\sum_{i=1}^M 4(x_i^T w_0 - y)^3 x_i^T \right]
 \end{aligned}$$

∴ the gradient of w_0 is $(\text{gradient}^T(-) = \text{derivative}(-))$

$$\nabla_w C(w_0) = \left[\frac{dC}{dw}(w_0) \right]^T = \sum_{i=1}^M 4(x_i^T w_0 - y) x_i$$

AT w_0 , we know

① $z_{0,i} = f(x_i, w_0)$
= $x_i^T w_0$.

② $c_0 = \text{sth}$

d

this is a vector in \mathbb{R}^d

Alternative Perspective:

Define $F(X, -) : \mathbb{R}^d \rightarrow \mathbb{R}^m$ ↗ no. of training samples

by $w_0 \mapsto \begin{pmatrix} x_1^T w_0 \\ x_2^T w_0 \\ \vdots \\ x_m^T w_0 \end{pmatrix}$

$L(-, Y) : \mathbb{R}^m \rightarrow \mathbb{R}^m$

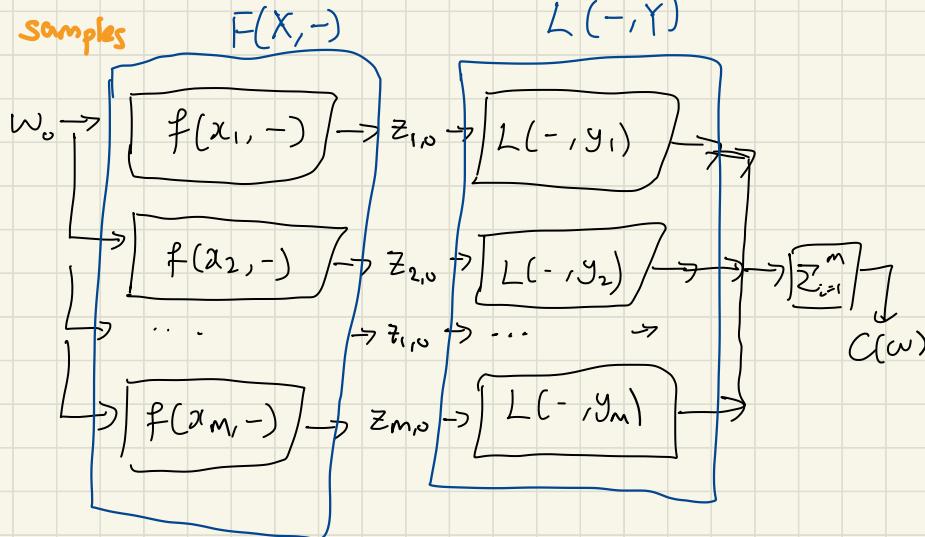
by $z_0 = \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_m \end{pmatrix} \mapsto \begin{pmatrix} (z_1 - y_1)^4 \\ (z_2 - y_2)^4 \\ \vdots \\ (z_m - y_m)^4 \end{pmatrix}$

$\sum_{i=1}^m : \mathbb{R}^m \rightarrow \mathbb{R}$

by $\begin{pmatrix} c_1 \\ \vdots \\ c_m \end{pmatrix} \mapsto c_1 + \dots + c_m$

Then, $\frac{dF(X, -)}{dw} : \mathbb{R}^d \rightarrow \mathbb{R}^{md}$ is

$$w_0 \mapsto \left[\frac{dF(X, -)}{dw}(w_0) \right] = \begin{bmatrix} \frac{\partial(x_1^T w)}{\partial w_1} |_{w=w_0} \\ \frac{\partial(x_2^T w)}{\partial w_1} |_{w=w_0} \\ \vdots \\ \frac{\partial(x_m^T w)}{\partial w_1} |_{w=w_0} \end{bmatrix}$$



described by $x_i^T w$, i mean the i^{th} component of $w \mapsto x_i^T w$. This is the component fn of F .

$$\begin{bmatrix} \frac{\partial(x_1^T w)}{\partial w_2} |_{w=w_0} & \frac{\partial(x_1^T w)}{\partial w_m} |_{w=w_0} \\ \frac{\partial(x_2^T w)}{\partial w_2} |_{w=w_0} & \vdots \\ \vdots & \vdots \\ \frac{\partial(x_m^T w)}{\partial w_2} |_{w=w_0} & \dots & \frac{\partial(x_m^T w)}{\partial w_m} |_{w=w_0} \end{bmatrix}^m = \begin{bmatrix} x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(d)} \\ x_2^{(1)} & x_2^{(2)} & \dots & x_2^{(d)} \\ \vdots & \vdots & \ddots & \vdots \\ x_m^{(1)} & x_m^{(2)} & \dots & x_m^{(d)} \end{bmatrix}^d = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_m^T \end{bmatrix}_{\mathbb{R}^{md}}$$

Similarly we get

$$\frac{dL(-, Y)}{dz} : \mathbb{R}^m \rightarrow \mathbb{R}^{m \times m}$$

$$\begin{pmatrix} z_{1,0} \\ \vdots \\ z_{m,0} \end{pmatrix} = z_0 \mapsto \left[\frac{dL(-, Y)}{dz}(z_0) \right]$$

$$M \begin{bmatrix} \frac{\partial(z_1 - y_1)^4}{\partial z_1} & \dots & \frac{\partial(z_1 - y_1)^4}{\partial z_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial(z_m - y_m)^4}{\partial z_1} & \dots & \frac{\partial(z_m - y_m)^4}{\partial z_m} \end{bmatrix}_{|z=z_0} = \begin{bmatrix} 4(z_{1,0} - y_1)^3 \\ 4(z_{2,0} - y_2)^3 \\ \vdots \\ 4(z_{m,0} - y_m)^3 \end{bmatrix}$$

diagonal matrix

(RECALL)

$$L(-, Y) : \mathbb{R}^m \rightarrow \mathbb{R}^m$$

$$\text{by } z_0 = \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_m \end{pmatrix} \mapsto \begin{pmatrix} (z_1 - y_1)^4 \\ (z_2 - y_2)^4 \\ \vdots \\ (z_m - y_m)^4 \end{pmatrix}$$

$$\text{Similarly, } \frac{d \sum_{i=1}^m}{dc} : \mathbb{R}^m \rightarrow \mathbb{R}^{1 \times m}$$

$$c_0 \mapsto \left[\frac{d \sum_{i=1}^m}{dc}(c_0) \right] = \begin{bmatrix} 1 & 1 & \dots & 1 \end{bmatrix}^T$$

$$\sum_{i=1}^m : \begin{pmatrix} c_1 \\ c_2 \\ \vdots \\ c_m \end{pmatrix} \mapsto c_1 + c_2 + \dots + c_m$$

RECALL

$$L(-, Y) : \mathbb{R} \rightarrow \mathbb{R}$$

$$z \mapsto (z - y_i)^4$$



$$\frac{dL(-, Y)}{dz} : \mathbb{R} \rightarrow \mathbb{R}^{1 \times 1}$$

$$z_0 \mapsto 4(z_{0,0} - y_i)^3$$

OBSERVE :

$$= \begin{bmatrix} \frac{dL(-, Y)}{dz}(z_{1,0}) \\ \vdots \\ \frac{dL(-, Y)}{dz}(z_{m,0}) \end{bmatrix}$$

$$TGT: C: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$\omega \mapsto C(\omega) = \sum_{i=1}^m \circ L(-, \gamma) \circ F(x, -)(\omega)$$

$$\rightarrow \boxed{F(x, -): \mathbb{R}^d \rightarrow \mathbb{R}^m} \rightarrow \boxed{L(-, \gamma): \mathbb{R}^m \rightarrow \mathbb{R}^m} \rightarrow \boxed{\sum_{i=1}^m: \mathbb{R}^m \rightarrow \mathbb{R}}$$

at position w_0 ,

get

$$\omega_0 \mapsto \begin{pmatrix} x_1^T w_0 \\ x_2^T w_0 \\ \vdots \\ x_m^T w_0 \end{pmatrix} \mapsto z_0 = \begin{pmatrix} z_{1,0} \\ z_{2,0} \\ \vdots \\ z_{m,0} \end{pmatrix} = \begin{pmatrix} (x_1^T w_0 - y_1)^4 \\ (x_2^T w_0 - y_2)^4 \\ \vdots \\ (x_m^T w_0 - y_m)^4 \end{pmatrix} \mapsto \sum_{i=1}^m (x_i^T w_0 - y_i)^4 \underbrace{c_0}_{C_0}$$

$$\frac{dC}{d\omega}: \mathbb{R}^d \rightarrow \mathbb{R}^{1 \times d}$$

$$\begin{aligned} \omega_0 \mapsto \left[\frac{dC}{d\omega}(\omega_0) \right] &= \underbrace{1 \left[\frac{d \sum_{i=1}^m}{dc} (c_0) \right]}_{m \times 1} \underbrace{m \left[\frac{dL(-, \gamma)}{dz}(z_0) \right]}_{m \times m} \underbrace{m \left[\frac{dF(x, -)}{dw}(w_0) \right]}_{m \times d} \\ &= \underbrace{1 \left[\begin{matrix} 1 & 1 & \dots & 1 \end{matrix} \right]}_{m \times 1} \underbrace{\left[\begin{matrix} 4(x_1^T w_0 - y_1)^3 & & & \\ & \ddots & & \\ & & 4(x_m^T w_0 - y_m)^3 & \end{matrix} \right]}_{m \times m} \underbrace{\left[\begin{matrix} x_1^T \\ x_2^T \\ \vdots \\ x_m^T \end{matrix} \right]}_{d \times 1} \\ &= \underbrace{1 \left[\begin{matrix} \frac{d}{m} \sum_{i=1}^m 4(x_i^T w_0 - y_i)^3 x_i^T \end{matrix} \right]}_{1 \times d} \end{aligned}$$

OBSERVE: each entry is raised by a scalar since this is diagonal

$$\therefore \nabla_w C(\omega_0) = \left[\frac{dC}{d\omega}(\omega_0) \right]^T = \sum_{i=1}^m 4(x_i^T w_0 - y_i) x_i \in \mathbb{R}^d$$

Question 4

Repeat Question 3 using $f(\mathbf{x}, \mathbf{w}) = \sigma(\mathbf{x}^T \mathbf{w})$, where $\sigma(a) = \frac{1}{1+\exp(-\beta a)}$

$$C: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$\downarrow \omega: \mathbb{R}^d \rightarrow \mathbb{R}^m$$

abuse notation

$$\begin{aligned} \omega & \mapsto \begin{pmatrix} \mathbf{x}_1^T \mathbf{w} \\ \mathbf{x}_2^T \mathbf{w} \\ \vdots \\ \mathbf{x}_m^T \mathbf{w} \end{pmatrix} \sigma: \mathbb{R}^m \rightarrow \mathbb{R}^m \mapsto \begin{pmatrix} \sigma(\mathbf{x}_1^T \mathbf{w}) \\ \sigma(\mathbf{x}_2^T \mathbf{w}) \\ \vdots \\ \sigma(\mathbf{x}_m^T \mathbf{w}) \end{pmatrix} L(-, Y): \mathbb{R}^m \rightarrow \mathbb{R}^d \mapsto \begin{pmatrix} (\sigma(\mathbf{x}_1^T \mathbf{w}) - y_1)^4 \\ (\sigma(\mathbf{x}_2^T \mathbf{w}) - y_2)^4 \\ \vdots \\ (\sigma(\mathbf{x}_m^T \mathbf{w}) - y_m)^4 \end{pmatrix} \sum_{i=1}^m: \mathbb{R}^d \rightarrow \mathbb{R} \mapsto \sum_{i=1}^m (\sigma(\mathbf{x}_i^T \mathbf{w}) - y_i)^4 \end{aligned}$$

$$\frac{dC}{d\mathbf{w}}: \mathbb{R}^d \rightarrow \mathbb{R}^{1 \times d}$$

$$\mathbf{w}_0 \mapsto \frac{1}{m} \left[\frac{d \sum_{i=1}^m}{dc} (c_i) \right] = \frac{1}{m} \left[\frac{dL(-, Y)}{dz} (z_i) \right] = \frac{1}{m} \begin{bmatrix} \frac{d\sigma(\mathbf{x}_1^T \mathbf{w}_0)}{da} \\ \vdots \\ \frac{d\sigma(\mathbf{x}_m^T \mathbf{w}_0)}{da} \end{bmatrix}$$

$$= \left[\begin{array}{c} \text{add everything up.} \\ \text{Differentiation is linear} \end{array} \right] \left[\begin{array}{c} \text{scalar mult.} \\ \text{the } i^{\text{th}} \text{ entry by} \\ \frac{dL(-, Y_i)}{dz} (\sigma(\mathbf{x}_i^T \mathbf{w}_0)) \end{array} \right] \left[\begin{array}{c} \text{scalar mult.} \\ \text{the } i^{\text{th}} \text{ entry by} \\ \frac{d\sigma}{da} (\mathbf{x}_i^T \mathbf{w}_0) \end{array} \right] \left[\begin{array}{c} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_m^T \end{array} \right]$$

vector in $\mathbb{R}^{1 \times d}$; m copies

$$= \sum_{i=1}^m 4(\sigma(\mathbf{x}_i^T \mathbf{w}_0) - y_i)^3 \cdot \beta \sigma(\mathbf{x}_i^T \mathbf{w}_0)(1 - \sigma(\mathbf{x}_i^T \mathbf{w}_0)) \mathbf{x}_i^T$$

$$\nabla_{\mathbf{w}} C(\mathbf{w}_0) = \left[\frac{dC}{d\mathbf{w}}(\mathbf{w}_0) \right]^T = \sum_{i=1}^m 4(\sigma(\mathbf{x}_i^T \mathbf{w}_0) - y_i)^3 \cdot \beta \sigma(\mathbf{x}_i^T \mathbf{w}_0)(1 - \sigma(\mathbf{x}_i^T \mathbf{w}_0)) \mathbf{x}_i$$

Question 5

$\sigma \in \mathbb{W}$.

Repeat Question 3 using $f(\mathbf{x}, \mathbf{w}) = \sigma(\mathbf{x}^T \mathbf{w})$, where $\sigma(a) = \max(0, a)$

Same story. what is derivative of σ ?

$$\frac{d\sigma}{da}(a_0) = \begin{cases} 1, & a_0 > 0 \\ 0, & a_0 < 0 \\ \text{undefined}, & a_0 = 0 \end{cases}$$

$$\sigma(a) = \begin{cases} a, & a > 0 \\ 0, & a \leq 0 \end{cases}$$

$$\frac{dc}{dw} : \mathbb{R}^d \rightarrow \mathbb{R}^{1 \times d}$$

$$w_0 \mapsto \frac{1}{m} \left[\frac{d \sum_{i=1}^m}{dc} (c_0) \right]$$

$$= \begin{bmatrix} \text{add everything up.} \\ \text{Differentiation is linear} \end{bmatrix}$$

$$= \begin{bmatrix} \text{scalar mult.} \\ \text{the } i\text{th entry by} \\ \frac{dL(-, Y)}{dz} (z_0) \end{bmatrix} \begin{bmatrix} \frac{d\sigma}{da}(a_{i,0}) \\ \vdots \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} 0 \\ \frac{d\sigma}{da}(a_m) \end{bmatrix} \begin{bmatrix} \frac{dF(X, -)}{dw} (w_0) \\ \vdots \\ x_i^T \end{bmatrix}$$

$$\text{what is this?} \\ \text{This is} \\ \begin{cases} 1, & \text{if } x_i^T w_0 > 0 \\ 0, & \text{if } x_i^T w_0 \leq 0 \end{cases}$$

$$= \sum_{i=1}^m 4 \left(\mathbb{1}_{\{x_i^T w_0 \geq 0\}} (x_i^T w_0) - y_i \right)^3 \cdot \begin{cases} 1 & \text{if } x_i^T w_0 > 0 \\ 0 & \text{if } x_i^T w_0 \leq 0 \end{cases} \cdot x_i^T$$

$$= \sum_{i=1}^m 4 \left(\mathbb{1}_{\{x_i^T w_0 \geq 0\}} (x_i^T w_0) - y_i \right)^3 \cdot \begin{cases} 1 & \text{if } x_i^T w_0 > 0 \\ 0 & \text{if } x_i^T w_0 \leq 0 \end{cases} \cdot x_i^T$$

|| save me.
 delta fn, indicator fn. Just a name

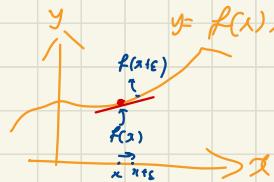
$$= \sum_{i=1, x_i^T w_0 > 0}^m 4 \left((x_i^T w_0) - y_i \right)^3 \cdot x_i^T$$

$$\therefore \nabla_w C(w_0) = \left[\frac{dc}{dw}(w_0) \right]^T$$

$$= \sum_{i=1, x_i^T w_0 > 0}^m 4 \left((x_i^T w_0) - y_i \right)^3 \cdot x_i$$

Recall

If $f: \mathbb{R} \rightarrow \mathbb{R}$ is a differentiable fn,



then $x \mapsto \frac{df}{dx}(x) = f'(x)$ is a function

$f': \mathbb{R} \rightarrow \mathbb{R}$ s.t.

$$f(x + \epsilon) - f(x) \approx f'(x) \cdot \epsilon.$$

If $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is a cts differentiable fn,

then $\vec{x} \mapsto \vec{\nabla} f(\vec{x}) = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{pmatrix}^*$

is a function

$$\vec{\nabla} f: \mathbb{R}^d \rightarrow \mathbb{R}^d$$

$$f(\vec{x} + \vec{\epsilon}) - f(\vec{x}) \approx \vec{\nabla} f(\vec{x}) \cdot \vec{\epsilon}$$

vector in \mathbb{R}^d s.t.
vector in \mathbb{R}^d .

Recall dot product of vectors.

$$\vec{v} \cdot \vec{w} = \|\vec{v}\| \|\vec{w}\| \cos \theta$$

$$\vec{v} \cdot \vec{w} = \vec{v}^T \vec{w}$$

dot prod. by
 $\vec{\nabla} f(\vec{x})$ s.t.
 $\mathbb{R}^d \rightarrow \mathbb{R}$.
and
is a
linear
map.

$$\begin{aligned} \vec{\nabla} f \begin{pmatrix} x_0 \\ y_0 \end{pmatrix} &= \begin{pmatrix} \frac{\partial f}{\partial x}(x_0, y_0) \\ \frac{\partial f}{\partial y}(x_0, y_0) \end{pmatrix} \\ &= \begin{pmatrix} 2x_0 \\ 2y_0 \end{pmatrix} \end{aligned}$$

