

---

---

---

---

---



## Lecture 12

08.10.21

### Admin

• Last class today. ~ GD for losses.

HW coming

• No class next week.

• New topics & projects following week

So regression / binary class

— analysis.

$$\nabla \hat{L}(w) = 0$$

— optimization

$$w_{k+1} = w_k - \eta \nabla \hat{L}(w_k)$$

to find  $w^*$   $h_w(x)$

$h_w(x)$  "best" model

|| Major topic Multiclass. ✓

|| Generalization

$h_w(x)$

$\hat{L}(w) \rightarrow$

fit all data  
perfectly

what about new data?

$\Rightarrow$  Statistical learning theory.

Defining other important ML-Deep AI  
problems.

contrastive loss pairs

$$s(x_1, x_2)$$

$$d(x_1, x_3)$$

learn  $f(x)$

$$\text{s.t. } f(x) \cdot f(x') \geq$$

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

when  $x \sim x'$   
 $x \not\sim x''$

NLP > 5 years

contrastive  
loss

< 2 years  
ImageNet

---

$$\text{pmi} = \log\left(\frac{p_{ij}}{p_i p_j}\right)$$

count word pair

freq  $p_{ij}$

ind freq  $p_i, p_j$

current

Solve

$$f(x_i) - f(x_j) = \log\left(\frac{p_{ij}}{p_i p_j}\right)$$

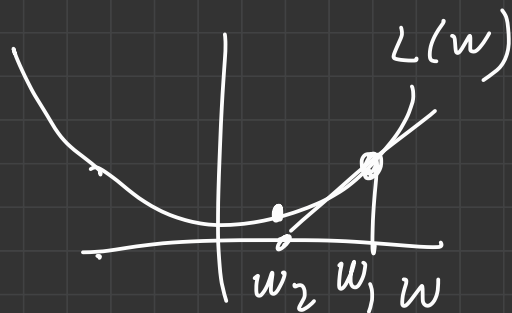
similarity loss

Chain Rule

$$h_w(x)$$

$$\ell(h, y)$$

$$\hat{L}(w) = \frac{1}{m} \sum_{i=1}^m \ell(h_w(x_i), y_i)$$



Standard  $\nabla \hat{L}(w) = 0$

Now GD  $w_{k+1} = w_k - h \nabla \hat{L}(w_k)$   
(with stepsize/  
learning rate  $h$ )

Chain Rule  $\nabla_w \hat{L}(w) = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial \ell}{\partial h}(h_w(x_i), y_i)}_{\text{loss}'}}_{\text{model gradient}} \nabla_w h_w(x_i)$

$$h_w(x) = w \cdot x \quad \nabla_w h_w = x$$

Abstract GD:

$\hat{L}(w)$   $\mu$ -convex &  $L$ -smooth

Defn

$$H(w) = D^2 L(w)$$

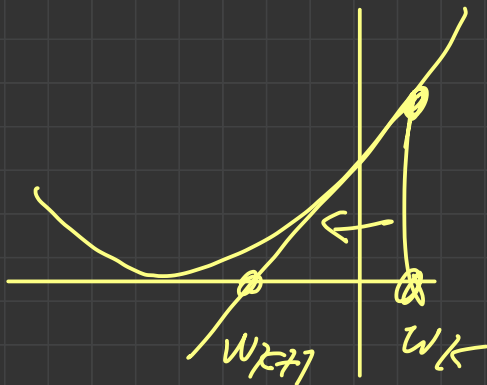
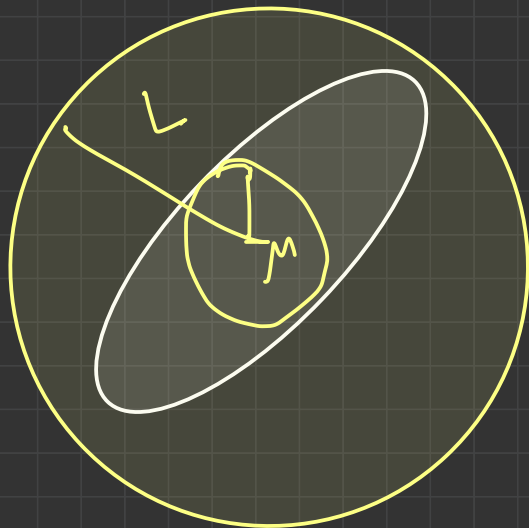
$$\mu I \leq H(w) \leq L I$$

$M, N$  p.s.d. matrix

$$M \leq N \iff x^T M x \leq x^T N x \quad \forall x$$

all eigenvals of  $H$  are between  
 $M$  and  $L$

$\hat{L}(w)$



$$w_{k+1} = w_k - h \nabla L(w_k)$$

cond number of convex fn

$$C = \frac{L}{\mu}$$

convergence rate of GD  
worst case analysis

$$h = \frac{1}{L}$$

rate.  $e_k \leq C^k e_0$

$$f^* = \min_x f(x)$$

$$\hat{L}^* = \min_w \hat{L}(w)$$

$$e_k = f(x_k) - f^*$$

$$e_k = \hat{L}(w_k) - \hat{L}^*$$

Line search ...

$$w_{k+1} = w_k - h \nabla \overset{1}{L}(w_k) \quad \checkmark$$

Look for better.

---



Loss gradients.

① Threshold models

$$S_m = \{y_1, \dots, y_m\}$$

$$h_w = w$$

$$= \ell_2(\underbrace{h-y}_e) = (h-y)^2/2 = e^2/2$$

$$\ell'_2(e)$$

$$\begin{aligned} \mathcal{L}'(w) &= \frac{1}{m} \sum_{i=1}^m \ell(h, y_i) \\ &= \frac{1}{m} \sum_{i=1}^m (w - y_i)^2 \end{aligned}$$

case 1  $h=1 \Rightarrow w_i = \bar{y}$

case 2  $h > 1$  diverges.

case 3  $h = \frac{1}{2}$

$$w_{k+1} - \bar{y} = \frac{1}{2}(w_k - \bar{y}) \Rightarrow \underline{\underline{e_k = \frac{1}{2^k} e_0}}}$$

$$w^* = \bar{y} = \frac{1}{m} \sum_{i=1}^m$$

$$\ell'(e) = e$$

$$\begin{aligned} \mathcal{L}'(w) &= \frac{1}{m} \sum_{i=1}^m w - y_i \\ &= w - \bar{y} \end{aligned}$$

G-D

$$\begin{aligned} w_{k+1} &= w_k - h(w_k - \bar{y}) \\ &= (1-h)w_k + h\bar{y} \end{aligned}$$

$$e_k = w_k - \bar{y}$$

GD2

Linear Model

$$x \in \mathbb{R}^d \quad w \in \mathbb{R}^d$$
$$h_w(x) = w \cdot x \quad \nabla_w h = x$$

Chain Rule

$$\nabla_w \hat{L}(w) = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial \ell}{\partial h}(h_w(x_i), y_i)}_{h(x_i) - y_i} \underbrace{\nabla_w h_w(x_i)}_{x_i}$$

$\ell_2$

$$\frac{\partial \ell}{\partial h} = e = h(x_i) - y_i$$

$$\nabla_w \hat{L}(w) = \frac{1}{m} \sum_{i=1}^m \underbrace{(h(x_i) - y_i)}_{w \cdot x_i} \underbrace{x_i}_{w^* \cdot x_i}$$

$$y = w^* \cdot x$$

$$(w - w^*) \cdot x_i$$

simplify  
...

$$\underline{d=2} \quad (x_1, x_2) = (S, 1)$$

$$y = mS + b$$

$$(w_1, w_2) = (w_m, w_b)$$

$$\underline{\underline{HW}} \quad \frac{\partial L}{\partial w_b} = (w_b - b) + (w_m - m) \bar{S}$$

$$\frac{\partial L}{\partial w_m} = (w_b - b) \bar{S} + (w_m - m) \bar{S}^2$$

$$\bar{S} = \frac{1}{m} \sum S_i \quad \bar{S}^2 = \frac{1}{m} \sum S_i^2$$

$$\begin{aligned} \nabla_w \mathcal{L}(w) &= H(w - w^*) \\ &= \begin{bmatrix} \bar{S}^2 & \bar{S} \\ \bar{S} & 1 \end{bmatrix} [w - w^*] \end{aligned}$$

