

## Lecture 2

$$h_1^* \in \underset{h \in H_1}{\operatorname{argmin}} L_T(h)$$

$$h_2^* \in \underset{h \in H_2}{\operatorname{argmin}} L_T(h)$$

$$H_1 \subset H_2$$

$$L_T(h_1^*) \geq \underline{\underline{L_T(h_2^*)}}$$



$$L_V(h_1^*) \leq L_V(h_2^*)$$

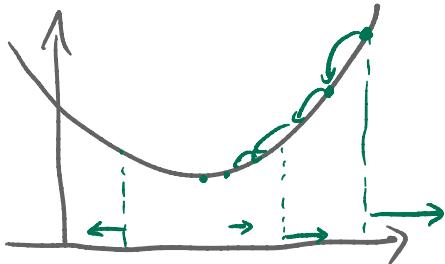
→ PICK  $h_2^*$

$h_2^*$  IS OVERFITTING

→ PICK  $h_1^*$

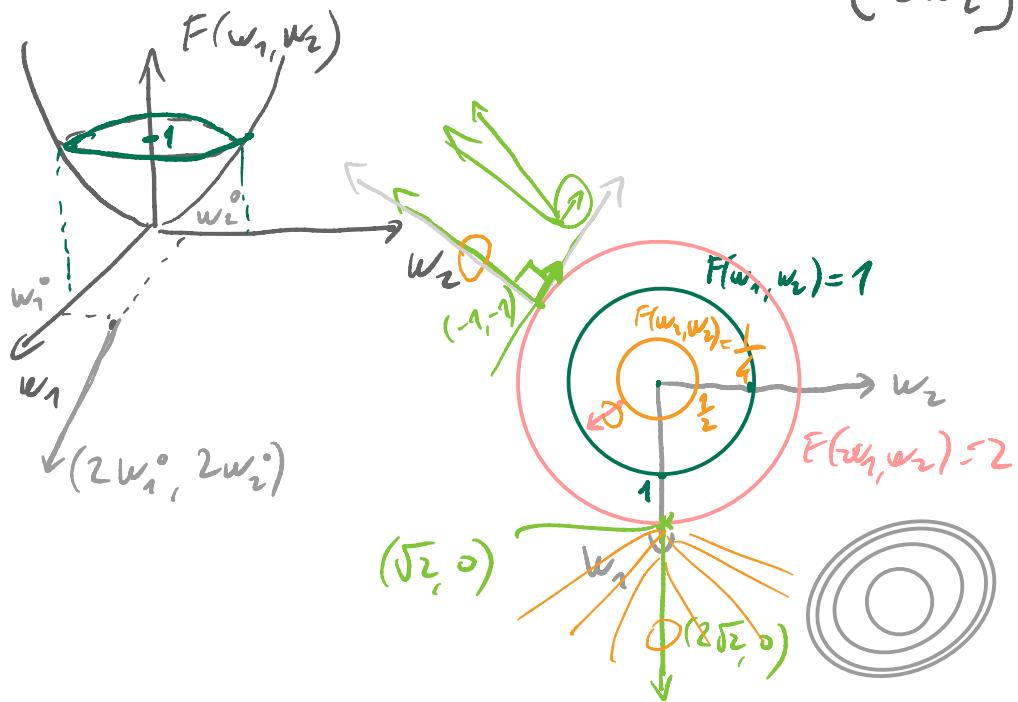
$h_1^*$  IS UNDERFITTING

$$\nabla F(\omega) = \begin{bmatrix} \frac{\partial F}{\partial \omega_1} \\ \vdots \\ \frac{\partial F}{\partial \omega_n} \end{bmatrix}$$

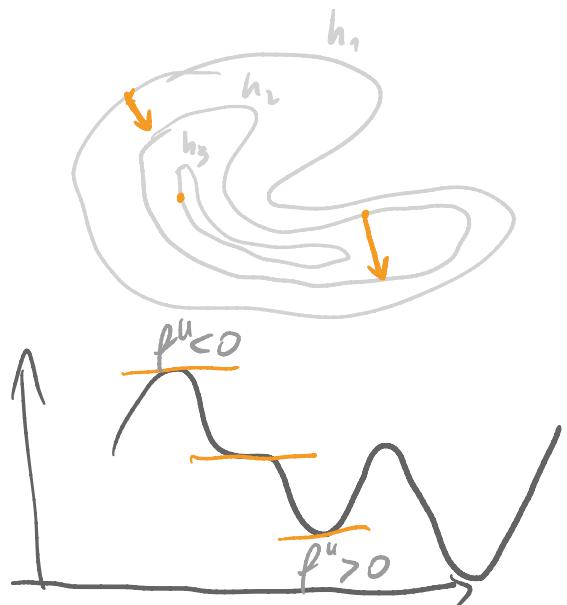


$$F(\omega_1, \omega_2) = \omega_1^2 + \omega_2^2$$

$$\nabla F = \begin{bmatrix} 2\omega_1 \\ 2\omega_2 \end{bmatrix}$$



$$h_3 > h_2 > h_1$$



$$f' \rightarrow \nabla F$$

$f'' \rightarrow \text{Hessian of } F \equiv H_F$   
 is a MATRIX

If  $F: \mathbb{R}^n \rightarrow \mathbb{R}$        $H_F$  is a  $n \times n$  MATRIX

$$(H_F)_{ij} = \frac{\partial}{\partial w_i} \frac{\partial}{\partial w_j} F$$

$$F(w_1, w_2) = w_1^2 + w_2^2$$

$$\begin{bmatrix} \frac{\partial^2 F}{\partial w_1^2} & \frac{\partial}{\partial w_1} \frac{\partial F}{\partial w_2} \\ \frac{\partial}{\partial w_2} \frac{\partial F}{\partial w_1} & \frac{\partial^2 F}{\partial w_2^2} \end{bmatrix} =$$

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

$$F(w_1, w_2) = w_1^2 + 2w_1 w_2$$

$f'' > 0 \iff H_F$  is a positive definite matrix.

$f'' \geq 0 \iff H_F$  is positive semidefinite matrix.

$A$  is positive semidefinite

$$\text{if } \forall x \in \mathbb{R}^n \quad x^T A x \geq 0$$

$$A \in \mathbb{R}^{1 \times 1} \quad x^T A x = x^2 A \geq 0 \Rightarrow x \geq 0$$

A positive definite

$$\forall x \in \mathbb{R}^n - \{0\} \quad x^T A x > 0$$

---

if  $\bar{w}$  :  $\nabla F(\bar{w}) = 0$  &  $H_F(\bar{w})$   
is POSITIVE  
DEFINITE

$\Rightarrow \bar{w}$  is a LOCAL MINIMUM

$$H_F \quad F''(w) \quad \nabla^2 F$$

---

A is diagonalizable, if  $\exists V$  (invertible)  
and  $\delta$  diagonal such that  
 $A = V \delta V^{-1}$

$\delta$  is diagonal

$$\delta = \begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{pmatrix} \quad \lambda_i: \text{eigenvalues}$$

IF  $\forall i \lambda_i \geq 0 \Rightarrow A$  is positive semidefinite

IF  $\forall i \lambda_i > 0 \Rightarrow A$  is positive definite

---

$$A = A^T \quad (A \text{ is symmetric})$$

$\Rightarrow A$  is diagonalizable

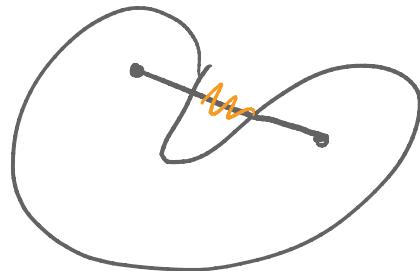
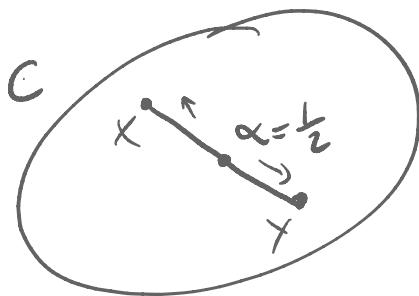
and  $V^{-1} = V^T \quad V$  is orthonormal

$$A = V \Delta V^{-1} = V \Delta V^T$$

$$V^{-1} = V^T \Rightarrow V^T V = I$$

# CONVEXITY

## CONVEX SET

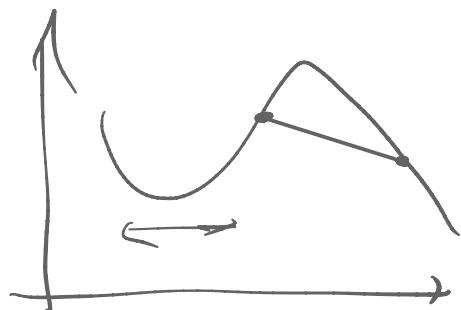
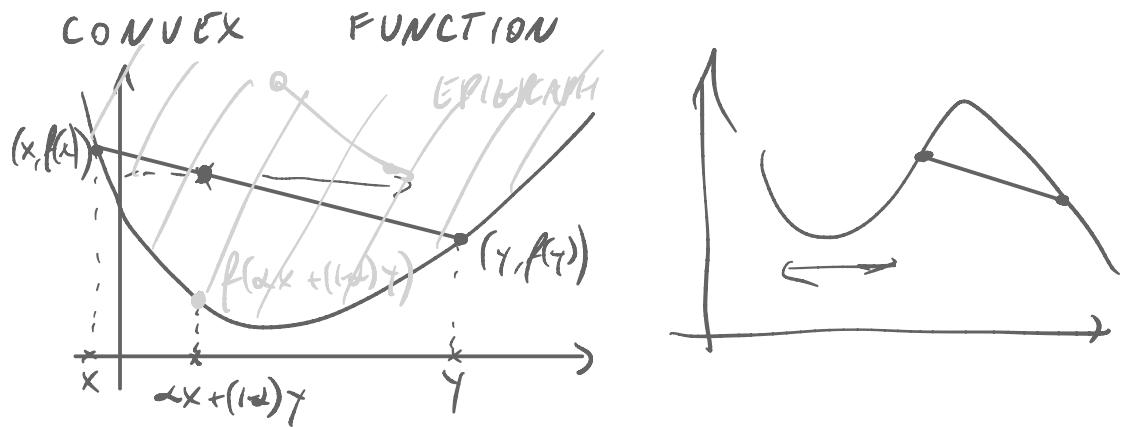


$C$  is a convex set if

$\forall x, y \in C$

$\forall \alpha \in [0, 1]$

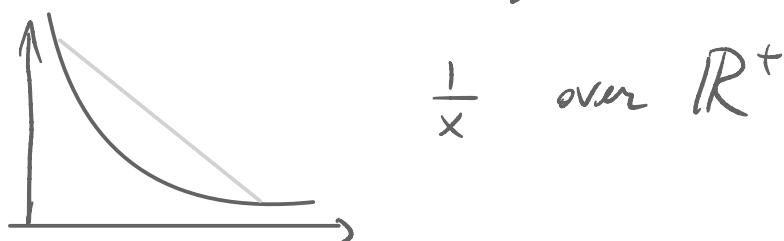
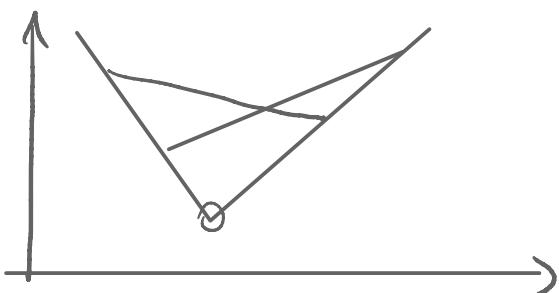
$$\alpha x + (1-\alpha) y \in C$$



$$f: B \rightarrow \mathbb{R}$$

$$\forall x, y \in B \quad \forall \alpha \in [0, 1] \quad B \text{ convex}$$

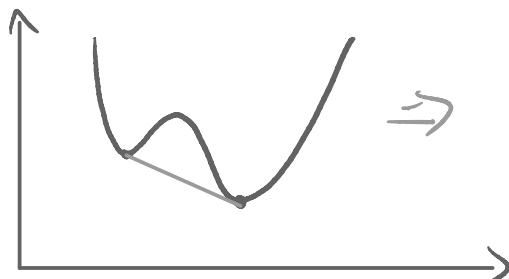
$$f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)$$



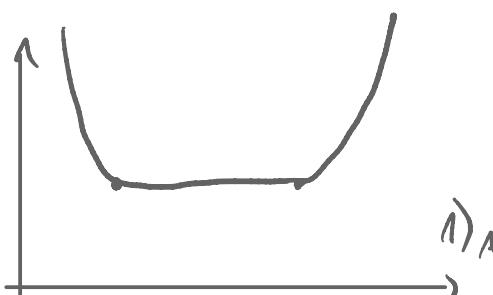
# CONVEX OPTIMIZATION

minimize  $f(x)$   
 $x \in B$

$B$  is convex  $f$  is convex

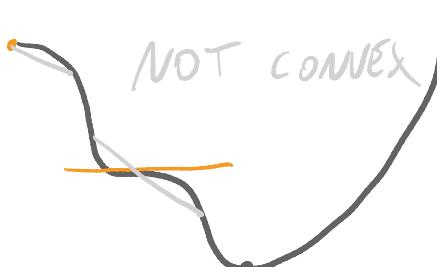


NO CONVEX



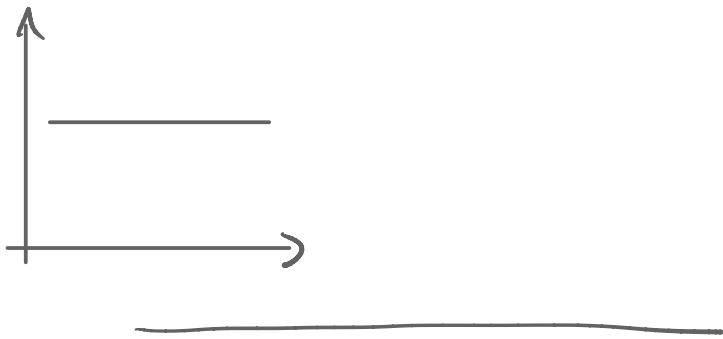
FOR CONVEX  
OPT PROBLEMS

1) A LOCAL MINIMIZER  
IS ALSO A GLOBAL  
ONE



NOT CONVEX / 2) A STATIONARY  
POINT ( $\nabla f = 0$ )

IS A GLOBAL  
MINIMIZER



## CLASSIFICATION PROBLEM

0-1 loss

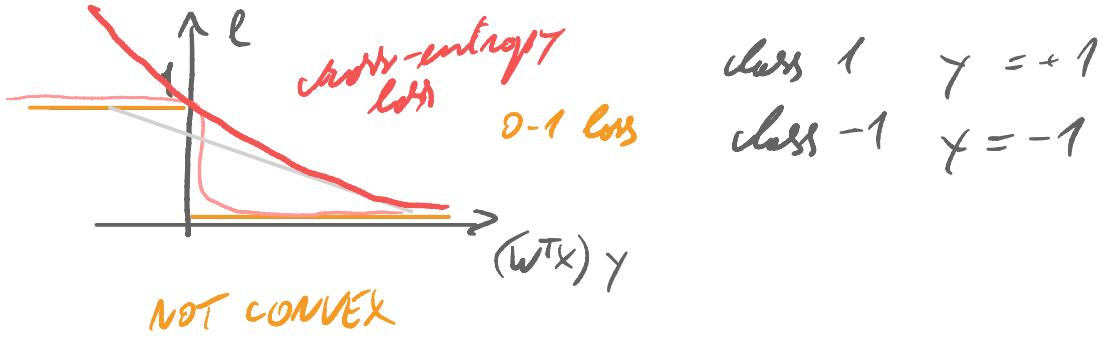
$$\ell(h, (x, \gamma)) = \mathbb{1}_{h(x) \neq \gamma} = \begin{cases} 1 & \text{if } h(x) \neq \gamma \\ 0 & \text{otherwise} \end{cases}$$

$$h(x) = w^T x$$

$$\ell(w, (x, \gamma)) = \mathbb{1}_{wx \neq \gamma}$$

CROSS-ENTROPY LOSS

$$\ell(w, (x, \gamma)) = -\gamma \ln wx - (1-\gamma) \ln(1-wx)$$



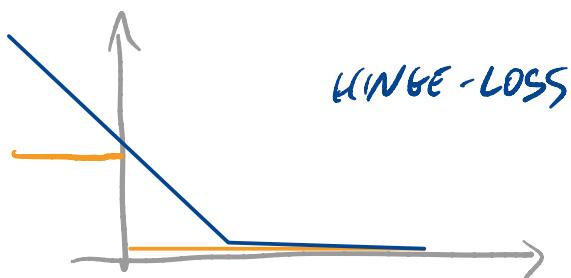
$$h(x) = \operatorname{sign} w^T x$$

$$w^T x > 0 \quad y=1 \quad l=0$$

$$w^T x < 0 \quad y=-1 \quad l=0$$

$$w^T x > 0 \quad y=-1 \quad l=1$$

$$w^T x < 0 \quad y=1 \quad l=1$$



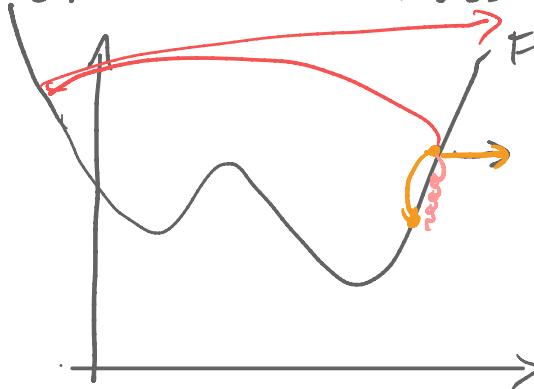
YOUR "NATURAL" LOSS  $l$  IS NON-CONVEX

YOU CONSIDER AN ALTERNATIVE  
"SURROGATE" LOSS  $\tilde{l}$

$$1) \tilde{l} \geq l$$

2)  $\tilde{l}$  IS CONVEX

## GRADIENT METHODS



minimize  
 $w$

$\alpha$  IS TOO LARGE

- 1) PICK AN INITIAL GUESS  $w_0$
- 2) FOR  $k = 1$  TO  $K$

$$w_k = w_{k-1} - \alpha_{k-1} \nabla F(w_{k-1})$$

LEARNING RATE

$$L_T(h) = \frac{1}{n} \sum_{i=1}^n \ell(h, (x_i, y_i))$$

$$T = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

$$h : X \rightarrow Y$$

EACH "h" WILL BE IDENTIFIED  
BY A VECTOR OF PARAMETERS  $w \in \mathbb{R}^n$

$$h(x) = w^T x$$

$L_T(h)$

$$L_T(w) = \frac{1}{n} \sum_{i=1}^n \ell(h_w, (x_i, y_i))$$

$f(w, i)$

$F(w)$

$$\underset{w}{\text{minimize}} \quad F(w) = \frac{1}{n} \sum_{i=1}^n f(w, i)$$

$w_0$

$$\begin{aligned}
 w_{k+1} &= w_k - \alpha_k \nabla F(w_k) = \\
 &= w_k - \alpha_k \nabla \left( \frac{1}{n} \sum_{i=1}^n f(w_k, i) \right) = \\
 &= w_k - \alpha_k \frac{1}{n} \sum_{i=1}^n \nabla f(w_k, i)
 \end{aligned}$$

FULL-BATCH GRADIENT DESCENT

STOCHASTIC GRADIENT DESCENT

START FROM  $w_0$

FOR  $k = 0$  TO  $K-1$

- PICK A RANDOM POINT  $\xi_k$  IN THE DATASET
- COMPUTE

$$w_{k+1} = w_k - \alpha_k \nabla f(w_k, \xi_k)$$

FB

SGD

Minibatch

PICK  $w_0$

FOR  $k = 0$  TO  $K-1$

PICK  $m$  POINTS FROM  $T$  AT RANDOM  
 $\xi_{k,1}, \xi_{k,2}, \dots, \xi_{k,m}$

$$w_{k+1} = w_k - \alpha_k \frac{1}{m} \sum_{i=1}^m \nabla f(w_k, \xi_{k,i})$$

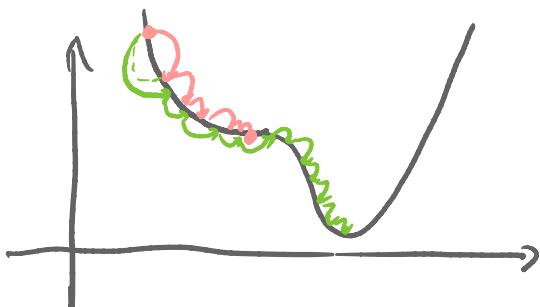
FB  $n \times$  MORE EXPENSIVE  
THAN SGD

IF YOUR GPU IS ABLE TO PROCESS  
B POINTS IN PARALLEL

THEN

FB IS ONLY  $\frac{n}{B}$  X MORE EXPENSIVE  
THAN SGD

GPU WITH  $B=n$



NOISE  
CAN HELP  
CONVERGENCE