

{ ECE - 6143 }  
 Lecture 3

✓ Recap: Linear regression

✗ Gradient descent

□ Application: linear regression (again) -

Recap: Data  $(x_i, y_i)$   $i=1, \dots, n$  -

- fit a linear model

$$y = w_0 + w_1 x$$

- Loss function

$$\text{MSE}(\omega) = \frac{1}{n} \sum_{i=1}^n [y_i - (w_0 + w_1 x_i)]^2$$

$$w_0^* = \bar{y} - w_1^* \bar{x}$$

$$w_1^* = \frac{\sum_i x_i y_i - \bar{x}^2}{\sum x_i^2 - \bar{x}^2}$$

Univariate  
Case

Multivariate

$$x = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \end{bmatrix}, y$$

$$\begin{bmatrix} \vdots \\ x^{(d)} \end{bmatrix}$$

## Linear model

$$y = w_0 + w_1 x^{(1)} + w_2 x^{(2)} + \dots + w_d x^{(d)}$$

↑                      ↑                      ↓  
 intercept            coefficients

Assumption 1 : Data & labels are zero mean.

⇒ Ignore  $w_0$ .

$$y = \sum_{i=1}^d w_i x_i$$

$w \cdot x$

Step 1

$$y = \langle w, x \rangle$$

Step 2

$$MSE = \frac{1}{2} \sum_{i=1}^n (y_i - \langle w, x_i \rangle)^2$$

Step 3

Minimize MSE with respect  
to  $w$ .

$$\begin{bmatrix} \langle w, x_1 \rangle \\ \langle w, x_2 \rangle \\ \vdots \end{bmatrix} = Xw = \begin{bmatrix} -x_1- \\ -x_2- \\ \vdots \\ -x_n- \end{bmatrix} \begin{bmatrix} w \end{bmatrix}$$

$X = \begin{bmatrix} -x_1- \\ -x_2- \\ \vdots \\ -x_n- \end{bmatrix}$

$\{w, \alpha_n\}$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$\text{MSE} = \frac{1}{2} \| y - Xw \|_2^2$$

Minimize with respect to  $w$ .

$$w^* = \left( \sum_{i=1}^n x_i x_i^T \right)^{-1} \sum_{i=1}^n y_i x_i$$

$$w^* = (X^T X)^{-1} X^T y$$

$\uparrow$        $\uparrow$        $\uparrow$        $\uparrow$   
 $d \times 1$      $d \times d$      $d \times n$      $n \times 1$

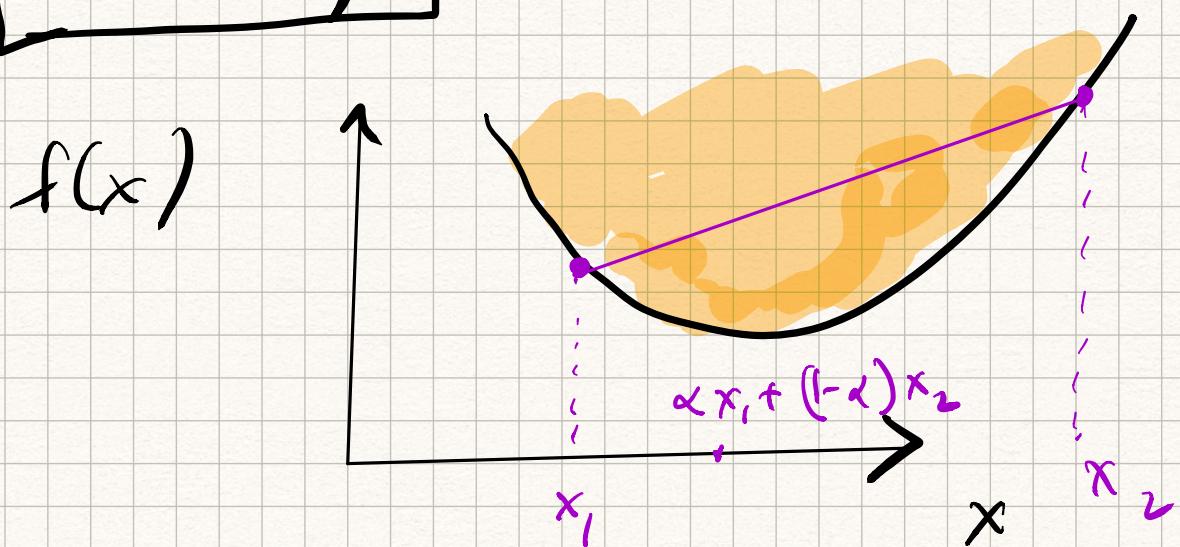
Multivariate  
case.

Bad! Inversion takes  $O(d^3)$  time.

Goal: Improve this running time

sufficiently. Using gradient descent.

# Convexity



words : Straight line b/w  $(x_1, f(x_1))$  &  $(x_2, f(x_2))$  should lie "above"  $f(x)$

---

Math :  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ . For any  $\underline{x}, \underline{y} \in \mathbb{R}^d$

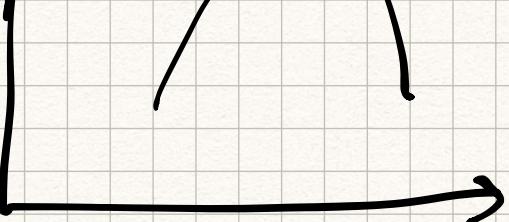
$$f(\alpha \underline{x} + (1-\alpha) \underline{y}) \leq \alpha f(\underline{x}) + (1-\alpha) f(\underline{y})$$

for all  $\alpha$  in  $[0, 1]$ .

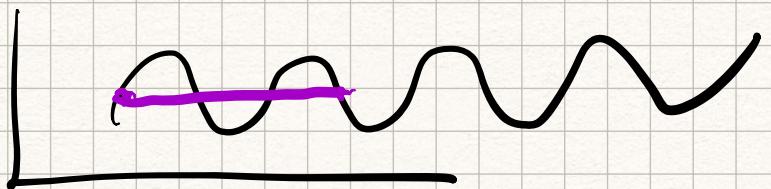
Then  $f$  is convex.

Opposite of convex: "concave"



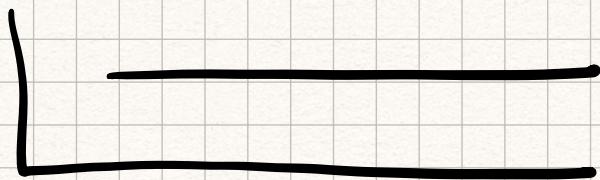


Sinusoid:



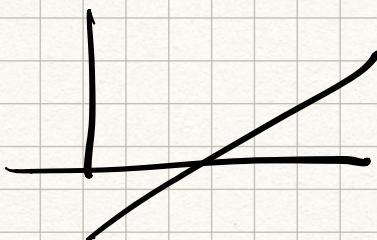
Neither concave nor convex.

Constant:



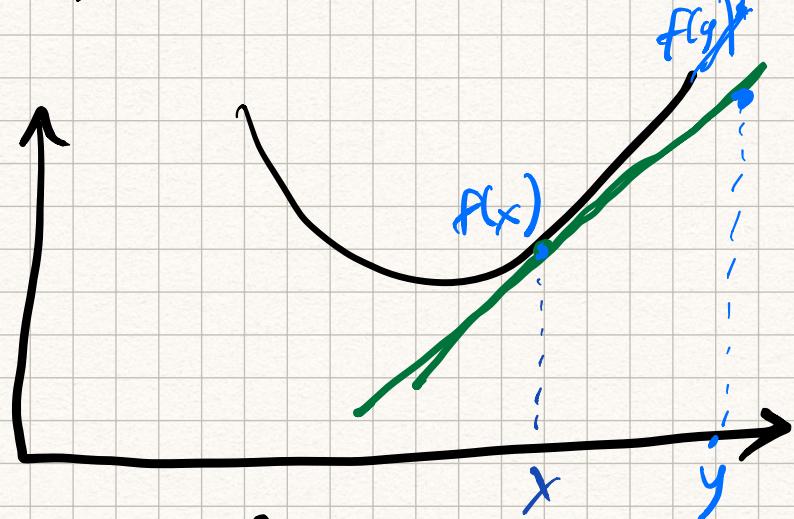
Both concave <sup>and</sup> ~~nor~~ convex

Straight



Both concave/  
convex

Property



Words

Tangent of  $f$  at any point  
 $x$  lies "below" the curve.

Math

For any  $x, y$ ,

Sobr  $f(y) \geq f(x) + \frac{df(x)}{dx} (y - x)$ .

Vectors

$$\begin{bmatrix} \frac{\partial f}{\partial x^{(1)}} \\ \frac{\partial f}{\partial x^{(2)}} \\ \vdots \\ \frac{\partial f}{\partial x^{(d)}} \end{bmatrix} = \nabla f(x)$$

For all  $x, y$ .

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle.$$

Observation

If  $\nabla f(x) = 0$ , then

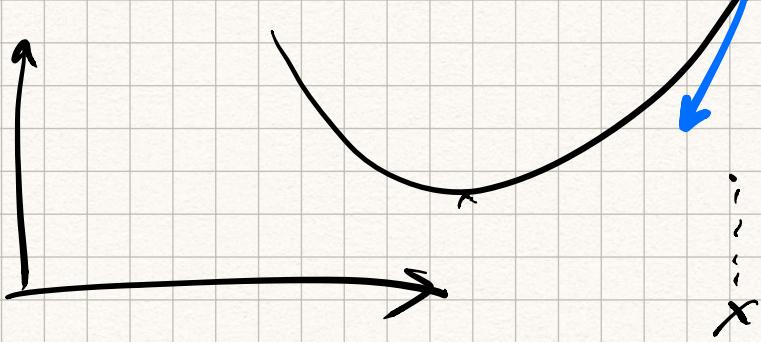
$f(y) \geq f(x)$  for all  $y$

i.e.  $x$  is a global minimum of  $f$ .

Caveat: Above discussion assumes  
differentiability / smoothness

Gradient descent | [GD]

$f(x)$



Goal: find  $x^*$  that minimizes  $f(x)$ .

Approach  
Find  $\Delta$

$$f(w + \Delta) < f(w).$$

- $w \leftarrow w + \Delta$
- Repeat

Principle of steepest descent:

$\nabla f(x)$  points in the direction of maximum change in  $f$ .

Mark

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

for  $k = 0, 1, 2, \dots, T$

"Step size"

"Learning rate"

"Number  
of  
epochs"

Terminate after  $T$  epochs

or when  $\sqrt{f}$  becomes small.

Applications of GD : Linear regression

$$L(\omega) = \frac{1}{2} \|y - X\omega\|_2^2$$

GD :

- Initialize  $\omega_0$  -

- Iterate -

$$\nabla L(\omega) = -X^T(y - X\omega)$$

$$\omega_{k+1} \leftarrow \omega_k - \alpha_k X^T(y - X\omega_k)$$

- Until  $T$  or  $\nabla L(\omega)$  becomes small.

Running time:  $O(nd) \cdot T$

Need to bound  $T$ , Set step size -

- Set error parameter  $\epsilon$ ,  
terminate when  $\|\nabla F(x)\|_2 \leq \epsilon$

- Can prove that this happens  
after  $T \geq \log \left[ \frac{1}{\epsilon} \|y - X\omega_0\|_2^2 \right]$

$$\log \gamma_p \left[ \frac{\lambda_{\min} - \lambda_2}{\Sigma} \right]$$

where  $\rho = \frac{1 - \frac{\lambda_1}{\lambda}}{1 + \frac{\lambda_1}{\lambda}}$

$\lambda_1, \lambda$  smallest & largest eigenvalues of  $X^T X$ .

Eigenvalue :

$$A v = \lambda v$$

$\uparrow$                      $\uparrow$   
eigenvector        eigenvalue