

## Lec 9 Gradient Descent

↳ How to minimize  $E_{in}(\underline{w}) = \frac{1}{N} \sum_{n=1}^N \log(1 + e^{-y_n \underline{w}^T \underline{x}_n})$ ?

### 1. Intro to gradient descent:

↳ Given differentiable function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ ,

we want to  $\min_{\underline{x} \in \mathbb{R}^n} f(\underline{x})$  // here  $\underline{x}$  is not data. Just a general term  
(unconstrained)

↳ Gradient descent is a numerical approach.

### ↳ example:

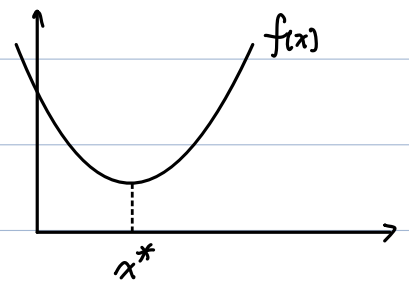
•  $x \in \mathbb{R}^n$ .  $n=1$  i.e.  $x \in \mathbb{R}$

// Assume  $f(x)$  is convex for now

• If  $x = x^*$ , then  $f'(x) = 0$

$x > x^*$ , then  $f'(x) > 0 \Rightarrow f(x)$  increases w.r.t  $x$

$x < x^*$ , then  $f'(x) < 0 \Rightarrow f(x)$  decreases w.r.t  $x$



### ↳ steps for $1D \in \mathbb{R}^1$ :

① initialize  $x = x_0$

② if  $f'(x) \approx 0$ , then stop

③ if  $f'(x) > 0$ , then  $x = x - \epsilon$ ; if  $f'(x) < 0$ , then  $x = x + \epsilon$

④ Go to step 2

//  $\epsilon$ : step size  $\left\{ \begin{array}{l} \text{large } \epsilon: \text{less accurate} \\ \text{small } \epsilon: \text{slow program} \end{array} \right.$

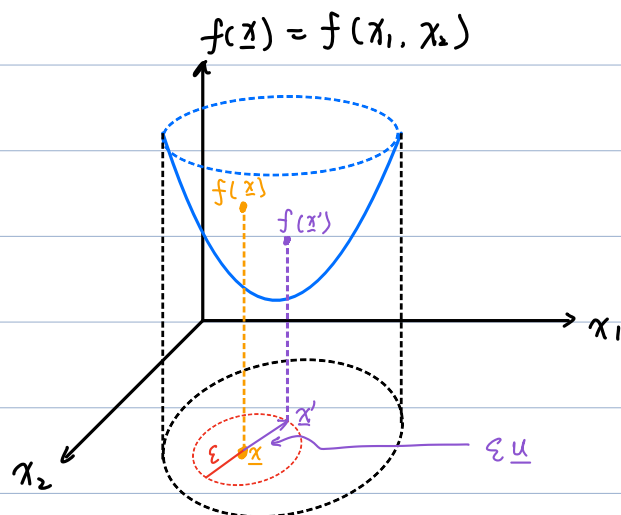
### ↳ setup for step size $\epsilon$ :

use  $\epsilon_k$ , where  $k$  is the iteration number.

e.g.  $\epsilon_k \sim \frac{1}{k}$

## 2. General Gradient descent:

↳ e.g.  $n=2$ :



Given current location  $\underline{x}$ , what's the next step to go?

$$\Rightarrow \underline{x}' \leftarrow \underline{x} + \epsilon \underline{u}$$

↳ unit vector, just for direction

↳ Idea: choose direction that maximize  $f(\underline{x}) - f(\underline{x}')$

$$\because f(\underline{x}') = f(\underline{x} + \epsilon \underline{u})$$

$$= f(\underline{x}) + \epsilon \underline{u}^T \nabla f(\underline{x}) + o(\epsilon^2)$$

// Taylor series expression

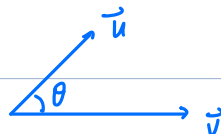
//  $o(\epsilon^2)$ : smaller term of  $\epsilon$ , negligible

$$\nabla f(\underline{x}) = \begin{bmatrix} \partial f / \partial x_1 \\ \vdots \\ \partial f / \partial x_n \end{bmatrix}$$

Pick  $\underline{u}$  to minimize  $\underline{u}^T \nabla f(\underline{x})$ , called  $\underline{u}^*$

// Recall:  $\underline{u}^T \underline{v} = \|\underline{u}\| \cdot \|\underline{v}\| \cos \theta$

$$\geq -\|\underline{u}\| \cdot \|\underline{v}\|$$



$$\therefore \underline{u}^* = - \frac{\nabla f(\underline{x})}{\|\nabla f(\underline{x})\|}$$

↑      ← normalize so that  $\|\underline{u}\| = 1$

\* direction of  $\underline{u}^*$  is the opposite of  $\nabla f(\underline{x})$

↳ Note:

- ①  $\nabla f(\underline{x})$  are points in the direction where  $f(\underline{x})$  has max ascent.
- ② Magnitude of  $\|f(\underline{x})\|$  in the denominator can be absorbed into  $\xi$ .

### 3. Gradient Descent Algorithm:

- Initialize  $\underline{x}_0$  (typically pick at random)
- For  $t = 0, 1, 2, \dots$ 
  - ① Compute  $\underline{g}_t = \nabla f(\underline{x}_t)$
  - ② Select direction  $\underline{v}_t = -\underline{g}_t$
  - ③ Update  $\underline{x}_{t+1} = \underline{x}_t + \xi_t \underline{v}_t$  // Normalization part absorbed.
  - ④ Go to step 1 until stopping criteria are reached

Note: • " $\xi_t$ " : learning rate

- Condition for stopping:  $\nabla f(\underline{x}_t) \approx 0$  is reasonable

### 4. GD of logistic regression

Recap:  $\underline{x}_{t+1} = \underline{x}_t - \xi_t \nabla f(\underline{x}_t)$   $\rightarrow \nabla_{\underline{x}} f(\underline{x})|_{\underline{x}=\underline{x}_t}$

e.g. logistic regression,

$$\begin{aligned} \text{how to minimize } E_{in}(\underline{w}) &= \frac{1}{N} \sum_{n=1}^N \log(1 + e^{-y_n \underline{w}^T \underline{x}_n}) ? \\ &= \frac{1}{N} \sum_{n=1}^N e_n(\underline{w}) \end{aligned}$$

↳ GD:  $\underline{w}_{k+1} = \underline{w}_k - \xi_k \nabla E_{in}(\underline{w}_k)$

$$= \underline{w}_k - \varepsilon_k \frac{1}{N} \sum_{n=1}^N \nabla \underline{e}_n(\underline{w}_k) ?$$

$$\nabla \underline{e}_n(\underline{w}_k) = \nabla_{\underline{w}} \log(1 + e^{-y_n \underline{w}^T \underline{x}_n})$$

$$= \frac{1}{1 + e^{-y_n \underline{w}^T \underline{x}_n}} \cdot \nabla_{\underline{w}} (1 + e^{-y_n \underline{w}^T \underline{x}_n})$$

$$= \frac{e^{-y_n \underline{w}^T \underline{x}_n}}{1 + e^{-y_n \underline{w}^T \underline{x}_n}} \cdot \nabla_{\underline{w}} (-y_n \underline{w}^T \underline{x}_n)$$

$\hookrightarrow = -y_n \cdot \underline{x}_n$

$$= \frac{-y_n \underline{x}_n}{1 + e^{y_n \underline{w}^T \underline{x}_n}}$$