CMPS 460 – Spring 2022

# MACHINE LEARNING

**Tamer Elsayed**
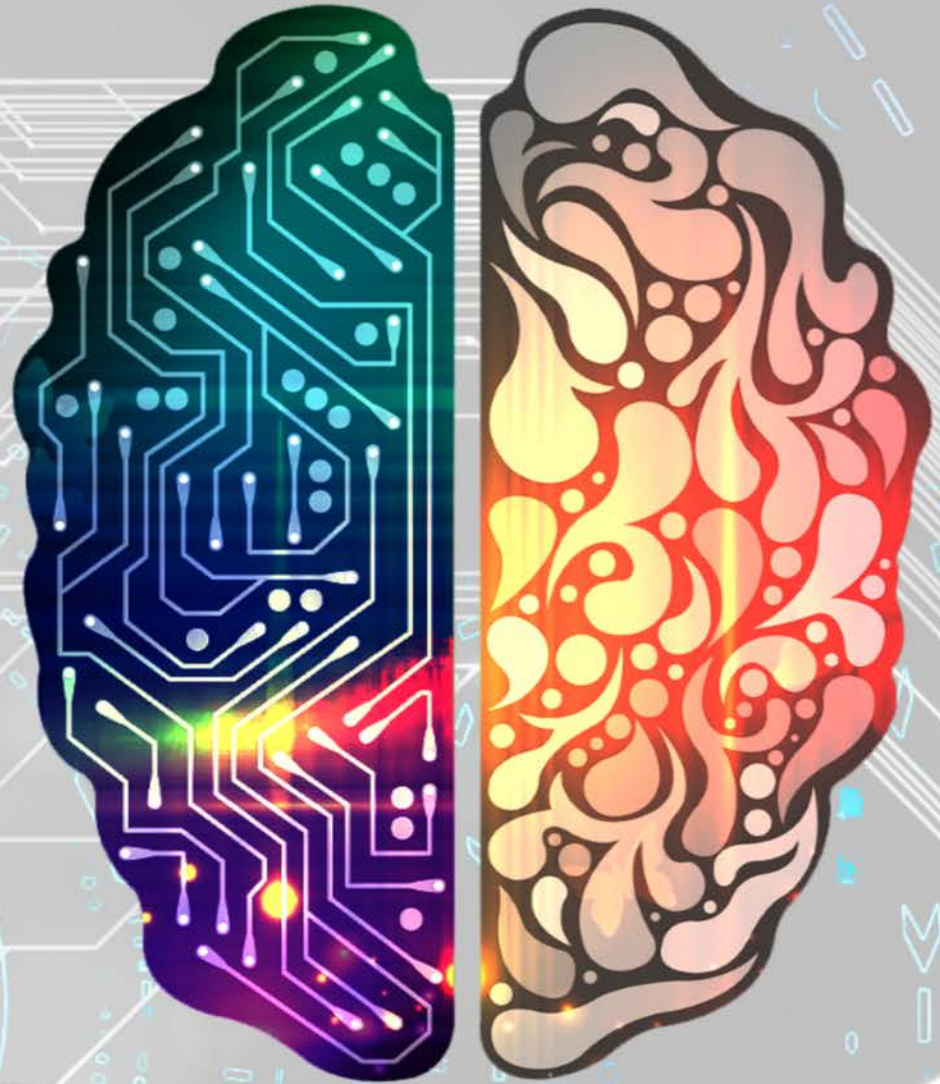
Image hosted by: WittySparks.com | Image source: Pixabay.com

**7.c**

# Linear Models: Gradient Descent

**Chapter 7: 7.4**

# Optimization Framework

**Objective function**

**Loss function** measures how well classifier fits training data

**Regularizer** prefers solutions that generalize well

$$\min_{\mathbf{w},b} L(\mathbf{w}, b) = \min_{\mathbf{w},b} \sum_{n=1}^{N} \mathbb{I}(y_n(\mathbf{w}^T x_n + b) < 0) + \lambda R(\mathbf{w}, b)$$

$l(y_n, \hat{y}_n)$

$\lambda$: parameter that controls the importance of the regularization term

- Different loss function approximations
  - easier to optimize

- Regularizer
  - prevents overfitting/prefers simple models.

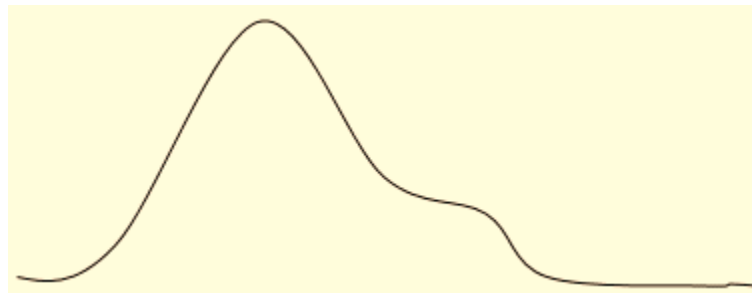# **Optimization with Gradient Descent**

# Gradient Descent

- A general solution for our optimization problem

$$\min_{\mathbf{w},b} L(\mathbf{w}, b) = \min_{\mathbf{w},b} \sum_{n=1}^{N} \mathbb{I}(y_n(\mathbf{w}^T \mathbf{x}_n + b) < 0) + \lambda R(\mathbf{w}, b)$$
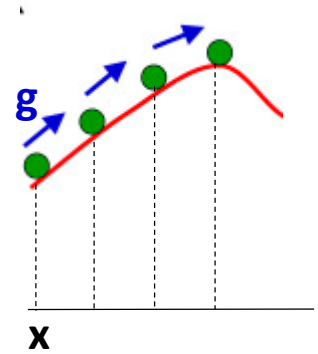
$$l(y_n, \hat{y}_n)$$

# Climbing a hill!

# Gradient-based Optimization

- Suppose the goal is to find maximum of f(**x**).

1. The optimizer maintains a current estimate of the parameter of interest, **x**.

2. At each step, measure the **gradient g** of f(**x**) at the current location, **x**.

3. Then take a step in the direction of the gradient, where the size of the step is controlled by $\eta$.

$$\textbf{x} \leftarrow \textbf{x} + \eta \, \textbf{g} \text{ (Gradient Ascent)}$$

- **Gradient Descent**: opposite of gradient ascent**.**

$$\textbf{x} \leftarrow \textbf{x} - \eta \, \textbf{g} \text{ (Gradient Descent)}$$

# Gradient Descent

- A general solution for our optimization problem

$$\min_{\mathbf{w},b} L(\mathbf{w}, b) = \min_{\mathbf{w},b} \sum_{n=1}^{N} \mathbb{I}(y_n(\mathbf{w}^T \mathbf{x}_n + b) < 0) + \lambda R(\mathbf{w}, b)$$

$$l(y_n, \hat{y}_n)$$

**_Idea_: _take iterative steps to update parameters in the direction of the gradient_**
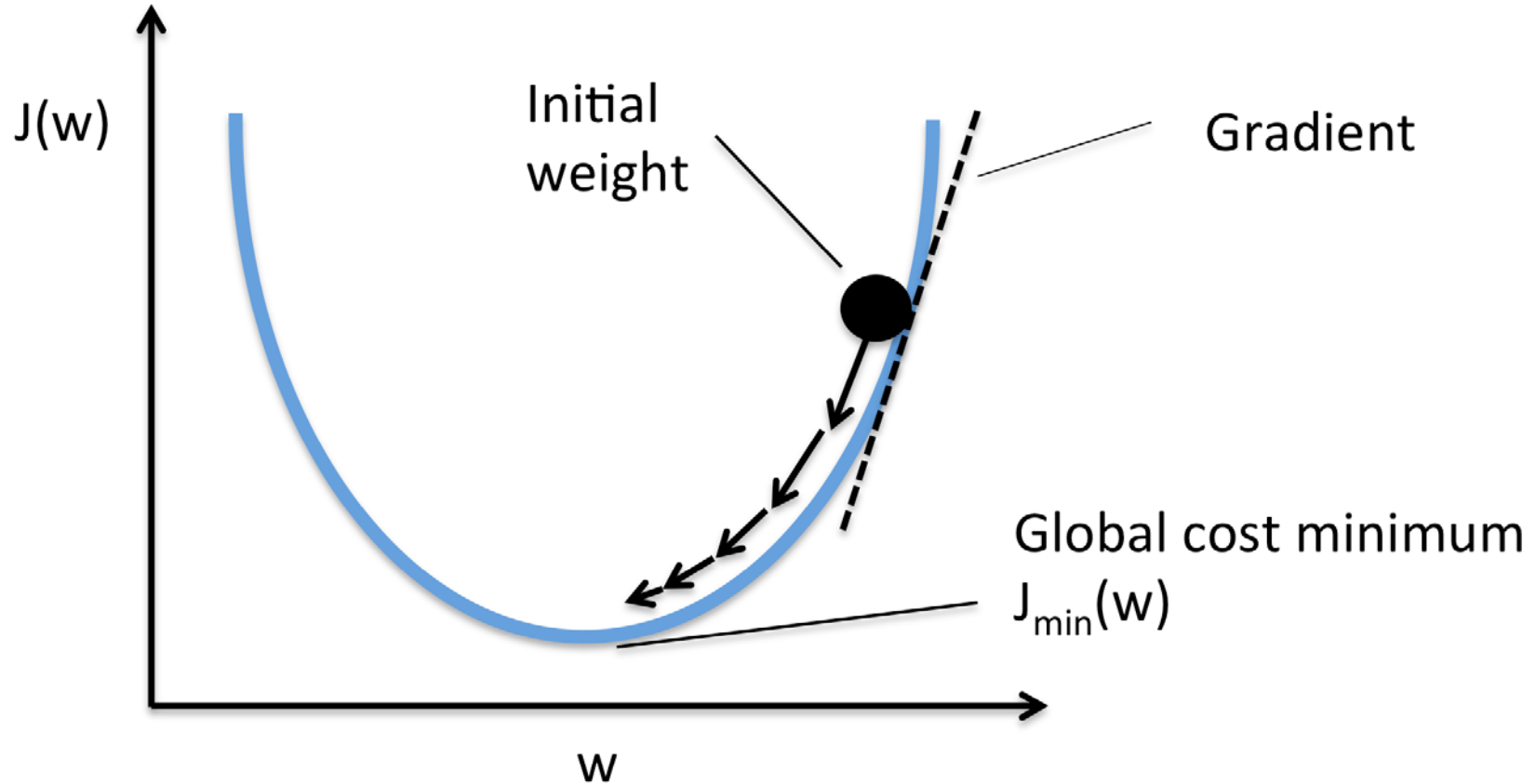
# Gradient Descent in 1-Dim



Image source: https://hackernoon.com/gradient-descent-aynk-7cbe95a778da

# Gradient Descent in 2-Dim



Image source: https://towardsdatascience.com/gradient- descent-in-a-nutshell-eaf8c18212f0
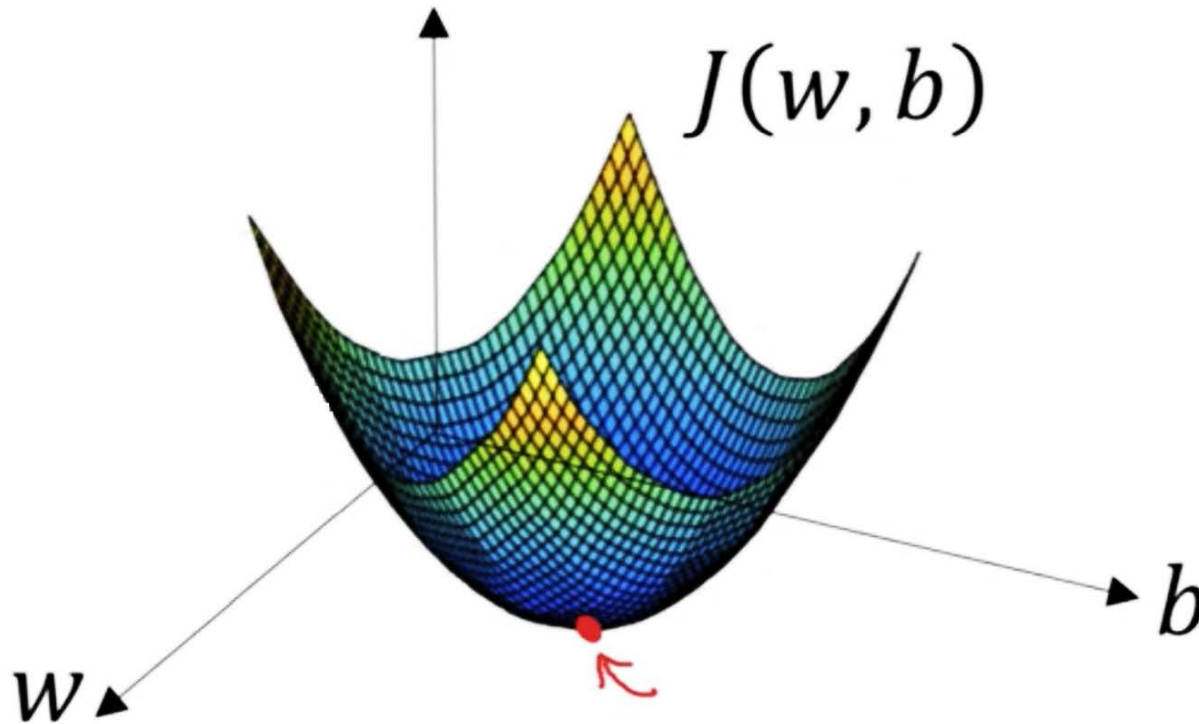
# Gradient Descent Algorithm

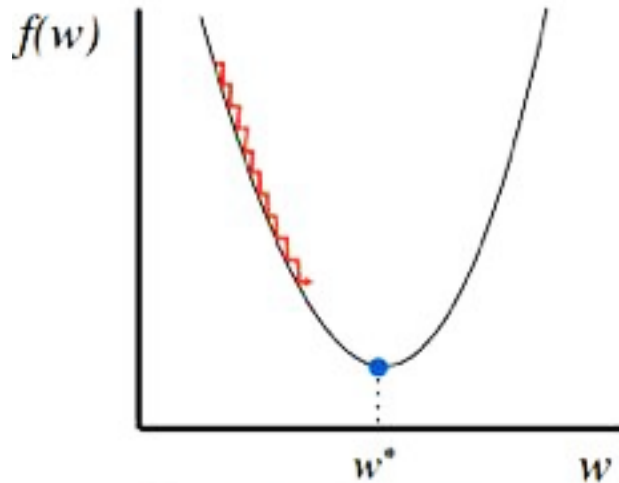Objective function to minimize

Number of steps

Step size

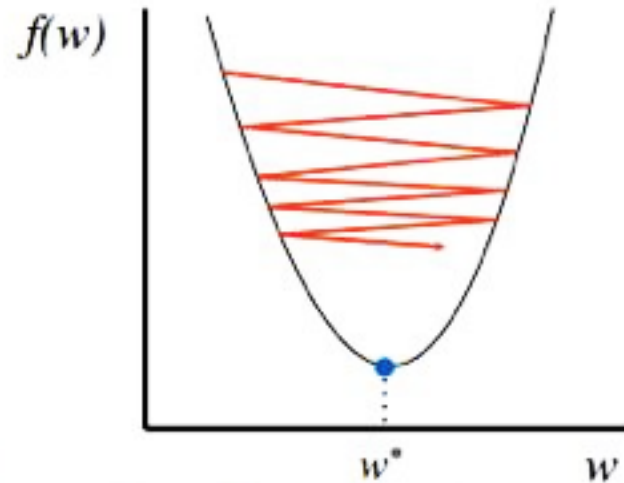**Algorithm 22** $\text{GRADIENTDESCENT}(\mathcal{F}, K, \eta_1, \dots)$

1: $z^{(0)} \leftarrow \langle 0, 0, \dots, 0 \rangle$    // initialize variable we are optimizing
2: **for** $k = 1 \dots K$ **do**
3:     $g^{(k)} \leftarrow \nabla_z \mathcal{F}\big|_{z^{(k-1)}}$    // compute gradient at current location
4:     $z^{(k)} \leftarrow z^{(k-1)} - \eta^{(k)} g^{(k)}$    // take a step down the gradient
5: **end for**
6: **return** $z^{(K)}$

# Impact of Step Size



Too small: converge very slowly

Too big: overshoot and even diverge

Image source: https://towardsdatascience.com/gradient- descent-in-a-nutshell-eaf8c18212f0
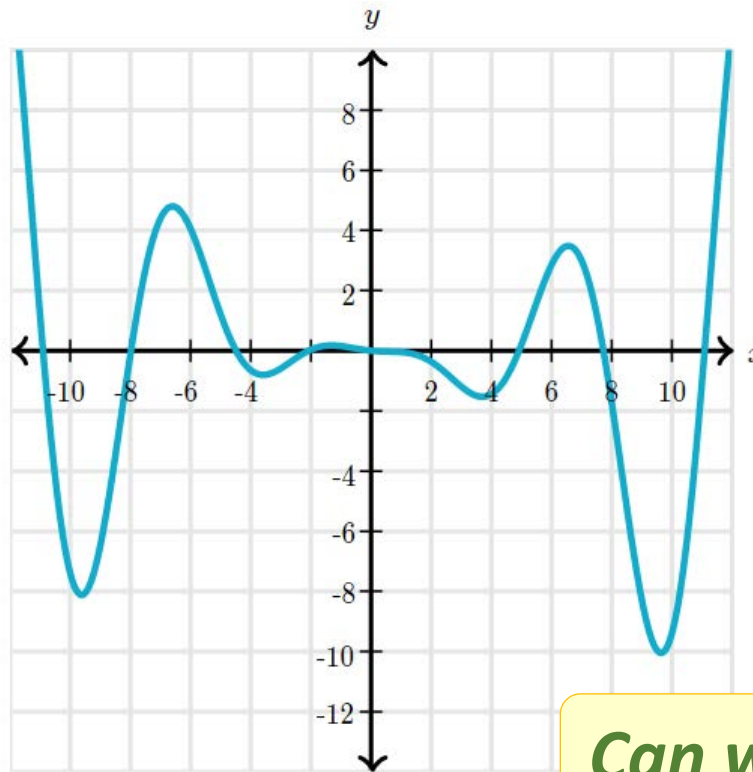
# Gradient Descent

**When to stop?**

- When the gradient gets close to zero
- When the objective stops changing much
- When the parameters stop changing much
- When performance on held-out dev set plateaus

**How to choose the step size?**

- Constant
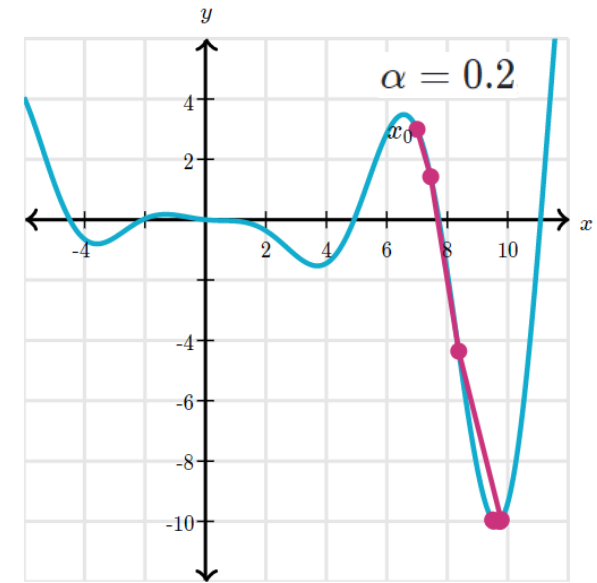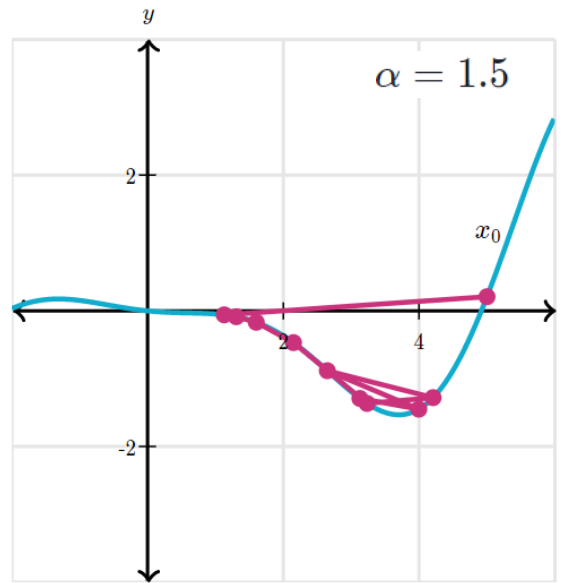- Better: start with large steps, then take smaller steps

# Example

$$f(x) = \frac{x^2 \cos(x) - x}{10}$$



**Can we apply GD?**

# Example

# Exercise

$$x^4 + 3x^3 - 2x$$

$$x_0 = 1, \propto = 0.2$$

$$x_0 = 1, \propto = 1.5$$

# Gradients for Multivariate Objectives

- Consider the following learning objective

$$\mathcal{L}(\boldsymbol{w}, b) = \sum_n \exp\left[-y_n(\boldsymbol{w} \cdot \boldsymbol{x}_n + b)\right] + \frac{\lambda}{2} \|\boldsymbol{w}\|^2$$

***What do we need to do to run gradient descent?***

# (1) Derivative with respect to $b$

$$\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial}{\partial b} \sum_n \exp\left[-y_n(\boldsymbol{w} \cdot \boldsymbol{x}_n + b)\right] + \frac{\partial}{\partial b} \frac{\lambda}{2} \|\boldsymbol{w}\|^2 \qquad (6.12)$$
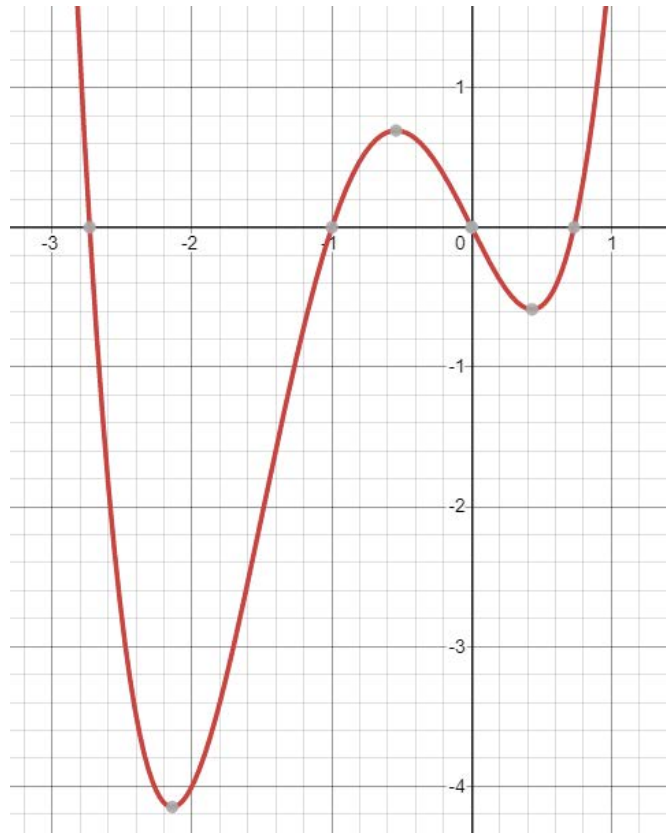
$$= \sum_n \frac{\partial}{\partial b} \exp\left[-y_n(\boldsymbol{w} \cdot \boldsymbol{x}_n + b)\right] + 0 \qquad (6.13)$$

$$= \sum_n \left(\frac{\partial}{\partial b} - y_n(\boldsymbol{w} \cdot \boldsymbol{x}_n + b)\right) \exp\left[-y_n(\boldsymbol{w} \cdot \boldsymbol{x}_n + b)\right] \qquad (6.14)$$

$$= -\sum_n y_n \exp\left[-y_n(\boldsymbol{w} \cdot \boldsymbol{x}_n + b)\right] \qquad (6.15)$$

# (2) Gradient with respect to $w$

$$\nabla_w \mathcal{L} = \nabla_w \sum_n \exp\left[-y_n(w \cdot x_n + b)\right] + \nabla_w \frac{\lambda}{2} \|w\|^2 \qquad (6.16)$$

$$= \sum_n \left(\nabla_w - y_n(w \cdot x_n + b)\right) \exp\left[-y_n(w \cdot x_n + b)\right] + \lambda w$$

$$\qquad (6.17)$$

$$= -\sum_n y_n x_n \exp\left[-y_n(w \cdot x_n + b)\right] + \lambda w \qquad (6.18)$$

# **Subgradients**

# Subgradients

- Problem: some objective functions are not differentiable everywhere
  - e.g., Hinge loss, $l_1$ norm

> **Solution: Subgradient Optimization**

- Let's ignore the problem, and just try to apply gradient descent anyway!!
- We will just differentiate by parts…

# Subgradients

- Generalizes the derivative to functions which are not differentiable.

- For any $x_0$ in the domain of the function, one can draw a line which goes through the point $(x_0, f(x_0))$ and which is everywhere either touching or below the graph of f.

- Set-valued

# Example: Subgradient of Hinge Loss

- For a given example n

$$\boldsymbol{\partial}_w \max\{0, 1 - y_n(\boldsymbol{w} \cdot \boldsymbol{x}_n + b)\} \tag{7.24}$$

$$= \boldsymbol{\partial}_w \begin{cases} 0 & \text{if } y_n(\boldsymbol{w} \cdot \boldsymbol{x}_n + b) > 1 \\ 1 - y_n(\boldsymbol{w} \cdot \boldsymbol{x}_n + b) & \text{otherwise} \end{cases} \tag{7.25}$$

$$= \begin{cases} \boldsymbol{\partial}_w 0 & \text{if } y_n(\boldsymbol{w} \cdot \boldsymbol{x}_n + b) > 1 \\ \boldsymbol{\partial}_w 1 - y_n(\boldsymbol{w} \cdot \boldsymbol{x}_n + b) & \text{otherwise} \end{cases} \tag{7.26}$$

$$= \begin{cases} \mathbf{0} & \text{if } y_n(\boldsymbol{w} \cdot \boldsymbol{x}_n + b) > 1 \\ -y_n\boldsymbol{x}_n & \text{otherwise} \end{cases} \tag{7.27}$$

# Subgradient Descent for Hinge Loss

**Algorithm 23** HINGEREGULARIZEDGD($\mathbf{D}, \lambda, MaxIter$)

1: $w \leftarrow \langle 0, 0, \ldots 0 \rangle$ , $b \leftarrow 0$       // initialize weights and bias
2: **for** $iter = 1 \ldots MaxIter$ **do**
3:     $g \leftarrow \langle 0, 0, \ldots 0 \rangle$ , $g \leftarrow 0$       // initialize gradient of weights and bias
4:     **for all** $(x,y) \in \mathbf{D}$ **do**
5:       **if** $y(w \cdot x + b) \leq 1$ **then**
6:         $g \leftarrow g + y\,x$       // update weight gradient
7:         $g \leftarrow g + y$       // update bias derivative
8:       **end if**
9:     **end for**
10:     $g \leftarrow g - \lambda w$       // add in regularization term
11:     $w \leftarrow w + \eta g$       // update weights
12:     $b \leftarrow b + \eta g$       // update bias
13: **end for**
14: **return** $w, b$

*Anything wrong here?*

# Perceptron ...

**Algorithm 5** PERCEPTRONTRAIN($\mathbf{D}$, *MaxIter*)

1: $w_d \leftarrow 0$, for all $d = 1 \ldots D$      // initialize weights
2: $b \leftarrow 0$      // initialize bias
3: **for** $iter = 1 \ldots MaxIter$ **do**
4:      **for all** $(x, y) \in \mathbf{D}$ **do**
5:          $a \leftarrow \sum_{d=1}^{D} w_d \, x_d + b$      // compute activation for this example
6:          **if** $ya \leq 0$ **then**
7:              $w_d \leftarrow w_d + yx_d$, for all $d = 1 \ldots D$      // update weights
8:              $b \leftarrow b + y$      // update bias
9:          **end if**
10:      **end for**
11: **end for**
12: **return** $w_0, w_1, \ldots, w_D, b$

*What is it optimizing?*

# Summary: Gradient Descent

- A generic algorithm to minimize objective functions

- Works well as long as functions are well behaved (i.e., convex)

- Subgradient descent can be used at points where derivative is not defined.


- Can we do better?
  - For some objectives, we can find closed form solutions (see CIML 7.6)