

Slides adapted from Prof Carpuat and Duraiswami



# CMSC 422 Introduction to Machine Learning

## **Lecture 14 (Sub)gradient Descent**

Furong Huang / [furongh@cs.umd.edu](mailto:furongh@cs.umd.edu)



UNIVERSITY OF  
MARYLAND

# Regrading deadline extended

---

March 29 11:00 am

Submit along with you exam sheets

Don't forget to put your name

# Recap: Linear Models

---

General framework for binary classification

Cast learning as optimization problem

Optimization objective combines 2 terms

loss function: measures how well classifier fits training data

Regularizer: measures how simple classifier is

- Does not assume data is linearly separable

Lets us separate model definition from training algorithm (**Gradient Descent**)

---

# Casting Linear Classification as an Optimization Problem

**Objective function**

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

**Regularizer**  
prefers solutions that generalize well

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

$\mathbb{I}(\cdot)$  Indicator function: 1 if  $(\cdot)$  is true, 0 otherwise

The loss function above is called the 0-1 loss



# Gradient descent

---

A general solution for our optimization problem

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

Idea: take iterative steps to update parameters in the direction of the gradient

# Gradient descent algorithm

Objective function  
to minimize

Number of  
steps

Step size

---

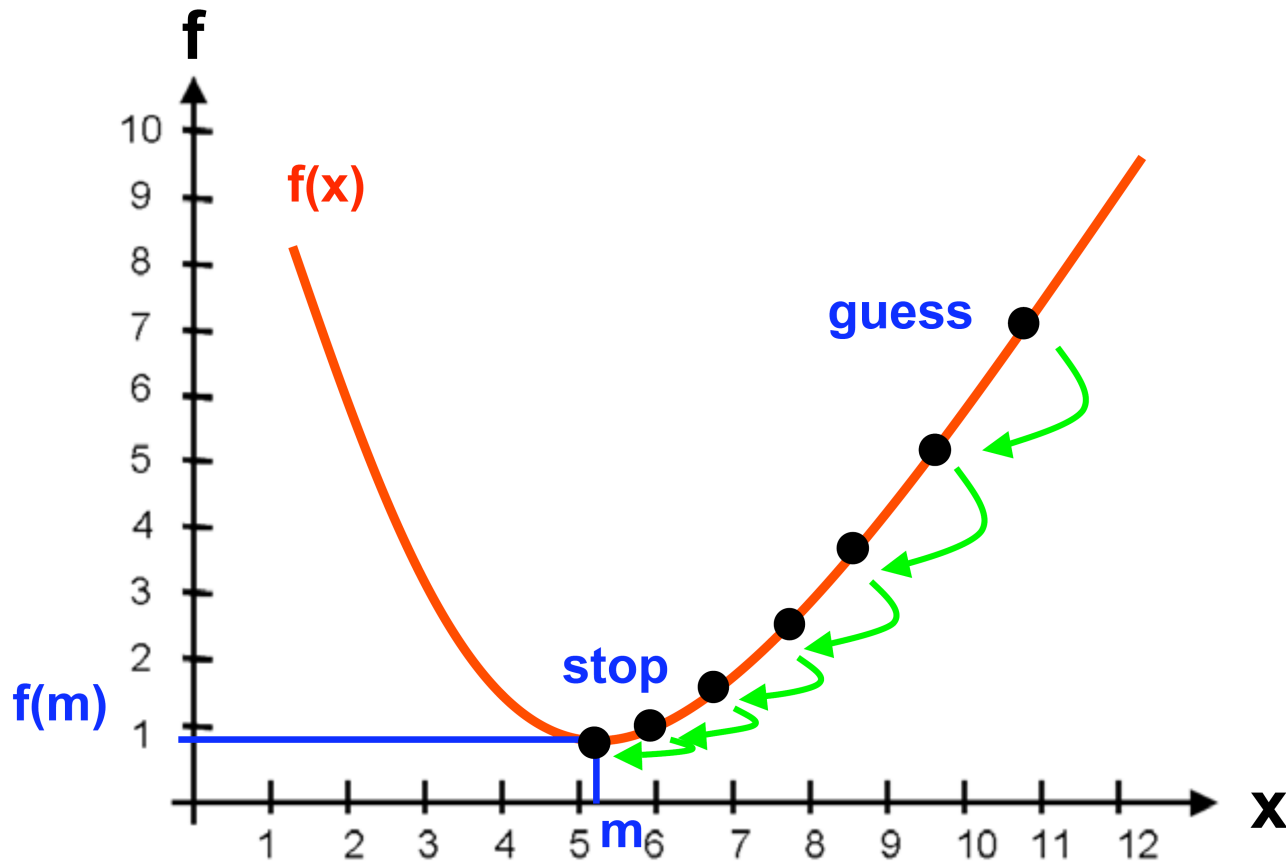
**Algorithm 22** GRADIENTDESCENT( $\mathcal{F}, K, \eta_1, \dots$ )

---

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

---

# Illustrating gradient descent in 1-dimensional case



# Gradient Descent

---

## 2 questions

### When to stop?

- When the gradient gets close to zero

- When the objective stops changing much

- When the parameters stop changing much

- Early

- When performance on held-out dev set plateaus

### How to choose the step size?

- Start with large steps, then take smaller steps

# Now let's calculate 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

---

$$\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial}{\partial b} \sum_n \exp [-y_n(\mathbf{w} \cdot \mathbf{x}_n + b)] + \frac{\partial}{\partial b} \frac{\lambda}{2} \|\mathbf{w}\|^2 \quad (6.12)$$

$$= \sum_n \frac{\partial}{\partial b} \exp [-y_n(\mathbf{w} \cdot \mathbf{x}_n + b)] + 0 \quad (6.13)$$

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

$$= - \sum_n y_n \exp [-y_n(\mathbf{w} \cdot \mathbf{x}_n + b)] \quad (6.15)$$

## (2) Gradient with respect to $w$

---

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

$$= \sum_n (\nabla_w - y_n(w \cdot x_n + b)) \exp [-y_n(w \cdot x_n + b)] + \lambda w \quad (6.17)$$

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

# Subgradients

---

Problem: some objective functions are not differentiable everywhere

Hinge loss,  $l_1$  norm

Solution: subgradient optimization

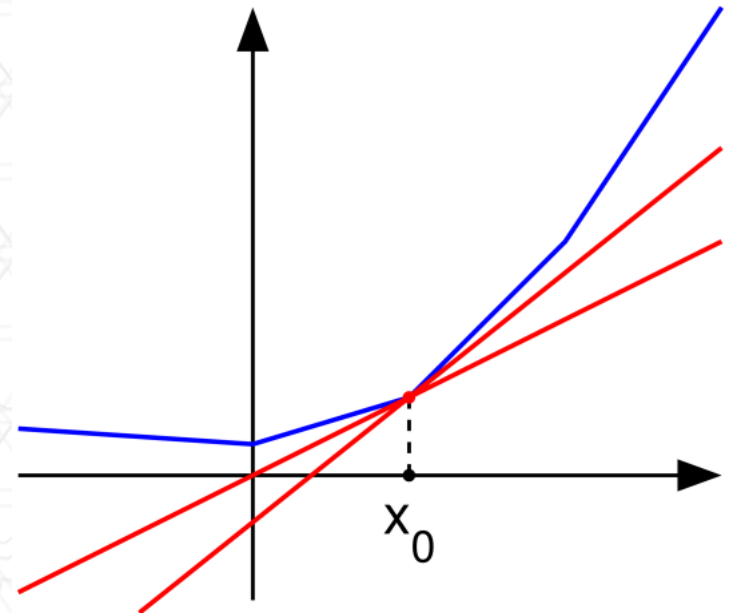
Let's ignore the problem, and just try to apply gradient descent anyway!!

we will just differentiate by parts...

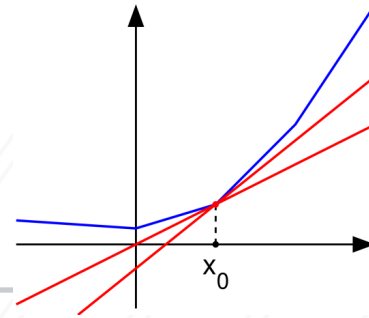


# Subgradient Review

- Subgradient generalized 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



# Subgradient Review



Rigorously, a *subgradient* of a convex function  $f: I \rightarrow \mathbf{R}$  at a point  $x_0$  in the open interval  $I$  is a real number  $c$  such that

$$f(x) - f(x_0) \geq c(x - x_0)$$

for all  $x$  in  $I$ . One may show that the set of subgradients at  $x_0$  for a convex function is a nonempty closed interval  $[a, b]$ , where  $a$  and  $b$  are the one-sided limits.

$$a = \lim_{x \rightarrow x_0^-} \frac{f(x) - f(x_0)}{x - x_0}$$

$$b = \lim_{x \rightarrow x_0^+} \frac{f(x) - f(x_0)}{x - x_0}$$

which are guaranteed to exist and satisfy  $a \leq b$ .

- The set  $[a, b]$  of all subgradients is called the subgradients of the function  $f$  at  $x_0$ . Since  $f$  is convex, if its subdifferential at  $x_0$  contains exactly one subgradient, then  $f$  is **differentiable** at  $x_0$ .

# Approximating the 0-1 loss with surrogate loss functions

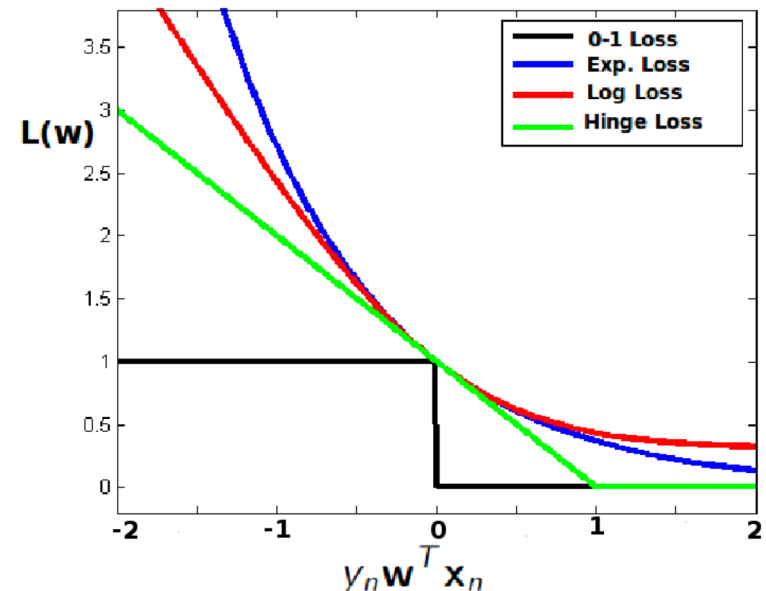
Examples (with  $b = 0$ )

Hinge loss  $[1 - y_n \mathbf{w}^T \mathbf{x}_n]_+ = \max\{0, 1 - y_n \mathbf{w}^T \mathbf{x}_n\}$

Log loss  $\log[1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n)]$

Exponential loss  $\exp(-y_n \mathbf{w}^T \mathbf{x}_n)$

What if  $b \neq 0$ ?



# Example: subgradient of hinge loss

---

For a given example  $n$

$$\partial_w \max\{0, 1 - y_n(\mathbf{w} \cdot \mathbf{x}_n + b)\} \quad (6.22)$$

$$= \partial_w \begin{cases} 0 & \text{if } y_n(\mathbf{w} \cdot \mathbf{x}_n + b) > 1 \\ y_n(\mathbf{w} \cdot \mathbf{x}_n + b) & \text{otherwise} \end{cases} \quad (6.23)$$

$$= \begin{cases} \partial_w 0 & \text{if } y_n(\mathbf{w} \cdot \mathbf{x}_n + b) > 1 \\ \partial_w y_n(\mathbf{w} \cdot \mathbf{x}_n + b) & \text{otherwise} \end{cases} \quad (6.24)$$

$$= \begin{cases} 0 & \text{if } y_n(\mathbf{w} \cdot \mathbf{x}_n + b) > 1 \\ y_n \mathbf{x}_n & \text{otherwise} \end{cases} \quad (6.25)$$

# Subgradient Descent for Hinge Loss

---

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

---

```
1:  $w \leftarrow \langle 0, 0, \dots, 0 \rangle$  ,  $b \leftarrow 0$  // initialize weights and bias
2: for  $iter = 1 \dots \text{MaxIter}$  do
3:    $g \leftarrow \langle 0, 0, \dots, 0 \rangle$  ,  $g \leftarrow 0$  // 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  $w, b$ 
```

---

# What is the perceptron optimizing?

---

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

---

```
1:  $w_d \leftarrow 0$ , for all  $d = 1 \dots D$  // initialize weights
2:  $b \leftarrow 0$  // initialize bias
3: for  $iter = 1 \dots MaxIter$  do
4:   for all  $(\mathbf{x}, y) \in \mathbf{D}$  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  $d = 1 \dots D$  // update weights
8:        $b \leftarrow b + y$  // update bias
9:     end if
10:   end for
11: end for
12: return  $w_0, w_1, \dots, w_D, b$ 
```

---

Loss function is a variant of the hinge loss

$$\max\{0, -y_n(\mathbf{w}^T \mathbf{x}_n + b)\}$$

# Recap: Linear Models

---

Lets us separate model definition from training algorithm (**Gradient Descent**)

# Summary

---

## Gradient descent

A generic algorithm to minimize objective functions

Works well as long as functions are well behaved (ie convex)

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

Choice of step size is important

## Optional: can we do better?

For some objectives, we can find closed form solutions (see CIML 6.6)





UNIVERSITY OF  
MARYLAND

**Furong Huang**

3251 A.V. Williams, College Park, MD 20740

301.405.8010 / [furongh@cs.umd.edu](mailto:furongh@cs.umd.edu)