

# What's the Perceptron Optimizing?

Machine Learning – CSE446

Carlos Guestrin

University of Washington

May 1, 2013

©Carlos Guestrin 2005-2013

1

## The Perceptron Algorithm

[Rosenblatt '58, '62]

- Classification setting:  $y$  in  $\{-1, +1\}$

- Linear model

- Prediction:  $\hat{y} = \text{Sign}(w \cdot x)$

- Training:  $w^{(0)} = 0$  or something smarter

- Initialize weight vector:

- At each time step:

- Observe features:  $x^{(t)} \leftarrow (\text{page}, \text{user}, \text{ad})$

- Make prediction:  $\hat{y} = \text{Sign}(w^{(t)} \cdot x^{(t)})$

- Observe true class:  $y^{(t)} \leftarrow \text{true label}$

- Update model:

- If prediction is not equal to truth

if  $\hat{y} \neq y^{(t)}$   
then  $w^{(t+1)} \leftarrow w^{(t)} + y^{(t)} x^{(t)}$   
else  $w^{(t+1)} \leftarrow w^{(t)}$

©Carlos Guestrin 2005-2013

2

I made a mistake:

e.g.  $y^{(t)} = +1$

$w^{(t)} \cdot x^{(t)} < 0$

but wanted  $> 0$

what to max  $w \cdot x^{(t)}$ ?

$x^{(t)}$  !!

by adding  $x^{(t)}$  to  $w$

I increase  $w^{(t+1)} \cdot x^{(t)}$

the most

if make a mistake!  
 $w^{(t+1)} \leftarrow w^{(t)} + \eta(\text{mistake})(y^{(t)} x^{(t)})$   
similarly when  $y^{(t)} = -1$

# What is the Perceptron Doing???

- When we discussed logistic regression:
  - Started from maximizing conditional log-likelihood

$$\max_w P(Y|X, w)$$

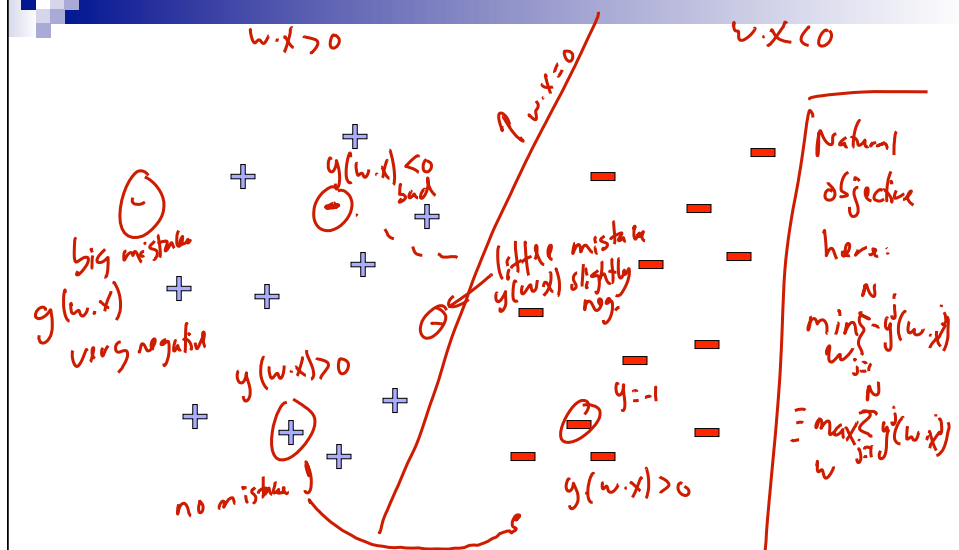
- When we discussed the Perceptron:
  - Started from description of an algorithm

- What is the Perceptron optimizing????

©Carlos Guestrin 2005-2013

3

## Perceptron Prediction: Margin of Confidence



©Carlos Guestrin 2005-2013

4

# Hinge Loss

- Perceptron prediction:  $\text{Sign}(w \cdot x)$

- Makes a mistake when:

$$y(w \cdot x) < 0 \Rightarrow$$

$\Rightarrow$

$$l(w, x) = \begin{cases} 0 & \text{if no mistake } y(w \cdot x) \geq 0 \\ -y(w \cdot x) & \text{otherwise } y(w \cdot x) < 0 \end{cases}$$

- Hinge loss (same as maximizing the margin used by SVMs)



©Carlos Guestrin 2005-2013

5

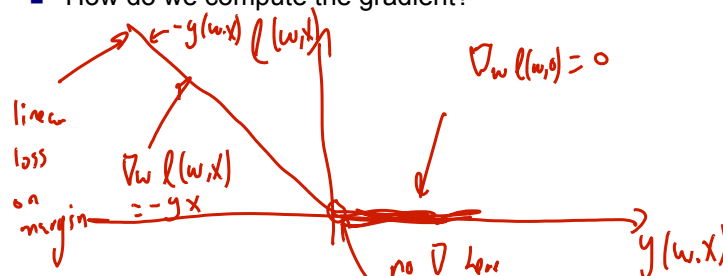
## Minimizing hinge loss in Batch Setting

- Given a dataset:  $(x^1, y^1) \dots (x^N, y^N)$

- Minimize average hinge loss:

$$\min_w \frac{1}{N} \sum_{i=1}^N l(w, x_i) = \begin{cases} 0 & \text{if } y(w \cdot x) \geq 0 \\ -y(w \cdot x) & \text{otherwise} \end{cases} \quad (-y(w \cdot x))_+$$

- How do we compute the gradient?

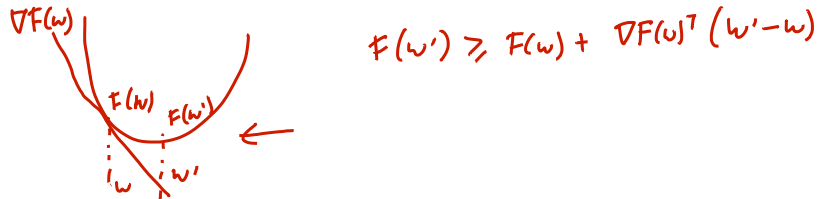


©Carlos Guestrin 2005-2013

6

# Subgradients of Convex Functions

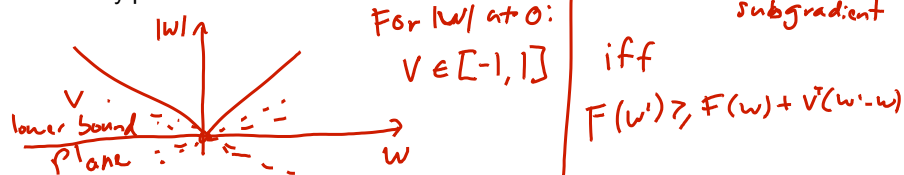
- Gradients lower bound convex functions:



- Gradients are unique at  $w$  iff function differentiable at  $w$

- Subgradients: Generalize gradients to non-differentiable points:

- Any plane that lower bounds function:



©Carlos Guestrin 2005-2013

7

# Subgradient of Hinge

- Hinge loss:



- Subgradient of hinge loss:

- If  $y^{(t)}(w, x^{(t)}) > 0$ :  $\partial l(w, x) = 0$
- If  $y^{(t)}(w, x^{(t)}) < 0$ :  $\partial l(w, x) = -y x$
- If  $y^{(t)}(w, x^{(t)}) = 0$ :  $\partial l(w, x) = [-y x, 0]$  e.g.,  $-y x$
- In one line:

$$\partial l(w, x) = \mathbb{1}(y(w, x) \leq 0) (-y x)$$

indicator of a mistake

©Carlos Guestrin 2005-2013

8

## Subgradient Descent for Hinge Minimization

- Given data:  $(x^1, y^1) \dots (x^N, y^N)$  I want  $\min_w$

- Want to minimize:

$$\frac{1}{N} \sum_{j=1}^N \ell(w, x^j) = \frac{1}{N} \sum_{j=1}^N (-y^j (w \cdot x^j))_+$$

- Subgradient descent works the same as gradient descent:

- But if there are multiple subgradients at a point, just pick (any) one:

$$w^{(t+1)} \leftarrow w^{(t)} - \eta \sum_{j=1}^N \partial \ell(w, x^j)$$

$\mathbb{1}_{\left\{ \begin{smallmatrix} \text{mistake} \\ y^j(w \cdot x^j) \leq 0 \end{smallmatrix} \right\}} (-y^j x^j)$

©Carlos Guestrin 2005-2013

9

## Perceptron Revisited

- Perceptron update:

$$\underline{w}^{(t+1)} \leftarrow w^{(t)} + \mathbb{1} \left[ y^{(t)} (w^{(t)} \cdot x^{(t)}) \leq 0 \right] y^{(t)} x^{(t)}$$

- Batch hinge minimization update:

$$w^{(t+1)} \leftarrow w^{(t)} + \eta \frac{1}{N} \sum_{i=1}^N \left\{ \mathbb{1} \left[ y^{(i)} (w^{(t)} \cdot x^{(i)}) \leq 0 \right] y^{(i)} x^{(i)} \right\}$$

- Difference?

Perceptron update is a stochastic gradient descent alg for hinge loss minimization with fixed step size ( $\eta=1$ )

©Carlos Guestrin 2005-2013

10

# What you need to know

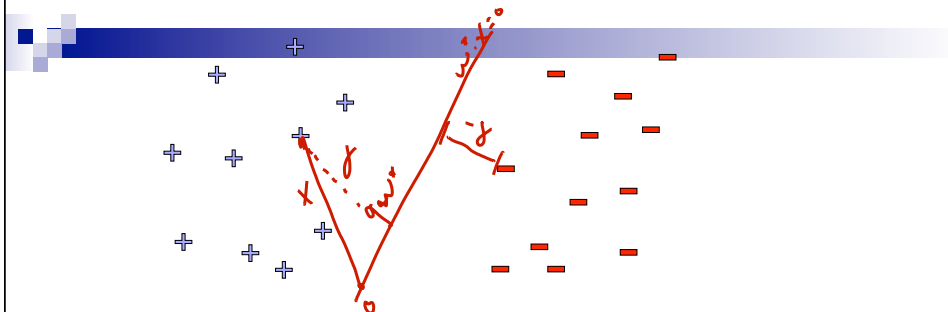
- Perceptron is optimizing hinge loss
- Subgradients and hinge loss
- (Sub)gradient decent for hinge objective

## Kernels

Machine Learning – CSE446  
Carlos Guestrin  
University of Washington

May 1, 2013

## Linear Separability: More formally, Using Margin



- Data linearly separable, if there exists

- a vector  $\exists w^*, \|w^*\|=1$

- a margin  $\gamma > 0$

- Such that

all points are  $\gamma$  far or more from  $w^* \cdot x = 0$

$\forall i$  if  $y^{(i)} = +1$   $w^* \cdot x^{(i)} > \gamma$

$y^{(i)} = -1$   $w^* \cdot x^{(i)} < -\gamma$

$y^{(i)} w^* \cdot x^{(i)} > \gamma$

linearly separable, margin  $\gamma$

©Carlos Guestrin 2005-2013

13

## Perceptron Analysis: Linearly Separable Case

- Theorem [Block, Novikoff]:

- Given a sequence of labeled examples:

$(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(T)}, y^{(T)})$

examples need not be i.i.d. or random...

- Each feature vector has bounded norm:

$$\forall t \quad \|x^{(t)}\| \leq R$$

$w^*$  is unknown!

- If dataset is linearly separable:

$$\exists w^*, \|w^*\|=1 \quad \forall t \quad y^{(t)} w^* \cdot x^{(t)} \geq \gamma, \text{ for } \gamma > 0$$

- Then the number of mistakes made by the online perceptron on this sequence is bounded by

$$\left(\frac{R}{\gamma}\right)^2$$

wow!!

constant, doesn't depend on  $T$

dimensionality of  $X$  !!

!

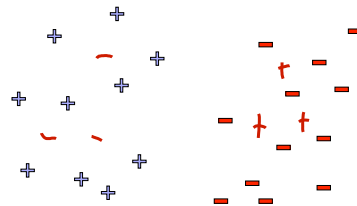
©Carlos Guestrin 2005-2013

14

# Beyond Linearly Separable Case

- Perceptron algorithm is super cool!
  - No assumption about data distribution!
    - Could be generated by an oblivious adversary, no need to be iid
  - Makes a fixed number of mistakes, and it's done for ever!
    - Even if you see infinite data

$$\left(\frac{R}{\gamma}\right)^2$$



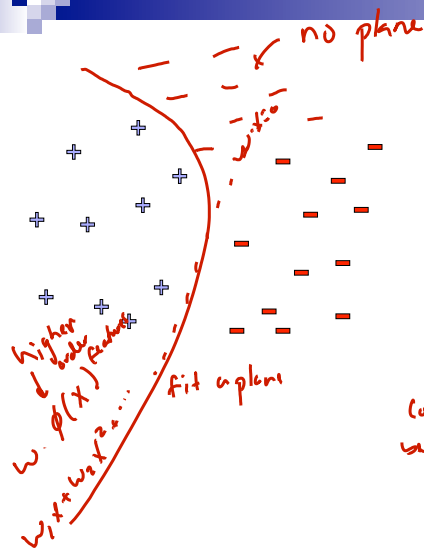
- However, real world not linearly separable
  - Can't expect never to make mistakes again
  - Analysis extends to non-linearly separable case
  - Very similar bound, see Freund & Schapire
  - Converges, but ultimately may not give good accuracy (make many many many mistakes)

we need features that make data as linearly separable as possible

©Carlos Guestrin 2005-2013

15

## What if the data is not linearly separable?



Use features of features of features of features....

$$\Phi(\mathbf{x}) : R^{n \times d} \mapsto F$$

$$\phi(x) = \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ x_1^2 \\ x_2^2 \\ x_1 x_2 \\ e^x \\ \sin x \\ -x^2 \\ \vdots \end{pmatrix}$$

Feature space can get really large really quickly!

©Carlos Guestrin 2005-2013

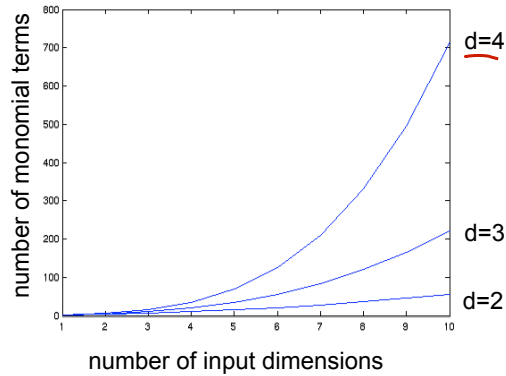
16



# Higher order polynomials

$$\text{num. terms} = \binom{d+m-1}{d} = \frac{(d+m-1)!}{d!(m-1)!}$$

m – input features  
d – degree of polynomial



Even though  
dims of  $\phi(x)$   
are huge, fit model  
~~very slowly~~  
very quickly  
grows fast!  
d = 6, m = 100  
about 1.6 billion terms

©Carlos Guestrin 2005-2013

17