

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

gradient

if make a mistake!

$w^{(t+1)} \leftarrow w^{(t)} + \mathbb{1}(\text{mistake}) (y^{(t)} x^{(t)})$

I made a mistake:
 e.g. $y^{(t)} = +1$
 $w^{(t)} \cdot x^{(t)} < 0$
 but wanted > 0
 what is max $w \cdot x^{(t)}$?
 $x^{(t)}$!!
 by adding $x^{(t)}$ to w
 I increase $w^{(t+1)} \cdot x^{(t)}$
 the most
 similarly when $y^{(t)} = -1$

©Carlos Guestrin 2005-2013

2

What is the Perceptron Doing???

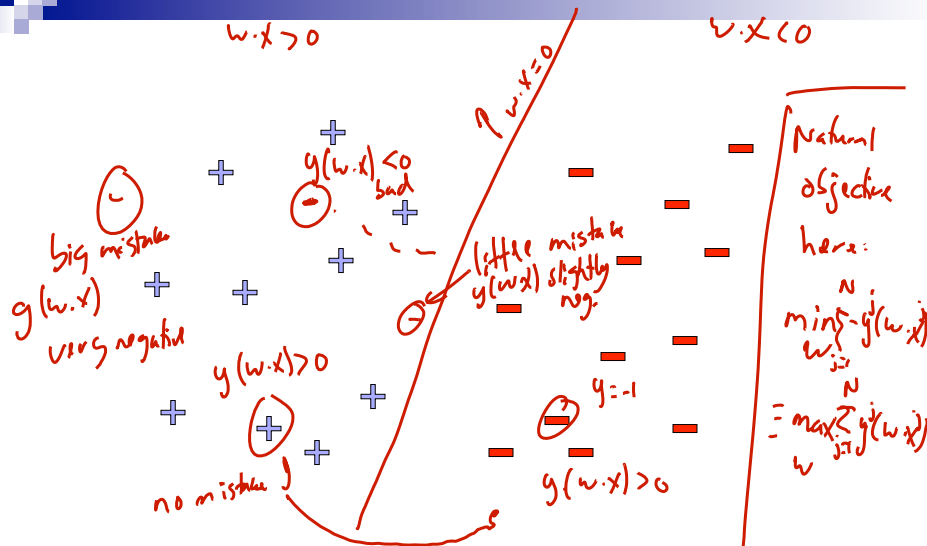
- When we discussed logistic regression:
 - Started from maximizing conditional log-likelihood

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

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

- What is the Perceptron optimizing????

Perceptron Prediction: Margin of Confidence



Hinge Loss

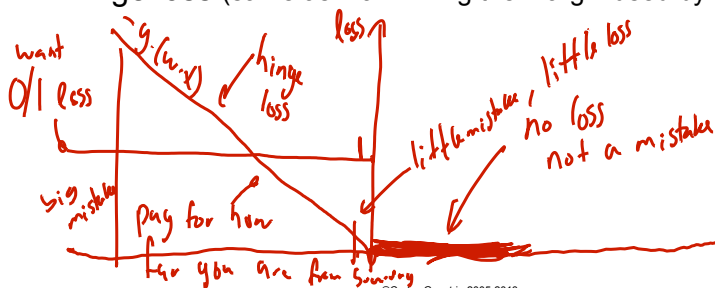
- Perceptron prediction: $\text{Sign}(w \cdot x)$

- Makes a mistake when:

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

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

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



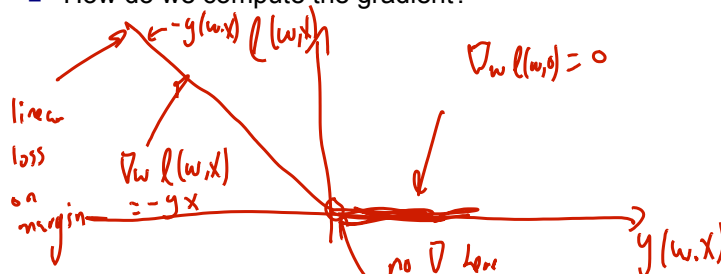
Minimizing hinge loss in Batch Setting

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

- Minimize average hinge loss:

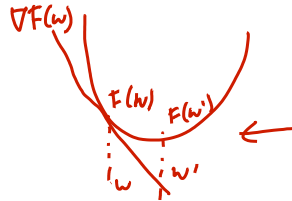
$$\min_w \frac{1}{N} \sum_{j=1}^N l(w, x_j) \equiv \begin{cases} 0 & \text{if } y(w \cdot x) \geq 0 \\ -y(w \cdot x) & \text{otherwise} \end{cases} \quad \left. \vphantom{\sum} \right\} (-y(w \cdot x))_+$$

- How do we compute the gradient?



Subgradients of Convex Functions

- Gradients lower bound convex functions:

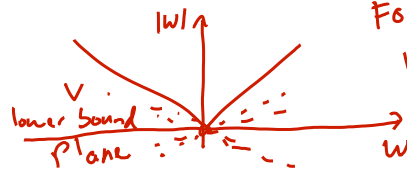


$$F(w') \geq F(w) + \nabla F(w)^T (w' - w)$$

- Gradients are unique at w iff function differentiable at w

- Subgradients: Generalize gradients to non-differentiable points:

- Any plane that lower bounds function:



For $w=0$ at 0 :
 $V \in [-1, 1]$

$V \in \partial_w F(w)$
 subgradient

iff
 $F(w') \geq F(w) + V^T (w' - w)$

©Carlos Guestrin 2005-2013

7

Subgradient of Hinge

- Hinge loss:



- Subgradient of hinge loss:

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

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

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:

$$w^{(t+1)} \leftarrow w^{(t)} - \eta \sum_{j=1}^N \underbrace{\partial \ell(w, x^j)}_{\mathbb{1}(\substack{\text{mistake} \\ (y^j (w \cdot x^j) \leq 0})} (-y^j x^j)}$$

©Carlos Guestrin 2005-2013

9

Perceptron Revisited

- Perceptron update: Subgradient

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

- Batch hinge minimization update: Subgradient

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

Sum over data points if mistake on point i
 param ... hopefully $\eta > 1$

- Difference?

Perceptron update is a stochastic gradient descent alg for hinge loss minimization with fixed step size ($\eta=2$)

©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 $\forall x \in \mathcal{X}$ if $y^{(i)} = +1$ $w^* \cdot x^{(i)} > \gamma$ and $y^{(i)} = -1$ $w^* \cdot x^{(i)} < -\gamma$
 - set γ for or more from $w^* \cdot x = 0$
 - } linearly separable; margin γ

©Carlos Guestrin 2005-2013 13

Perceptron Analysis: Linearly Separable Case

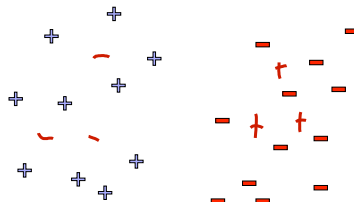
- Theorem [Block, Novikoff]:
 - Given a sequence of labeled examples: $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(T)}, y^{(T)})$
examples need not be i.i.d. or random...
 - Each feature vector has bounded norm: $\forall t \ \|x^{(t)}\| \leq R$ *w^* is unknown!*
 - If dataset is linearly separable: $\exists w^*, \|w^*\|=1 \ \forall t \ y^{(t)} w^* \cdot x^{(t)} \geq \gamma$, for $\gamma > 0$
- Then the number of mistakes made by the online perceptron on ~~the~~ ^{any} 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

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



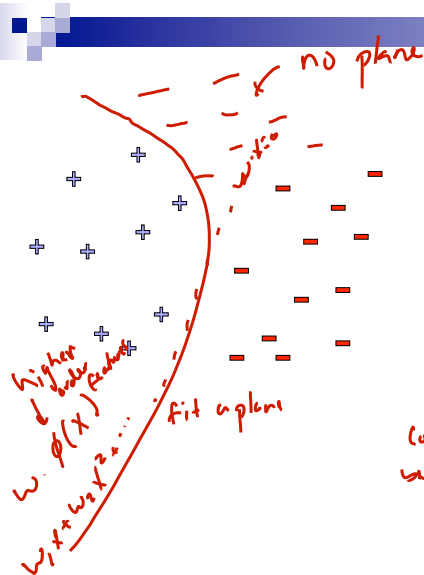
- However, real world not linearly separable
 - Can't expect never to make mistakes again
 - Analysis extends to non-linearly separable case
 - Very similar bound, see Freund & Schapire
 - Converges, but ultimately may not give good accuracy (make many many many mistakes)

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....

$\Phi(\mathbf{x}) : R^{k \text{ dim}} \mapsto F$

$\phi(x) = \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ x_3 \\ x_1^2 \\ x_2^2 \\ x_1 x_2 \\ e^x \\ \sin x \\ x_1^2 x_2 \\ \dots \end{pmatrix}$

could even be ∞ dim

spec here

Feature space can get really large really quickly!

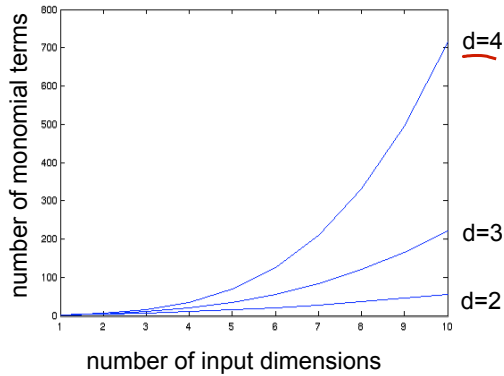
©Carlos Guestrin 2005-2013

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