

Online Learning Perceptron Algorithm

Machine Learning – CSE546
Carlos Guestrin (taught by Sameer)
University of Washington
October 23, 2014

©Carlos Guestrin 2005-2013

1

Challenge 1: Complexity of Computing Gradients

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \left\{ -\lambda w_i^{(t)} + \underbrace{\sum_j x_i^j [y^j - \hat{P}(Y^j = 1 | \mathbf{x}^j, \mathbf{w}^{(t)})]}_{\substack{\text{over data} \\ N}} \right\}$$

$\underbrace{\quad}_{\substack{\text{each feature} \\ K}}$

For each t , $O(NK)$

large datasets \rightarrow problem

SGD

©Carlos Guestrin 2005-2013

2

Challenge 2: Data is streaming

- Assumption thus far: **Batch data**

All the data points are available

- But, e.g., in click prediction for ads is a streaming data task:

- User enters query, and ad must be selected:
 - Observe x^t , and must predict y^t



$\vec{x}^t = \langle \text{query, user, ad} \rangle$
 \hat{y} = click or not
 $y = \begin{matrix} +1 & -1 \end{matrix}$

- User either clicks or doesn't click on ad:

- Label y^t is revealed afterwards
 - Google gets a reward if user clicks on ad

if ($y^t = \hat{y}$) Google becomes rich

$y^t = \{+1, -1\}$
if $\hat{y} \neq y^t$

- Weights must be updated for next time:

$w^{t+1} \leftarrow w^t + \Delta_t$

©Carlos Guestrin 2005-2013

3

Online Learning Problem

- At each time step t:

- Observe features of data point:

- Note: many assumptions are possible, e.g., data is iid, data is adversarially chosen... details beyond scope of course

$X_t = \langle \text{query, user, ad} \rangle \langle 1010111000 \dots \rangle$

- Make a prediction:

- Note: many models are possible, we focus on linear models
- For simplicity, use vector notation

$\hat{y} = \begin{cases} +1 & w^t \cdot x^t \geq 0 \\ -1 & w^t \cdot x^t < 0 \end{cases}$

$w^t \cdot x^t = \sum_{j=0}^k w_j^t x_j^t = w_0 + \sum_{j=1}^k w_j^t x_j^t$

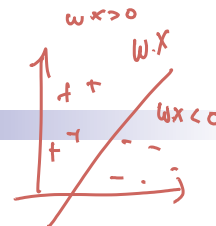
- Observe true label:

- Note: other observation models are possible, e.g., we don't observe the label directly, but only a noisy version... Details beyond scope of course

Observe $y^t = \begin{cases} +1 \\ -1 \end{cases}$ Mistake if $\hat{y} \neq y^t$

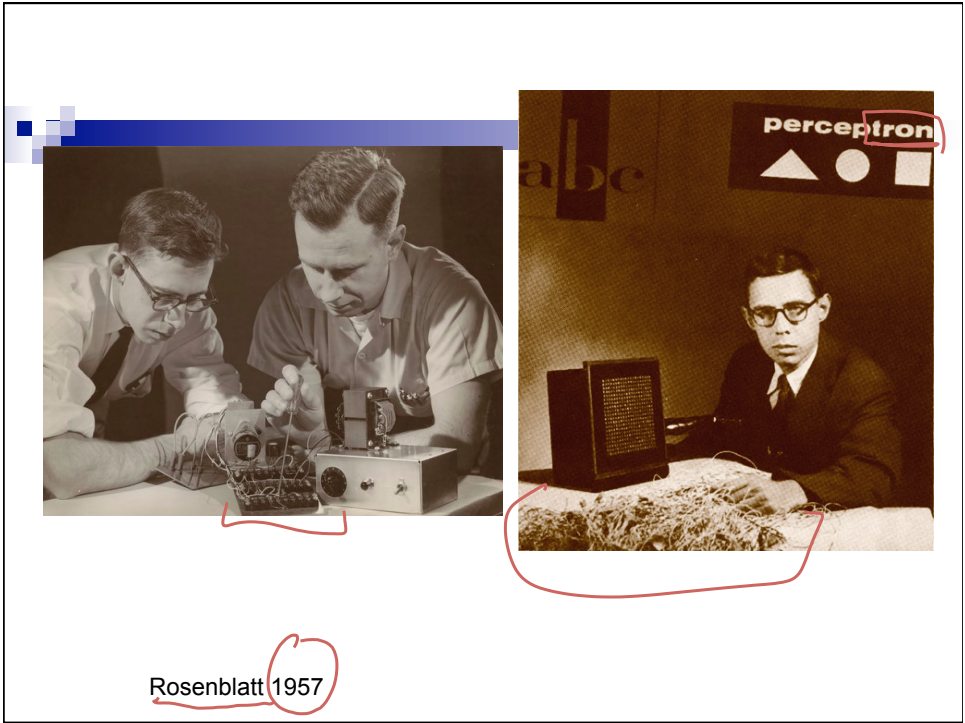
- Update model:

$w^{t+1} \leftarrow w^t + \Delta_t$



©Carlos Guestrin 2005-2013

4



The Perceptron Algorithm [Rosenblatt '58, '62]

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

- Linear model

□ Prediction: $\hat{y} = \text{sign}(w \cdot x)$

- Training:

□ Initialize weight vector: $w^0 = 0$ or random

□ At each time step:

- Observe features: x^t
- Make prediction: $\hat{y} = \text{sign}(w \cdot x)$
- Observe true class: $y^t \leftarrow +1$ or -1

- Update model:

□ If prediction is not equal to truth

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

$\text{sign}(w \cdot x) \neq y^t$
 if $y^t = 1$ but $\hat{y} = -1$
 $w^t \cdot x^t < 0$
 $\rightarrow y^t w^t \cdot x^t < 0$
 if $y^t = -1$ but $\hat{y} = 1$
 $w^t \cdot x^t \geq 0$
 $\rightarrow y^t w^t \cdot x^t \leq 0$
 Mistake when $y^t w^t \cdot x^t \leq 0$

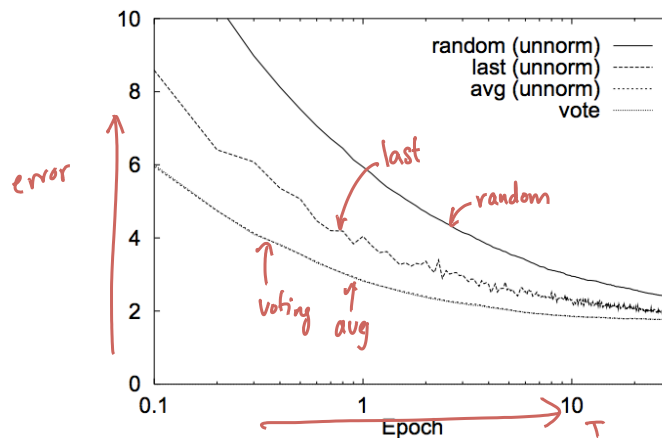
Fundamental Practical Problem for All Online Learning Methods: **Which weight vector to report?**

- Perceptron prediction: $\hat{y} = \text{sign}(w \cdot x)$
- Suppose you run online learning method and want to sell your learned weight vector... Which one do you sell???
- Last one? $\hat{w} \leftarrow w^T \leftarrow$ too noisy
- Random One - $\hat{w} \leftarrow w^t \leftarrow$ too noisy
- Average Weight $\hat{w} = \frac{1}{T+1} \sum_{t=0}^T w^t \leftarrow$ Good!
- Voting \leftarrow we won't cover

©Carlos Guestrin 2005-2013

7

Choice can make a huge difference!!



[Freund & Schapire '99]

©Carlos Guestrin 2005-2013

8

Mistake Bounds

- Algorithm "pays" every time it makes a mistake:

Google

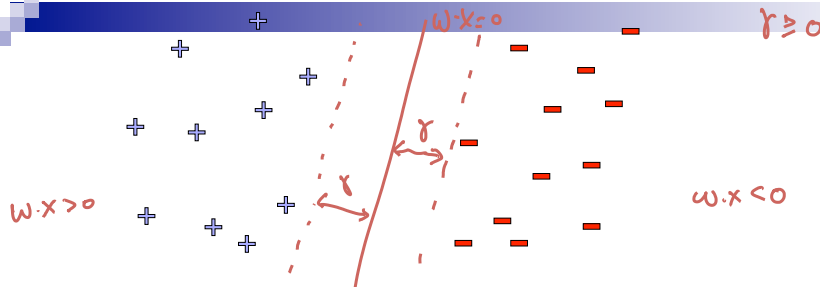
loss function

mistakes

- How many mistakes is it going to make?

Mistake Bound

Linear Separability: More formally, Using Margin



- Data linearly separable, if there exists

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

- a margin $\gamma > 0$

- Such that all points are at least γ away from $w^* x$

$$\forall t \quad \begin{array}{l} \text{if } y^t = +1 \quad w^* x^t \geq \gamma \\ y^t = -1 \quad w^* x^t \leq -\gamma \end{array} \quad \left| \quad \begin{array}{l} \text{Linear Sep:} \\ \forall t \quad y^t w^* x^t \geq \gamma \end{array} \right.$$

Perceptron Analysis: Linearly Separable Case

- Theorem [Block, Novikoff]:
 - Given a sequence of labeled examples: $(x^1, y^1), \dots, (x^T, y^T)$
 - Each feature vector has bounded norm: $\forall t \ \|x^t\| \leq R$
 - If dataset is linearly separable:
 - $\exists w^*, \|w^*\|=1 \ (y^t w^* x^t \geq \gamma) \ \forall t \ \gamma \geq 0$
- Then the number of mistakes made by the online perceptron on any such sequence is bounded by

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

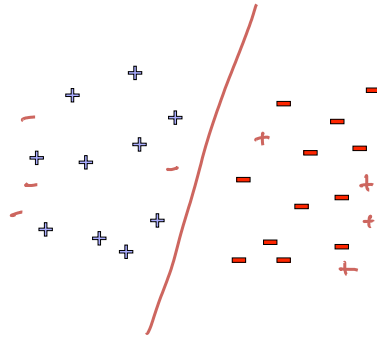
Doesn't depend on T
 Constant number of mistakes
 Independent of data size

Perceptron Proof for Linearly Separable case

- ✓ Every time we make a mistake, we get gamma closer to w^* :
 - Mistake at time t : $w^{(t+1)} = w^{(t)} + y^{(t)} x^{(t)}$
 - Taking dot product with w^* : $w^* \cdot w^{t+1} = w^* \cdot (w^t + y^t x^t) = w^* \cdot w^t + y^t w^* \cdot x^t$
 - Thus after m mistakes: $w^* \cdot w^{t+1} - (w^* \cdot w^0) \geq m \gamma$
- ✓ Similarly, norm of $w^{(t+1)}$ doesn't grow too fast:
 - $\|w^{(t+1)}\|^2 = \|w^{(t)}\|^2 + 2y^{(t)}(w^{(t)} \cdot x^{(t)}) + \|x^{(t)}\|^2$
 - Thus, after m mistakes: $\|w^{t+1}\|^2 - \|w^0\|^2 \leq m R^2$
- Putting all together:
 - $m \gamma \leq w^* \cdot w^{t+1} \leq (\|w^*\|) \|w^{t+1}\| \leq \sqrt{m} R$
 - $m \gamma \leq \sqrt{m} R \implies m \leq \left(\frac{R}{\gamma}\right)^2$

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
- 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)
(degree of non-linearity)



©Carlos Guestrin 2005-2013

13

What you need to know

- Notion of online learning → *Loss: # mistakes*
- Perceptron algorithm
- Mistake bounds and proof
- In online learning, report averaged weights at the end

©Carlos Guestrin 2005-2013

14

What's the Perceptron Optimizing?

Machine Learning – CSE546

Carlos Guestrin

University of Washington

October 23, 2013

©Carlos Guestrin 2005-2013

15

What is the Perceptron Doing???

- When we discussed logistic regression:

- Started from maximizing conditional log-likelihood

$$\max_{\omega} \ell(y, X, \omega) = \max_{\omega} \log \prod_{i=1}^n P(y_i | x_i, \omega) \rightarrow \text{Gradient} \\ \rightarrow \text{LR algorithm}$$

- When we discussed the Perceptron:

- Started from description of an algorithm

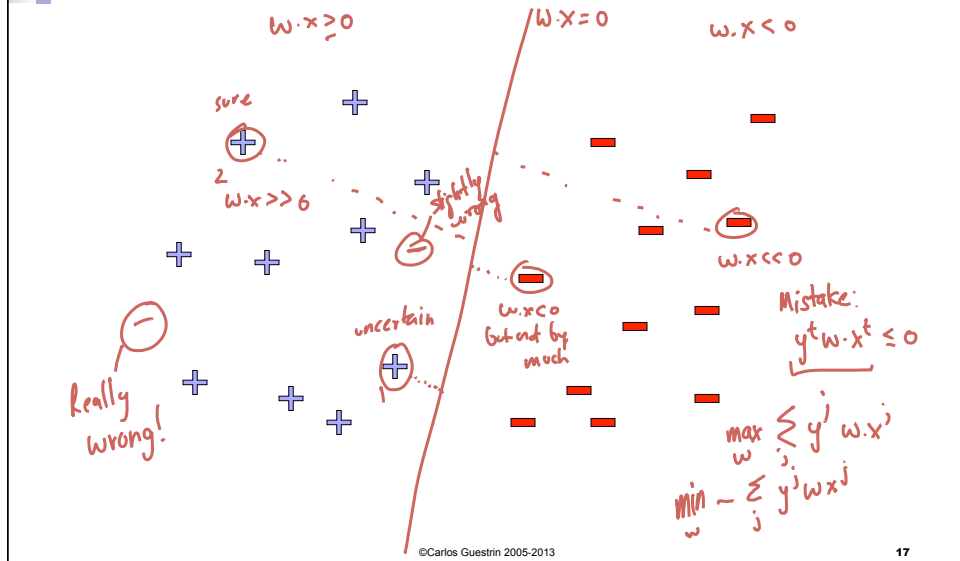
- What is the Perceptron (optimizing)???

$$\min_{\omega} \sum_t \text{loss}(y^t, x^t, \omega)$$

©Carlos Guestrin 2005-2013

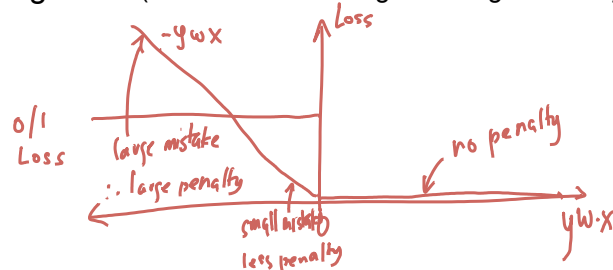
16

Perceptron Prediction: Margin of Confidence



Hinge Loss

- Perceptron prediction: $\text{sign}(w \cdot x)$
- Makes a mistake when: $yw \cdot x \leq 0$
- Hinge loss (same as maximizing the margin used by SVMs)



©Carlos Guestrin 2005-2013

18

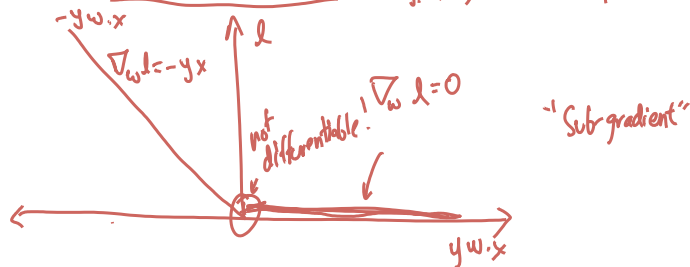
Minimizing hinge loss in Batch Setting

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

- Minimize average hinge loss:

$$\min_{\omega} \frac{1}{N} \sum_{j=1}^N l(y^j, x^j, \omega) \quad \begin{cases} 0 & \text{if } y^j \omega x^j \geq 0 \\ -y^j \omega x^j & \text{if } y^j \omega x^j < 0 \end{cases}$$

- How do we compute the gradient? $l(y^j, x^j, \omega) = (-y^j \omega x^j)_+ \quad (a)_+ = \begin{cases} a & \text{if } a > 0 \\ 0 & \text{o.w.} \end{cases}$

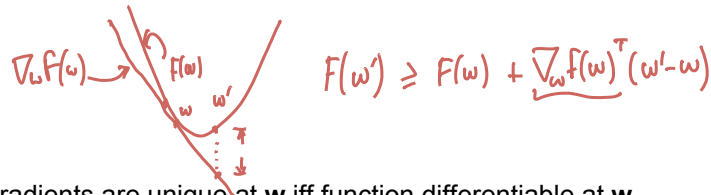


©Carlos Guestrin 2005-2013

19

Subgradients of Convex Functions

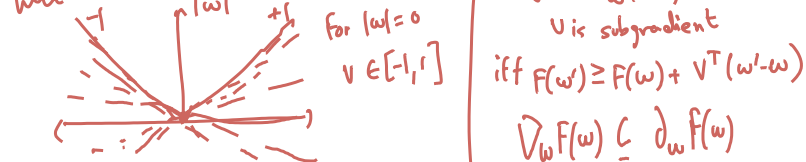
- Gradients lower bound convex functions:



- Gradients are unique at w iff function differentiable at w

- Subgradients: Generalize (gradients) to non-differentiable points:

- Any plane that lower bounds function:

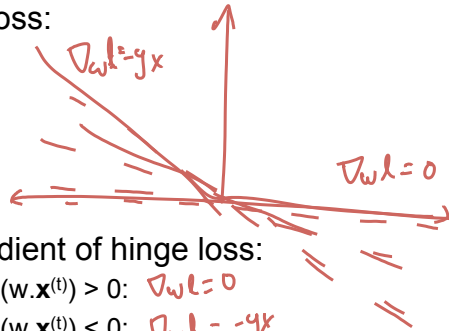


©Carlos Guestrin 2005-2013

20

Subgradient of Hinge

- Hinge loss:



- Subgradient of hinge loss:

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

$$\partial_w l(w, x, y) = \underbrace{\mathbb{1}(y w \cdot x \leq 0)}_{\text{mistake}} (-y \cdot x)$$

©Carlos Guestrin 2005-2013

21

Subgradient Descent for Hinge Minimization

- Given data:

- Want to minimize:

- Subgradient descent works the same as gradient descent:

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

©Carlos Guestrin 2005-2013

22

Perceptron Revisited

- Perceptron update:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \underbrace{\mathbb{1} \left[y^{(t)} (\mathbf{w}^{(t)} \cdot \mathbf{x}^{(t)}) \leq 0 \right]}_{\text{mistake}} \underbrace{y^{(t)} \mathbf{x}^{(t)}}_{\text{add } yx}$$

no step size ($\epsilon=1$)

- Batch hinge minimization update:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \underbrace{\eta}_{\text{all data}} \frac{1}{N} \sum_{i=1}^N \left\{ \underbrace{\mathbb{1} \left[y^{(i)} (\mathbf{w}^{(t)} \cdot \mathbf{x}^{(i)}) \leq 0 \right]}_{\text{mistake}} \underbrace{y^{(i)} \mathbf{x}^{(i)}}_{\text{add } yx} \right\}$$

- Difference?

Perceptron is hinge loss minimization using SGD and $\eta=1$

©Carlos Guestrin 2005-2013

23

What you need to know

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

©Carlos Guestrin 2005-2013

24