

Online Learning Perceptron Algorithm

Machine Learning – CSE446

Carlos Guestrin

University of Washington

April 29, 2013

©Carlos Guestrin 2005-2013

1

Challenge 1: Complexity of Computing Gradients in LR

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \left\{ -\lambda w_i^{(t)} + \sum_{j=1}^N x_i^j [y^j - \hat{P}(Y^j = 1 \mid \mathbf{x}^j, \mathbf{w}^{(t)})] \right\}$$

$O(Nk)$

if lots of data ... N is very large

→ very slow.

We talked about SGD instead

small change after each data point

©Carlos Guestrin 2005-2013

2

Challenge 2: Data is streaming

- Assumption thus far: **Batch data**

Have data, will machine learn

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

- User enters query, and ad must be selected:

- Observe x^i , and must predict y^i

*Q → [] → x^i → predict \hat{y} will the user click
 webpage (page, user, ad) → which ads have high click prob.*

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

- Label y^i is revealed afterwards

- Google gets a reward if user clicks on ad

reward is like $0/1$ loss in classification

- Weights must be updated for next time:

$w^{(t+1)} \leftarrow w^{(t)} + \Delta$ what's Δ ?

©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

$x^{(t)} \in (\text{page}, \text{user}, \text{ad})$

- Make a prediction: \hat{y} *hopefully $\approx y^{(t)} \leftarrow \text{unknown}$*

- Note: many models are possible, we focus on linear models

- For simplicity, use vector notation

*$w_0 + \sum_{i=1}^n w_i x_i^{(t)} > 0?$
 \Rightarrow click*

*$w^{(t)} \cdot x^{(t)} > 0?$
 $\sum_{i=1}^n w_i x_i^{(t)} > 0?$*

$x^{(t)} = \begin{pmatrix} 1 \\ \text{page features} \\ \text{user features} \\ \text{ad features} \end{pmatrix}$

$x_0^{(t)} = 1 \forall t$

- 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 $y^{(t)} \rightarrow$ clicked
 \rightarrow not clicked*

what's Δ ?

- Update model:

*$w^{(t+1)} \leftarrow w^{(t)} + \Delta^{(t)}$ \checkmark
 something*

©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: $w^{(0)} = 0$ 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 $\hat{y} \neq y^{(t)}$
 $w^{(t+1)} \leftarrow w^{(t)}$
 else
 $w^{(t+1)} \leftarrow w^{(t)} + y^{(t)} x^{(t)}$

I made a mistake:
 e.g. $y^{(t)} = +1$
 $w^{(t)} \cdot x^{(t)} < 0$
 but wanted > 0
 what u max $w \cdot x$?
 $x^{(t)}$!!
 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

5

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

- Suppose you run online learning method and want to sell your learned weight vector... Which one do you sell???

- Last one? $w^{(T)}$? No... very noisy, influenced by last mistake

- Random time step? $\frac{1}{T} \sum_{t=1}^T w^{(t)}$ (easier to keep track)

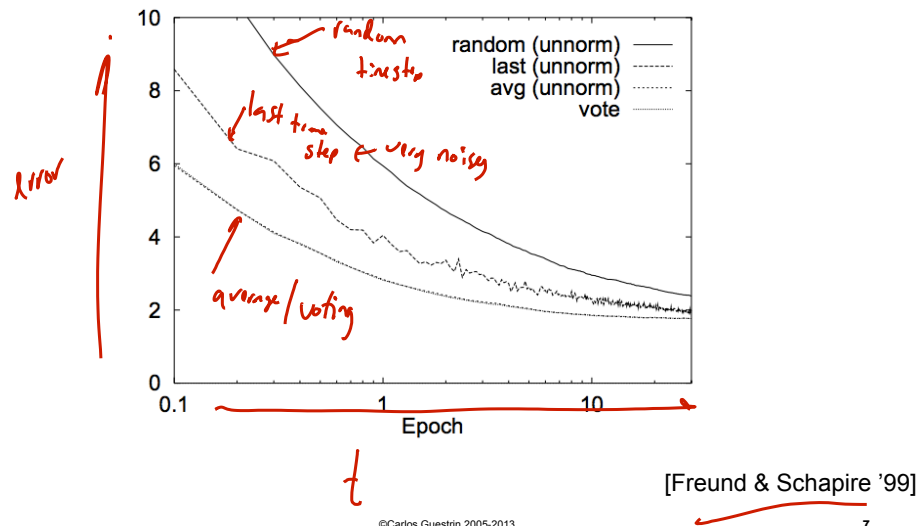
- average!! $\bar{w} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$

- Voting & more advanced methods: how long has this param been good

©Carlos Guestrin 2005-2013

6

Choice can make a huge difference!!



Mistake Bounds

why does it work??

- Algorithm “pays” every time it makes a mistake:
loss function for online setting: number of mistakes up to time T

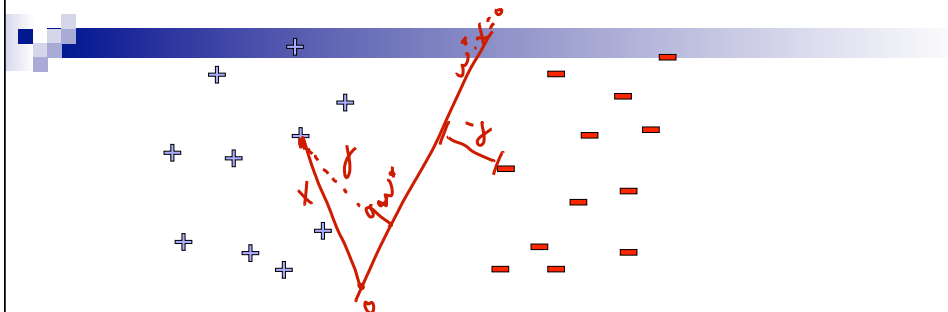
\Rightarrow Google pays for its mistake

- How many mistakes is it going to make?

in its lifetime

Mistake bound

Linear Separability: More formally, Using Margin



- Data linearly separable, if there exists

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

- a margin $\gamma > 0$

- Such that

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

γ for or more from $w^* \cdot x = 0$

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

linearly separable, margin γ

©Carlos Guestrin 2005-2013

9

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

10

$$a \cdot b \leq \|a\| \cdot \|b\|$$

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^* : $w^* \cdot w^{(t+1)} = w^* \cdot w^{(t)} + y^{(t)} (w^* \cdot x^{(t)})$
 - Thus after m mistakes: $w^* \cdot w^{(m+1)} \geq m \gamma$ (by induction)
- Similarly, norm of $w^{(t+1)}$ doesn't grow too fast:
 - $\|w^{(t+1)}\|^2 = \|w^{(t)}\|^2 + 2y^{(t)}(w^{(t)} \cdot x^{(t)}) + \|x^{(t)}\|^2 \leq R^2$ (since $y^{(t)}(w^{(t)} \cdot x^{(t)}) < 0$ for a mistake)
 - Thus, after m mistakes: $\|w^{(m+1)}\|^2 \leq m R^2$
- Putting all together:

$$m \gamma \leq w^* \cdot w^{(m+1)} \leq \|w^*\| \cdot \|w^{(m+1)}\| \leq \sqrt{m} R$$

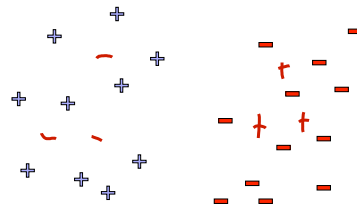
$$\Rightarrow m \gamma \leq \sqrt{m} R \Rightarrow m \leq \left(\frac{R}{\gamma}\right)^2$$

©Carlos Guestrin 2005-2013

11

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



We need features that make data as linearly separable as possible

©Carlos Guestrin 2005-2013

12

What you need to know

- Notion of online learning
- Perceptron algorithm
- Mistake bounds and proof
- In online learning, report averaged weights at the end