# Online Learning
# Perceptron Algorithm

Machine Learning – CSE546

Carlos Guestrin

University of Washington

October 23, 2013

©Carlos Guestrin 2005-2013

---

# Challenge 1: Complexity of Computing Gradients

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

$\forall i$, cost is $O(Nk)$

if $N$ is huge, we have a problem

$\rightarrow$ SGD, looks at one data point at a time
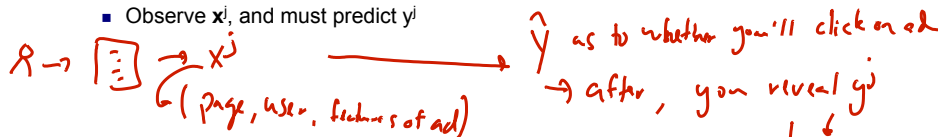
# Challenge 2: Data is streaming

- Assumption thus far: **Batch data**

  *Have all data before you learn*

- But, e.g., in click prediction for ads is a streaming data task:
  - User enters query, and ad must be selected:
    - Observe $x^j$, and must predict $y^j$

  $R \rightarrow \boxed{\vdots} \rightarrow x^j$

  $\hookleftarrow$ ( page, user, features of ad )

  $\hat{y}$ *as to whether you'll click on ad*

  $\rightarrow$ *after, you reveal* $y^j$

  *either you click or you don't*

  - User either clicks or doesn't click on ad:
    - Label $y^j$ is revealed afterwards
      - Google gets a reward if user clicks on ad, *lose money if* $\mathbb{I}(\hat{y} \neq y^j)$

  - Weights must be updated for next time: *what's* $\Delta$?

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

  *update model*

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 *constant*

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

  - Make a prediction: $\hat{y} \leftarrow \text{sign}(w^{(t)} \cdot x^{(t)})$

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

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

  $w_0^{(t)} + \sum_i w_i^{(t)} x_i^{(t)} \geq 0 \Rightarrow 1$

  $\text{otherwise} \Rightarrow -1$

  $w^{(t)} \cdot x^{(t)} = \sum_{i=0}^{k} w_i^{(t)} x_i^{(t)}$

  $= w_0^{(t)} + \sum_i w_i^{(t)} x_i^{(t)}$

  - 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
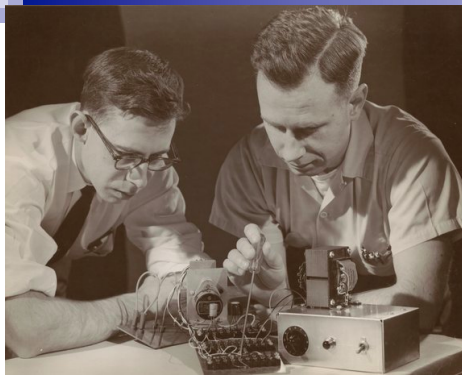
  *observe* $y^{(t)} \rightarrow$ *clicked* $\searrow$ *didn't click*

  *Mistake* $\hat{y} \neq y^{(t)}$

  - Update model:

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

4

Rosenblatt 1957

# The Perceptron Algorithm [Rosenblatt '58, '62]

$$\text{sign}\left(w^{(t)}\cdot x^{(t)}\right) \neq y^{(t)}$$

- Classification setting: y in {-1,+1}
- Linear model
  - Prediction: $\hat{y} = \text{sign}\left(w^{(t)}\cdot x^{(t)}\right)$

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

- Training: $w^{(0)}$ e.g. 0 or random motor settings
  - Initialize weight vector:
  - At each time step:
    - Observe features: $x^{(t)}$
    - Make prediction: $\hat{y} = \text{sign}\left(w^{(t)}\cdot x^{(t)}\right)$
    - Observe true class:
      $y^{(t)} \leftarrow$ true label
    - Update model:
      - If prediction is not equal to truth

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

if $y^{(t)} = +1$
$\quad w^{(t)}\cdot x^{(t)} < 0$
$\Rightarrow$ mistake
$\Rightarrow y^{(t)} w^{(t)}\cdot x^{(t)} < 0$
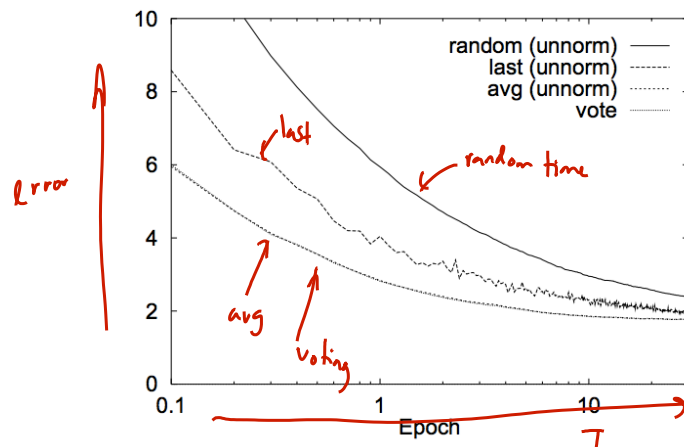similarly for $y^{(t)} = -1$

6

3

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

- Perceptron prediction: $\text{sign}(w \cdot x)$
- Suppose you run online learning method and want to sell your learned weight vector… Which one do you sell???

- Last one? $w^{(T)}$ ? ← very noisy

- Random time step? ← very noisy

- average!! $\hat{w} = \frac{1}{T+1} \sum_{t=0}^{T} w^{(t)}$

- Voting or more advanced avg, see readings

7


## Choice can make a huge difference!!



random (unnorm)
last (unnorm)
avg (unnorm)
vote

*last*

*random time*

*error*

*avg*

*voting*

Epoch

$T$

[Freund & Schapire '99]

8

4

# Mistake Bounds

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

  Loss function for this online setting is # mistakes
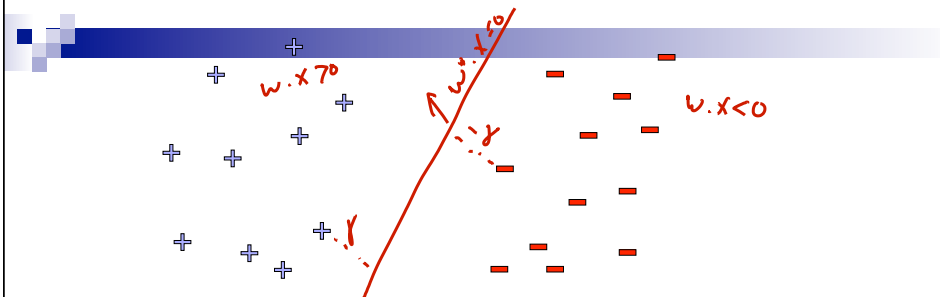
  $\Rightarrow$ Google pays for every mistake

- How many mistakes is it going to make?

  Mistake bound

9

---

# Linear Separability: More formally, Using Margin



$w \cdot x > 0$

$w^* \cdot x = 0$

$w \cdot x < 0$

- Data linearly separable, if there exists
  - a vector $\exists \; w^*, \; \|w^*\| = 1$
  - a margin $\gamma > 0$
- Such that  all points are  at least $\gamma$ away from $w \cdot x = 0$

  $\forall t : \text{if } y^{(t)} = +1, \; w^* \cdot x^{(t)} \geq \gamma$

  $\quad\quad\quad\quad y^{(t)} = -1, \; w^* \cdot x^{(t)} \leq -\gamma$

  Linearly Separable:

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

10

5

# Perceptron Analysis: Linearly Separable Case

- Theorem [Block, Novikoff]:
  - Given a sequence of labeled examples: $(x^{(1)}, y^{(1)})$ .... $(x^{(T)}, y^{(T)})$
    - not iid
  - Each feature vector has bounded norm:
    $$\forall t \quad \|x^{(t)}\| \leq R$$
  - If dataset is linearly separable:
    $$\exists w^*, \|w^*\| = 1, \forall t \quad y^{(t)} w^* \cdot x^{(t)} \geq \gamma \quad \text{for } \gamma > 0$$
- Then the number of mistakes made by the online perceptron on any such sequence is bounded by
  $$\left(\frac{R}{\gamma}\right)^2$$

wow!!

constant, doesn't grow with $T$ !!

---

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

# Perceptron Proof for Linearly Separable case

Assume $\|w^{(0)}\| = 0$

- 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 $w^*$:  $w^* \cdot w^{(t+1)} = w^* (w^{(t)} + y^{(t)} x^{(t)})$
  - Thus after m mistakes:  $= w^* \cdot w^{(t)} + y^{(t)} w^* x^{(t)}$
    by induction ... $w^* \cdot w^{(t+1)} \geq m \gamma$
    $$\underbrace{w^* \cdot w^{(t+1)}}_{\geq w^* \cdot w^{(t)} + \gamma} \geq \gamma$$
- Similarly, norm of $w^{(t+1)}$ doesn't grow too fast:
  - $\|w^{(t+1)}\|^2 = \|w^{(t)}\|^2 + 2y^{(t)}(w^{(t)} \cdot x^{(t)}) + \|x^{(t)}\|^2$
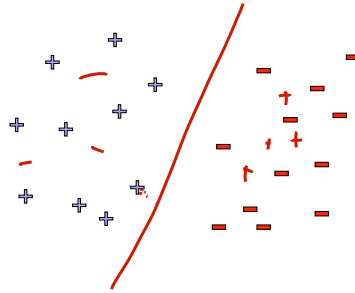    mistake: $< 0$  $\leq R^2$  $\|w^{(t+1)}\|^2 \leq \|w^{(t)}\|^2 + R^2$
  - Thus, after m mistakes:
    $$\|w^{(t+1)}\|^2 \leq m R^2$$
- Putting all together:
  $$m\gamma \leq w^* \cdot w^{(t+1)} \leq \|w^*\| \|w^{(t+1)}\| \leq \sqrt{m} R$$
  $$\Rightarrow m\gamma \leq \sqrt{m} R \Rightarrow m \leq \left(\frac{R}{\gamma}\right)^2 \quad \text{wow!!}$$

# 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)

*d ; f data is very very non-linearly separable*

13

# What you need to know

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

14

# What's the Perceptron Optimizing?

Machine Learning – CSE546

Carlos Guestrin

University of Washington

October 23, 2013

---

# What is the Perceptron Doing???

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

$$\max_w \; \log \; \prod_j P(y^{(j)} \mid x^{(j)}, w)$$

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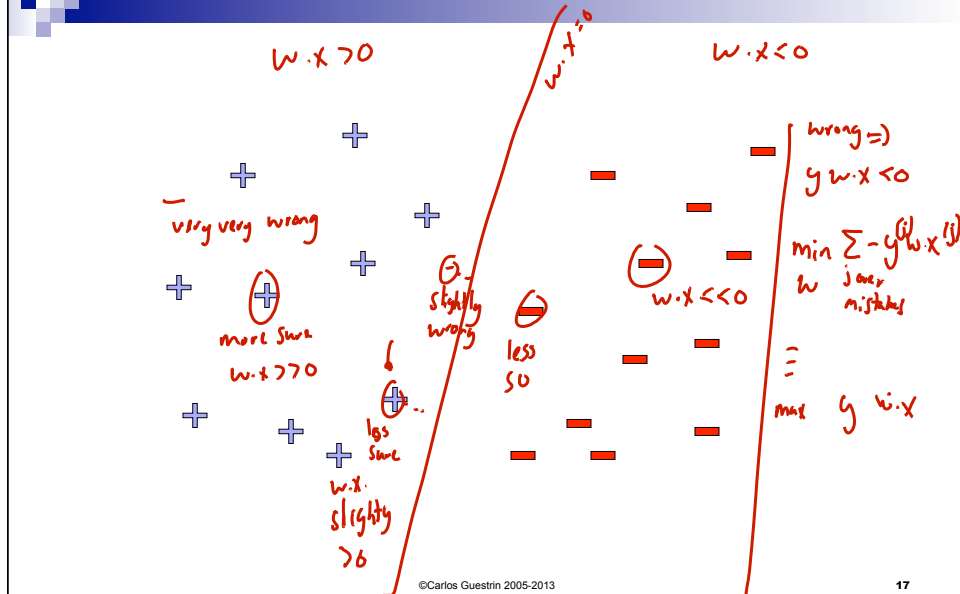  - Started from description of an algorithm

- **What is the Perceptron optimizing????**

# Perceptron Prediction: Margin of Confidence

min-f (=) max f

$w \cdot x > 0$ $\quad w \cdot t = 0$ $\quad w \cdot x < 0$

wrong =)
$y\, w \cdot x < 0$

$\min_{w} \sum -y^{(i)} w \cdot x^{(i)}$ over mistakes

$=$

$\max \; y \; w \cdot x$

very very wrong

more sure
$w \cdot t > 0$

slightly wrong

less so

$w \cdot x << 0$

less sure

$w \cdot x$ slightly $> 0$

©Carlos Guestrin 2005-2013                          17

---

# Hinge Loss

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

- Makes a mistake when:

  $y\, w \cdot x < 0 \;\Rightarrow\;$

  $\ell(w,x) : \begin{cases} 0 & \text{if } y\, w \cdot x \geq 0 \\ -y\, w \cdot x & \text{otherwise } (y\, w \cdot x < 0) \end{cases}$

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

  0/1 loss mistakes

  $-y\, w \cdot x$

  loss $\ell(w,x)$

  don't pay here

  pay linearly here

  $y\, w \cdot x$

©Carlos Guestrin 2005-2013                          18

9

# Minimizing hinge loss in Batch Setting

$(a)_+ \begin{cases} a & \text{if } a \geq 0 \\ 0 & \text{otherwise} \end{cases}$
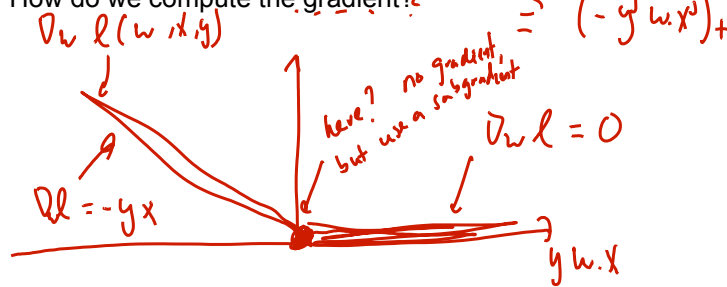
- Given a dataset: $(x^1, y^1) \dots (x^N, y^N)$

- Minimize average hinge loss:

$$\min_w \frac{1}{N} \sum_{j=1}^{N} \ell(w, x^j, y^j) \quad \begin{cases} 0 & \text{if } y^j w x^j \geq 0 \\ -y^j w x^j & \text{otherwise} \end{cases} = (-y^j w \cdot x^j)_+$$
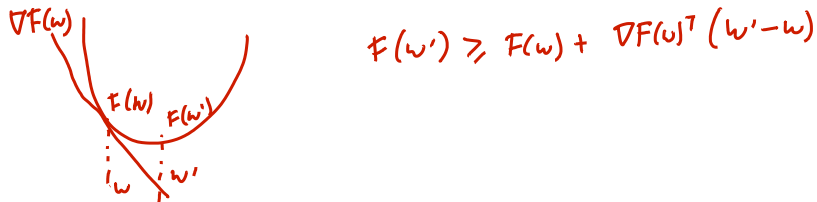
- How do we compute the gradient?

$$\nabla_w \ell(w, x, y)$$

no gradient here? but use a subgradient

$\nabla_w \ell = 0$

$\nabla \ell = -yx$

$y \, w \cdot x$

19

# Subgradients of Convex Functions

- Gradients lower bound convex functions:

$\nabla F(w)$

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

$F(w) \quad F(w')$

$w \quad w'$

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

- Subgradients: Generalize gradients to non-differentiable points:
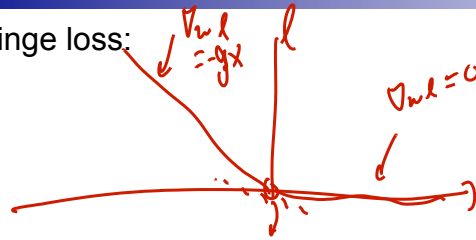  - Any plane that lower bounds function:

  |w|

  For |w| at 0:

  $V \in [-1, 1]$

  lower bound plane

  $w$

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

  iff

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

20

10

# Subgradient of Hinge

- Hinge loss:



- Subgradient of hinge loss:
  - If $y^{(t)}(w.\mathbf{x}^{(t)}) > 0$: $\nabla_w \ell = 0$
  - If $y^{(t)}(w.\mathbf{x}^{(t)}) < 0$: $\nabla_w \ell = -yx$
  - If $y^{(t)}(w.\mathbf{x}^{(t)}) = 0$: $\partial_w \ell = [-yx, 0]$  ←
  - In one line:

$$\partial_w \ell(w,x,y) = \underbrace{\mathbb{1}(y \, w.x \le 0)}_{\text{mistake}} (-yx)$$

in subgradient descent, you can pick __any__ of these, e.g. $-yx$

21

---

# Subgradient Descent for Hinge Minimization

- Given data: $(x^1, y^1) \ldots (x^N, y^N)$

  I want $\min_w$

- Want to minimize:

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

- Subgradient descent works the same as gradient descent:
  - But if there are multiple subgradients at a point, just pick (any) one:

$$w^{(t+1)} \leftarrow w^{(t)} - \underbrace{\eta}_{\text{step size}} \sum_{j=1}^{N} \underbrace{\partial \ell(w^{(t)}, x^j, y^j)}_{\underbrace{\mathbb{1}(y^j w^{(t)}.x^j \le 0)}_{\text{mistake?}} \left(-y^j x^j\right)}$$

22

# Perceptron Revisited

*[handwritten: step size 1] [handwritten: if mistake] [handwritten: add $y^{(i)}x^{(t)}$]*

- Perceptron update:

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

*[handwritten: loss] [handwritten: sum over all points] [handwritten: mistake on point i?] [handwritten: add $g^i \cdot x^i$]*

- Batch hinge minimization update:

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

*[handwritten: step size]*

*[handwritten: Perception as SGD for hinge loss minimization with step size constant $\eta = 1$, and no regularization]*

- Difference?

©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**