

# Case Study 1: Estimating Click Probabilities

## Perceptron Algorithm Kernels (continued)

Machine Learning/Statistics for Big Data  
CSE599C1/STAT592, University of Washington

Carlos Guestrin  
January 15<sup>th</sup>, 2013

©Carlos Guestrin 2013

1

## 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

#### □ Make a prediction:

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

$$w_0 + \sum_i w_i x_i^{(t)} > 0? \Rightarrow w \cdot x^{(t)} > 0$$

$$\sum_{i=0}^d w_i x_i^{(t)} \rightarrow x_0^{(t)} = 1$$

$$x^{(t)} = \begin{pmatrix} x^{(t)} \\ 1 \end{pmatrix}$$

#### □ 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  
or  
not clicked

#### □ Update model:

$$w^{(t+1)} \leftarrow w^{(t)} + \Delta^{(t)} \epsilon_i \dots \text{what is } \Delta^{(t)}? \\ \text{something}$$

©Carlos Guestrin 2013

2

# 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:  $w^{(0)} = 0$
  - At each time step:
    - Observe features:  $x^{(t)} \leftarrow$  user, page, and features
    - Make prediction:  $\hat{y} = \text{Sign}(w^{(t)} \cdot x^{(t)})$
    - Observe true class:  $y^{(t)} \leftarrow$  true label
    - Update model:
      - If prediction is not equal to truth, if make a mistake
        - if  $\hat{y} \neq y^{(t)}$
        - else:  $w^{(t+1)} \leftarrow w^{(t)} + y^{(t)} x^{(t)}$

©Carlos Guestrin 2013

3

## What if the data is not linearly separable?

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

$\Phi(x) : R^m \mapsto F$

$\Phi(x) = \begin{pmatrix} x \\ x^2 \\ x^3 \\ x^4 \\ e^{-x} \\ \sin \log \cos x \\ \vdots \end{pmatrix}$

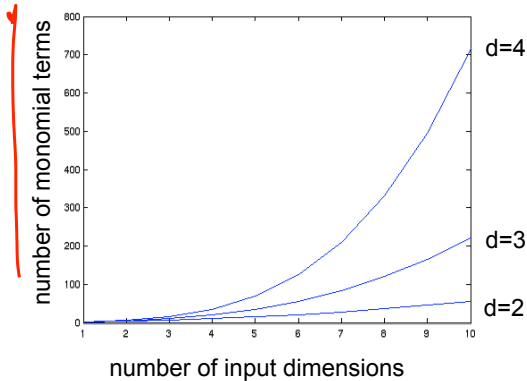
**Feature space can get really large really quickly!**

©Carlos Guestrin 2013

# Higher order polynomials

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

*dim of  $\phi$*



$m$  – input features  
 $d$  – degree of polynomial

$\phi(u) \cdot \phi(v) = (u \cdot v + 1)^d$

grows fast!  
 $d = 6, m = 100$   
about 1.6 billion terms

©Carlos Guestrin 2013

5

# Perceptron Revisited

- Given weight vector  $w^{(t)}$ , predict point  $x$  by:

$$\hat{y} = \text{sign}(w^{(t)} \cdot x)$$

- Mistake at time  $t$ :  $w^{(t+1)} = w^{(t)} + y^{(t)} x^{(t)}$
- Thus, write weight vector in terms of mistaken data points only:
  - Let  $M^{(t)}$  be time steps up to  $t$  when mistakes were made:

$$w^{(t)} = \sum_{i \in M^{(t)}} y^{(i)} x^{(i)}$$

- Prediction rule now:

$$\text{sign}(w^{(t)} \cdot x) = \text{sign}\left(\sum_{i \in M^{(t)}} y^{(i)} x^{(i)}\right) \cdot x = \text{sign}\left(\sum_{i \in M^{(t)}} y^{(i)} x^{(i)} \cdot x\right)$$

- When using high dimensional features:

$$\text{sign}\left(\sum_{i \in M^{(t)}} y^{(i)} \phi(x^{(i)}) \cdot \phi(x)\right) \leftarrow \text{list of mistakes ever made}$$

dot product between  $x$  and mistake  $i$

©Carlos Guestrin 2013

6

# Dot-product of polynomials

$$u = (u_1, u_2)$$

$$v = (v_1, v_2)$$

$\Phi(u) \cdot \Phi(v) =$  polynomials of degree exactly  $d$

$$d=1 \quad \phi(u) \cdot \phi(v) = \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} \cdot \begin{pmatrix} v_1 \\ v_2 \end{pmatrix} = u_1 v_1 + u_2 v_2 = u \cdot v$$

$$d=2 \quad \phi(u) \cdot \phi(v) = \begin{pmatrix} u_1^2 \\ u_1 u_2 \\ u_2 u_1 \\ u_2^2 \end{pmatrix} \cdot \begin{pmatrix} v_1^2 \\ v_1 v_2 \\ v_2 v_1 \\ v_2^2 \end{pmatrix} = u_1^2 v_1^2 + 2u_1 u_2 v_1 v_2 + u_2^2 v_2^2 = (u_1 v_1 + u_2 v_2)^2 = (u \cdot v)^2$$

proof by one step of induction

if  $\phi(\cdot)$  is poly of degree exactly  $d$ , compute in time of basically  $u \cdot v$

$$\underbrace{\phi(u)} \cdot \underbrace{\phi(v)} = \underbrace{(u \cdot v)^d}$$

©Carlos Guestrin 2013

7

# Finally the Kernel Trick!!! (Kernelized Perceptron)

- Every time you make a mistake, remember  $(\mathbf{x}^{(t)}, y^{(t)})$

↳ Keep indices  $M(t)$  of mistakes up to  $t$   
 &  $\mathbf{x}^{(i)}, y^{(i)}$  for these mistakes

- Kernelized Perceptron prediction for  $\mathbf{x}$ :

$$\text{sign}(\mathbf{w}^{(t)} \cdot \phi(\mathbf{x})) = \sum_{i \in M^{(t)}} y^{(i)} \phi(\mathbf{x}^{(i)}) \cdot \phi(\mathbf{x})$$

predict at any  $\mathbf{x}$

$$= \sum_{i \in M^{(t)}} y^{(i)} k(\mathbf{x}^{(i)}, \mathbf{x})$$

"  $\phi(\mathbf{x}^{(i)}) \cdot \phi(\mathbf{x})$

©Carlos Guestrin 2013

8

## Polynomial kernels

- All monomials of degree  $d$  in  $O(d)$  operations:

$$\Phi(\mathbf{u}) \cdot \Phi(\mathbf{v}) = (\mathbf{u} \cdot \mathbf{v})^d = \text{polynomials of degree exactly } d$$

- How about all monomials of degree up to  $d$ ?

□ Solution 0:  $\phi(\mathbf{u}) \cdot \phi(\mathbf{v}) = \sum_{i=0}^d \binom{d}{i} (\mathbf{u} \cdot \mathbf{v})^i \leftarrow \text{about } O(d^2)$

□ Better solution:  $(\mathbf{u} \cdot \mathbf{v})^1 + (\mathbf{u} \cdot \mathbf{v})^2 + (\mathbf{u} \cdot \mathbf{v})^0 + (\mathbf{v} \cdot \mathbf{u})^1 = (\mathbf{u} \cdot \mathbf{v} + 1)^2$

Kernels for all monomials of degree up to and including  $d$

$$K(\mathbf{u}, \mathbf{v}) = \phi(\mathbf{u}) \cdot \phi(\mathbf{v}) = (\mathbf{u} \cdot \mathbf{v} + 1)^d \leftarrow \text{about } O(d)$$

©Carlos Guestrin 2013

9

## Common kernels

- Polynomials of degree exactly  $d$  ✓

$$K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v})^d$$

- Polynomials of degree up to  $d$  ✓

$$K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v} + 1)^d$$

- Gaussian (squared exponential) kernel ✓ radial basis functions...

$$K(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{u} - \mathbf{v}\|^2}{2\sigma^2}\right)$$

- Sigmoid

$$K(\mathbf{u}, \mathbf{v}) = \tanh(\eta \mathbf{u} \cdot \mathbf{v} + \nu)$$

? on strings  
↓ on graphs ...

infinite dimensional feature spaces

©Carlos Guestrin 2013

10

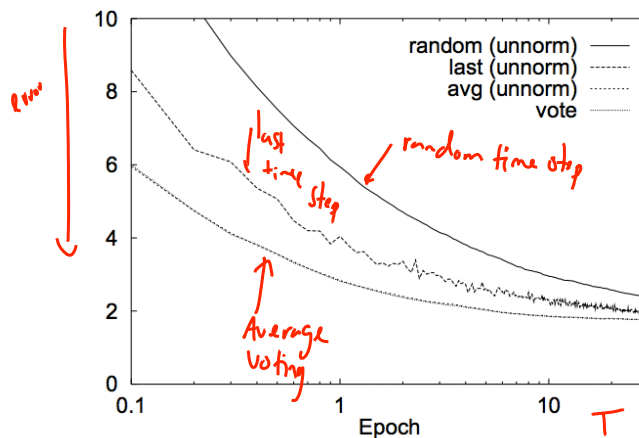
## Fundamental Practical Problem for All Online Learning Methods: **Which weight vector to report?**

- Suppose you run online learning method and want to sell your learned weight vector... Which one do you sell???
- Last one?  $w(t)$  ← can be noisy, highly affected by last mistake
- Random time step?
- average!!  $\hat{w} = \frac{1}{T} \sum_{t=0}^T w(t)$  (easy to keep track of sum)
- Voting, see Freund, Schapire Very important from readings

©Carlos Guestrin 2013

11

## Choice can make a huge difference!!



[Freund & Schapire '99]

©Carlos Guestrin 2013

12

## What you need to know

- Notion of online learning
- Perceptron algorithm
- Mistake bounds and proofs
- The kernel trick
- Kernelized Perceptron
- Derive polynomial kernel
- Common kernels
- In online learning, report averaged weights at the end

*Ⓟ in online settings, often user preferences change over time. can address this by averaging biased towards recent w(t)*

©Carlos Guestrin 2013

13

## Case Study 1: Estimating Click Probabilities

### Stochastic Gradient Descent

Machine Learning/Statistics for Big Data  
CSE599C1/STAT592, University of Washington

Carlos Guestrin  
January 15<sup>th</sup>, 2013

©Carlos Guestrin 2013

14

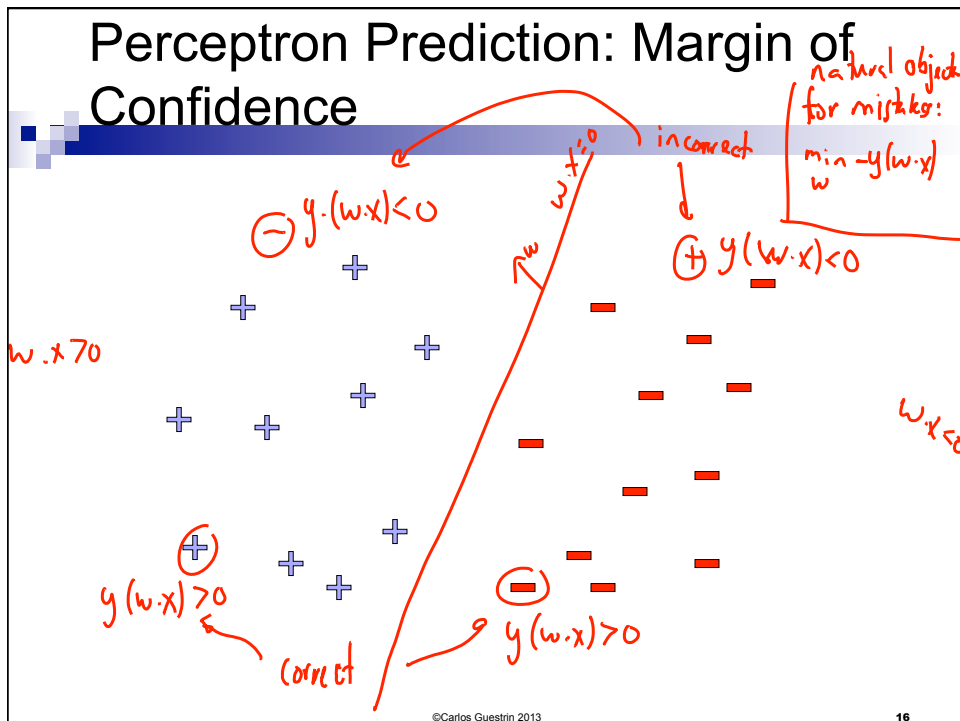
# What is the Perceptron Doing???

- When we discussed logistic regression:
  - Started from maximizing conditional log-likelihood
- When we discussed the Perceptron:
  - Started from description of an algorithm
- What is the Perceptron optimizing????

©Carlos Guestrin 2013

15

## Perceptron Prediction: Margin of Confidence

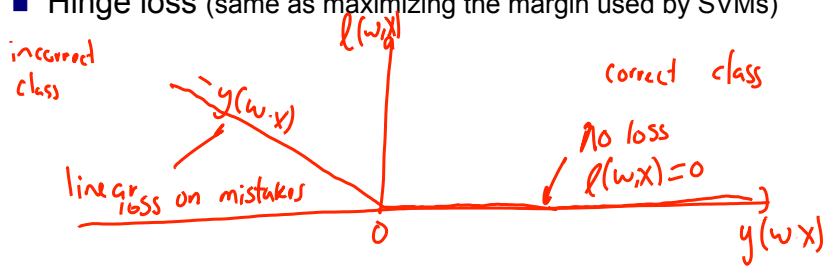


16



# Hinge Loss

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



©Carlos Guestrin 2013

17

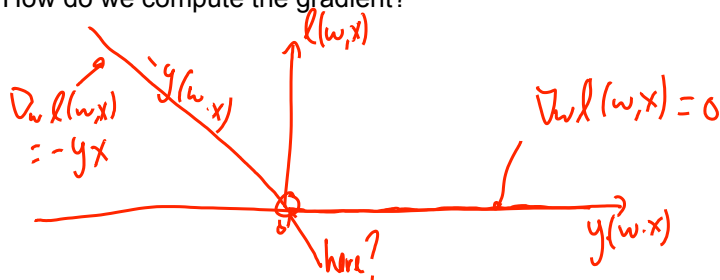
# Minimizing hinge loss in Batch Setting

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

■ Minimize average hinge loss:

$$\min_w \frac{1}{N} \sum_{i=1}^N \underbrace{l(w, x^{(i)})}_{(-y^{(i)}(w \cdot x^{(i)}))_+}$$

- How do we compute the gradient?  $\nabla_w l(w, x) ?$

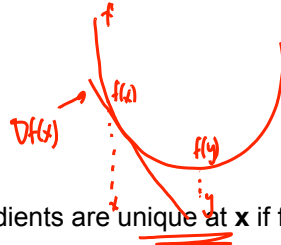


©Carlos Guestrin 2013

18

# Subgradients of Convex Functions

- Gradients lower bound convex functions:

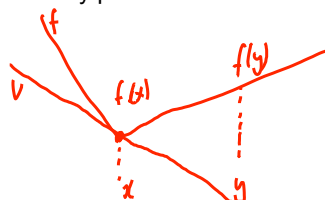


$$f(y) \geq f(x) + \nabla f(x) \cdot (y-x)$$

- Gradients are unique at  $x$  if function differentiable at  $x$

- Subgradients: Generalize gradients to non-differentiable points:

- Any plane that lower bounds function:



$$v \in \partial f(x) \text{ subgradient}$$

if

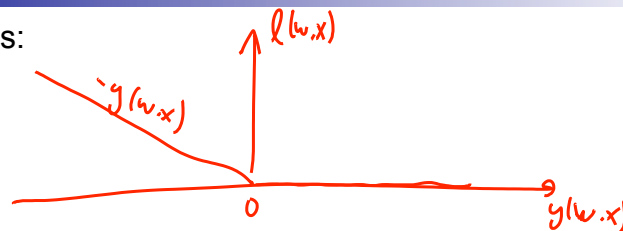
$$f(y) \geq f(x) + v \cdot (y-x)$$

©Carlos Guestrin 2013

19

# Subgradient of Hinge

- Hinge loss:



- Subgradient of hinge loss:

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

$$\partial l(w, x^{(t)}) = \mathbb{I}(y^{(t)}(w \cdot x^{(t)}) \leq 0) [-y x]$$

↑ think of this as gradient of hinge loss

©Carlos Guestrin 2013

20

# Announcements

- No recitation this week
- Comments on readings:
  - Material in readings are superset of what you need
  - Read foundations, e.g., from Kevin Murphy's book, before class
  - Fill in details after class
- Homework out today
  - Start early, start early, start early, start early, start early, start early, start early, start early, start early, start early, start early, start early...
  - Warm-up part of programming due on 1/22
  - Full homework due on 1/29, beginning of class

©Carlos Guestrin 2013

21

# Subgradient Descent for Hinge Minimization

- Given data:  $(x^{(1)}, y^{(1)}) \dots (x^{(N)}, y^{(N)})$
- Want to minimize:
$$\frac{1}{N} \sum_{i=1}^N \ell(w, x^{(i)}) = \frac{1}{N} \sum_{i=1}^N \max(0, 1 - y^{(i)} (w \cdot x^{(i)}))$$
- Subgradient descent works the same as gradient descent:
  - But if there are multiple subgradients at a point, just pick (any) one:

$$w^{(t+1)} \leftarrow w^{(t)} - \eta \frac{1}{N} \sum_{i=1}^N \underbrace{\partial \ell(w, x^{(i)})}_{\mathbb{1}(y^{(i)}(w^{(t)} \cdot x^{(i)}) \leq 0) [-y^{(i)} x^{(i)}]}$$

©Carlos Guestrin 2013

22

# Perceptron Revisited

- Perceptron update:

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

Step size = 1

- Batch hinge minimization update:

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

Step size

average over points

perceptron update

- Difference?

©Carlos Guestrin 2013

23

# Learning Problems as Expectations

- Minimizing loss in training data:

- Given dataset:  $\mathbf{x}^{(1)} \dots \mathbf{x}^{(N)}$

- Sampled iid from some distribution  $p(\mathbf{x})$  on features:

$\mathbf{x}^{(i)} \sim_{\text{iid}} p(\mathbf{x})$

- Loss function, e.g., hinge loss, logistic loss,...

- We often minimize loss in training data:

$$\min_{\mathbf{w}} \left\{ \ell_{\mathcal{D}}(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \ell(\mathbf{w}, \mathbf{x}^{(i)}) \right\}$$

monte carlo integration

- However, we should really minimize expected loss on all data:

$$\min_{\mathbf{w}} \left\{ \ell(\mathbf{w}) = E_{\mathbf{x}} [\ell(\mathbf{w}, \mathbf{x})] = \int p(\mathbf{x}) \ell(\mathbf{w}, \mathbf{x}) d\mathbf{x} \right\}$$

expected loss

- So, we are approximating the integral by the average on the training data

©Carlos Guestrin 2013

24

## Gradient descent in Terms of Expectations

- “True” objective function:

$$\min_{\mathbf{w}} \left\{ \ell(\mathbf{w}) = E_{\mathbf{x}} [\ell(\mathbf{w}, \mathbf{x})] = \int p(\mathbf{x}) \ell(\mathbf{w}, \mathbf{x}) d\mathbf{x} \right\}$$

- Taking the gradient:

$$\nabla_{\mathbf{w}} \ell(\mathbf{w}) = \nabla_{\mathbf{w}} \left[ E_{\mathbf{x}} [\ell(\mathbf{w}, \mathbf{x})] \right] = E_{\mathbf{x}} [\nabla_{\mathbf{w}} \ell(\mathbf{w}, \mathbf{x})]$$

- “True” gradient descent rule:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \underbrace{E_{\mathbf{x}} [\nabla_{\mathbf{w}} \ell(\mathbf{w}, \mathbf{x})]}_{\text{estimate expected gradient to solve true loss}}$$

- How do we estimate expected gradient?

estimate expected gradient to solve true loss

©Carlos Guestrin 2013

25

## SGD: Stochastic Gradient Descent (or Ascent)

- “True” gradient:  $\nabla_{\mathbf{w}} \ell(\mathbf{w}) = E_{\mathbf{x}} [\nabla \ell(\mathbf{w}, \mathbf{x})]$

- Sample based approximation: take iid samples  $\mathbf{x}^{(i)}$

$$\nabla \ell(\mathbf{w}) = E_{\mathbf{x}} [\nabla \ell(\mathbf{w}, \mathbf{x})] \approx \hat{\nabla} \ell(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \nabla \ell(\mathbf{w}, \mathbf{x}^{(i)})$$

the bigger  $N$ , the closer  $\hat{\nabla} \ell$  is to  $\nabla \ell$

- What if we estimate gradient with just one sample???

- Unbiased estimate of gradient

- Very noisy!

- Called stochastic gradient descent

- Among many other names

- VERY useful in practice!!!

$$\nabla \ell(\mathbf{w}) \approx \hat{\nabla} \ell(\mathbf{w}) = \nabla \ell(\mathbf{w}, \mathbf{x}^{(i)})$$

$$E[\hat{\nabla} \ell(\mathbf{w})] = \nabla \ell(\mathbf{w})$$

$\mathbf{x}^{(i)}$  iid sample

©Carlos Guestrin 2013

26

# Perceptron & Stochastic Gradient descent

- Perceptron update:

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

Complexity in terms of N:  $O(1)$

$O(N)$

1 sample estimate of  $E[\nabla L(\mathbf{w})]$   
Stochastic Gradient descent!!  
with  $\eta=1$

- Batch hinge minimization update:

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

N sample estimate of  $E[\nabla L(\mathbf{w})]$

©Carlos Guestrin 2013

27

# Stochastic Gradient Descent: general case

after sufficiently large  $t$ , running avg of loss approx. true loss, use for stopping.

- Given a stochastic function of parameters:  $f(\mathbf{w}) = E_{\mathbf{x}} [f(\mathbf{w}, \mathbf{x})]$ 
  - Want to find minimum

$$\mathbf{w}^* = \underset{\mathbf{w}}{\text{argmin}} E_{\mathbf{x}} [f(\mathbf{w}, \mathbf{x})]$$

- Start from  $\mathbf{w}^{(0)}$ ; e.g.,  $\mathbf{w}^{(0)} = \mathbf{0}$
- Repeat until convergence:

- Get a sample data point  $\mathbf{x}^{(t)}$
- Update parameters:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta_t \nabla f(\mathbf{w}^{(t)}, \mathbf{x}^{(t)})$$

Stopping criteria  
very hard...  
→ theory bounds # iterations in terms of uncomputable constants

- Works on the online learning setting!
- Complexity of gradient computation is constant in number of examples!

→ in practice, validation set

- In general, step size changes with iterations, e.g.,  $\eta_t = \frac{\kappa}{t}$  for  $\kappa > 0$

©Carlos Guestrin 2013

28

## Stochastic Gradient Ascent for Logistic Regression

- Logistic loss as a stochastic function:

$$E_{\mathbf{x}} [\ell(\mathbf{w}, \mathbf{x})] = E_{\mathbf{x}} [\ln P(y|\mathbf{x}, \mathbf{w}) - \lambda \|\mathbf{w}\|_2^2]$$

- Batch gradient ascent updates:

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

- Stochastic gradient ascent updates:

- Online setting:

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta_t \left\{ -\lambda w_i^{(t)} + x_i^{(t)} [y^{(t)} - P(Y=1|\mathbf{x}^{(t)}, \mathbf{w}^{(t)})] \right\}$$

↑ 1 point at a time