

## Case Study 1: Estimating Click Probabilities

# L2 Regularization for Logistic Regression

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

Carlos Guestrin  
January 10<sup>th</sup>, 2013

©Carlos Guestrin 2013

1

## Logistic Regression

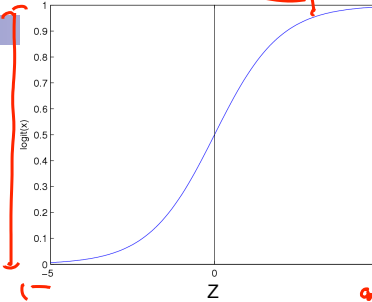
Logistic  
function  
(or Sigmoid):

$$\frac{1}{1 + \exp(-z)}$$

### ■ Learn $P(Y|X)$ directly

- Assume a particular functional form
- Sigmoid applied to a linear function of the data:

$$P(Y = 0|X, W) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$



linear function of the features

$w_0 + \sum_i w_i X_i$  ← not bounded  
could be negative

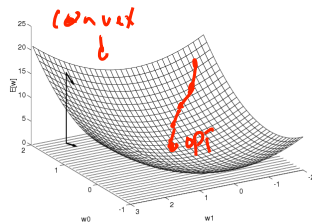
**Features can be discrete or continuous!**

©Carlos Guestrin 2013

2

## Optimizing concave function – Gradient ascent

- Conditional likelihood for Logistic Regression is concave. Find optimum with gradient ascent



Gradient:  $\nabla_{\mathbf{w}} l(\mathbf{w}) = \left[ \frac{\partial l(\mathbf{w})}{\partial w_0}, \dots, \frac{\partial l(\mathbf{w})}{\partial w_n} \right]'$

Step size,  $\eta > 0$

Update rule:  $\Delta \mathbf{w} = \eta \nabla_{\mathbf{w}} l(\mathbf{w})$

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \frac{\partial l(\mathbf{w})}{\partial w_i}$$

- Gradient ascent is simplest of optimization approaches
  - e.g., Conjugate gradient ascent much better (see reading)

©Carlos Guestrin 2013

3

## Gradient Ascent for LR

Gradient ascent algorithm: iterate until change  $< \epsilon$

$$w_0^{(t+1)} \leftarrow w_0^{(t)} + \eta \sum_j [y^j - \hat{P}(Y^j = 1 | \mathbf{x}^j, \mathbf{w}^{(t)})]$$

For  $i=1, \dots, d$

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

repeat

©Carlos Guestrin 2013

4

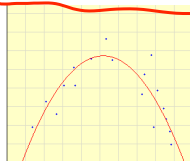
## Test set error as a function of model complexity



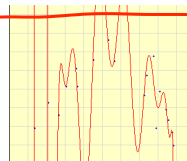
## Regularization in linear regression

- Overfitting usually leads to very large parameter choices, e.g.:

$$-2.2 + 3.1 X - 0.30 X^2$$



$$-1.1 + 4,700,910.7 X - 8,585,638.4 X^2 + \dots$$



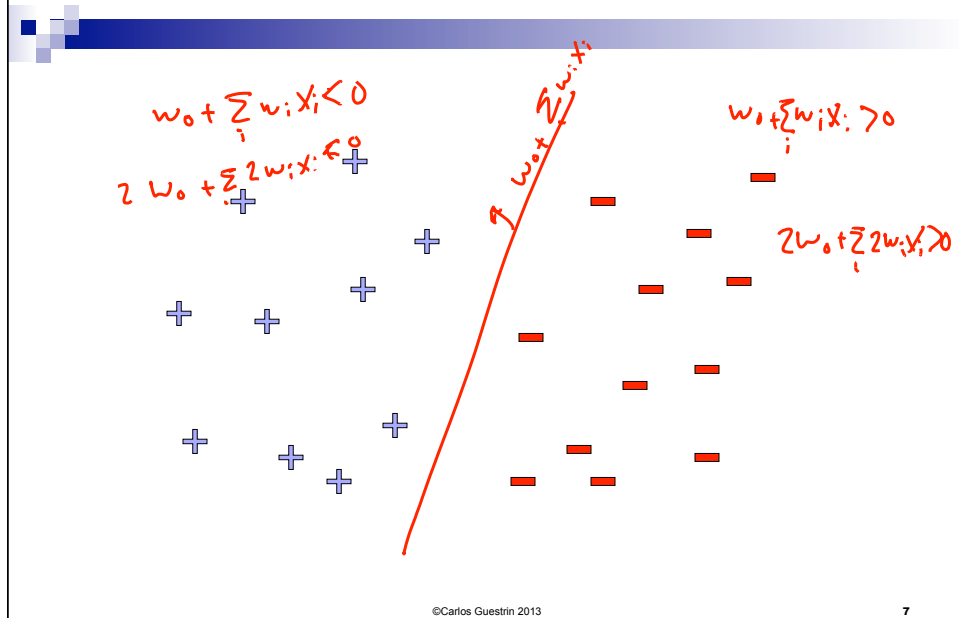
- Regularized least-squares (a.k.a. ridge regression), for  $\lambda > 0$ :

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \underbrace{\sum_j \left( t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2}_{\text{Squared error on train data}} + \underbrace{\lambda \sum_{i=1}^k w_i^2}_{\text{regularization term}} \quad \leftarrow \| \mathbf{w} \|_2^2$$

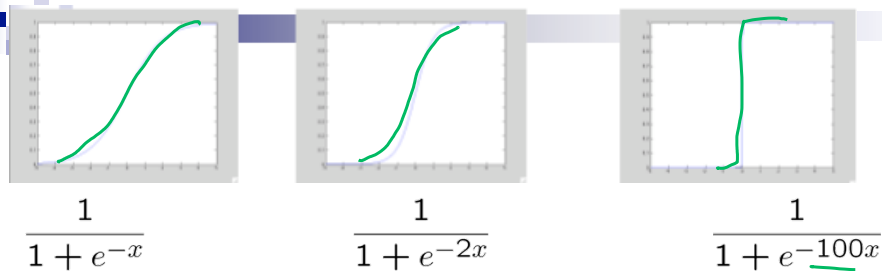
©Carlos Guestrin 2013

6

# Linear Separability



## Large parameters $\rightarrow$ Overfitting



- If data is linearly separable, weights go to infinity
- Leads to overfitting:

- Penalizing high weights can prevent overfitting...

©Carlos Guestrin 2013

8

## Regularized Conditional Log Likelihood

- Add regularization penalty, e.g.,  $L_2$ :

$$\max_{\mathbf{w}} \ell(\mathbf{w}) = \ln \prod_i P(y^j | \mathbf{x}^j, \mathbf{w}) - \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

Conditional log likelihood
regularization

$-\frac{\lambda}{2} \sum_{i=1}^n w_i^2$

- Practical note about  $w_0$ :

not regularize  $w_0$

- Gradient of regularized likelihood:

$$\frac{\partial \ell(\mathbf{w})}{\partial w_i} = \frac{\partial}{\partial w_i} \ln \prod_j P(y^j | \mathbf{x}^j, \mathbf{w}) - \lambda w_i$$

©Carlos Guestrin 2013

9

## Standard v. Regularized Updates

- Maximum conditional likelihood estimate

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \ln \left[ \prod_{j=1}^N P(y^j | \mathbf{x}^j, \mathbf{w}) \right]$$

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

- Regularized maximum conditional likelihood estimate

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \ln \left[ \prod_j P(y^j | \mathbf{x}^j, \mathbf{w}) \right] - \lambda \sum_{i>0} w_i^2$$

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

©Carlos Guestrin 2013

10

## Stopping criterion

for LR, strong concavity  
constant  $\gamma \geq \lambda$

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

- Regularized logistic regression is strongly concave
  - Negative second derivative bounded away from zero:

$f(x)$ :

concave  $\Leftrightarrow -f''(x) \geq 0$  ,  $-f''(x) \geq \gamma$  ;  $\gamma > 0$   
diff strongly concave

- Strong concavity (convexity) is super helpful!!

- For example, for strongly concave  $\ell(\mathbf{w})$ :

$$\underbrace{\ell(\mathbf{w}^*)}_{\text{true}} - \underbrace{\ell(\mathbf{w})}_{\text{estimate}} \leq \frac{1}{2\lambda} \|\nabla \ell(\mathbf{w})\|_2^2$$

norm of gradient  $\|\nabla \ell(\mathbf{w})\|_2^2 \leq 2\lambda \epsilon$

©Carlos Guestrin 2013

11

## Convergence rates for gradient descent/ascent

- Number of iterations to get to accuracy

$$\ell(\mathbf{w}^*) - \ell(\mathbf{w}) \leq \epsilon$$

- If func Lipschitz:  $O(1/\epsilon^2)$   
 $\|\ell(x) - \ell(y)\| \leq K \|x - y\|$
- If gradient of func Lipschitz:  $O(1/\epsilon)$
- If func is strongly convex:  $O(\ln(1/\epsilon))$

exponentially  
fewer  
iterations

constant  
step size

©Carlos Guestrin 2013

12

## What you should know about Logistic Regression (LR) and Click Prediction

- Click prediction problem:
  - Estimate probability of clicking
  - Can be modeled as logistic regression
- Logistic regression model: Linear model
- Optimize conditional likelihood
- Gradient computation
- Overfitting
- Regularization
- Regularized optimization
- Convergence rates and stopping criterion

©Carlos Guestrin 2013

13

### Case Study 1: Estimating Click Probabilities

Online Learning  
Perceptron Algorithm  
Kernels

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

Carlos Guestrin  
January 10<sup>th</sup>, 2013

©Carlos Guestrin 2013

14

## Challenge 1: Complexity of Computing Gradients

*N - data points*  
*d → dims*

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

*O(d)*

complexity is  $O(Nd)$   
 in big data  $d$  is large  
 $N$  can be huge

just to take a small  $\eta$  step

©Carlos Guestrin 2013

15

## Challenge 2: Data is streaming

- Assumption thus far: **Batch data**

- But, click prediction is a streaming data task:

- User enters query, and ad must be selected:

- Observe  $\mathbf{x}^i$ , and must predict  $y^i$

$Q \rightarrow [ ] \rightarrow \mathbf{x}^i$

( $Q$  features,  $[ ]$  features, ad features)

must predict  $y^i$  ← click?  
 online

- 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

- Weights must be updated for next time:

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

©Carlos Guestrin 2013

16



# 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

$x^{(t)}$   
 $\hat{y}^{(t)} \approx y$   
 $x^{(t)} = \begin{pmatrix} x^{(t)} \\ 1 \end{pmatrix}$   
 $w_0 + \sum_i w_i x_i^{(t)} > 0 ? \Rightarrow w^{(t)} \cdot x^{(t)} > 0$   
 $\sum_{i=0}^n w_i x_i^{(t)} \rightarrow x_0^{(t)} = 1$

- 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:

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

©Carlos Guestrin 2013

17

# 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  $\hat{y} \neq y^{(t)}$

$w^{(t+1)} \leftarrow w^{(t)} + y^{(t)} x^{(t)}$   
 else:  $w^{(t+1)} \leftarrow w^{(t)}$

©Carlos Guestrin 2013

18

# Mistake Bounds

- Algorithm “pays” every time it makes a mistake:

loss function. number of mistakes made  
 $\Rightarrow$  Google pays for its mistakes

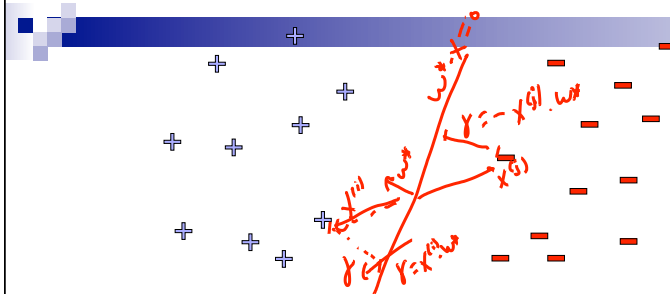
- How many mistakes is it going to make?

$\rightarrow$  Mistake Bound

©Carlos Guestrin 2013

19

## Linear Separability: More formally, Using Margin



- Data linearly separable, if there exists

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

- a margin  $\gamma > 0$

- Such that

$$\forall t \quad \left. \begin{array}{l} \text{if } y^{(t)} = +1, w^* \cdot x^{(t)} \geq \gamma \\ \text{if } y^{(t)} = -1, -w^* \cdot x^{(t)} \geq \gamma \end{array} \right\} \quad \forall t \quad y^{(t)} (w^* \cdot x^{(t)}) \geq \gamma$$

©Carlos Guestrin 2013

20

## Perceptron Analysis: Linearly Separable Case

- Theorem [Block, Novikoff]:

- Given a sequence of labeled examples:  $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})$

- Each feature vector has bounded norm:  $\forall x \quad \|x^{(i)}\| \leq R$

- If dataset is linearly separable:

$\exists w^*, \|w^*\| = 1, \forall x \quad y^{(i)}(w^* \cdot x^{(i)}) \geq \gamma, \gamma > 0$

- Then the number of mistakes made by the online perceptron on this sequence is bounded by

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

wow!!

constant, independent of  $T$  and dim of data  $d$

©Carlos Guestrin 2013

21

## 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  $w^*$ :  $w^* \cdot w^{(t+1)} = w^* \cdot w^{(t)} + \underbrace{y^{(t)}(w^* \cdot x^{(t)})}_{\geq \gamma}$

- Thus after  $k$  mistakes:

$$w^* \cdot w^{(k+1)} \geq k\gamma$$

- Similarly, norm of  $w^{(t+1)}$  doesn't grow too fast:

$$\|w^{(t+1)}\|^2 = \|w^{(t)}\|^2 + 2y^{(t)}(w^{(t)} \cdot x^{(t)}) + \|x^{(t)}\|^2 \leq R^2$$

$$\dots \leq \|w^{(1)}\|^2 + R^2$$

- Thus, after  $k$  mistakes:

$$\|w^{(k+1)}\|^2 \leq kR^2$$

- Putting all together:

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

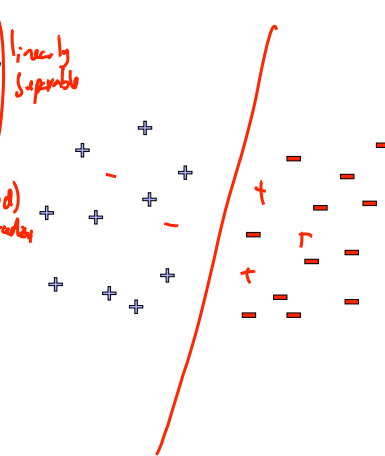
$$\Rightarrow k \leq \left(\frac{R}{\gamma}\right)^2 \leftarrow \text{wow!!}$$

©Carlos Guestrin 2013

22

# Beyond Linearly Separable Case

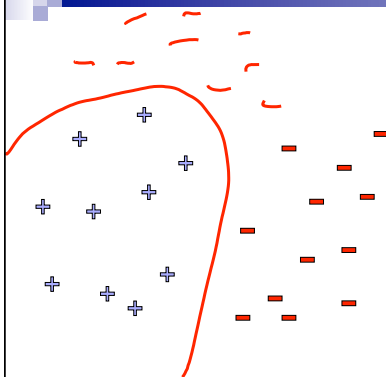
- Perceptron algorithm is super cool!
  - No assumption about data distribution!
    - Could be generated by an oblivious adversary, no need to be iid
  - Makes a fixed number of mistakes, and it's done for ever!
    - Even if you see infinite data
  - Constant cost per iteration
    - Converges in  $O(1/\epsilon)$
- 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 from Readings
  - Converges, but ultimately may not give good accuracy (make many many mistakes)



©Carlos Guestrin 2013

23

## What if the data is not linearly separable?



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

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

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

Feature space can get really large really quickly!

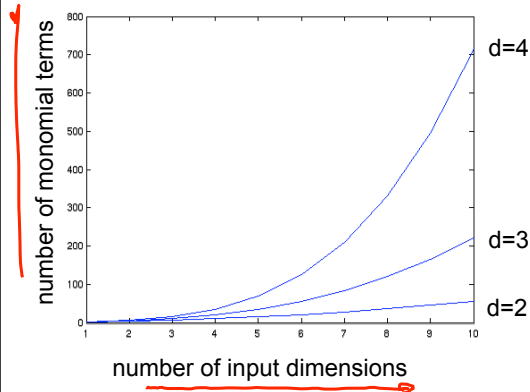
©Carlos Guestrin 2013

24

# 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

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

©Carlos Guestrin 2013

25

# 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:

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

- When using high dimensional features:

$$\text{sign}\left(\sum_{i \in M^{(t)}} y^{(i)} \underbrace{\phi(x^{(i)}) \cdot \phi(x)}_{\text{dot product between } x \text{ and mistake } i}\right) \leftarrow \text{list of mistakes ever made}$$

©Carlos Guestrin 2013

26

## 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  $d$ ,  
 $\phi(u) \cdot \phi(v) = (u \cdot v)^d$  ← compute in time of basicaly  $u \cdot v$

©Carlos Guestrin 2013

27

## Finally the Kernel Trick!!!

### (Kernelized Perceptron

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

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

$$\begin{aligned} \text{sign}(\mathbf{w}^{(t)} \cdot \phi(\mathbf{x})) &= \sum_{i \in M^{(t)}} \phi(\mathbf{x}^{(i)}) \cdot \phi(\mathbf{x}) \\ &= \sum_{i \in M^{(t)}} k(\mathbf{x}^{(i)}, \mathbf{x}) \end{aligned}$$

©Carlos Guestrin 2013

28

## 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:
  - Better solution:

©Carlos Guestrin 2013

29

## 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  
$$K(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{u} - \mathbf{v}\|^2}{2\sigma^2}\right)$$
- Sigmoid  
$$K(\mathbf{u}, \mathbf{v}) = \tanh(\eta \mathbf{u} \cdot \mathbf{v} + \nu)$$

©Carlos Guestrin 2013

30

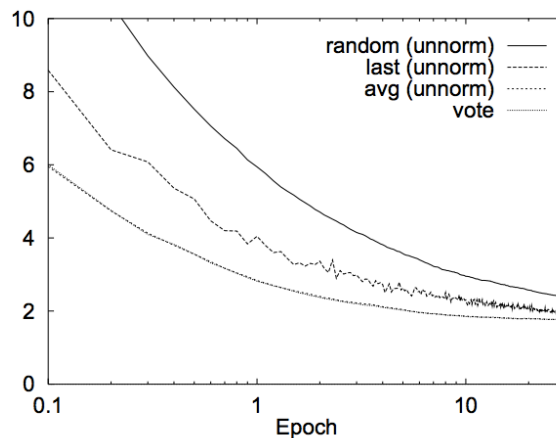
## Fundamental Practical Problem for All Online Learning Methods: **Which weight vector to report?**

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

©Carlos Guestrin 2013

31

## Choice can make a huge difference!!



[Freund & Schapire '99]

©Carlos Guestrin 2013

32



# 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