

Kernels

Machine Learning – CSE546

Carlos Guestrin

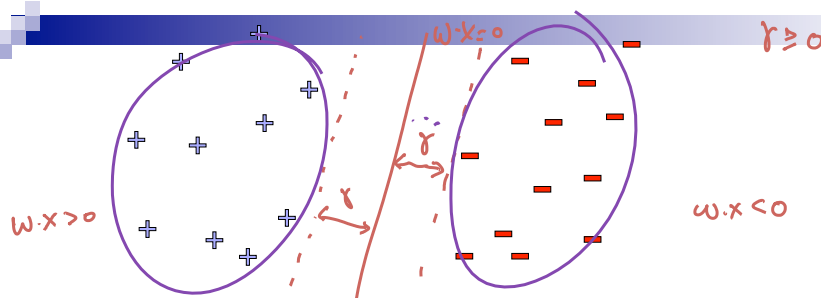
University of Washington

October 28, 2014

©Carlos Guestrin 2005-2014

1

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 $w^* x$

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

©Carlos Guestrin 2005-2014

2

Perceptron Analysis: Linearly Separable Case

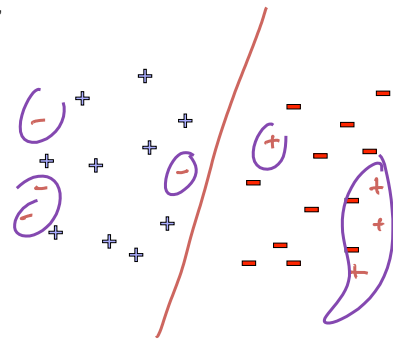
- Theorem [Block, Novikoff]:
 - Given a sequence of labeled examples: $(x^1, y^1), \dots, (x^T, y^T)$
 - Each feature vector has bounded norm: $\forall t, \|x^t\| \leq R$
 - If dataset is linearly separable:
 - $\exists w^*, \|w^*\| = 1, (y^t w^* x^t \geq \gamma) \text{ for } \gamma \geq 0$
- Then the 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

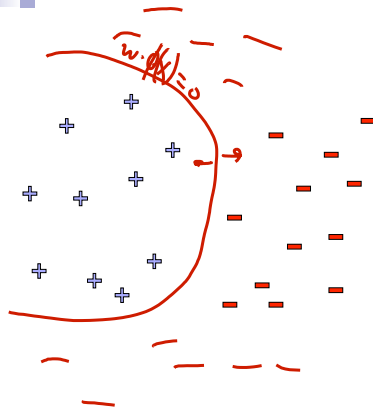
Beyond Linearly Separable Case

- 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)

What if the data is not linearly separable?



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

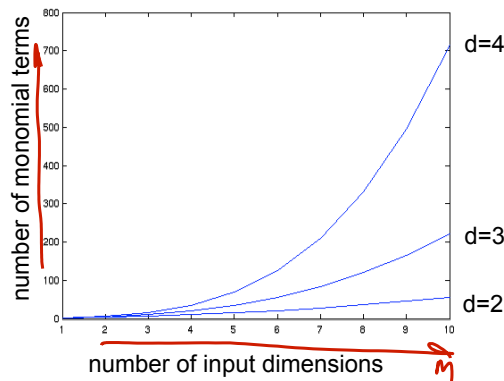
$$\Phi(x) : R^m \mapsto F$$

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

Feature space can get really large really quickly!

Higher order polynomials

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



m – input features
d – degree of polynomial

*Kernels:
deal with a lot
high dim very
efficiently*

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

instead x , use high dim features $\phi(x)$

Perceptron Revisited

$$x \cdot x^{(j)} = \sum_{i=1}^m x_i x_i^{(j)}$$

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

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

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

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

- Prediction rule now:

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

- When using high dimensional features:

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

prediction only depends on $\phi(x)$ and $\phi(x^{(j)})$

©Carlos Guestrin 2005-2014

Dot-product of polynomials

$$u = \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} \quad v = \begin{pmatrix} v_1 \\ v_2 \end{pmatrix}$$

$\Phi(u) \cdot \Phi(v)$ ~~is~~ ^{for} polynomials of degree exactly d

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

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

proof by single step
of induction
for poly of degree exactly d

$$\Phi(u) \cdot \Phi(v) = (u \cdot v)^d = k(u, v)$$

kernel trick !!

©Carlos Guestrin 2005-2014

8

Finally the Kernel Trick!!! (Kernelized Perceptron)

- Every time you make a mistake, remember $(\mathbf{x}^{(t)}, y^{(t)})$
↳ keep list of all mistakes ever made

- Kernelized Perceptron prediction for \mathbf{x} :

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

©Carlos Guestrin 2005-2014

9

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: $\phi(u) \cdot \phi(v) = \sum_{i=0}^d \binom{d}{i} (u \cdot v)^i$

- Better solution:

$$k=2 \quad (u \cdot v)^1 + (u \cdot v)^2 + (v \cdot u)^1 + (u \cdot v)^0 = (u \cdot v + 1)^2$$

proof by "induction"

For polynomials of degree d :

$$\begin{aligned} \phi(u) \cdot \phi(v) &= k(u, v) = \\ &= (u \cdot v + 1)^d \quad \text{!!} \\ &\quad \vdots \\ &\quad \text{etc} \end{aligned}$$

©Carlos Guestrin 2005-2014

10

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

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

↙ Radial basis functions
RBF kernel
↓
Projecting into
infinite dim
space.

©Carlos Guestrin 2005-2014

11

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

©Carlos Guestrin 2005-2014

12

Your Midterm...

- Content: Everything up to last Tuesday (nearest neighbors/decision trees)...
- Only 80mins, so arrive early and settle down quickly, we'll start and end on time
- "Open book"
 - Textbook, Books, Course notes, Personal notes
- Bring a calculator that can do log ☺
- No:
 - Computers, tablets, phones, other materials, internet devices, wireless telepathy or wandering eyes...
- The exam:
 - Covers key concepts and ideas, work on understanding the big picture, and differences between methods

©Carlos Guestrin 2005-2014

13

Support Vector Machines

Machine Learning – CSE546
Carlos Guestrin
University of Washington

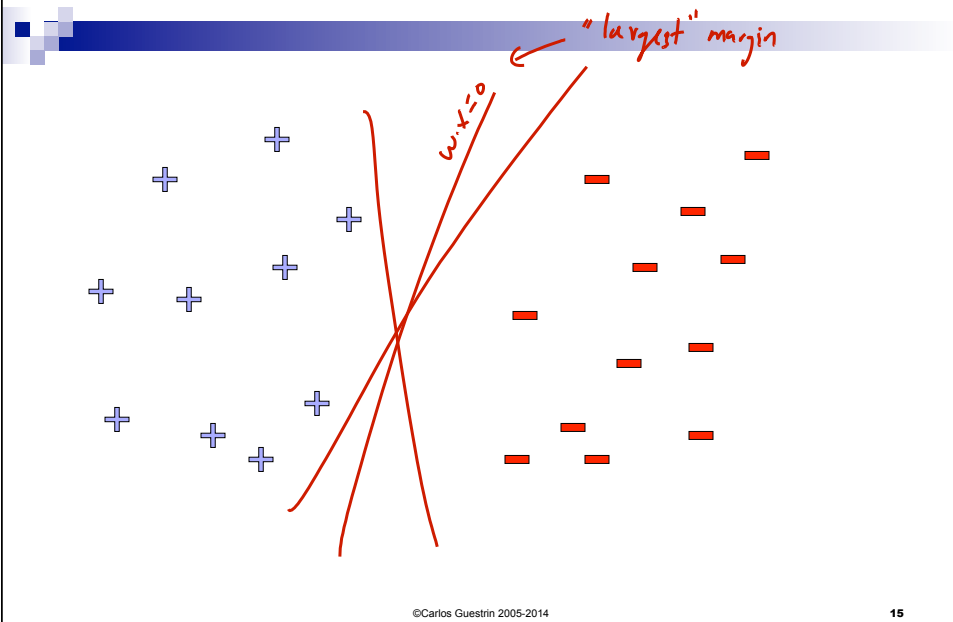
October 28, 2014

©Carlos Guestrin 2005-2014

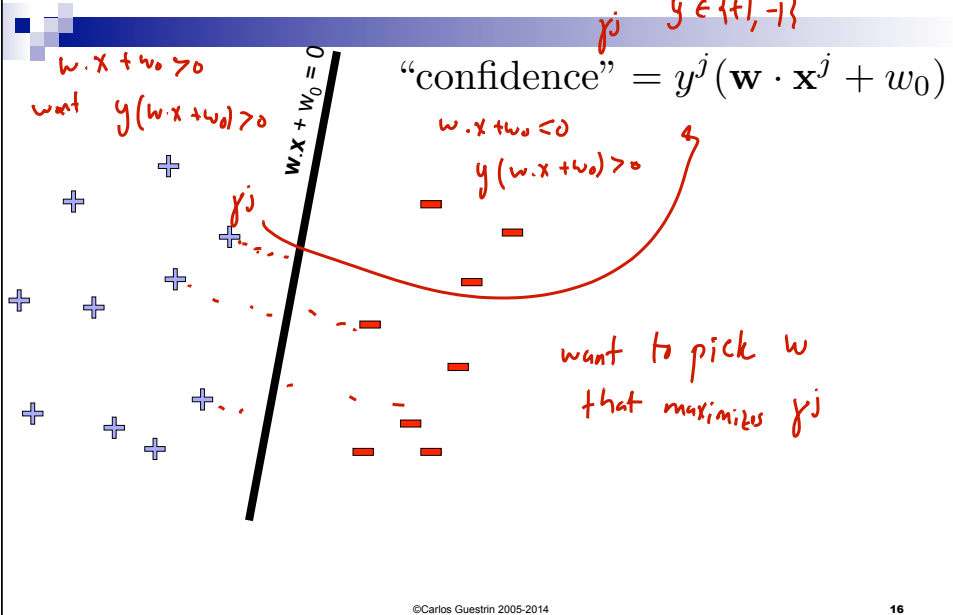
14

perception + L2 regularization

Linear classifiers – Which line is better?

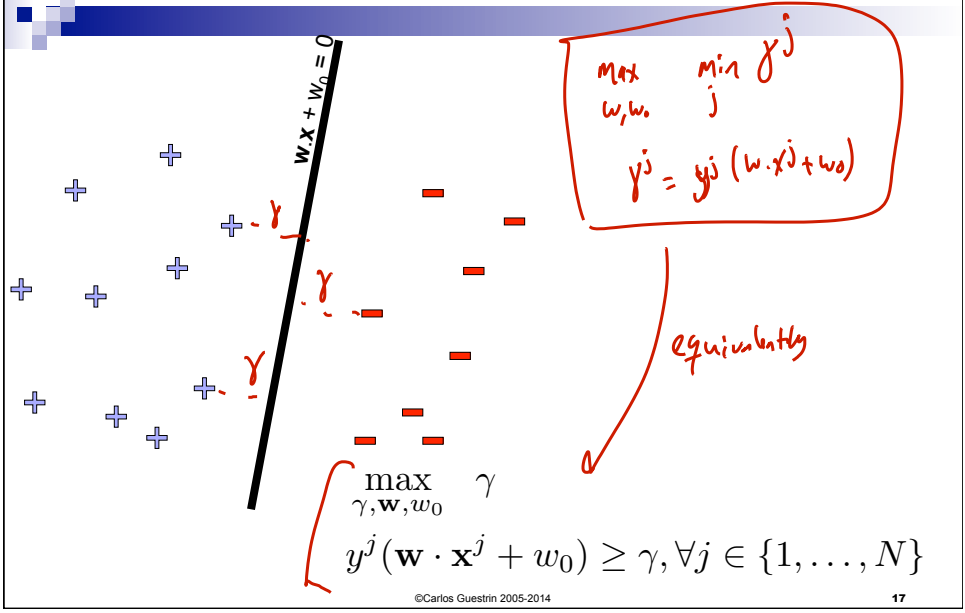


Pick the one with the largest margin!

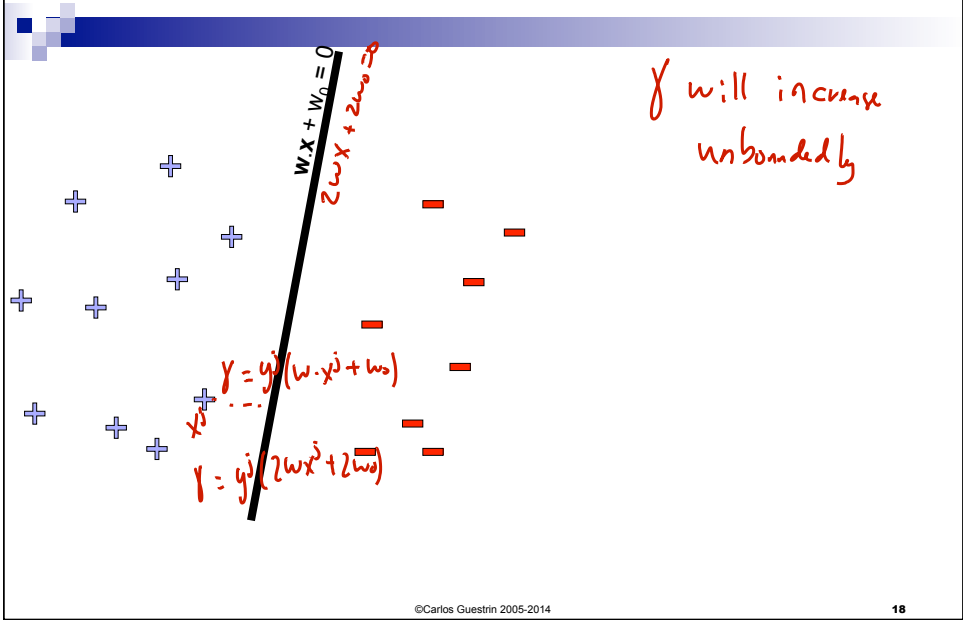


Maximize the margin

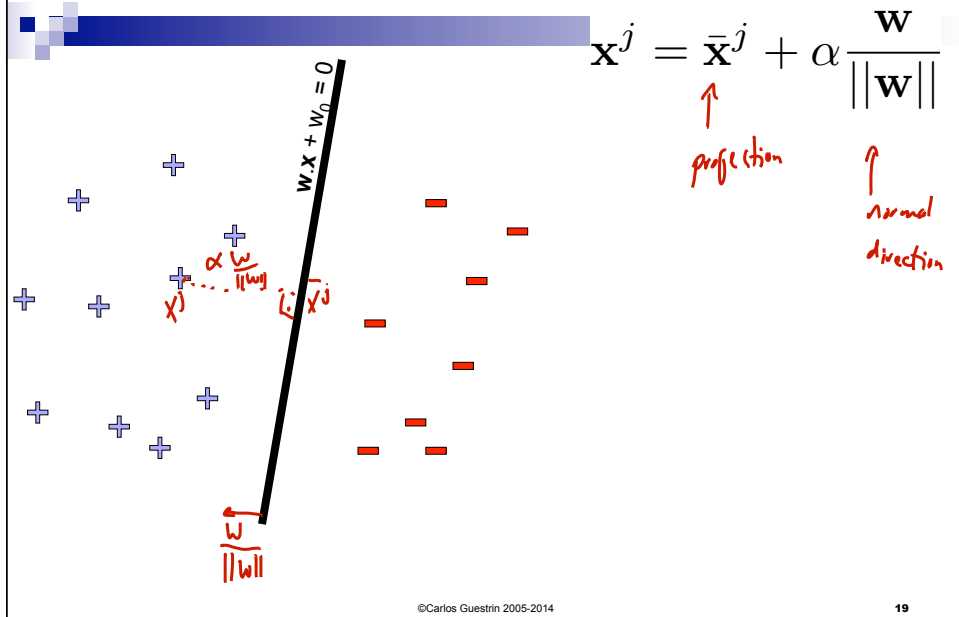
SMvs.
max worst-case margin



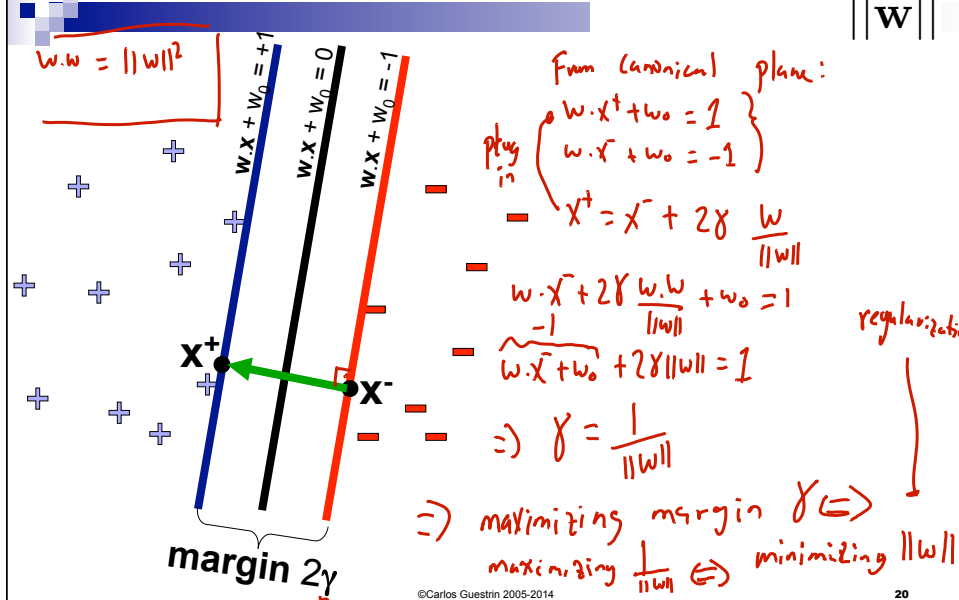
But there are many planes...



Review: Normal to a plane



A Convention: Normalized margin – Canonical hyperplanes



Margin maximization using canonical hyperplanes

Unnormalized problem: $\max_{\gamma, \mathbf{w}, w_0} \gamma$
 $y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0) \geq \gamma, \forall j \in \{1, \dots, N\}$

Normalized Problem:
 $\max \gamma \equiv \max \frac{1}{\|\mathbf{w}\|} \equiv \min \|\mathbf{w}\| \equiv \min \|\mathbf{w}\|^2$
 $y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0) \geq 1$

$\min_{\mathbf{w}, w_0} \|\mathbf{w}\|_2^2$
 $y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0) \geq 1, \forall j \in \{1, \dots, N\}$

Support vector machine

©Carlos Guestrin 2005-2014 21

Support vector machines (SVMs)

$\min_{\mathbf{w}, w_0} \|\mathbf{w}\|_2^2$
 $y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0) \geq 1, \forall j \in \{1, \dots, N\}$

- Solve efficiently by many methods, e.g.,
 - quadratic programming (QP)
 - Well-studied solution algorithms
 - Stochastic gradient descent
- Hyperplane defined by support vectors

what if I forget all data but the support vectors => same solution!

©Carlos Guestrin 2005-2014 22