

\hat{w}
 \hat{w}

perception + L2 regularization

Support Vector Machines

Machine Learning – CSE546
Carlos Guestrin
University of Washington

October 28, 2014
©Carlos Guestrin 2005-2014

1

Linear classifiers – Which line is better?

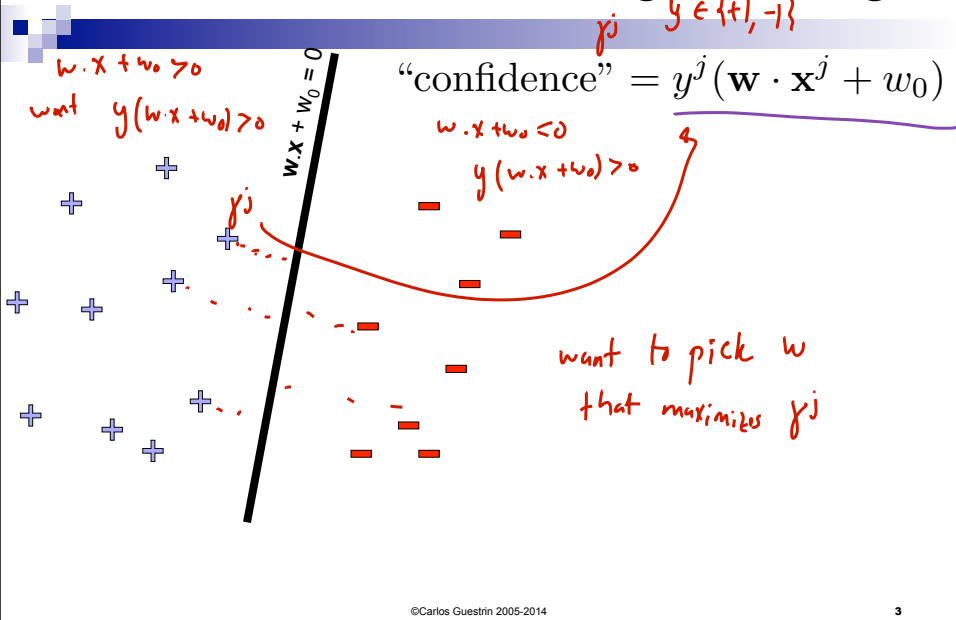
$w \cdot x = 0$

"largest" margin

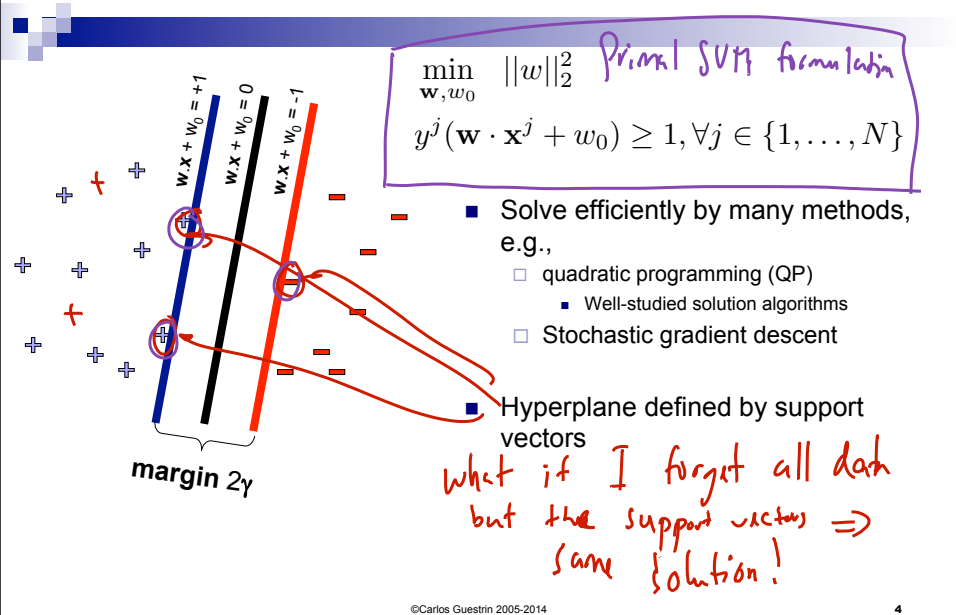
©Carlos Guestrin 2005-2014

2

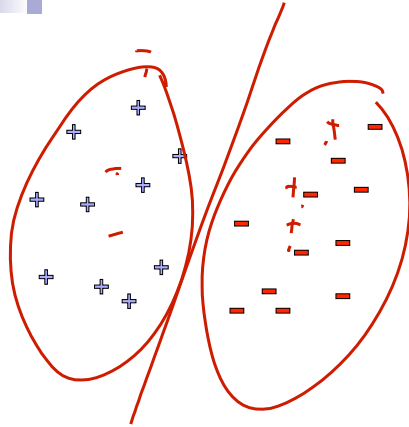
Pick the one with the largest margin!



Support vector machines (SVMs)



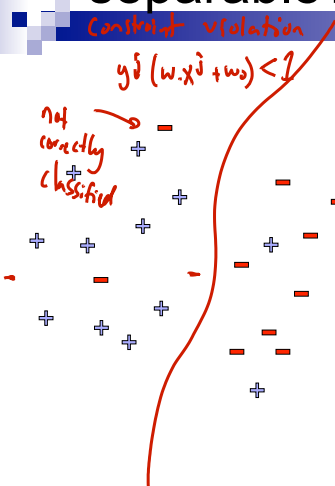
What if the data is not linearly separable?



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

→ Kernels

What if the data is still not linearly separable?



or outliers

$$\min_{\mathbf{w}, w_0} \|\mathbf{w}\|_2^2$$

$$(z)_+ = \begin{cases} 0 & \text{if } z \leq 0 \\ z & \text{otherwise} \end{cases}$$

$$\rightarrow y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0) \geq 1, \forall j \leftarrow$$

- If data is not linearly separable, some points don't satisfy margin constraint:

$$\exists j \quad y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0) < 1 \Rightarrow 1 - y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0) > 0$$

- How bad is the violation?

$$\text{constraint violation} = \begin{cases} 0 & \text{if } 1 - y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0) \leq 0 \\ 1 - y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0) & \text{otherwise} \end{cases} = (1 - y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0))_+$$

- Tradeoff margin violation with $\|\mathbf{w}\|$:

$$\min_{\mathbf{w}, w_0} \underbrace{\|\mathbf{w}\|_2^2}_{\text{margin}} + C \sum_{j=1}^N \underbrace{(1 - y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0))_+}_{\text{stacks}} \quad \text{penalty } C > 0$$

SVMs for Non-Linearly Separable meet my friend the Perceptron...

- Perceptron was minimizing the hinge loss:

$$\min_{w, w_0} \sum_{j=1}^N (-y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0))_+$$

not convex

train error for classifier

hinge loss

confidence

0 1

$y^j(\mathbf{w} \cdot \mathbf{x}^j + w_0)$

- SVMs minimizes the regularized hinge loss!!

$$\|\mathbf{w}\|_2^2 + C \sum_{j=1}^N (1 - y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0))_+$$

regularization term

like $\frac{1}{\lambda}$

SVMs convention

convex upper bound on train error \Rightarrow "easy" to optimize

©Carlos Guestrin 2005-2014

7

Stochastic Gradient Descent for SVMs

- Perceptron minimization:

$$\sum_{j=1}^N (-y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0))_+$$

- SGD for Perceptron:

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

mistake

$$\eta = 1 \quad \nabla (-y^{(t)}(\mathbf{w}^{(t)} \cdot \mathbf{x}^{(t)} + w_0))$$

update weight

- SVMs minimization:

$$\|\mathbf{w}\|_2^2 + C \sum_{j=1}^N (1 - y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0))_+$$

sum over mistakes $x^{(j)}$

- SGD for SVMs:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \eta \left[-2\mathbf{w}^{(t)} + C \mathbb{1}_{[y^{(t)}(\mathbf{w}^{(t)} \cdot \mathbf{x}^{(t)} + w_0) \leq 1]} y^{(t)} \mathbf{x}^{(t)} \right]$$

margin violation

step size is important

©Carlos Guestrin 2005-2014

8

What you need to know

- Maximizing margin
- Derivation of SVM formulation
- Non-linearly separable case
 - Hinge loss
 - A.K.A. adding slack variables
- SVMs = Perceptron + L2 regularization
- Can also use kernels with SVMs
- Can optimize SVMs with SGD
 - Many other approaches possible