



# Kernels and Support Vector Machines

Machine Learning – CSE446

Sham Kakade

University of Washington

November 1, 2016

©2016 Sham Kakade

1

## Announcements:



- Project Milestones coming up
- HW2
  - You've implemented GD, SGD, etc...
- HW3 posted this week.
  - Let's get state of the art on MNIST!
  - It'll be collaborative
- Today:
  - Review: the perceptron, margins, and separability
  - Kernels & SVMs

©2016 Sham Kakade

2

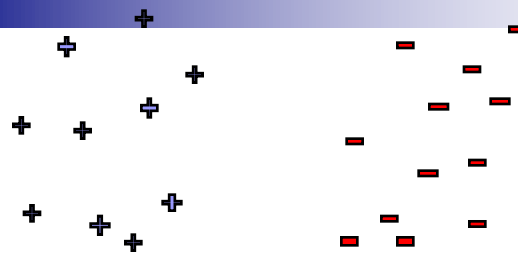
## Support Vector Machines (Two Ideas Mixed up)

- 1) An attempt to better optimize the classification loss?
  - Questionable?
  - Latent SVMs are interesting.
- 2) Kernels
  - Warp the feature space
  - This idea is actually more general
- The success of SVMs?

©2016 Sham Kakade

3

## Linear Separability: More formally, Using Margin



- Data linearly separable, if there exists
  - a vector
  - a margin
- Such that

©2016 Sham Kakade

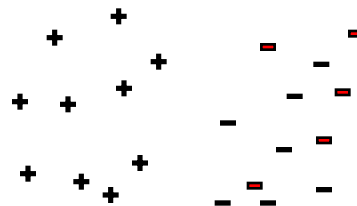
4

## Perceptron Analysis: Linearly Separable Case

- Theorem [Block, Novikoff]:
  - Given a sequence of labeled examples:
  - Each feature vector has bounded norm:
  - If dataset is linearly separable:
- Then the number of mistakes made by the online perceptron on any such sequence is bounded by

## 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



# Kernels

Machine Learning – CSE446

Sham Kakade

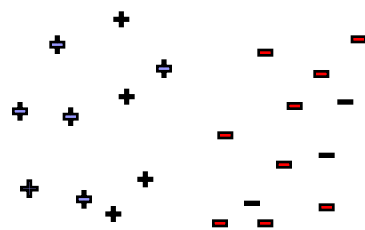
University of Washington

November 1, 2016

©2016 Sham Kakade

7

What if the data is not linearly separable?



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

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

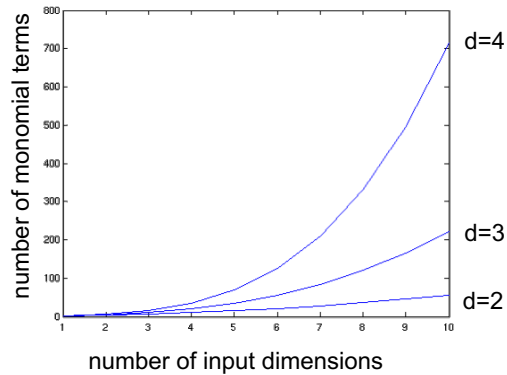
**Feature space can get really large really quickly!**

©2016 Sham Kakade

# 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



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

©2016 Sham Kakade

9

# Perceptron Revisited

- Given weight vector  $w^{(t)}$ , predict point  $x$  by:
- 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:
- Prediction rule now:
- When using high dimensional features:

©2016 Sham Kakade

10

## Dot-product of polynomials

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

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

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

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

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

## Polynomial kernels

- All monomials of degree  $d$  in  $O(d)$  operations:  
 $\Phi(\mathbf{u}) \cdot \Phi(\mathbf{v}) = (\mathbf{u} \cdot \mathbf{v})^d =$  polynomials of degree exactly  $d$
- How about all monomials of degree up to  $d$ ?
  - Solution 0:
  - Better solution:

©2016 Sham Kakade

13

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

©2016 Sham Kakade

14

# Support Vector Machines

Machine Learning – CSE446

Sham Kakade

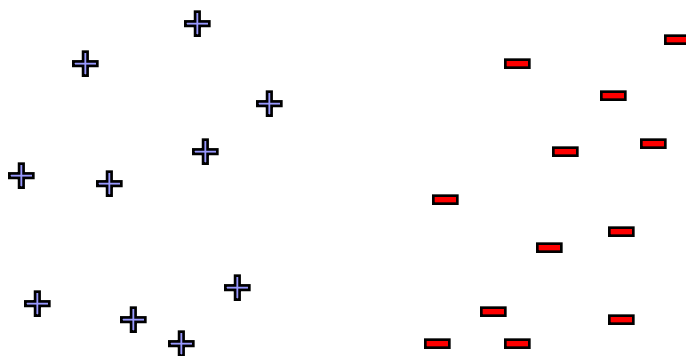
University of Washington

November 1, 2016

©2016 Sham Kakade

15

## Linear classifiers – Which line is better?

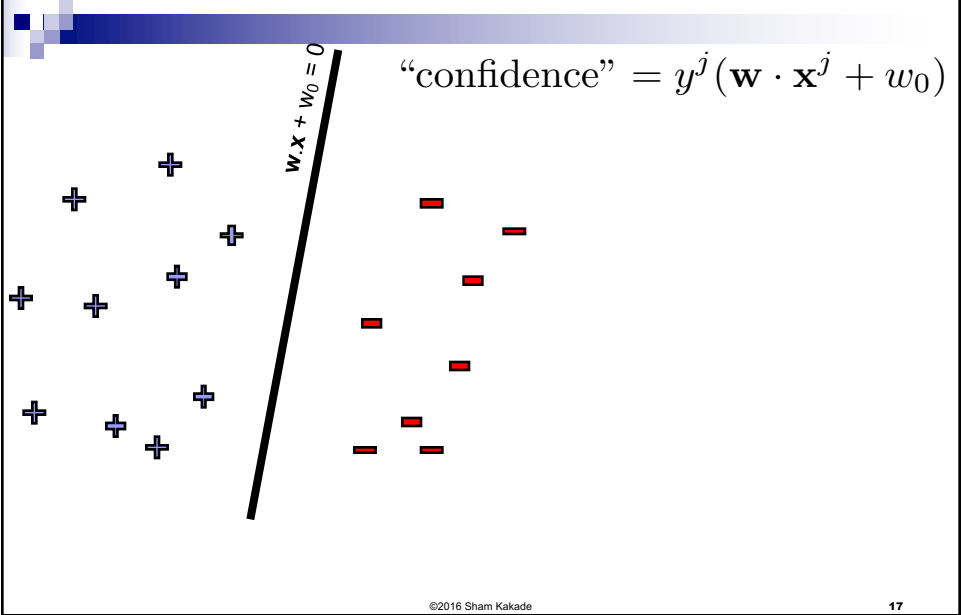


©2016 Sham Kakade

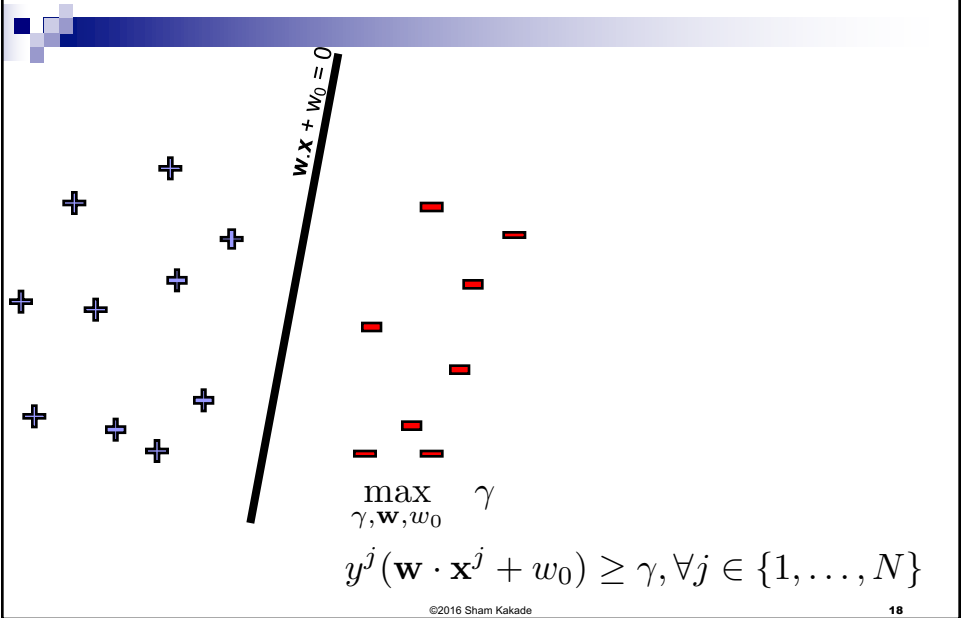
16



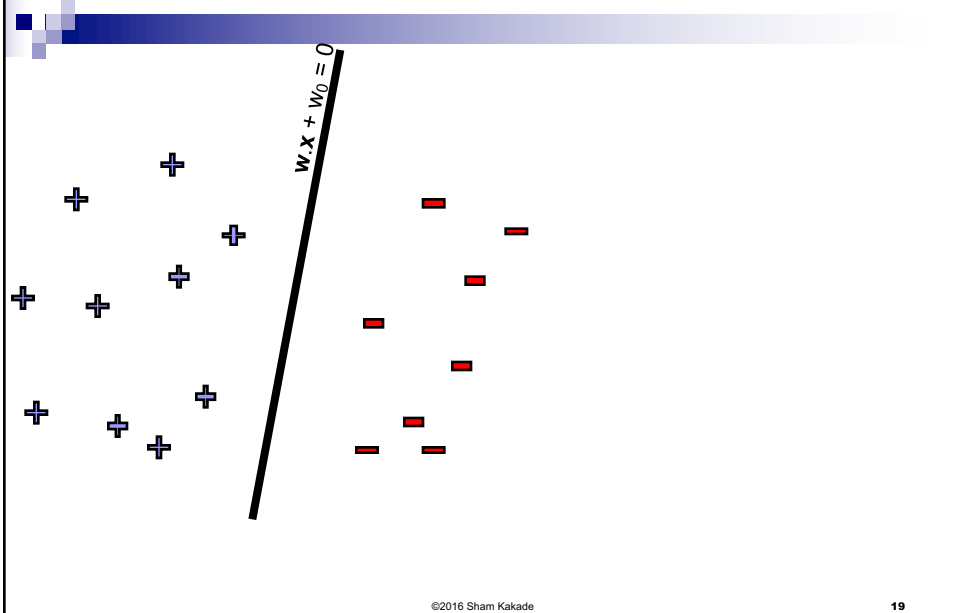
Pick the one with the largest margin!



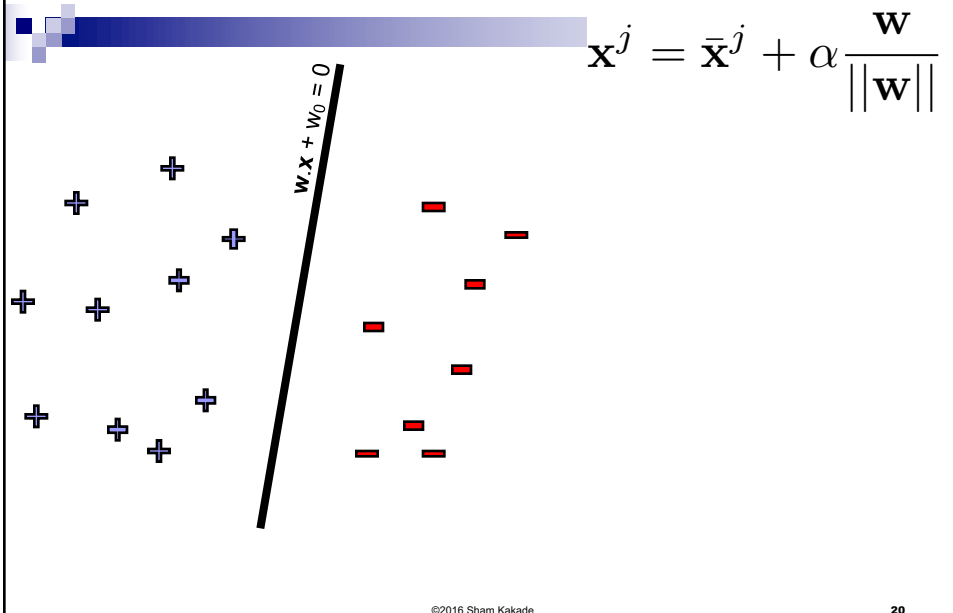
Maximize the margin



# But there are many planes...

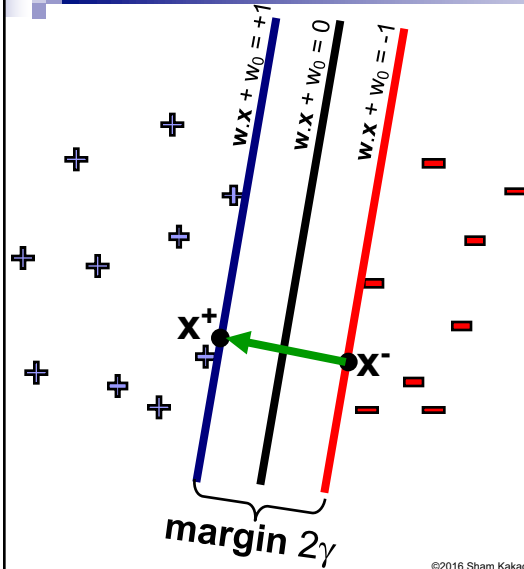


# Review: Normal to a plane



## A Convention: Normalized margin – Canonical hyperplanes

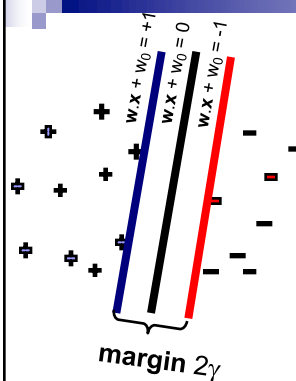
$$\mathbf{x}^j = \bar{\mathbf{x}}^j + \alpha \frac{\mathbf{w}}{\|\mathbf{w}\|}$$



©2016 Sham Kakade

21

## 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:

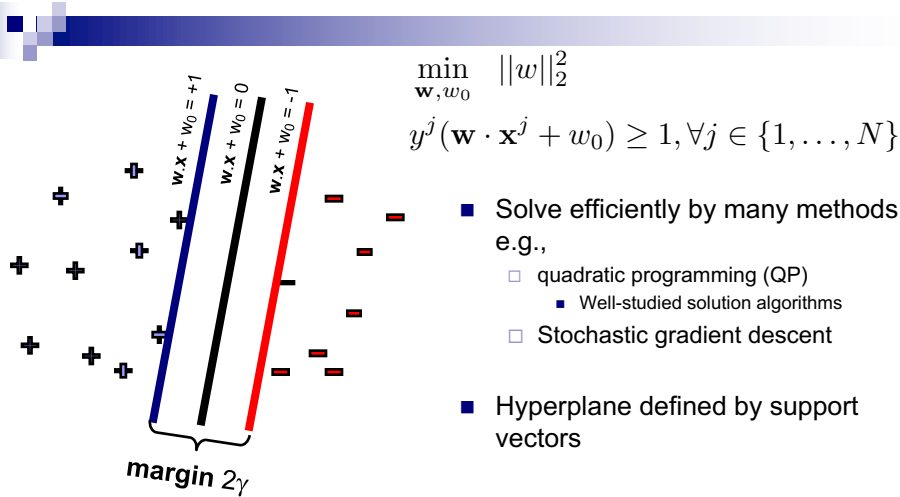
$$\min_{\mathbf{w}, w_0} \|\mathbf{w}\|_2^2$$

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

©2016 Sham Kakade

22

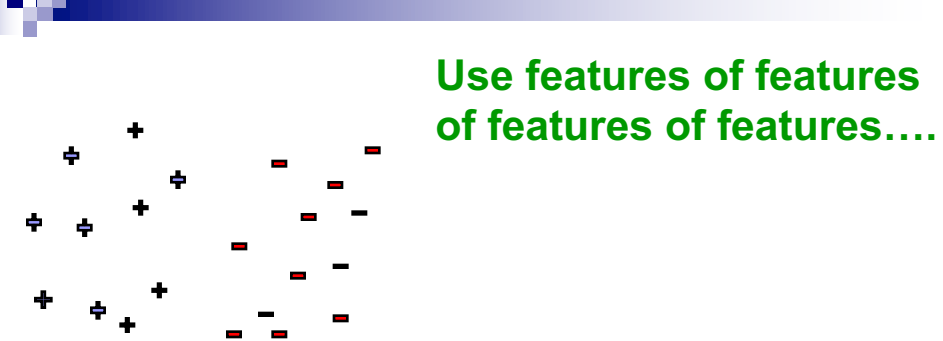
# Support vector machines (SVMs)



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

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

# What if the data is not linearly separable?

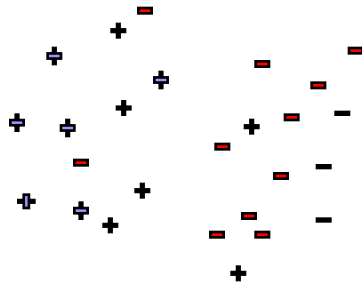


**Use features of features of features of features....**

## What if the data is still not linearly separable?

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

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



- If data is not linearly separable, some points don't satisfy margin constraint:
- How bad is the violation?
- Tradeoff margin violation with  $\|\mathbf{w}\|$ :

©2016 Sham Kakade

25

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

- Perceptron was minimizing the hinge loss:

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

- SVMs minimizes the regularized hinge loss!!

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

©2016 Sham Kakade

26

## Stochastic Gradient Descent for SVMs

- Perceptron minimization:

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

- SGD for Perceptron:

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

- SVMs minimization:

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

- SGD for SVMs:

©2016 Sham Kakade

27

## SVMs vs logistic regression

- We often want probabilities/confidences (logistic wins here)
- For classification loss, they are comparable
- Multiclass setting:
  - Softmax naturally generalizes logistic regression
  - SVMs have
- What about good old least squares?

©2016 Sham Kakade

28

# Multiple Classes

- One can generalize the hinge loss
  - If no error (by some margin) -> no loss
  - If error, penalize what you said against the best
- SVMs vs logistic regression
  - We often want probabilities/confidences (logistic wins here)
  - For classification loss, they are
- Latent SVMs
  - When you have many classes it's difficult to do logistic regression
- 2) Kernels
  - Warp the feature space

# Slack variables – Hinge loss

$$\text{minimize}_{w,b} \quad w \cdot w$$

$$(\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1 - \xi_j, \quad \forall j$$

- If margin  $\geq 1$ , don't care
- If margin  $< 1$ , pay linear penalty

©2016 Sham Kakade 31

# Side note: What's the difference between SVMs and logistic regression?

**SVM:**

$$\text{minimize}_{w,b} \quad w \cdot w + C \sum_j \xi_j$$

$$(\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1 - \xi_j, \quad \forall j$$

$$\xi_j \geq 0, \quad \forall j$$

**Logistic regression:**

$$P(Y = 1 | x, w) = \frac{1}{1 + e^{-(w \cdot x + b)}}$$

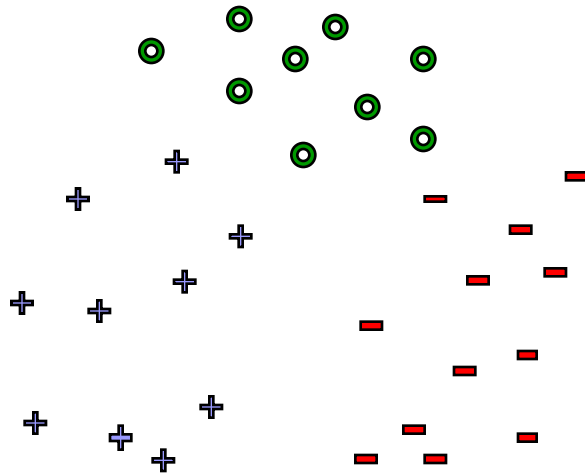
**Log loss:**

$$-\ln P(Y = 1 | x, w) = \ln(1 + e^{-(w \cdot x + b)})$$

©2016 Sham Kakade 32



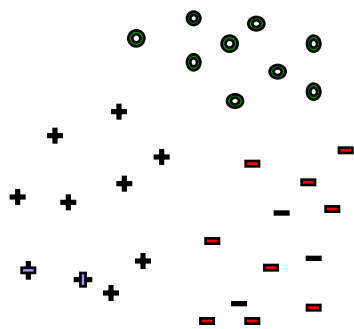
# What about multiple classes?



©2016 Sham Kakade

33

# One against All



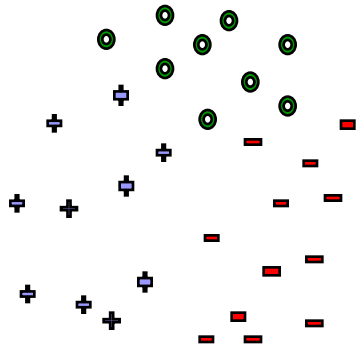
Learn 3 classifiers:

©2016 Sham Kakade

34

## Learn 1 classifier: Multiclass SVM

Simultaneously learn 3 sets of weights



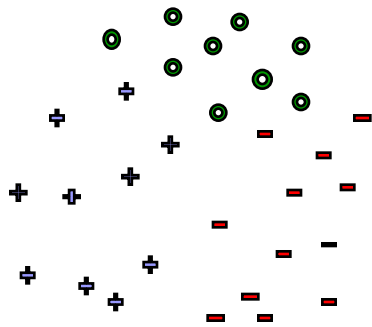
$$w^{(y_j)} \cdot x_j + b^{(y_j)} \geq w^{(y')} \cdot x_j + b^{(y')} + 1, \quad \forall y' \neq y_j, \quad \forall j$$

©2016 Sham Kakade

35

## Learn 1 classifier: Multiclass SVM

$$\begin{aligned} & \text{minimize}_{w,b} \quad \sum_y w^{(y)} \cdot w^{(y)} + C \sum_j \xi_j \\ & w^{(y_j)} \cdot x_j + b^{(y_j)} \geq w^{(y')} \cdot x_j + b^{(y')} + 1 - \xi_j, \quad \forall y' \neq y_j, \quad \forall j \\ & \xi_j \geq 0, \quad \forall j \end{aligned}$$



©2016 Sham Kakade

36