# Support Vector Machines (SVMs)
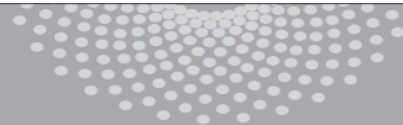
Aarti Singh

Machine Learning 10-315
Oct 21, 2020

# Discriminative Classifiers

Optimal Classifier:

$$f^*(x) = \arg \max_{Y=y} P(Y = y | X = x)$$
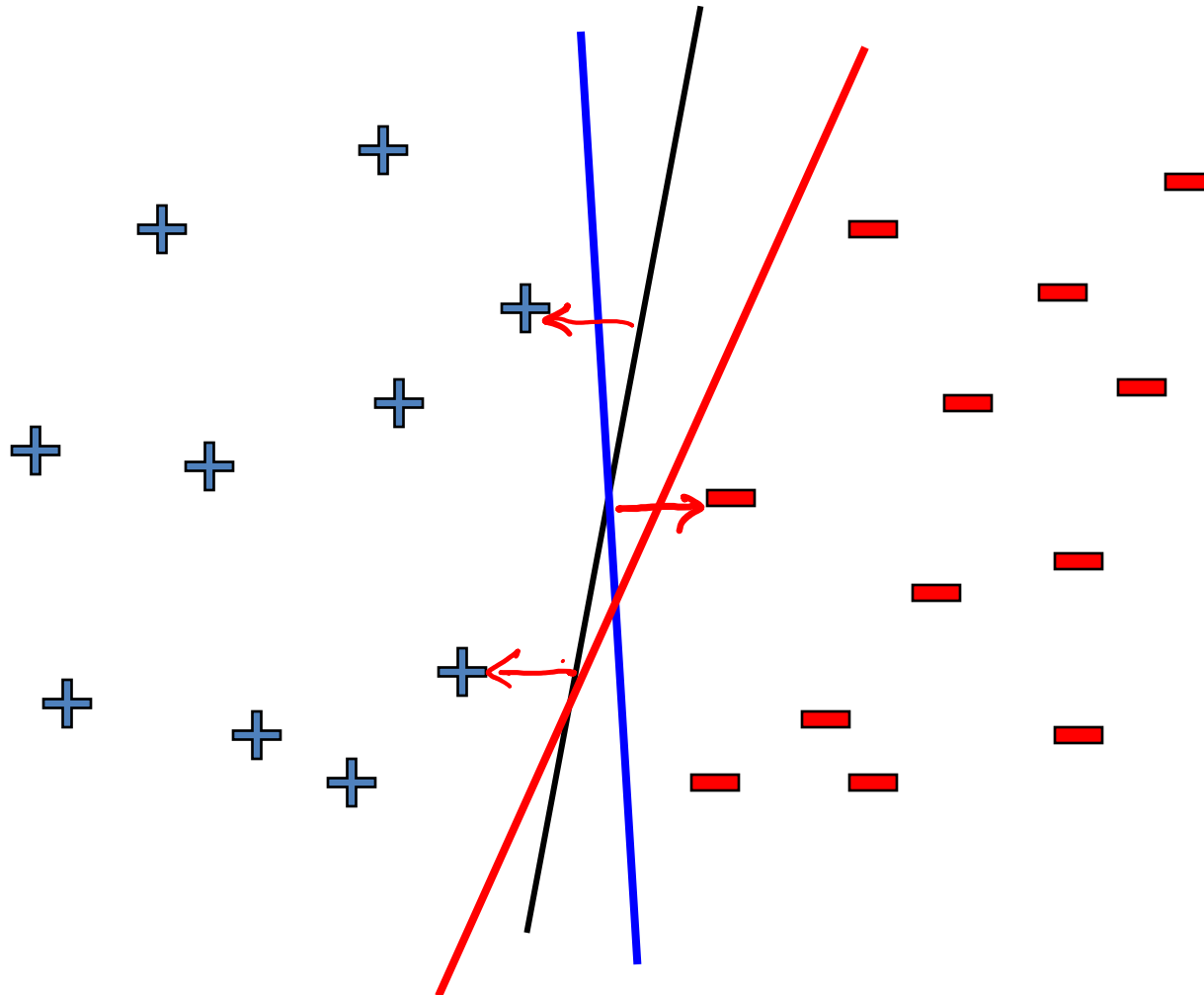
$$= \arg \max_{Y=y} P(X = x | Y = y) P(Y = y)$$

Why not learn P(Y|X) directly? Or better yet, why not learn the decision boundary directly?

- Assume some functional form for P(Y|X) (e.g. Logistic Regression) or for the decision boundary (e.g. Neural nets, SVMs)

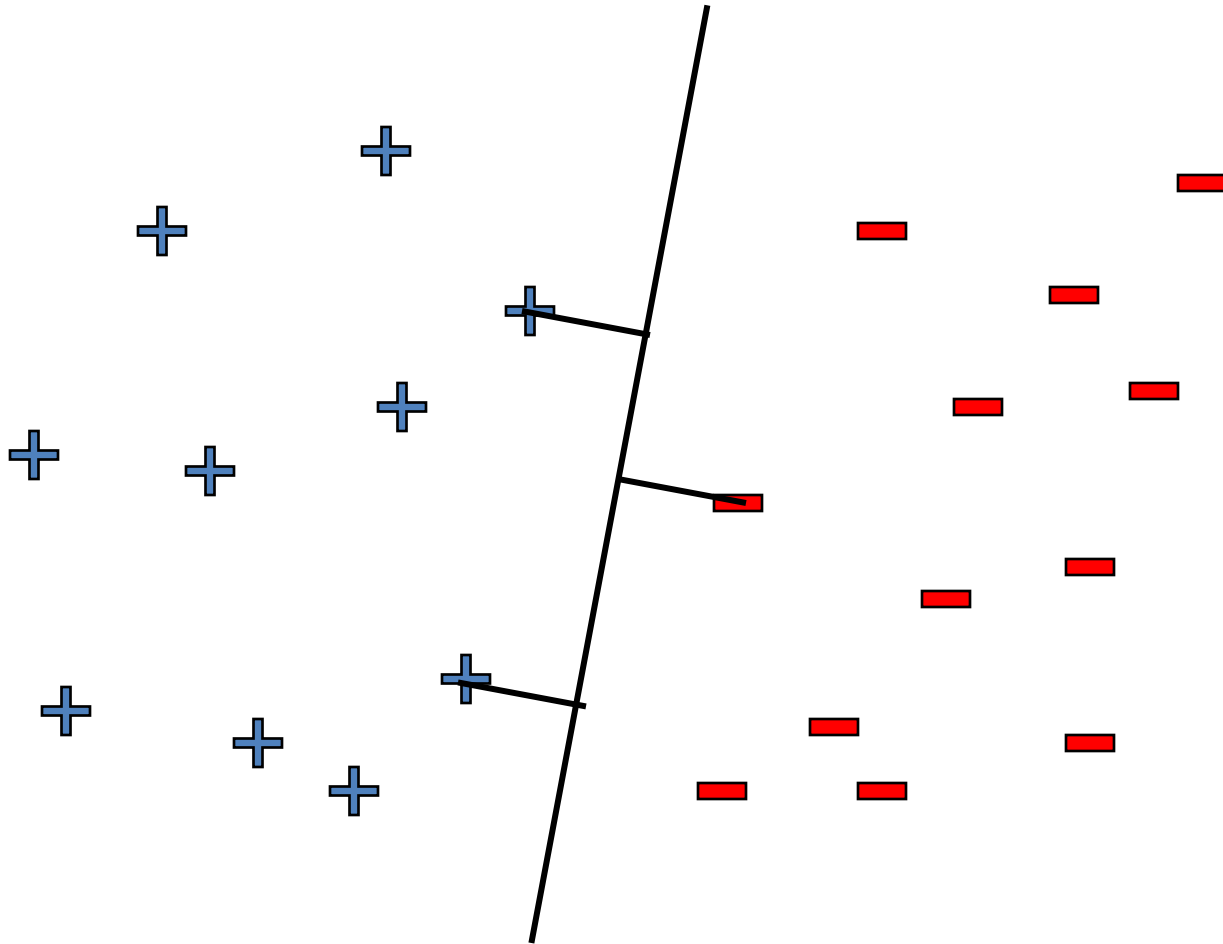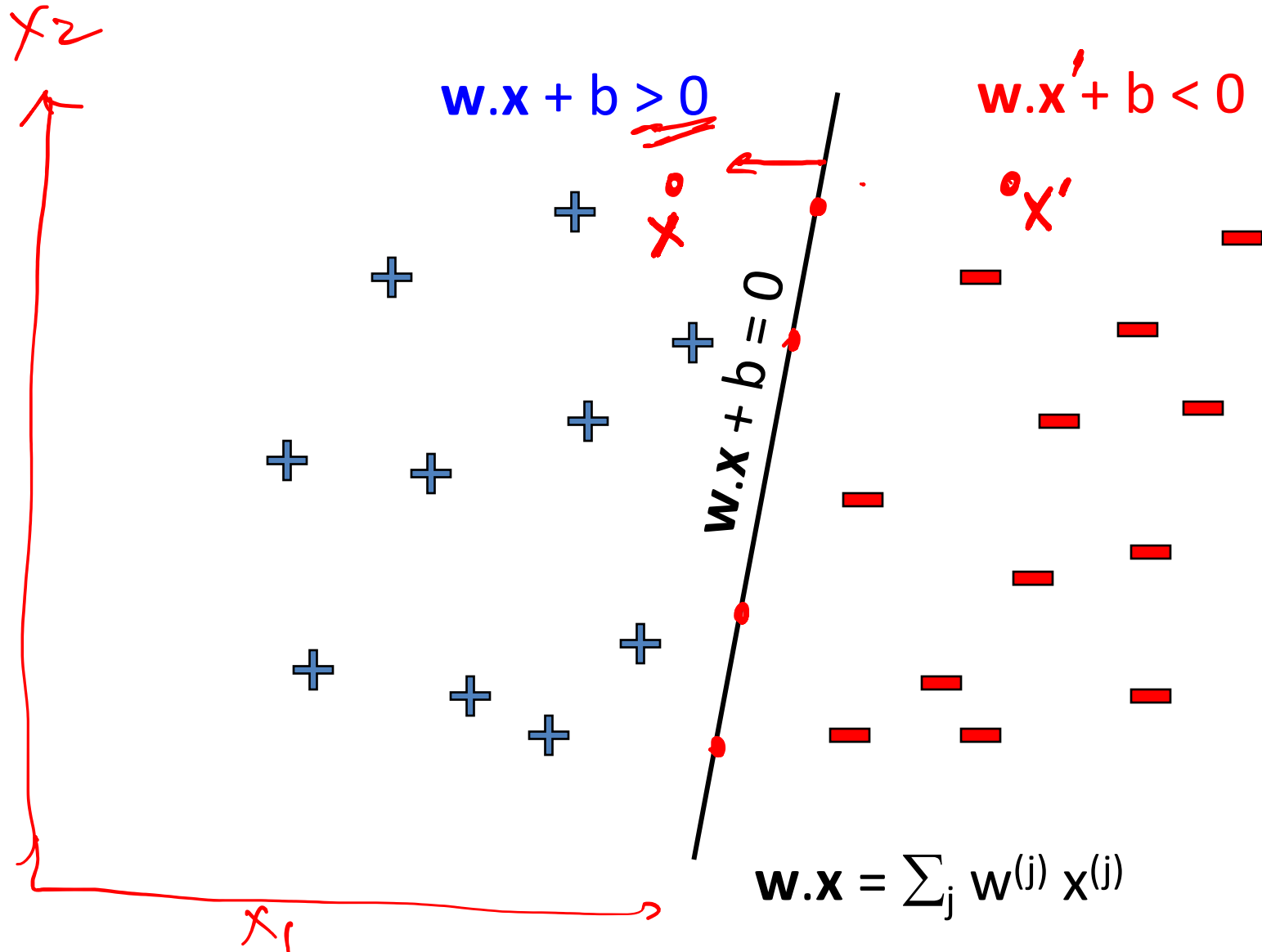- Estimate parameters of functional form directly from training data

# At Pittsburgh G-20 summit …

# Linear classifiers – which line is better?

# Pick the one with the largest margin!

# Parameterizing the decision boundary

$x_2$

$\mathbf{w}.\mathbf{x} + b > 0$

$\mathbf{w}.\mathbf{x}' + b < 0$

$\mathbf{w}.\mathbf{x} + b = 0$

$\mathbf{w}.\mathbf{x} = \sum_j w^{(j)} x^{(j)}$

$x_1$

# Parameterizing the decision boundary

**w.x** + b > 0　　　　**w.x** + b < 0

**w.x** + b = 0

$$y_j \in \{-1, +1\} \quad \text{—— class}$$

$$\text{"confidence"} = \left(\mathbf{w}.\mathbf{x}_j + b\right) y_j \quad \pm 1$$

# Maximizing the margin

**w.x** + b > 0          **w.x** + b < 0

**w.x**₊ + b = a

**w.x** + b = 0

**w.x**₋ + b = -a

Distance of closest examples from the line/hyperplane

$$\text{margin} = \gamma = 2a/\|w\|$$

Step 1: **w** is perpendicular to lines since for any $x_1$, $x_2$ on line  **w**.$(x_1 - x_2) = 0$

$x_1$
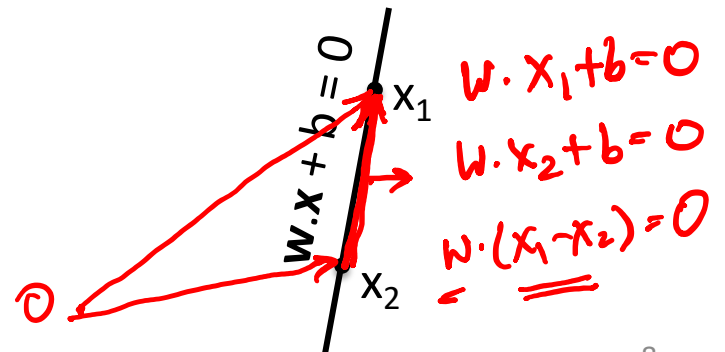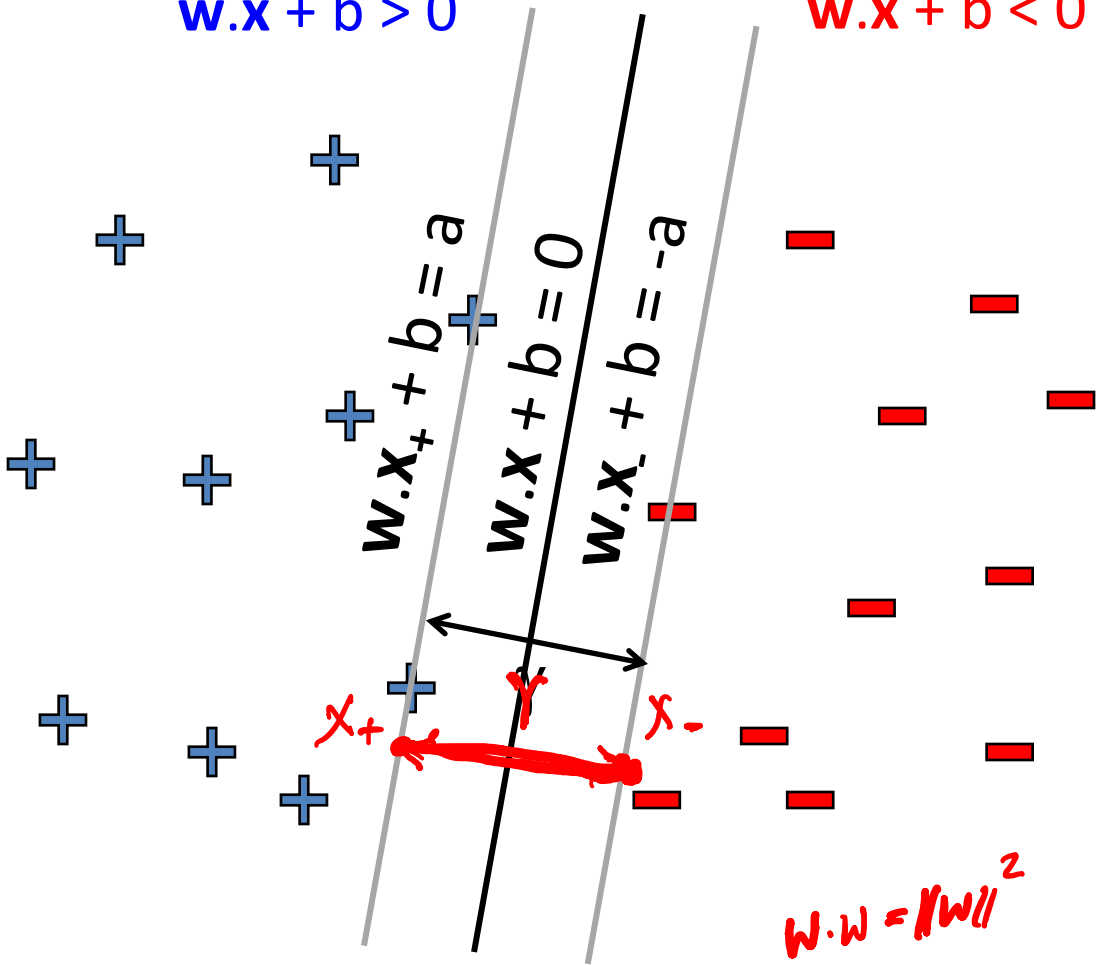
$x_2$

w.x + b = 0

$W \cdot X_1 + b = 0$

$W \cdot X_2 + b = 0$

$W \cdot (X_1 - X_2) = 0$

O

8

# Maximizing the margin

$\gamma = 2a/\|w\|$

**w.x** + b > 0          **w.x** + b < 0

$$\boxed{\text{margin} = \gamma = 2a/\|w\|}$$

**w.x$_+$** + b = a

**w.x** + b = 0

**w.x$_-$** + b = -a

$x_+$    $\gamma$    $x_-$

$w.w = \|w\|^2$

Step1: **w** is perpendicular to lines

Step 2: Take a point x$_-$ on w.x$_-$ +b = -a and move to point x$_+$ that is $\gamma$ away on line w.x+b = a

$$\mathbf{x}_+ = \mathbf{x}_- + \gamma \mathbf{w}/\|w\|$$

$$\mathbf{w.x}_+ = \mathbf{w.x}_- + \gamma \mathbf{w. w}/\|w\|$$

$$a - b = -a - b + \gamma \|w\|$$

$$2a = 2a/\|w\| \;\gamma\; \|w\|$$

9

# Maximizing the margin

$\mathbf{w.x} + b > 0$    $\mathbf{w.x} + b < 0$

$\mathbf{w.x_+} + b = a$

$\mathbf{w.x} + b = 0$

$\mathbf{w.x_-} + b = -a$

$\gamma$

Distance of closest examples from the line/hyperplane

margin = $\gamma$ = 2a/‖w‖

Smaller margin ⇔ larger ‖w‖

# Maximizing the margin

$w \cdot x + b = 0$

$\frac{w}{a} \cdot x + \frac{b}{a} = 0$

$\mathbf{w.x} + b > 0$

$(w \cdot x + b) y_a \geq a$

$+$

$\mathbf{w.x} + b < 0$

$(w \cdot x + b) \cdot y \geq a$

$i$

$\mathbf{w.x}_+ + b = a$

$\mathbf{w.x} + b = 0$

$\mathbf{w.x}_- + b = -a$

$\gamma$

Distance of closest examples
from the line/hyperplane

margin = $\gamma$ = 2a/‖w‖

$\max_{\mathbf{w},b} \; \gamma = 2a/‖w‖$

s.t. $(\mathbf{w.x}_j + b) \, y_j \geq a \quad \forall j$

Note: 'a' is arbitrary (can normalize
equations by a)

11

# Support Vector Machines

$\mathbf{w}.\mathbf{x} + b > 0$     $\mathbf{w}.\mathbf{x} + b < 0$

$\mathbf{w}.\mathbf{x}_+ + b = 1$

$\mathbf{w}.\mathbf{x} + b = 0$

$\mathbf{w}.\mathbf{x}_- + b = -1$

$\gamma$

$$\min_{\mathbf{w},b} \ \mathbf{w}.\mathbf{w}$$

$\|w\|^2$

$$\text{s.t. } (\mathbf{w}.\mathbf{x}_j + b) \ y_j \geq 1 \ \ \forall j$$

Solve efficiently by quadratic programming (QP)

– Quadratic objective, linear constraints

– Well-studied solution algorithms

# Support Vectors

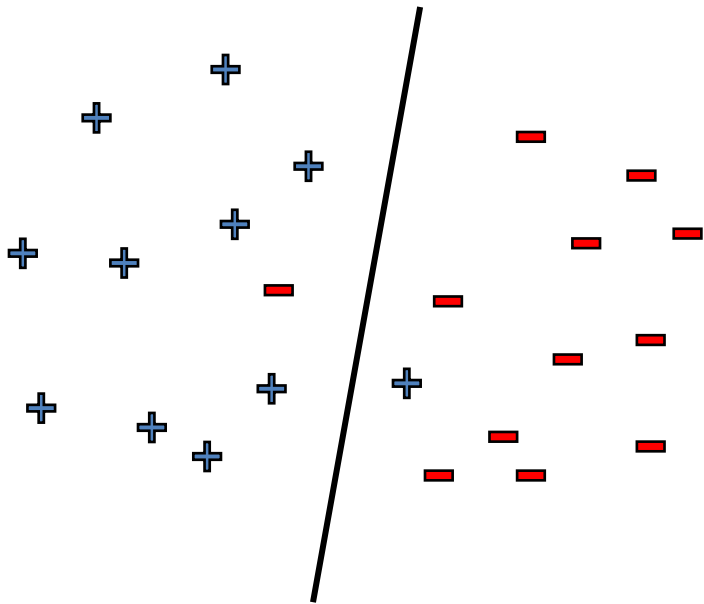$\mathbf{w.x} + b > 0$          $\mathbf{w.x} + b < 0$

Linear hyperplane defined by "support vectors"

Moving other points a little doesn't effect the decision boundary

only need to store the support vectors to predict labels of new points

For support vectors
$(\mathbf{w.x}_j + b)\, y_j = 1$
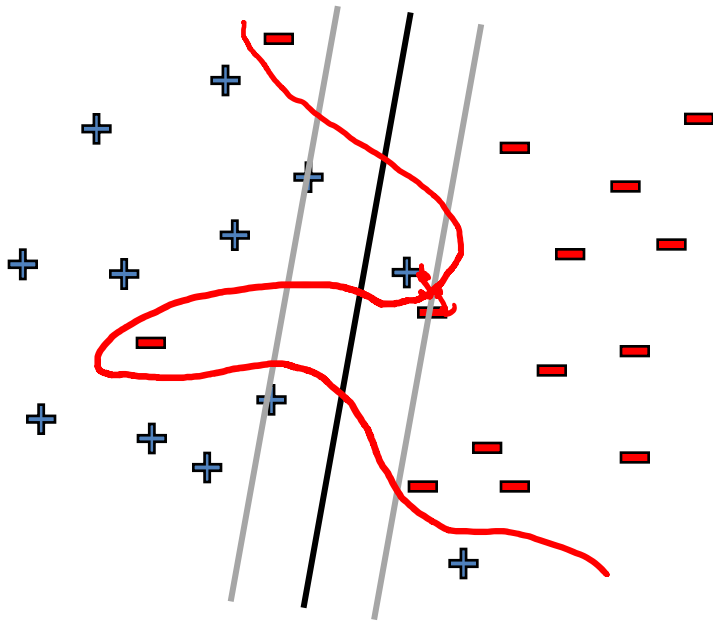
# What if data is not linearly separable?

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

$x_1^2, x_2^2, x_1x_2, ...., \exp(x_1)$

But run risk of overfitting!

# What if data is still not linearly separable?

Allow "error" in classification

$$\sum_j 1_{(w \cdot x_j + b) \, y_j > 0}$$

$$\min_{w,b} \; \mathbf{w}.\mathbf{w} + C \, \#\text{mistakes}$$

$$\text{s.t. } (\mathbf{w}.\mathbf{x}_j + b) \, y_j \geq 1 \quad \forall j$$
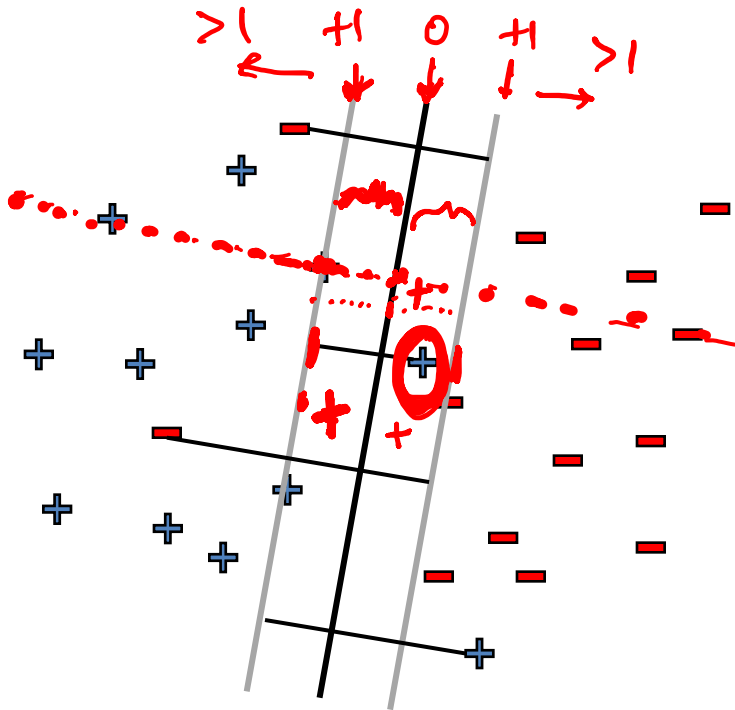
Maximize margin and minimize # mistakes on training data
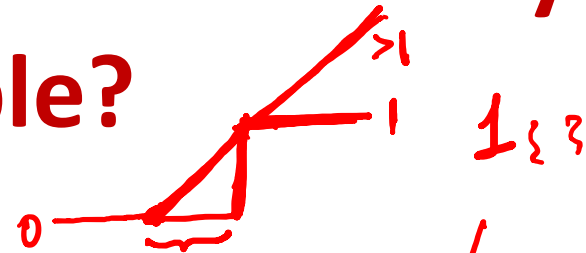
C - tradeoff parameter

Not QP ☹

0/1 loss (doesn't distinguish between near miss and bad mistake)

Smaller margin ⇔ larger ‖w‖

# What if data is still not linearly separable?

Allow "error" in classification



**Soft margin approach**

$$\min_{\mathbf{w},b,\{\xi_j\}} \mathbf{w}.\mathbf{w} + C \sum_j \xi_j$$

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

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

$\xi_j$  - "slack" variables
    $= (>1$ if $x_j$ misclassifed)
pay linear penalty if mistake

C  - tradeoff parameter (chosen by cross-validation)

Still QP ☺

# Soft-margin SVM

$1-\xi_j < 0$

$\xi_j > 1$

$(\mathbf{w}\cdot x + b)\, y \geq 1$

$0 < \xi_j < 1$

$\xi_j = 0$

$\xi_j > 1$

$\xi_j > 1$

$\xi_j = 0$

$\xi_j = 0$

$\xi_j = 0$

$\xi_j > 1$

**w.x + b = 1**

**w.x + b = -1**

Soften the constraints:

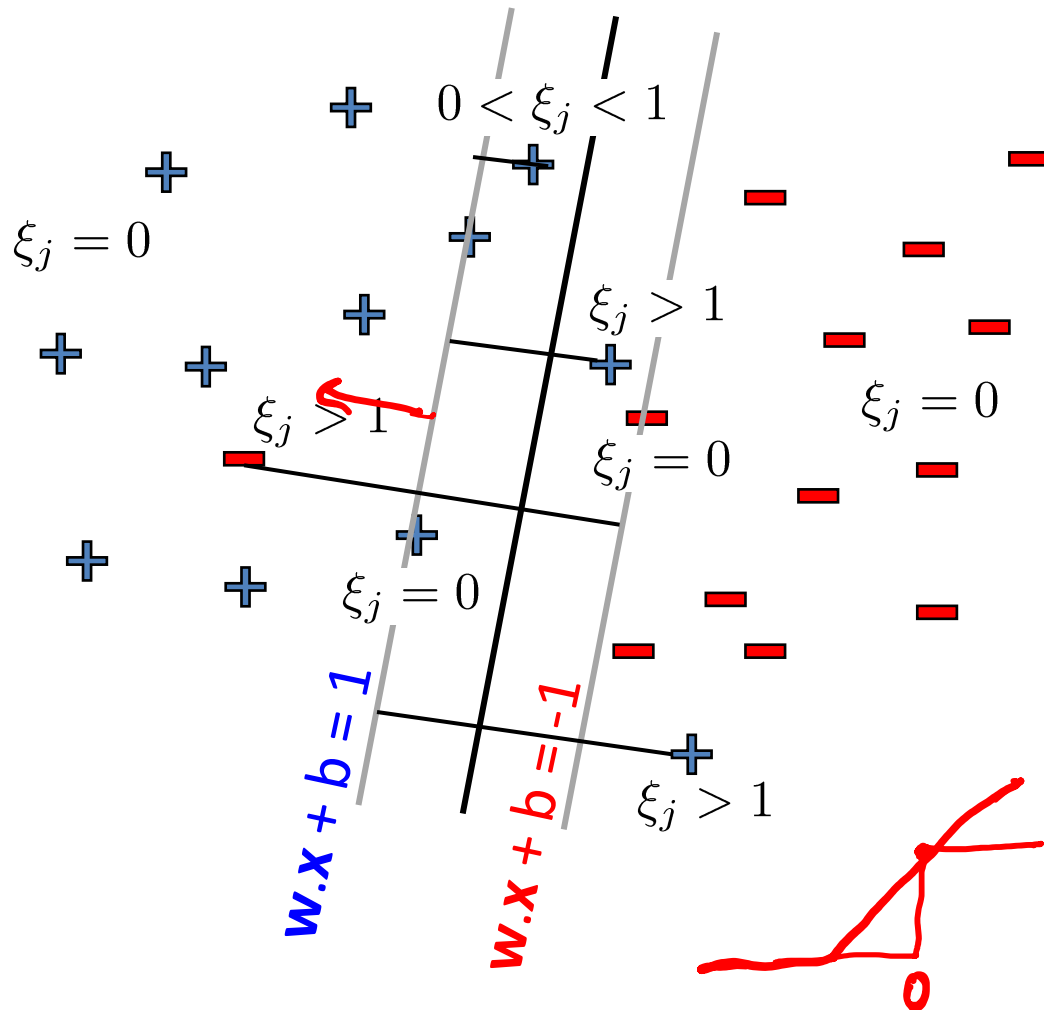$$(\mathbf{w}\cdot\mathbf{x}_j + b)\, y_j \geq 1 - \xi_j \quad \forall j$$

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

Penalty for misclassifying:

$$C\,\xi_j$$

How do we recover hard margin SVM?

Set C = ∞

17

# Slack variables – Hinge loss

$0 < \xi_j < 1$

$\xi_j = 0$

$\xi_j > 1$

$\xi_j > 1$

$\xi_j = 0$

$\xi_j = 0$

$\xi_j = 0$

$\xi_j = 0$

$\mathbf{w}.\mathbf{x} + b = 1$

$\mathbf{w}.\mathbf{x} + b = -1$

$\xi_j > 1$

$\|w\|^2 + C \sum_j \xi_j$

$B_+ = \begin{cases} B & \text{if } B > 0 \\ 0 & \text{o.w.} \end{cases}$

Notice that

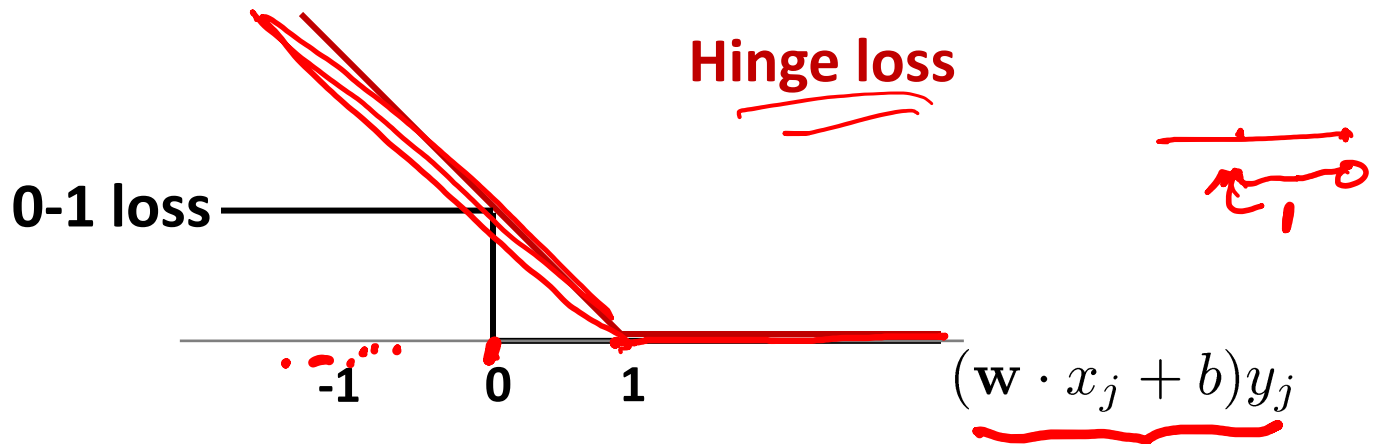$$\xi_j = \overbrace{(1 - \underbrace{(\mathbf{w} \cdot x_j + b)y_j}_{\text{Confidence}}))}_{}{}_+$$

$$= \begin{cases} 1 & \text{''} 0 \\ 1 + \cdots & < 0 \\ 0 & > 1 \end{cases}$$

18

# Slack variables – Hinge loss

$$\xi_j = (1 - (\mathbf{w} \cdot x_j + b)y_j))_+$$

**Hinge loss**

**0-1 loss**

**-1**     **0**     **1**     $(\mathbf{w} \cdot x_j + b)y_j$

$$\min_{\mathbf{w},b,\{\xi_j\}} \mathbf{w}.\mathbf{w} + C \sum_j \xi_j$$

s.t. $(\mathbf{w}.\mathbf{x}_j + b)\, y_j \geq 1 - \xi_i \quad \forall j$
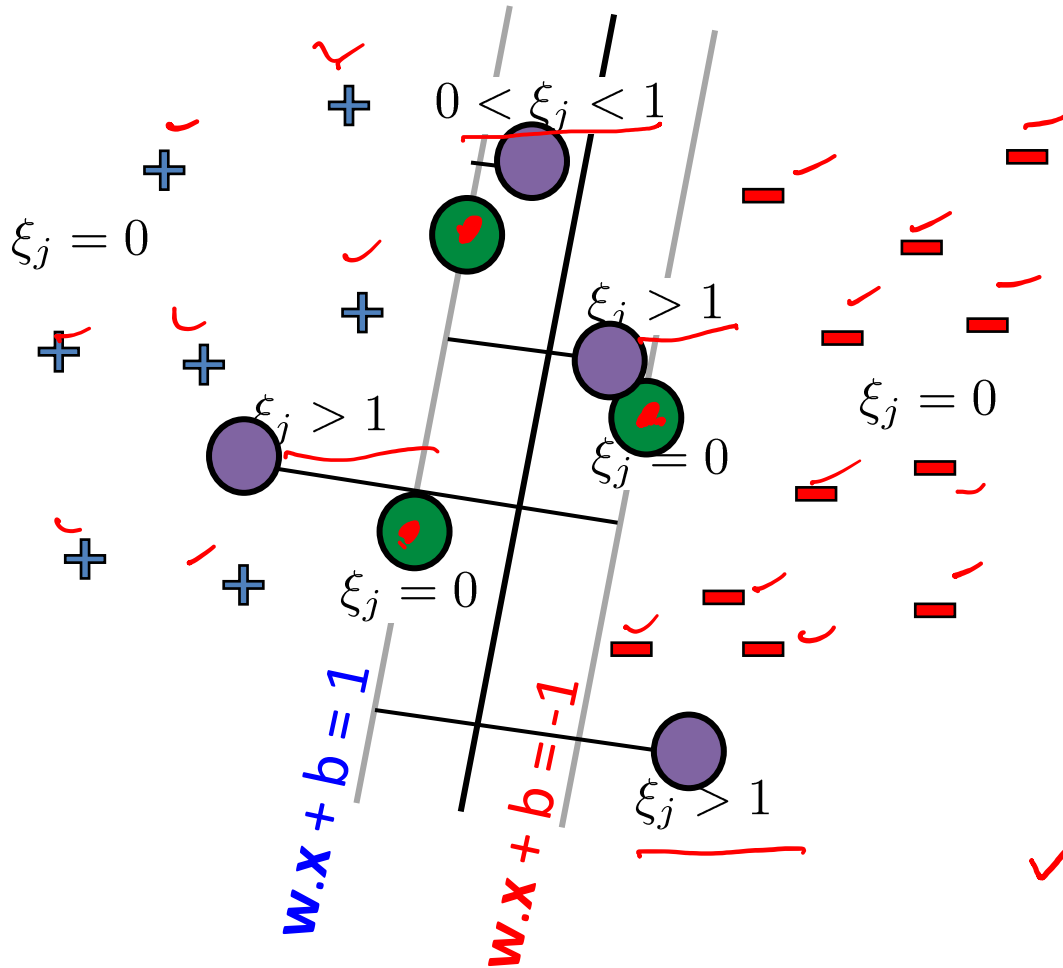
$\xi_j \geq 0 \quad \forall j$

⟺

Regularized hinge loss

$$\min_{\mathbf{w},b} \mathbf{w}.\mathbf{w} + C \sum_j (1 - (\mathbf{w}.x_j + b)y_j)_+$$

regularization parameter

$\xi_j$

# Support Vectors

← any training points that affect decision boundary $(w, b)$



$0 < \xi_j < 1$

$\xi_j = 0$

$\xi_i > 1$

$\xi_j > 1$

$\xi_j = 0$

$\xi_j = 0$

$\xi_j = 0$

$\xi_j > 1$

$\mathbf{w}.\mathbf{x} + b = 1$

$\mathbf{w}.\mathbf{x} + b = -1$

**Margin support vectors**

$\xi_j = 0$, $(\mathbf{w}.\mathbf{x}_j + b)\, y_j = 1$ ✓

(don't contribute to objective but enforce constraints on solution)

Correctly classified but on margin

**Non-margin support vectors**

$\xi_j > 0$    $(w \cdot x_j + b)\, y_j < 1$

(contribute to both objective and constraints)

✓→ $1 > \xi_j > 0$   Correctly classified but inside margin

✓→ $\xi_j > 1$ Incorrectly classified

20

# SVM vs. Logistic Regression

$$\|w\|^2 + C \text{ hingeloss}$$

SVM : **Hinge loss**

$$\text{loss}(f(x_j), y_j) = (1 - (\mathbf{w} \cdot x_j + b)y_j))_+$$

Logistic Regression : **Log loss** ( -ve log conditional likelihood)

$$\text{loss}(f(x_j), y_j) = -\log P(y_j \mid x_j, \mathbf{w}, b) = \log(1 + e^{-(\mathbf{w} \cdot x_j + b)y_j})$$

**Log loss**    **Hinge loss**

*Surrogate losses*

**0-1 loss**

**-1    0    1**    $(\mathbf{w} \cdot x_j + b)y_j$