

Support Vector Machines (SVMs) Recap...

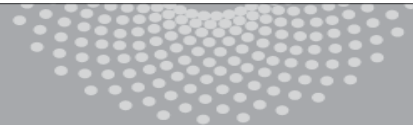
Aarti Singh

Machine Learning 10-315

Oct 26, 2020



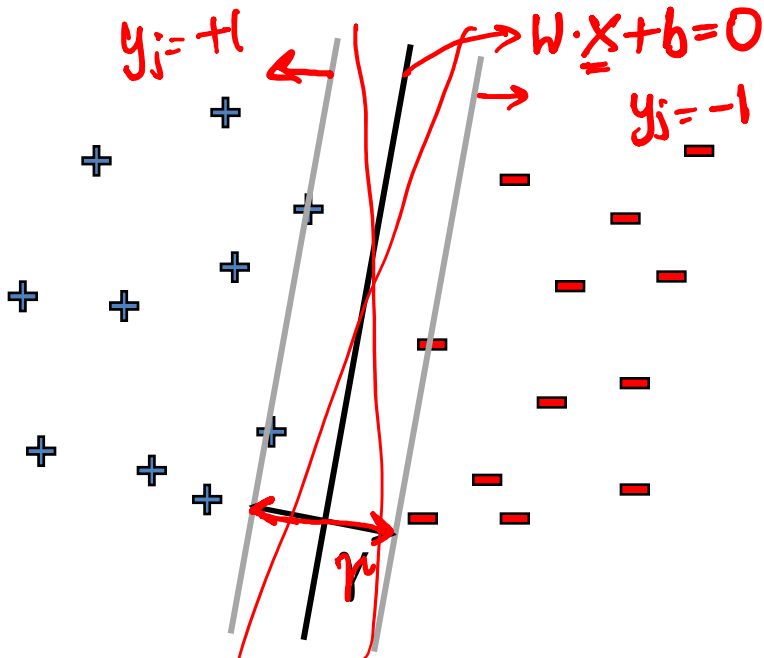
MACHINE LEARNING DEPARTMENT



Carnegie Mellon.
School of Computer Science

Hard-margin SVM

Data perfectly separable by a linear decision boundary



Hard margin approach

$$\min_{w,b} w \cdot w$$

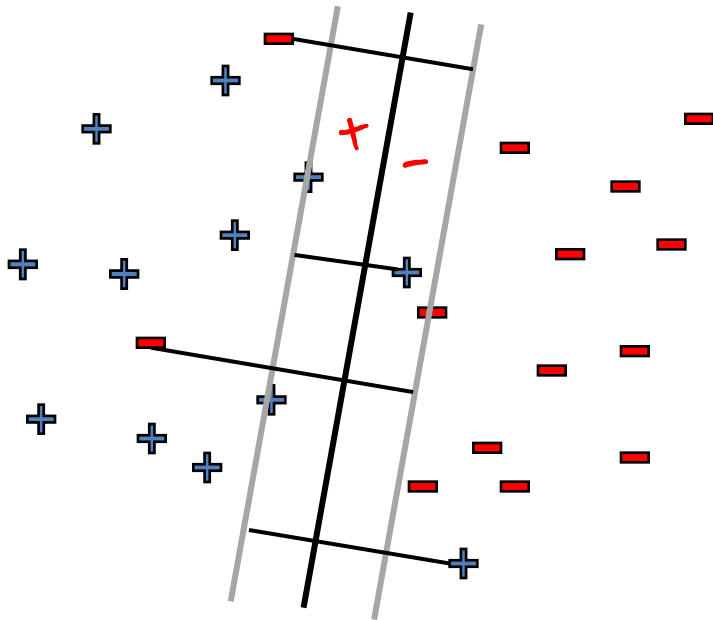
$$\text{s.t. } \underbrace{(w \cdot x_j + b)}_{\text{confidence}} y_j \geq 1 \quad \forall j$$

Solve using Quadratic Programming (QP)

$$\text{Margin, } \gamma \quad \propto \quad \frac{1}{\|w\|}$$

Soft-margin SVM

Allow “error” in classification



Soft margin approach

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_j\}} \quad & \mathbf{w} \cdot \mathbf{w} + C \sum_j \xi_j \\ \text{s.t.} \quad & (\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1 - \xi_j \quad \forall j \\ & \xi_j \geq 0 \quad \forall j \end{aligned}$$

penalty for misclassified

ξ_j - “slack” variables
= (>1 if x_j misclassified)
pay linear penalty if mistake

$C \rightarrow \infty$
C - tradeoff parameter (chosen by cross-validation)

Still QP 😊

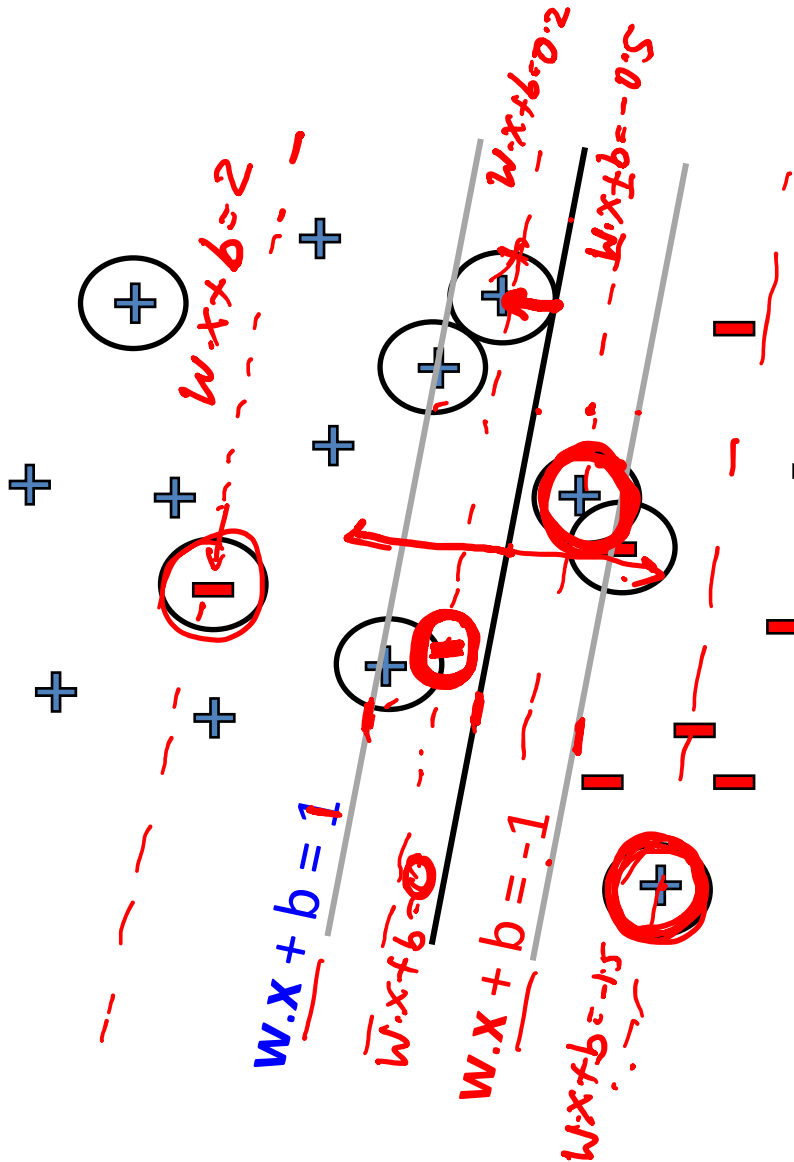
Slack variables – Hinge loss

$\rightarrow \min_{w, b, \xi} \|w\|^2 + C \sum_j \xi_j$

$\rightarrow \xi_j \geq 0$

$\rightarrow (w \cdot x_j + b) y_j \geq 1 - \xi_j \quad \forall j$
 confidence

What is the slack ξ_j for the following points?



Confidence | Slack

1 | 0 ✓

$\rightarrow > 1$ | 0 ✓

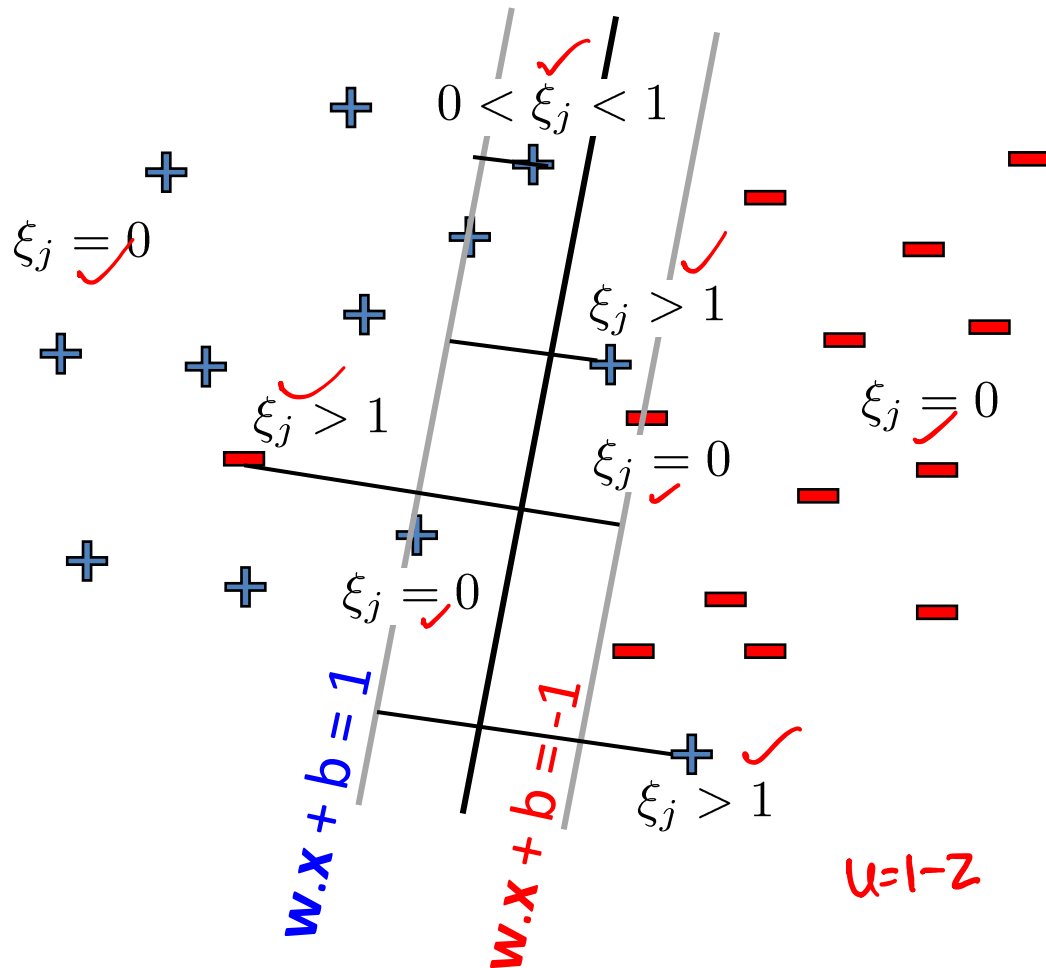
0.7
 $0.5 \geq 1 - \xi_j$ | 0 - 1 (C)

$\xi_j \geq 0.5$ 0.3 $\rightarrow < 0$ (C) | > 1 ✓ (1-C)

$(C) - 0.2$
 $-0.5 \geq 1 - \xi_j$ | 1.2

$\xi_j \geq 1 - (-0.5) = 1.5$ 1.2

Slack variables – Hinge loss



Notice that

$$\xi_j = \underbrace{(1 - (\mathbf{w} \cdot \mathbf{x}_j + b)y_j)}_c$$

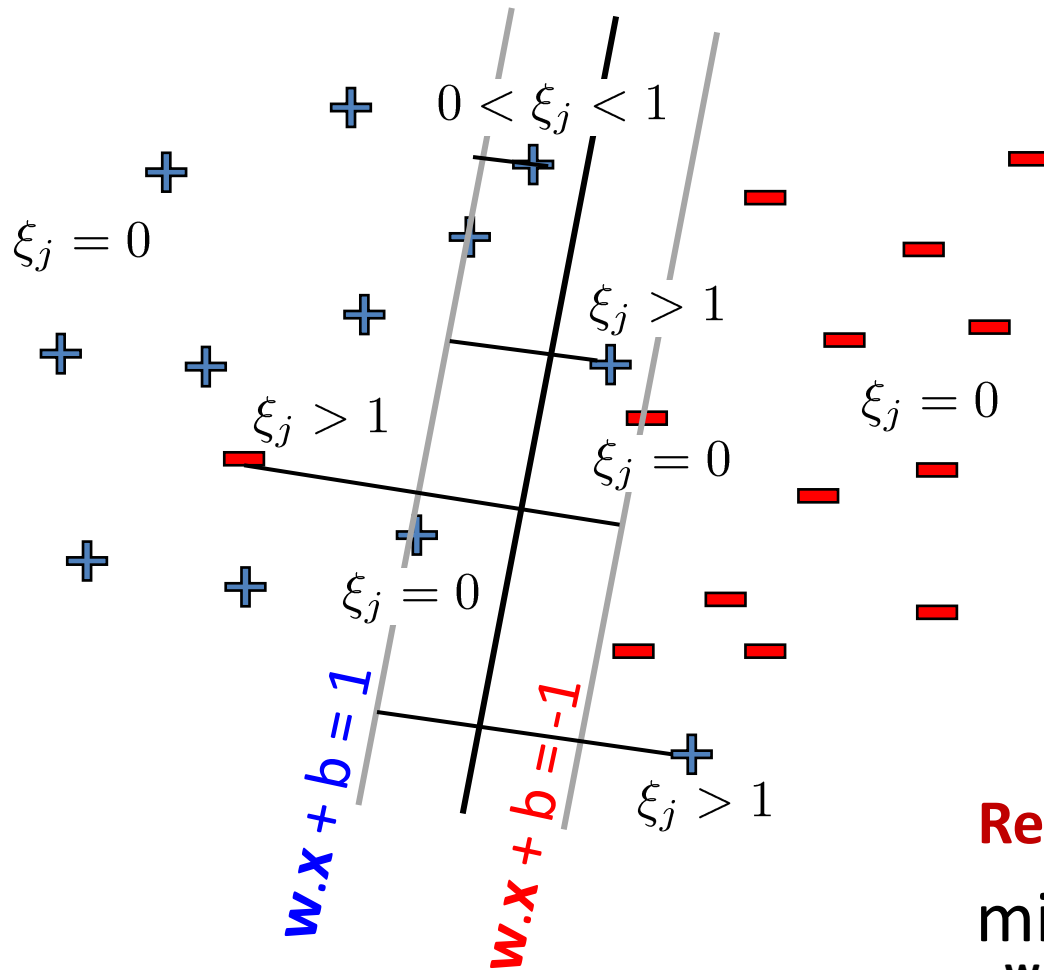
$$= \begin{cases} 1 - (\mathbf{w} \cdot \mathbf{x}_j + b)y_j & \text{if } c \leq 1 \\ 0 & \text{if } c > 1 \end{cases}$$

$$(1-z)_+ = \begin{cases} 1-z & \text{if } z \leq 1 \\ 0 & \text{if } z > 1 \end{cases}$$

$$u_+ = \begin{cases} u & \text{if } u \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

$$u = 1 - z$$

Slack variables – Hinge loss

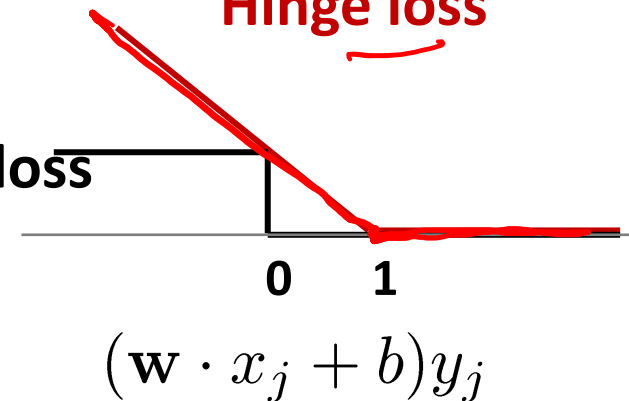


Notice that

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

Hinge loss

0-1 loss



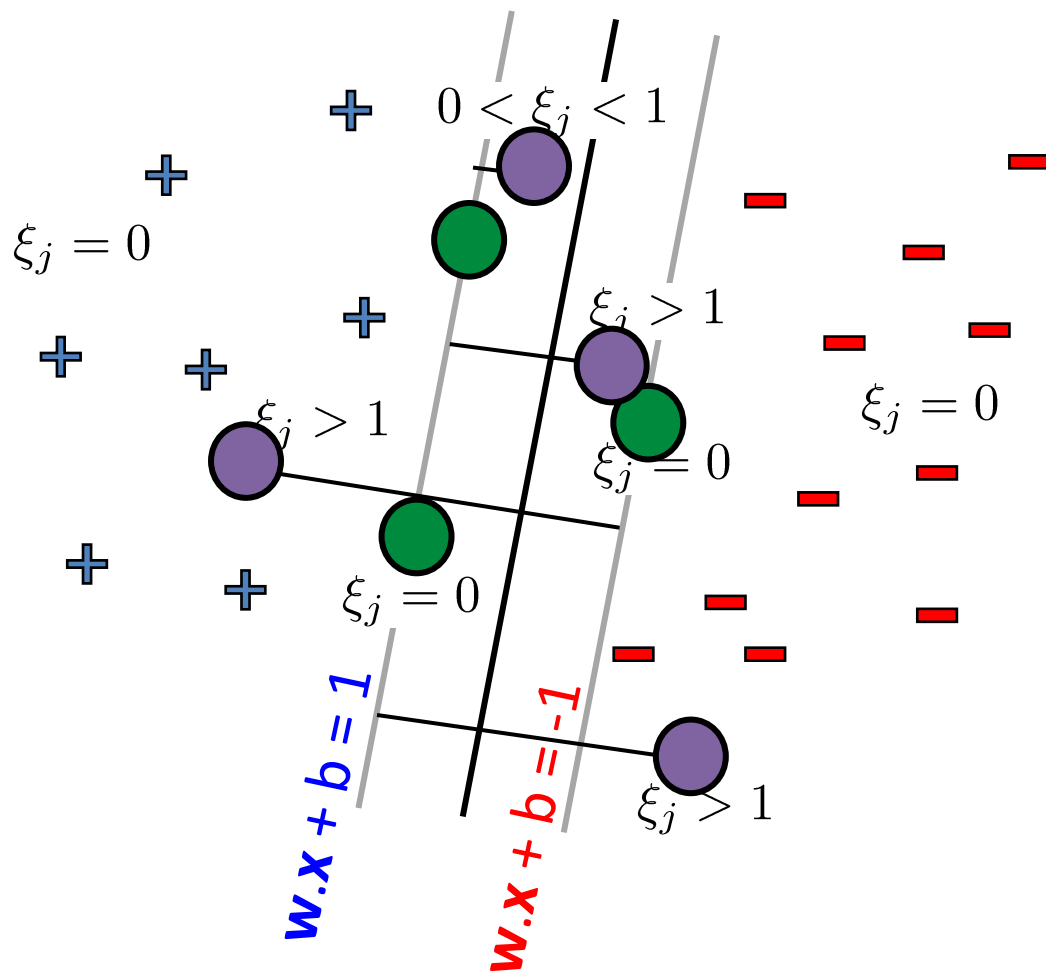
Regularized Hinge loss

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

$w \cdot w + C \sum_j \xi_j$

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_j\}} \quad & \mathbf{w} \cdot \mathbf{w} + C \sum \xi_j \\ \text{s.t.} \quad & (\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1 - \xi_j \quad \forall j \\ & \xi_j \geq 0 \quad \forall j \end{aligned}$$

Support Vectors



Margin support vectors

→ $\xi_j = 0$, $(\mathbf{w} \cdot \mathbf{x}_j + b) y_j = 1$
 (don't contribute to objective but enforce constraints on solution)

Correctly classified but on margin

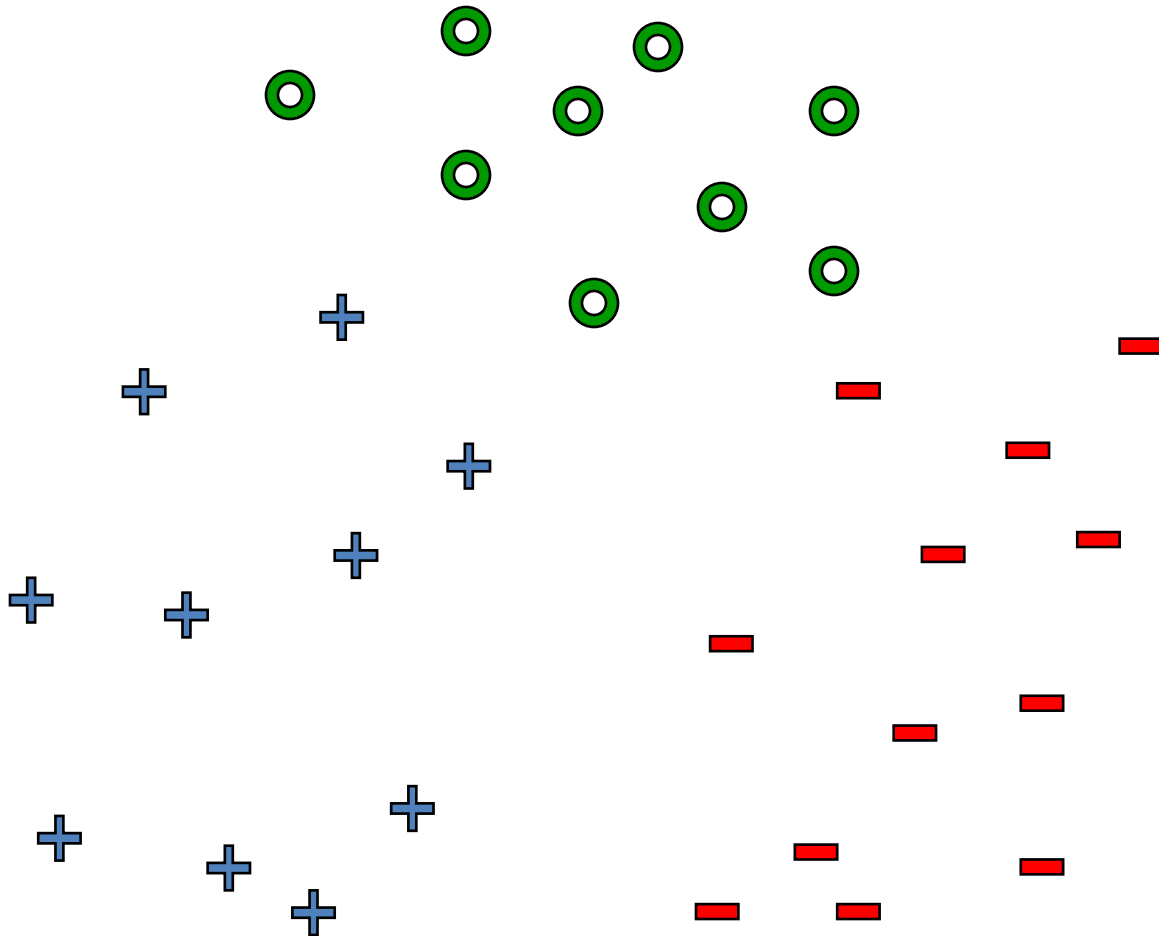
Non-margin support vectors

$\xi_j > 0$
 (contribute to both objective and constraints)

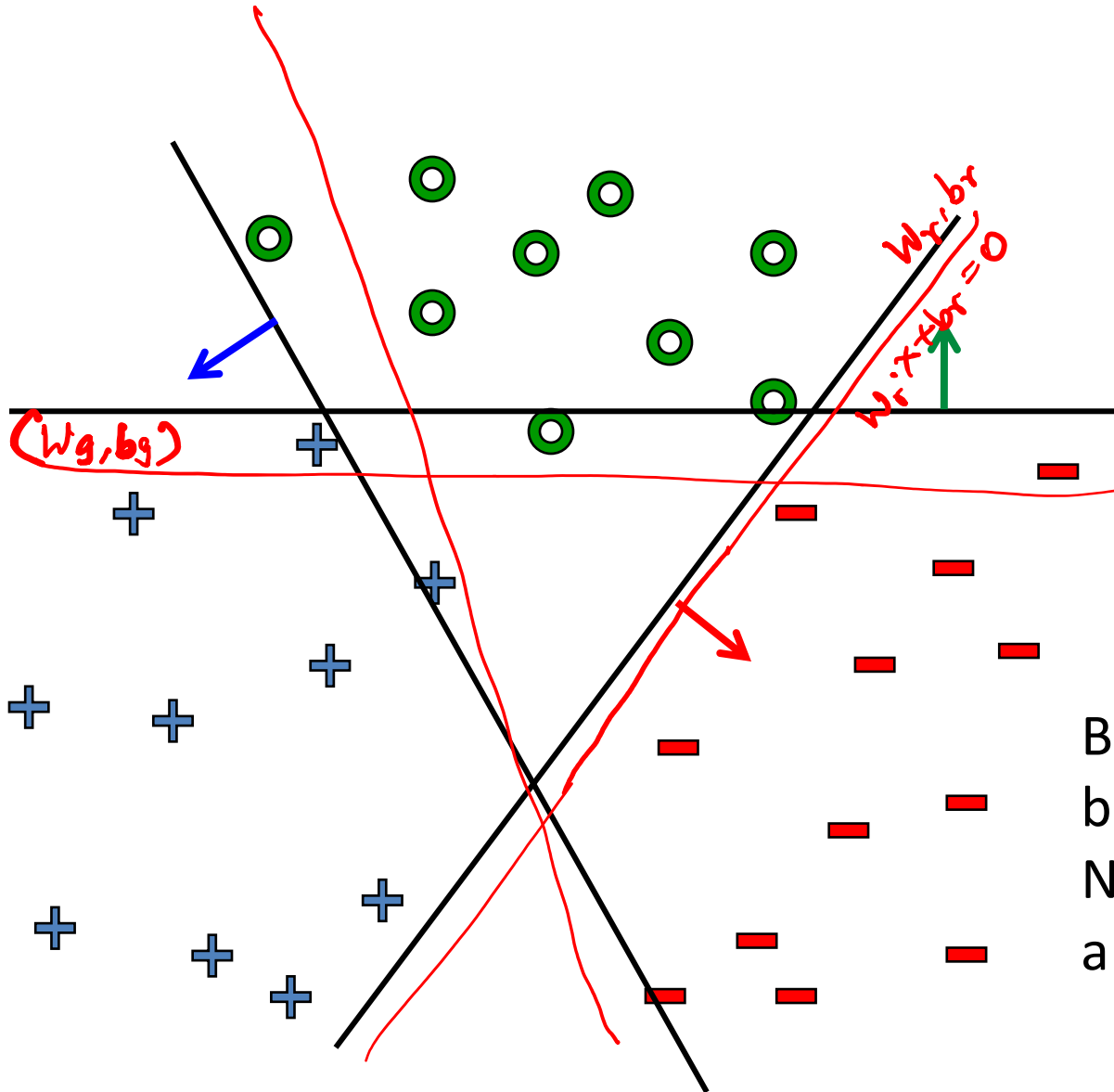
→ $1 > \xi_j > 0$ Correctly classified but inside margin

→ $\xi_j > 1$ Incorrectly classified

What about multiple classes?



One vs. rest



Learn 3 classifiers separately:

Class k vs. rest

→ $(\mathbf{w}_k, b_k)_{k=1,2,3}$

$$y = \arg \max_k \underbrace{\mathbf{w}_k \cdot \mathbf{x}} + \underbrace{b_k}$$

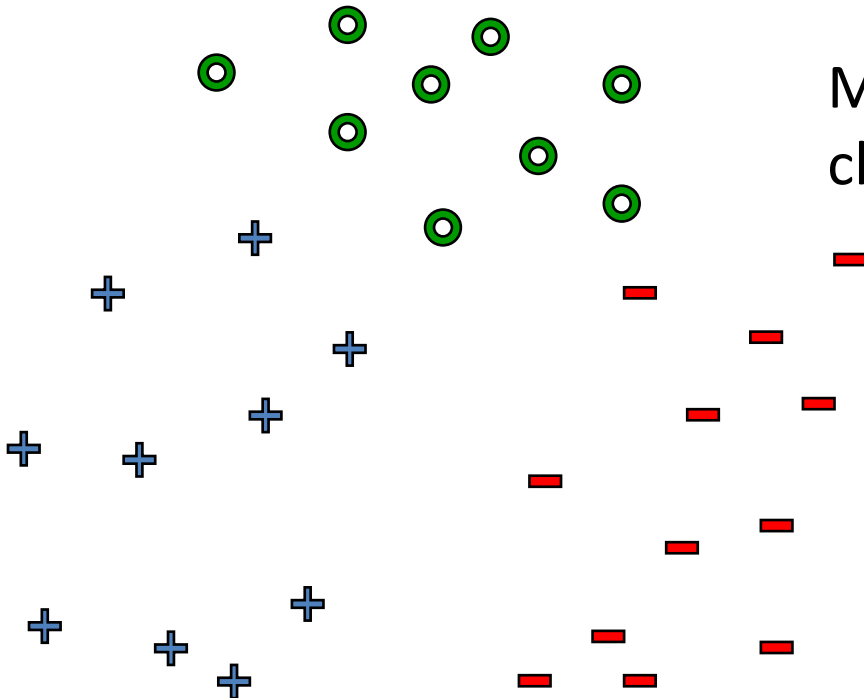
But \mathbf{w}_k s may not be based on the same scale.
Note: $(a\mathbf{w}) \cdot \mathbf{x} + (ab)$ is also a solution

Learn 1 classifier: Multi-class SVM

Simultaneously learn 3 sets of weights

$$\min_{\{w^{(y)}\}, \{b^{(y)}\}} \sum_y w^{(y)} \cdot w^{(y)}$$

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



Margin - gap between correct class and nearest other class

$$y = \arg \max_k w^{(k)} \cdot x + b^{(k)}$$

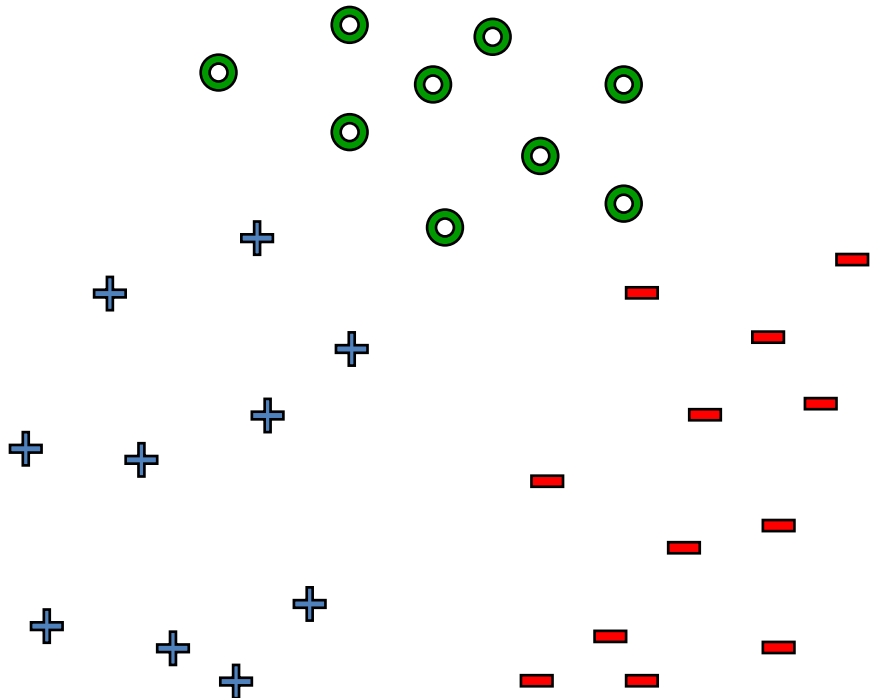
Learn 1 classifier: Multi-class SVM

Simultaneously learn 3 sets of weights

$$\text{minimize } \sum_y \mathbf{w}^{(y)} \cdot \mathbf{w}^{(y)} + C \sum_j \sum_{y \neq y_j} \xi_j^{(y)} \quad \text{over } \{\mathbf{w}^{(y)}\}, \{b^{(y)}\}, \{\xi_j^{(y)}\}$$

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

$$\underline{\xi_j^{(y)}} \geq 0, \quad \forall y \neq y_j, \quad \forall j$$



$$y = \arg \max \mathbf{w}^{(k)} \cdot \mathbf{x} + b^{(k)}$$

Joint optimization: \mathbf{w}_k s have the same scale.

Support Vector Machines - Dual formulation

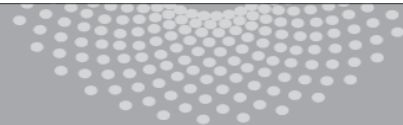
Aarti Singh

Machine Learning 10-315

Oct 26, 2020



MACHINE LEARNING DEPARTMENT



Carnegie Mellon.
School of Computer Science

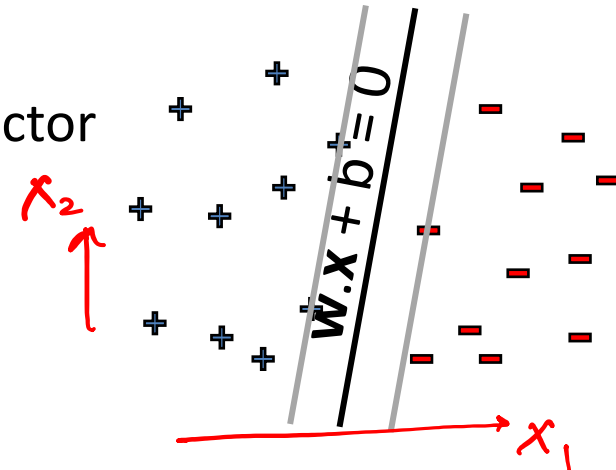
SVM – linearly separable case

→ n training points
d features

$(\mathbf{x}_1, \dots, \mathbf{x}_n)$

\mathbf{x}_j is a d-dimensional vector

- Primal problem: minimize w, b $\frac{1}{2} \mathbf{w} \cdot \mathbf{w}$
 $(\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1, \forall j$



w - weights on features (d-dim problem)

- Convex quadratic program – quadratic objective, linear constraints
- But expensive to solve if d is very large
- Often solved in dual form (n-dim problem)

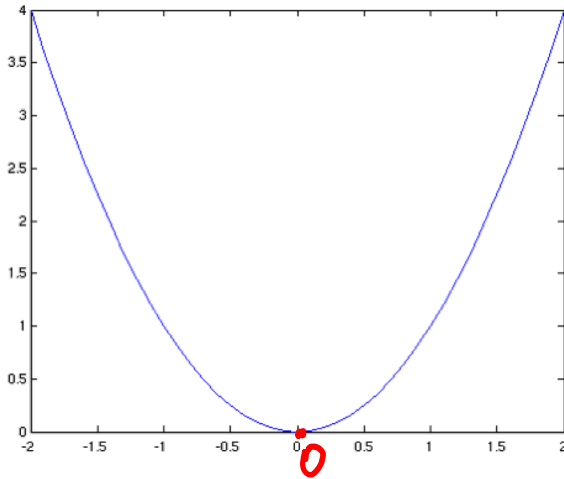
↳ # training points

Detour - Constrained Optimization

$$\begin{aligned} \min_x \quad & x^2 \leftarrow \\ \text{s.t.} \quad & \underline{x \geq b} \end{aligned}$$

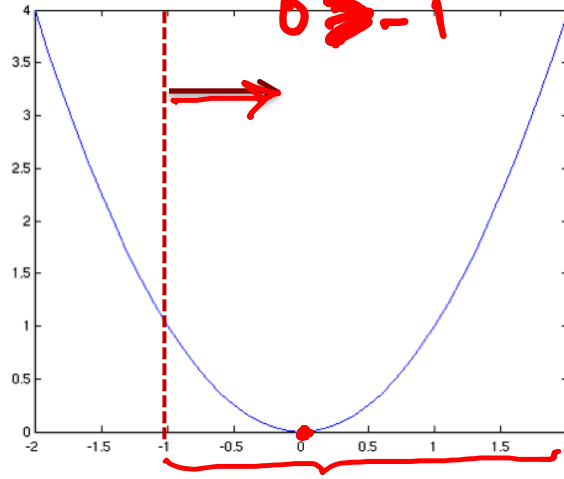
$$x^* = \max(\underline{b}, 0)$$

$$\min_x x^2$$



$$x^* = 0$$

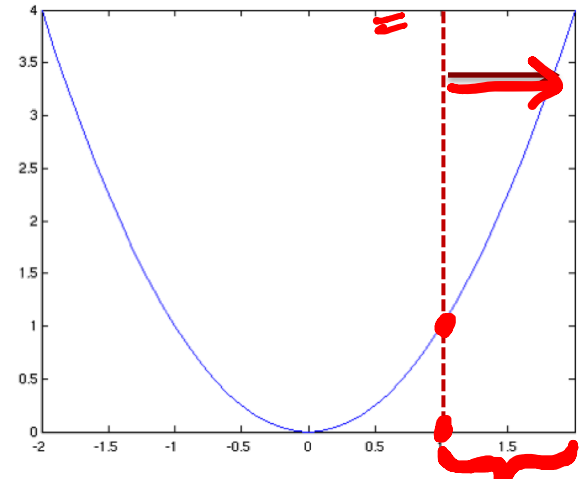
$$\begin{aligned} \min_x \quad & x^2 \\ \text{s.t.} \quad & x \geq -1 \leftarrow \end{aligned}$$



$$x^* = 0$$

Constraint inactive

$$\begin{aligned} \min_x \quad & x^2 \\ \text{s.t.} \quad & x \geq 1 \leftarrow \end{aligned}$$

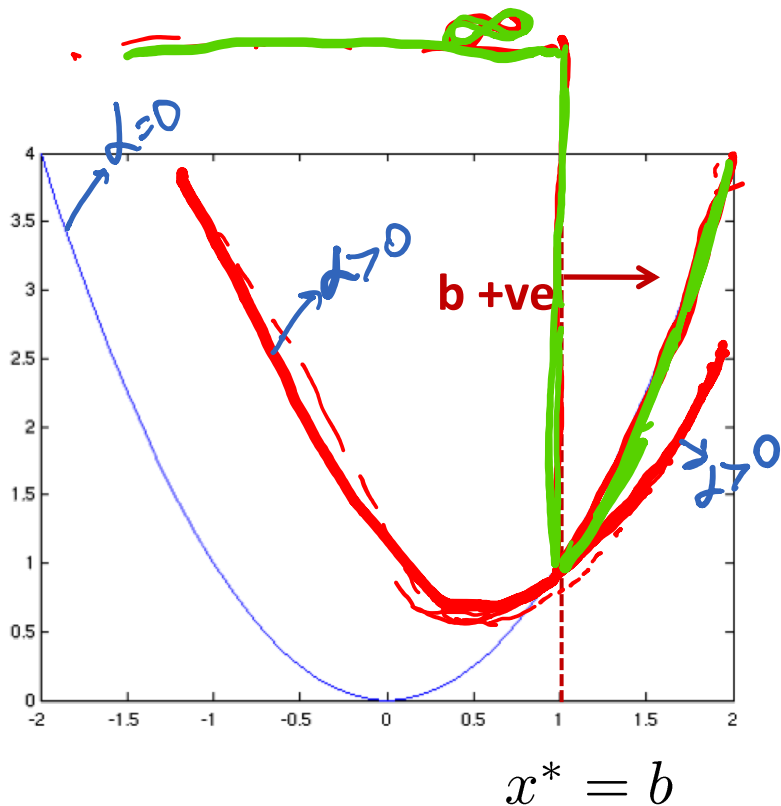


$$x^* = 1$$

Constraint active

(tight)

Constrained Optimization



$$\begin{aligned}
 &x < 2 \\
 &-x > -2 \\
 &\underline{-x + 2 > 0}
 \end{aligned}$$

$$\begin{aligned}
 &\min_x x^2 \\
 &\text{s.t. } \underline{x \geq b} \quad \left. \begin{array}{l} \text{---} \\ \text{---} \end{array} \right\} x - b \geq 0
 \end{aligned}$$

Equivalent unconstrained optimization:

$$\begin{aligned}
 &\min_x x^2 + \begin{cases} \infty & x < b \\ 0 & x \geq b \end{cases} \quad \underline{I(x-b)} \\
 &\min_x x^2 + \underline{I(x-b)}
 \end{aligned}$$

Replace with lower bound ($\alpha \geq 0$)

$$\underline{x^2 + I(x-b)} \geq \underline{x^2 - \alpha(x-b)}$$

$\alpha > 0 \quad x < b$

Primal and Dual Problems

Notice that

$$L(x, \alpha) = \underbrace{x^2 - \alpha(x-b)}_{-\alpha'(x-c)} \quad \alpha \geq 0$$

Lagrangian *Lagrange multiplier*

Primal problem: $p^* = \min_x x^2 \quad \text{s.t. } x \geq b$ $= \min_x \max_{\alpha \geq 0} L(x, \alpha)$

$x-b \geq 0$ $x-c \geq 0$

Why? $L(x, \alpha) = x^2 - \alpha(x - b)$

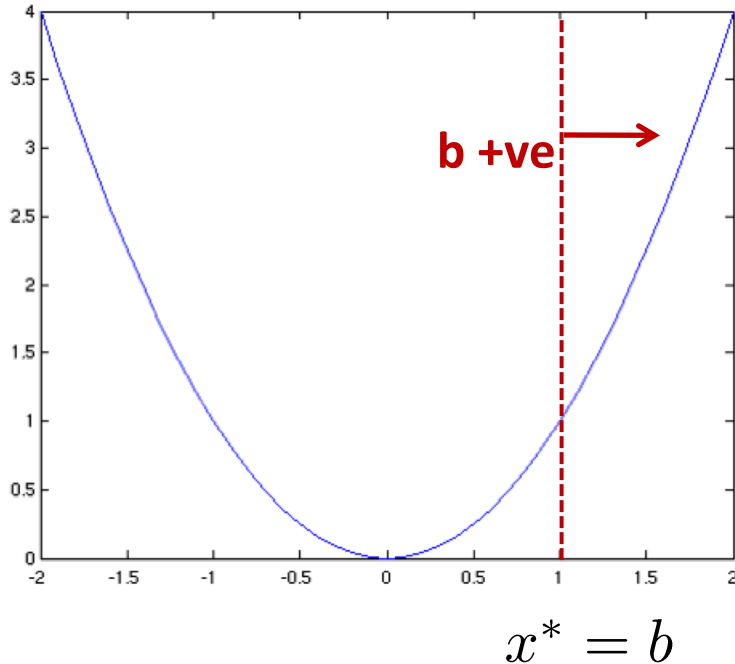
$$\max_{\alpha \geq 0} L(x, \alpha) = x^2 - \min_{\alpha \geq 0} \alpha(x - b) = \begin{cases} x^2 & x \geq b \\ \infty & x < b \end{cases}$$

$\downarrow \infty$ $\downarrow -ve$

Dual problem: $d^* = \max_{\alpha} d(\alpha) \quad \text{s.t. } \alpha \geq 0$ $= \max_{\alpha} \min_x L(x, \alpha) \quad \text{s.t. } \alpha \geq 0$

hard margin SVM
 $(w \cdot x_j + b) y_j \geq 1 \quad \forall j$

Constrained Optimization – Dual Problem



$\alpha = 0$ constraint is inactive
 $\alpha > 0$ constraint is active

Primal problem:

$$\begin{aligned} \min_x \quad & x^2 \\ \text{s.t.} \quad & x \geq b \quad \rightarrow \alpha \geq 0 \\ & \underline{x - b \geq 0} \end{aligned}$$

Moving the constraint to objective function
 Lagrangian:

$$\begin{aligned} L(x, \alpha) &= x^2 - \alpha(x - b) \\ \text{s.t.} \quad & \underline{\alpha \geq 0} \end{aligned}$$

Dual problem:

$$\begin{aligned} \max_{\alpha} \quad & d(\alpha) \rightarrow \min_x L(x, \alpha) \\ \text{s.t.} \quad & \alpha \geq 0 \end{aligned}$$

Connection between Primal and Dual

$$\text{Primal problem: } p^* = \min_x x^2 \quad \downarrow$$
$$\text{s.t. } x \geq b$$

$$= \min_x \max_{\alpha \geq 0} L(x, \alpha)$$

$$\text{Dual problem: } d^* = \max_{\alpha} d(\alpha)$$
$$\text{s.t. } \alpha \geq 0$$

$$= \max_{\alpha} \min_x L(x, \alpha)$$
$$\text{s.t. } \alpha \geq 0 \quad d(\alpha)$$

- **Dual problem (maximization) is always concave even if primal is not convex**

Why? Pointwise infimum of concave functions is concave.

[Pointwise supremum of convex functions is convex.]

$$L(x, \alpha) = x^2 - \alpha(x - b) \quad \text{linear in } \alpha$$

- **As many dual variables α as constraints, helpful if fewer constraints than dimension of primal variable x**

Connection between Primal and Dual

Primal problem: $p^* = \min_x x^2$
s.t. $x \geq b$

Dual problem: $d^* = \max_{\alpha} d(\alpha)$
s.t. $\alpha \geq 0$

➤ **Weak duality:** The dual solution d^* lower bounds the primal solution p^* i.e. $d^* \leq p^*$

To see this, recall $L(x, \alpha) = x^2 - \alpha(x - b)$

For every feasible x' (i.e. $x' \geq b$) and feasible α' (i.e. $\alpha' \geq 0$), notice that

$$d(\alpha) = \min_x L(x, \alpha) \leq x'^2 - \underbrace{\alpha'(x' - b)}_{\substack{\text{+ve} \\ \text{+ve}}} \leq x'^2$$

Since above holds true for every feasible x' , we have $d(\alpha) \leq x^{*2} = p^*$

Connection between Primal and Dual

Primal problem: $p^* = \min_x x^2$
s.t. $x \geq b$

Dual problem: $d^* = \max_{\alpha} d(\alpha)$
s.t. $\alpha \geq 0$

- **Weak duality:** The dual solution d^* lower bounds the primal solution p^* i.e. $d^* \leq p^*$
- **Strong duality:** $d^* = p^*$ holds often for many problems of interest e.g. if the primal is a feasible convex objective with linear constraints