# Support Vector Machines (SVMs) contd...
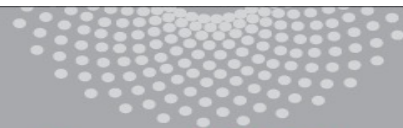
Aarti Singh

Machine Learning 10-315
Mar 23, 2022
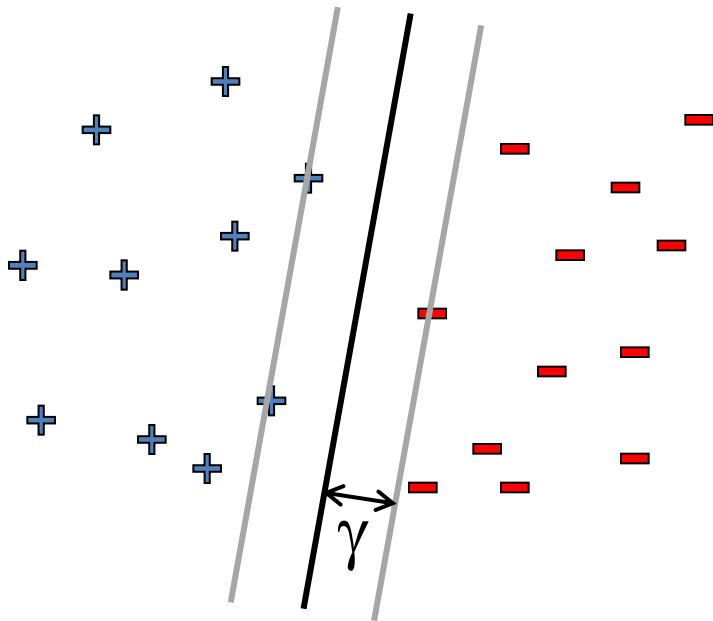
# Hard-margin SVM

Data perfectly separable by a linear decision boundary



**Hard margin approach**

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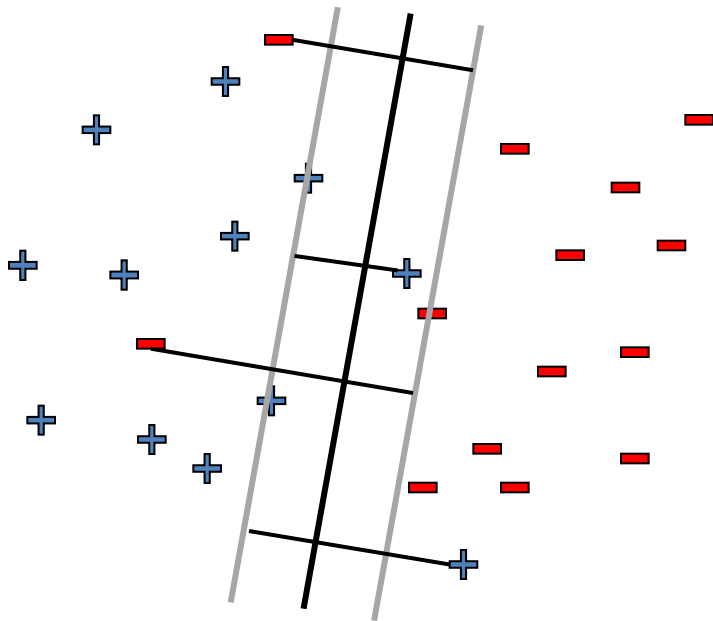
$$\text{s.t. } (\mathbf{w}.\mathbf{x}_j + b) \ y_j \geq 1 \quad \forall j$$

Solve using Quadratic Programming (QP)

Margin, $\gamma \ \alpha \ 1/\|w\|$

# Soft-margin SVM

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

$\quad$ = (>1 if $x_j$ misclassifed)

pay linear penalty if mistake

C  -  tradeoff parameter (chosen by

$\quad$ cross-validation)

Still QP ☺

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

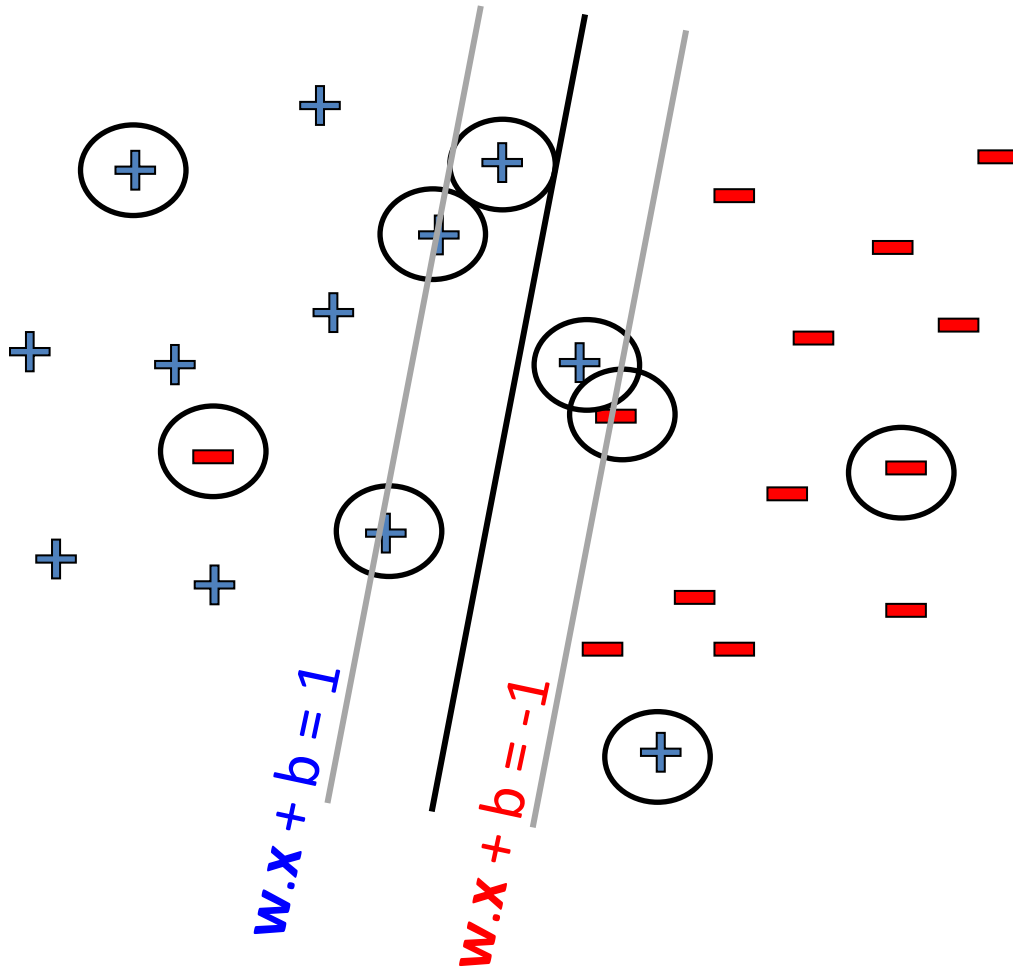$$\text{s.t. } (\mathbf{w}.\mathbf{x}_j+b) \, y_j \geq 1-\xi_j \quad \forall j$$
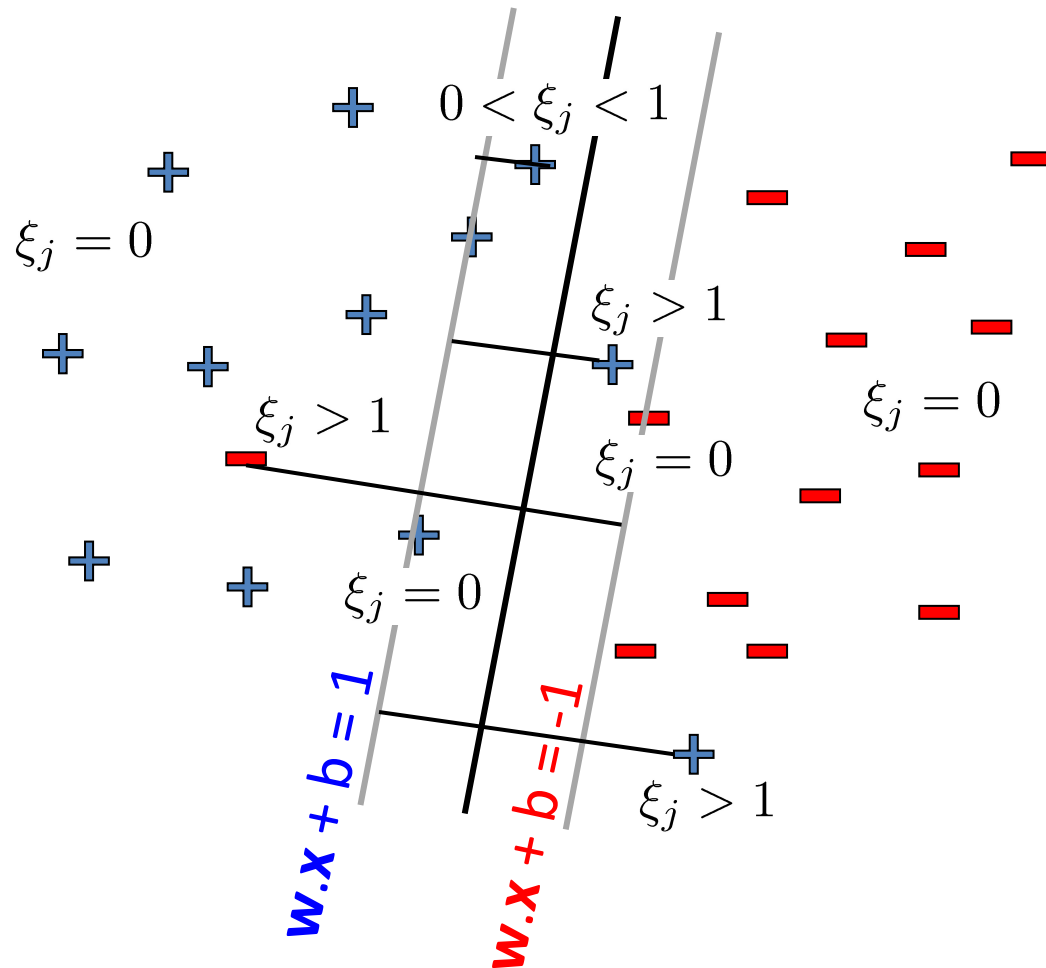
$$\xi_j \geq 0 \quad \forall j$$

# Slack variables

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

What is the slack $\xi_j$ for the following points?

Confidence | Slack



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

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

# Slack variables – Hinge loss



$0 < \xi_j < 1$

$\xi_j = 0$

$\xi_j > 1$

$\xi_j > 1$

$\xi_j = 0$

$\xi_j = 0$

$\xi_j = 0$

$\xi_j > 1$

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

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

Notice that

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

**Hinge loss**

**0-1 loss**

0    1

$$(\mathbf{w} \cdot x_j + b)y_j$$

5

# Slack variables – Hinge loss

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

**Hinge loss**

**0-1 loss**

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

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

Regularized hinge loss

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

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

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

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



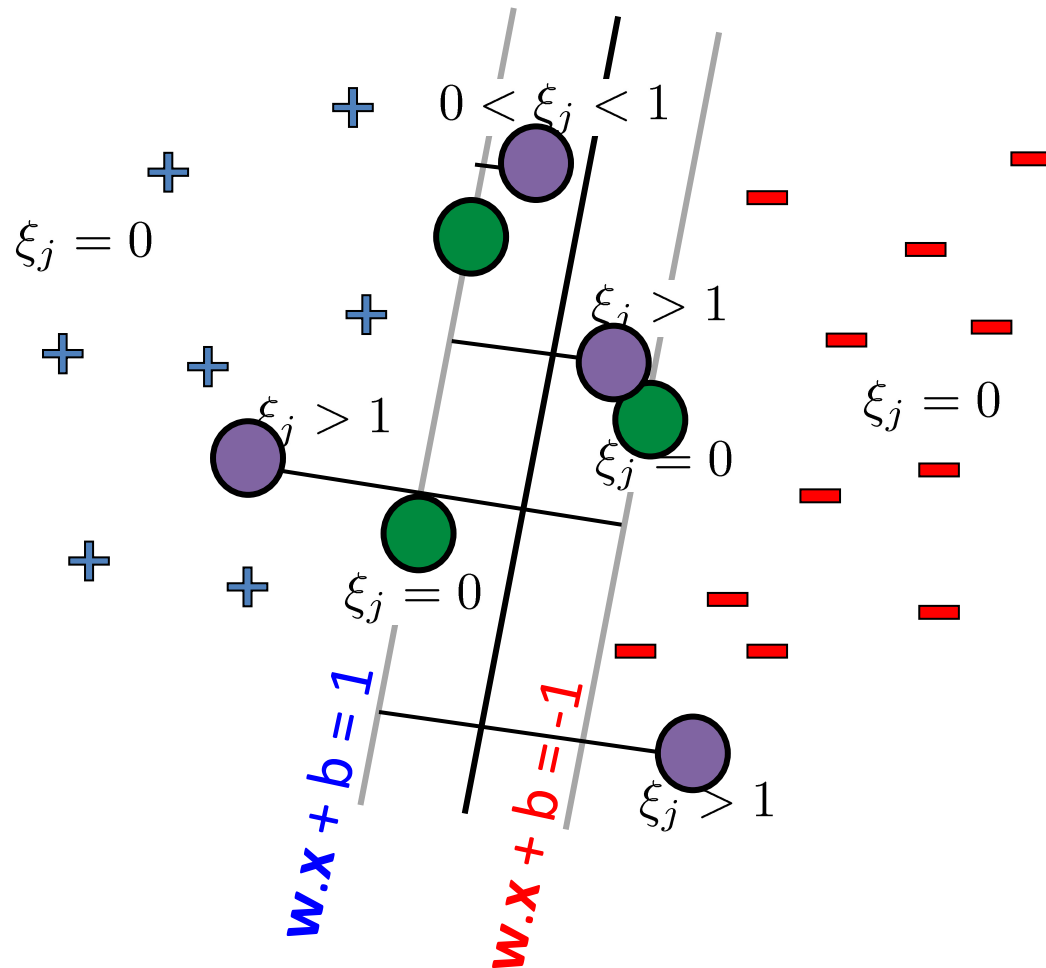$0 < \xi_j < 1$

$\xi_j = 0$

$\xi_j > 1$

$\xi_j > 1$

$\xi_j = 0$

$\xi_j = 0$

$\xi_j = 0$

$\xi_j > 1$

w.x + b = 1

w.x + b = -1

**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**

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

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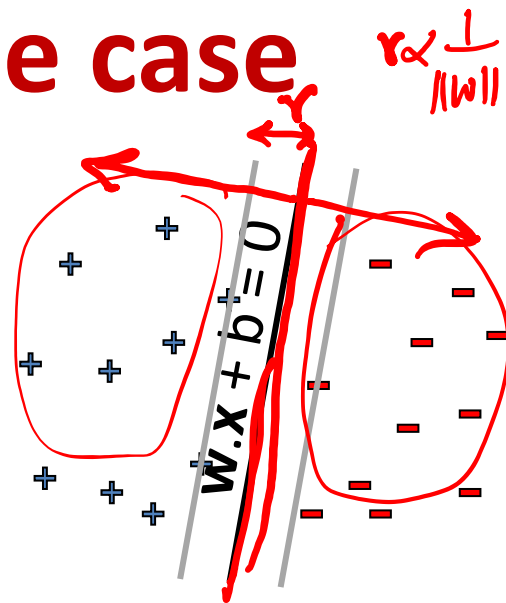$\xi_j > 1$ Incorrectly classified

# SVM – linearly separable case

$r \propto \frac{1}{\|w\|}$

n training points   $(\mathbf{x}_1, ..., \mathbf{x}_n)$

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

- Primal problem: $\text{minimize}_{\mathbf{w},b}\quad \frac{1}{2}\mathbf{w}.\mathbf{w}$

$$\left(\mathbf{w}.\mathbf{x}_j + b\right) y_j \geq 1, \;\; \forall j$$

**w – weights on features (d-dim problem)**

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

- Convex quadratic program – quadratic objective, linear constraints

- But expensive to solve if d is very large

- Often solved in dual form (n-dim problem)

8

# Detour - Constrained Optimization

quadratic objective $\longrightarrow$ $\min_x \ x^2$

linear constraint $\longrightarrow$ s.t. $\quad x \geq b$

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

$\min_x \ x^2$



$x^* = 0$

$\min_x \ x^2$
s.t. $\quad x \geq -1$



$x^* = 0$
Constraint inactive

$\min_x \ x^2$
s.t. $\quad x \geq 1$



$x^* = 1$
Constraint active
(tight)

# Constrained Optimization



$$\min_x \ x^2$$
$$\text{s.t.} \ \ x \geq b$$

Equivalent unconstrained optimization:
$\min_x \ x^2 + I(x-b)$

$$I(x-b) = \begin{cases} \infty & x < b \\ 0 & x \geq b \end{cases}$$

$x^* = b$

Replace with lower bound ($\alpha >= 0$)
$x^2 + I(x-b) \ >= \ \underbrace{x^2 - \alpha(x-b)}_{\max_{\alpha \geq 0} L(x,\alpha)}$

LHS

**b +ve**

$x$

10

# Primal and Dual Problems

**Primal problem:** p* = $\min_x x^2$
s.t. $x \geq b$     $\leftrightarrow \alpha \geq 0$

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

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

where Lagrangian $L(x, \alpha) = x^2 - \alpha(x - b)$

How to form the Lagrangian?

For each constraint, introduce a positive Lagrange multiplier
Fold constraints into objective

$\alpha \geq 0$

$x^2 - \alpha(x - b)$

$\min_{x_1, x_2} x_1^2 + x_2^2$    s.t. $x_1 \geq b_1$    $\alpha_1 \geq 0$
$x_2 \geq b_2$    $\alpha_2 \geq 0$

$x_1^2 + x_2^2 - \alpha_1(x_1 - b_1) - \alpha_2(x_2 - b_2)$

11

# Why solve the Dual problem?

$w, b$    $(d+1) dim$

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

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

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

= $\max_\alpha \min_x L(x, \alpha)$
s.t.  $\alpha \geq 0$

➢ **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.] ✓

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

➢ **As many dual variables $\alpha$ as constraints, helpful if fewer** ✓
**constraints than dimension of primal variable x**

12

# Connection between Primal and Dual

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

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

➢ **Weak duality:** The dual solution d* lower bounds the primal solution p* i.e. d* ≤ p*

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

For every feasible x' (i.e. x' ≥ b) and feasible α' (i.e. α' ≥ 0) , notice that

$$d(\alpha) = \min_x L(x, \alpha) \ \leq \ x'^2 - \alpha'(x'-b) \ \leq \ x'^2$$

$$L(x', \alpha')$$

Since above holds true for every feasible x', we have d(α) ≤ x*² = p*

13

# 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* ≤ 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

*Primal variables* $w, b^*$ ⟷ *dual variables* $\alpha_1^* \ldots \alpha_n^*$

✓ SVM

# Connection between Primal and Dual

What does strong duality say about $\alpha^*$ (the $\alpha$ that achieved optimal value of dual) and $x^*$ (the $x$ that achieves optimal value of primal problem)?

$$\min_x \max_\alpha L(x,\alpha) \qquad\qquad \max_\alpha \min_x L(x,\alpha)$$

Whenever strong duality holds, the following conditions (known as KKT conditions) are true for $\alpha^*$ and $x^*$:

- 1. $\nabla L(x^*, \alpha^*) = 0$ i.e. Gradient of Lagrangian at $x^*$ and $\alpha^*$ is zero. ✔

- 2. $x^* \geq b$ i.e. $x^*$ is primal feasible ✔

  $\alpha \geq 0$ ✔

- 3. $\alpha^* \geq 0$ i.e. $\alpha^*$ is dual feasible ✔

  $x \geq b$ ✔

  $\max \|w\|^2$

- 4. $\alpha^*(x^* - b) = 0$ (called as complementary slackness)

  s.t. $(wx_i + b) y_i \geq 1$

  dual const ↑         ↑ primal constraint

    We use the first one to relate $x^*$ and $\alpha^*$. We use the last one (complimentary slackness) to argue that $\alpha^* = 0$ if constraint is inactive and $\alpha^* > 0$ if constraint is active and tight.

# Primal and Dual Problems

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

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

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

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

where Lagrangian $\;L(x, \alpha) = x^2 - \alpha(x - b)$

How to form the Lagrangian?

       For each constraint, introduce a positive Lagrange multiplier
       Fold constraints into objective

# Dual SVM – linearly separable case

n training points, d features        ($\mathbf{x}_1$, …, $\mathbf{x}_n$) where $x_i$ is a d-dimensional vector

- <u>Primal problem</u>:         $\underset{\mathbf{w},b}{\text{minimize}} \quad \frac{1}{2}\mathbf{w}.\mathbf{w}$

$$\left(\mathbf{w}.\mathbf{x}_j + b\right) y_j \geq 1, \quad \forall j_{=1,...n}$$

$\alpha_1 .. \alpha_n \geq 0$

**w – weights on features (d-dim problem)**

- <u>Dual problem</u> (derivation):

$\max_{\alpha_1 .. \alpha_n \geq 0} d(\alpha) = \min_{w,b} L(w,b,\alpha)$

$$L(\mathbf{w},b,\alpha) = \frac{1}{2}\mathbf{w}.\mathbf{w} - \sum_j \alpha_j \left[\left(\mathbf{w}.\mathbf{x}_j + b\right) y_j - 1\right]$$
$$\alpha_j \geq 0, \quad \forall j$$

**$\alpha$ – weights on training pts (n-dim problem)**

# Dual SVM – linearly separable case

- Dual problem:

$$\max_\alpha \min_{\mathbf{w},b} L(\mathbf{w}, b, \alpha) = \tfrac{1}{2}\mathbf{w}.\mathbf{w} - \sum_j \alpha_j \left[ \left( \mathbf{w}.\mathbf{x}_j + b \right) y_j - 1 \right]$$

$$\alpha_j \geq 0, \ \forall j$$

$d(\alpha)$

$\frac{\partial L}{\partial w} = w - \sum_j \alpha_j x_j y_j = 0$

$$\frac{\partial L}{\partial \mathbf{w}} = 0 \qquad \Rightarrow \mathbf{w} = \sum_j \alpha_j y_j \mathbf{x}_j$$

If we can solve for $\alpha$s (dual problem), then we have a solution for **w** (primal problem)

$$\frac{\partial L}{\partial b} = 0 \qquad \Rightarrow \sum_j \alpha_j y_j = 0$$

$\frac{\partial L}{\partial b} = \sum_j \alpha_j y_j = 0$

# Dual SVM – linearly separable case

- Dual problem:

$$\max_\alpha \min_{\mathbf{w},b} L(\mathbf{w}, b, \alpha) = \tfrac{1}{2}\mathbf{w}.\mathbf{w} - \sum_j \alpha_j \left[ \left( \mathbf{w}.\mathbf{x}_j + b \right) y_j - 1 \right]$$

$$\underbrace{\min_{\mathbf{w},b} L(\mathbf{w}, b, \alpha)}_{d(\alpha)}$$

$$\alpha_j \geq 0, \ \forall j$$

$$\Rightarrow \mathbf{w} = \sum_j \alpha_j y_j \mathbf{x}_j \qquad \Rightarrow \sum_j \alpha_j y_j = 0$$

$$L(\mathbf{w}^*, b^*, \alpha) = \frac{1}{2} \sum_j \alpha_j y_j x_j \cdot \sum_i \alpha_i y_i x_i - \sum_j \alpha_j \left[ \left( \sum_i \alpha_i y_i x_i \cdot x_j + b \right) y_j - 1 \right]$$

$$= \frac{1}{2} \sum_j \alpha_j y_j x_j \cdot \sum_i \alpha_i y_i x_i - \sum_i \alpha_i x_i y_i \cdot \sum_j \alpha_j y_j x_j - b \sum_j \alpha_j y_j + \sum_j \alpha_j$$

$$= -\frac{1}{2} \sum_i \alpha_i y_i x_i \cdot \sum_j \alpha_j y_j x_j + \sum_j \alpha_j \quad = d(\alpha)$$

# Dual SVM – linearly separable case

$$\text{maximize}_{\alpha} \quad \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i . \mathbf{x}_j$$

$\alpha_1 \ldots \alpha_n$

$$\sum_i \alpha_i y_i = 0$$
$$\alpha_i \geq 0$$

n-dim

Dual problem is also QP

Solution gives $\alpha_j$s

$\longrightarrow$

$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$$

What about b?

# Dual SVM: Sparsity of dual solution

$$\alpha_j \left[ ((\mathbf{w} \cdot \mathbf{x}_j + b) y_j - 1 \right] = 0$$

$$\mathbf{w} = \sum_j \alpha_j y_j \mathbf{x}_j$$



$\alpha_j = 0$

$\alpha_j > 0$

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

$\alpha_j = 0$

$\alpha_j > 0$

$\alpha_j > 0$

$\alpha_j = 0$

**Complementary slackness implies**
Only few $\alpha_j$s can be non-zero : where constraint is active and tight

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

$$(\mathbf{w} \cdot \mathbf{x}_j + b) y_j > 1$$

**Support vectors** – training points j whose $\alpha_j$s are non-zero

21