# Constrained Optimization – Dual Problem



**Primal problem:**

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

**Moving the constraint to objective function**

**Lagrangian:**

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

**Dual problem:**

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

# 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*

**Duality gap = p*-d***

➤ **Strong duality:** d* = p* holds often for many problems of interest e.g. if the primal is a feasible convex objective with linear constraints (Slater's condition)

# Solving the dual

**Solving:**

$$\max_\alpha \min_x \underbrace{x^2 - \alpha(x - b)}_{L(x, \alpha)}$$

$$\text{s.t.} \quad \alpha \geq 0$$

<u>Find the dual</u>: Optimization over x is unconstrained.

$$\frac{\partial L}{\partial x} = 2x - \alpha = 0 \Rightarrow x^* = \frac{\alpha}{2} \qquad L(x^*, \alpha) = \frac{\alpha^2}{4} - \alpha\left(\frac{\alpha}{2} - b\right)$$

$$= -\frac{\alpha^2}{4} + b\alpha$$

<u>Solve</u>: Now need to maximize $L(x^*, \alpha)$ over $\alpha \geq 0$

Solve unconstrained problem to get $\alpha'$ and then take max($\alpha', 0$)

$$\frac{\partial}{\partial \alpha} L(x^*, \alpha) = -\frac{\alpha}{2} + b \Rightarrow \alpha' = 2b$$

$$\Rightarrow \alpha^* = \max(2b, 0) \qquad \Rightarrow x^* = \frac{\alpha^*}{2} = \max(b, 0)$$

$\alpha = 0$ **constraint is inactive,** $\alpha > 0$ **constraint is active (tight)**

# Dual SVM – linearly separable case

n training points, d features

$(\mathbf{x}_1, ..., \mathbf{x}_n)$ where $x_i$ is a d-dimensional vector

- Primal problem:

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

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

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

- Dual problem (derivation):

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

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

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

# Dual SVM – linearly separable case

- Dual problem (derivation):

$$\max_\alpha \min_{w,b} L(w, b, \alpha) = \frac{1}{2}w \cdot w - \sum_j \alpha_j \left[ \left( w \cdot x_j + b \right) y_j - 1 \right]$$

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

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

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

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

# Dual SVM – linearly separable case

- Dual problem:

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

$$\alpha_j \geq 0, \ \forall_j$$

$$\Rightarrow w = \sum_j \alpha_j y_j x_j$$

$$\Rightarrow \sum_j \alpha_j y_j = 0$$

# Dual SVM – linearly separable case

maximize$_\alpha$ $\sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i . x_j$

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

Dual problem is also QP

Solution gives $\alpha_i$s

$$w = \sum_i \alpha_i y_i x_i$$

What about b?

43

# Dual SVM – linearly separable case

maximize$_\alpha$ $\sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i \cdot x_j$

$$\sum_i \alpha_i y_i = 0$$

$$\alpha_i \geq 0$$

Dual problem is also QP

Solution gives $\alpha_i$s

Use any one of support vectors with

$\alpha_k > 0$ to compute b since constraint is

**tight** ($w.x_k + b)y_k = 1$

$$w = \sum_i \alpha_i y_i x_i$$

$$b = y_k - w \cdot x_k$$

for any k where $\alpha_k > 0$

1. $x_i$ with non-zero $\alpha_i$
are called SV.

2. The decision boundary is
determined only by the SV.

3. Let $t_j$ ($j=1,2,...,s$) be the
indices of the SV $\delta$. Then
$w = \sum_{j=1}^{s} \alpha_{t_j} y_{t_j} x_{t_j}$

$y = w^T x + b$
$\Rightarrow b = y - w^T x$
for any $x_k$, $y_k > 0$,
$(x_k, y_k) \Rightarrow b = y_k - w^T x_k$
$b = y_k - w^T x_k$

If $y_k=1$ then $w.x_k + b = 1$
& If $y_k=-1$ then $w.x_k + b = -1$
$\Rightarrow b = y_k - w \cdot x_k$ for avg k where
$\Rightarrow w.x_k + b = y_k \Rightarrow b = y_k - w.x_k$ for avg k where $\alpha_k > 0$
$\frac{1}{12} \alpha_k > 0$

# Dual formulation only depends on dot-products, not on w!

maximize$_\alpha$ $\sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i \cdot x_j$

$\sum_i \alpha_i y_i = 0$
$C \geq \alpha_i \geq 0$ $\longrightarrow$ $\begin{cases} \text{Regularization} \\ \text{Parameter} \end{cases}$

$\Phi(x_i)$ map $x_i$ to $\Phi(x_i)$

$K(x_i, x_j) = \Phi(x_i) \cdot \Phi(x_j)$

maximize$_\alpha$ $\sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \underline{K(x_i, x_j)}$

$\underline{K(x_i, x_j) = \Phi(x_i) \cdot \Phi(x_j)}$

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

$\Phi(x)$ – High-dimensional feature space, but never need it explicitly as long as we can compute the dot product fast using some Kernel K

# Dot Product of Polynomials

$\Phi(\mathbf{x}) = $ polynomials of degree exactly d

$$K(x,z) = \Phi(x) \cdot \Phi(z)$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \qquad \mathbf{z} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$

d=1 $\quad \Phi(\mathbf{x}) \cdot \Phi(\mathbf{z}) = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \cdot \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} = x_1 z_1 + x_2 z_2 = \mathbf{x} \cdot \mathbf{z}$

$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}; \quad \Phi(x) = \begin{bmatrix} x_1^2 \\ \sqrt{2}\, x_1 x_2 \\ x_2^2 \end{bmatrix}$

$x \xrightarrow{\Phi} \phi(x)$

d=2 $\quad \Phi(\mathbf{x}) \cdot \Phi(\mathbf{z}) = \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1 x_2 \\ x_2^2 \end{bmatrix} \cdot \begin{bmatrix} z_1^2 \\ \sqrt{2}z_1 z_2 \\ z_2^2 \end{bmatrix} = x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1 x_2 z_1 z_2$

$$= (x_1 z_1 + x_2 z_2)^2$$

$$= (\mathbf{x} \cdot \mathbf{z})^2$$

d $\quad \Phi(\mathbf{x}) \cdot \Phi(\mathbf{z}) = K(\mathbf{x}, \mathbf{z}) = (\mathbf{x} \cdot \mathbf{z})^d$

20

for many mappings from a low-D space to a high-D space, there is a simple operation on two vectors in the low-D space that can be used to compute the scalar product of their two images in the High-D space. ④8

$K(x^a, x^b) = \Phi(x^a) \cdot \Phi(x^b)$

Letting the kernel do the work

doing the scalar product in the obvious way.

# Finally: The Kernel Trick!

maximize$_\alpha$  $\sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K(x_i, x_j)$

$$K(x_i, x_j) = \Phi(x_i) \cdot \Phi(x_j)$$

$$\sum_i \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0$$

- Never represent features explicitly
  - Compute dot products in closed form

- Constant-time high-dimensional dot-products for many classes of features

$$\mathbf{w} = \sum_i \alpha_i y_i \Phi(x_i)$$

$$b = y_k - \mathbf{w} \cdot \Phi(x_k)$$

for any $k$ where $C > \alpha_k > 0$

# Common Kernels

- Polynomials of degree d

$$K(u, v) = (u \cdot v)^d$$

- Polynomials of degree up to d

$$K(u, v) = (u \cdot v + 1)^d$$

- Gaussian/Radial kernels (polynomials of all orders – recall series expansion of exp)

$$K(u, v) = \exp\left(-\frac{\|u - v\|^2}{2\sigma^2}\right)$$

- Sigmoid

$$K(u, v) = \tanh(\eta \, u \cdot v + v)$$

$$(\beta \, \overrightarrow{u \cdot v} + \gamma)$$

22

# Example

**Q:** Consider the following dataset;

| X | Y |
|---|---|
| 1 | 1 |
| 2 | 1 |
| 4 | -1 |
| 5 | -1 |
| 6 | 1 |

Where, X is the conditional feature and Y is the decision feature (class) of the objects. Answer the following:

a) Graphically demonstrate that the objects are not linearly separable.

b) After by the SVM and kernel function $K(u,v) = (uv+1)^2$ to generate the discriminant function. Assume that:

- Suppose we have 5 one-dimensional data points

- $x_1=1, x_2=2, x_3=4, x_4=5, x_5=6$, with 1, 2, 6 as class 1 and 4, 5 as class 2 $\Rightarrow y_1=1, y_2=1, y_3=-1, y_4=-1, y_5=1$

- We use the polynomial kernel of degree 2

- $K(x,y) = (xy+1)^2$

- C is set to 100 ✓

- We first find $\alpha_i$ (i=1, ..., 5) by

$$\max. \quad \sum_{i=1}^{5} \alpha_i - \frac{1}{2}\sum_{i=1}^{5}\sum_{i=1}^{5} \alpha_i \alpha_j y_i y_j (x_i x_j + 1)^2$$

$$\text{subject to } \boxed{100} \geq \alpha_i \geq 0, \sum_{i=1}^{5} \alpha_i y_i = 0$$

| X | Y |
|---|---|
| 1 | 1 |
| 2 | 1 |
| 4 | -1 |
| 5 | -1 |
| 6 | 1 |

$\alpha_1 = x_1$ $\alpha_1$
$\alpha_2 = x_2$ $\alpha_2$
$\alpha_3 = x_3, 4$ $\alpha_3$
$0 = x_4, 5$ $\alpha_4$
$0 = x_5, 6$ $\alpha_5$

the Lagrangian multipliers corresponding to the objects are
$\alpha_1 = 0, \alpha_2 = 2.5, \alpha_3 = 0, \alpha_4 = 7.3, \alpha_5 = 4.8$

1) Use the discriminant function to predict the class label of
2) Use the discriminant function to predict the class label of
3) Use the discriminant function to predict the class label of object with X = 3.

$$f(\bar{z}) = wz + b = \sum \alpha_i y_i \, K(x_i, \bar{z}) + b$$
$$= \sum \alpha_i y_i \, \phi(x_i)$$
$$= \sum \alpha_i y_i \, (x_i z + 1)^2 + b$$
$$= 2.5(1)(2z+1)^2 + 7.333(-1)(5z+1)^2 + 4.833(1) \cdot$$
$$\frac{(6z+1)^2}{(6z+1)^2} + b$$

| x | y | α |
|---|---|---|
| 1 | 1 → 0 |
| 2 | 1 → 2.5 |
| 4 | -1 → 0 |
| 5 | -1 → 7.333 |
| 6 | 1 → 4.833 |

# Example

- By using a QP solver, we get
  - $\alpha_1=0$, $\alpha_2=2.5$, $\alpha_3=0$, $\alpha_4=7.333$, $\alpha_5=4.833$
  - Note that the constraints are indeed satisfied
  - The support vectors are $\{x_2=2, x_4=5, x_5=6\}$
  - The discriminant function is

$$w = \sum \alpha_i y_i \, \phi(x_i)$$

$$f(z)$$

$$= 2.5 \, \phi(2)$$
$$\underset{\alpha_5}{\overset{y_5}{\downarrow}} \quad \overset{K(z, x_5)}{\downarrow}$$
$$-7.333 \, \phi(5)$$
$$+ 4.833 \, \phi(6)$$

$$= 2.5(1)(2z+1)^2 + 7.333(-1)(5z+1)^2 + 4.833(1)(6z+1)^2 + b$$

$$= 0.6667 z^2 - 5.333 z + b$$

- $b$ is recovered by solving $f(2)=1$ or by $f(5)=-1$ or by $f(6)=1$,
  as $x_2$ and $x_5$ lie on the line $\phi(w)^T \phi(x) + b = 1$ and $x_4$
  lies on the line $\phi(w)^T \phi(x) + b = -1$

$$w = \alpha_2 y_2 \phi(x_2)$$
$$= 2.5 f(1) \phi(x_1)$$

$$b = y_k - w \phi(x_k)$$

$$\boxed{ \begin{array}{l} x_2=2, \ y_2=1 \\ \alpha_2=2.5 \end{array} }$$

$$\phi(w)^T \phi(x) + b$$

$$f(2)=1 \Rightarrow x=2, y=1$$

$$w \Rightarrow \phi \quad \therefore b = 1 - \phi(w)^T \phi(x)$$

- All three give $b=9 \implies$ $\quad f(z) = 0.6667 z^2 - 5.333 z + 9$

$$b = 1 - [2.5 \, \phi(x_2) \cdot \phi(x_2) - 7.333 \, \phi(x_4) + 4.833 \, \phi(x_5)] \cdot \phi(x_2)$$

$$b = 1 - [2.5 \, \phi(x_2) \cdot \phi(x_2) - 7.333 \, \phi(x_4) \cdot \phi(x_2) + 4.833 \, \phi(x_5) \cdot \phi(x_2)]$$

$$= 1 - [2.5 \times K(x_2, x_2) - 7.333 \, K(x_4, x_2) + 4.833 \, K(x_5, x_2)]$$

$$= 1 - [2.5 \times (x_2^2+1)^2 - 7.333 \, (x_4 x_2 + 1)^2 + 4.833 \, (x_5 x_2 + 1)^2]$$

$\Rightarrow b = 1 - [2.5(5^2) - 7.333(11)^2 + 4.833(13)^2]$

$= 1 - [62.5 - 887.293 + 816.777]$

$= 1 - [879.277 - 887.293] = 1 - [-8.016] = 9.016 \approx 9$

$\dfrac{9}{6} \rightarrow 9$

$2.5$
$0$
$7.333$
$4.833$

# Example

Value of discriminant function



$f(z) = 0.6667\, z^2 - 5.333\, z + 9$

$f(3) = 0.6667(3)^2 - 5.333(3) + 9$

$= 6.0003 - 15.999 + 9$

$= 15.0003 - 15.999 < 0$

$\Rightarrow x = 3$ lies in class 2