# Constrained Optimization – Dual Problem



$x^* = b$

$$\frac{\alpha w^T x + b|}{||w||} \geq \gamma$$

$$2) \geq -\gamma$$

$$2) \quad \frac{w^T x + b}{||w||} \leq -\gamma$$

$$w^T x + b \leq -\gamma ||w|| \leq -1$$

$$\therefore -\gamma ||w|| \leq -1$$

$$\therefore \gamma ||w|| \geq 1$$

$$\gamma \geq \frac{1}{||w||}$$

**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:**

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

2

# 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*

**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:**

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

s.t. $\alpha \geq 0$

<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)

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

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

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

6

# 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>

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

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

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

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

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

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

  **α – weights on training pts (n-dim problem)**

# Dual SVM – linearly separable case

- Dual problem (derivation):

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

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

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

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

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:

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

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

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

# Dual SVM – linearly separable case

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

$$\sum_i \alpha_i y_i = 0$$

$$\alpha_i \geq 0$$

Dual problem is also QP

Solution gives $\alpha_j$s $\longrightarrow$

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

What about b?

# Dual SVM – linearly separable case

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

$$\sum_i \alpha_i y_i = 0$$

$$\alpha_i \geq 0$$

$y = w^T x + b$

$\rightarrow b = y - w^T x$

for any $\alpha_k > 0$,

$(x_k, y_k) \Rightarrow -w^T x_k$

$b = y_k - w^T x_k$

Dual problem is also QP

Solution gives $\alpha_j$s $\longrightarrow$

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

$$b = y_k - \mathbf{w}.\mathbf{x}_k$$

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

**Use any one of support vectors with $\alpha_k > 0$ to compute b since constraint is tight $(w.x_k + b)y_k = 1$**

If $y_k = 1$ then $w.x_k + b = 1$ & if $y_k = -1$ then $w.x_k + b = -1$
$\Rightarrow w.x_k + b = y_k \Rightarrow b = y_k - w.x_k$ for any $k$ where $\alpha_k > 0$

12

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

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

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

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

$$\sum_i \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0 \longrightarrow \begin{cases} Regularazition \\ Parameter \, C \end{cases}$$

⇓

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

$$K(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)$$

$$\sum_i \alpha_i y_i = 0$$

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

$\Phi(\mathbf{x})$ – 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

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

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

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

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

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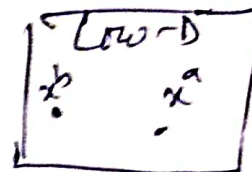
$$= (x_1 z_1 + x_2 z_2)^2$$
$$= (\mathbf{x} \cdot \mathbf{z})^2$$
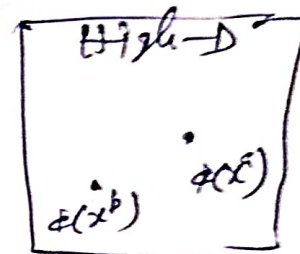
d  $\Phi(\mathbf{x}) \cdot \Phi(\mathbf{z}) = K(\mathbf{x}, \mathbf{z}) = (\mathbf{x} \cdot \mathbf{z})^d$

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)

# Finally: The Kernel Trick!

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

↑ *Letting the Kernel do the work*

↑ *doing the scalar product in the obvious way.*

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

$$K(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)$$

$$\sum_i \alpha_i y_i = 0$$

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

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

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

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

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

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

21

# Common Kernels

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

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

- Sigmoid

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

$$\beta \, \vec{u} \cdot \vec{v} + \gamma$$

22

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:

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

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

# Example

- 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 $\Rightarrow Y_1=1$, $Y_2=1$, $Y_3=-1$, $Y_4=-1$, $Y_5=1$

| X | Y |
|---|---|
| 1 | 1, $\alpha_1$ |
| 2 | 1, $\alpha_2$ |
| 4 | -1, $\alpha_3$ |
| 5 | -1, $\alpha_4$ |
| 6 | 1, $\alpha_5$ |

$O_1 = x_1 \quad$ 
$O_2 = x_2 = 2$
$O_3 = x_3 = 4$
$O_4 = x_4 = 5$
$O_5 = x_5 = 6$

- 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

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

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

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

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

30

$$f(z) = wz + b = \sum_i \alpha_i Y_i K(x_i, z) + b$$
$$= \sum \alpha_i Y_i (x_i z + 1)^2 + b$$

## Example
$$= 2.5(1)(2z+1)^2 + 7.333(-1)(5z+1)^2 + \frac{4.833(1)\cdot}{(6z+1)^2} + b$$

| x | y | $\alpha$ |
|---|---|---|
| 1 | 1 | 0 |
| 2 | 1 | 2.5 |
| 4 | -1 | 0 |
| 5 | -1 | 7.333 |
| 6 | 1 | 4.833 |

- 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

$W = \sum \alpha_i Y_i \phi(K_i)$

$= 2.5 \phi(2)$
$- 7.333 \phi(5)$
$+ 4.833 \phi(6)$

$y_5$ $\alpha_5$ $K(z, x_5)$

$$f(z)$$
$$= 2.5(1)(2z+1)^2 + 7.333(-1)(5z+1)^2 + 4.833(1)(6z+1)^2 + b$$
$$= 0.6667z^2 - 5.333z + b$$

$W = \alpha_2 Y_2 \phi(x_2)$
$= 2.5(1)\phi(2)$

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

$x_2 = 2, Y_2 = 1$
$\alpha_2 = 2.5$

$\phi(w)^T \phi(x) + b = 1$

for $f(2) = 1 \Rightarrow x = 2, Y = 1$

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

- $b$ 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 $\phi(w)^T\phi(x) + b = -1$
- All three give b=9 $\Longrightarrow$ $f(z) = 0.6667z^2 - 5.333z + 9$

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

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

# Example

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

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

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

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

$= 6.0003 - 15.999 + 9$

$= 15.0003 - 15.999 < 0$

$\Rightarrow x = 3$ lies in class 2

Value of discriminant function



class 1    class 2    class 1

1   2   3   4   5   6

2022/10/18