

# Support Vector Machines - Dual formulation and Kernel Trick

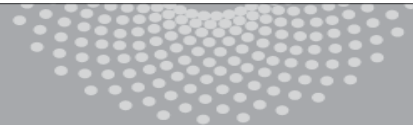
Aarti Singh

Machine Learning 10-315

Oct 28, 2020



MACHINE LEARNING DEPARTMENT



**Carnegie Mellon.**  
School of Computer Science



# 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.  $d^* \leq p^*$

$$\text{Duality gap} = p^* - d^*$$

- **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)

$$(w^*, b^*) \equiv \underline{x^*} \quad \underline{\alpha^*}$$

# 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)?

*Karush-Kuhn-Tucker*

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 ✓
  - 3.  $\alpha^* \geq 0$  i.e.  $\alpha^*$  is dual feasible ✓
  - ✓ • 4.  $\alpha^*(x^* - b) = 0$  (called as complementary slackness)
- $\alpha_i \quad (w \cdot x_i + b) y_i \geq 1$
- ⇒

$\alpha^* = 0$	$x^* > b$
$\alpha$	↓
$x^* = b$	$\alpha^* = 0$
	⇒
	$\alpha^* > 0 \Rightarrow x^* = b$

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.

# Solving the dual

Solving:

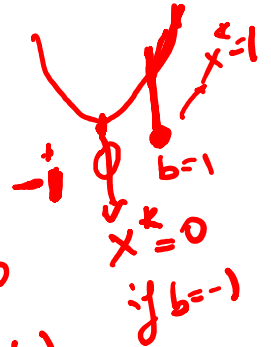
$$\text{Dual: } \max_{\alpha} \left( \min_x \overbrace{x^2 - \alpha(x-b)}^{L(x, \alpha)} \right)$$

s.t.  $\alpha \geq 0$

$$\min_x x^2$$

s.t.  $x \geq b, \alpha \geq 0$

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



Dual derivation:

$$\frac{\partial L(x, \alpha)}{\partial x} = 2x - \alpha = 0 \Rightarrow x = \alpha/2$$

$$\min_x L(x, \alpha) = \left(\frac{\alpha}{2}\right)^2 - \alpha\left(\frac{\alpha}{2} - b\right) = \frac{\alpha^2}{4} - \frac{\alpha^2}{2} + \alpha b = -\frac{\alpha^2}{4} + \alpha b$$

$$\text{Dual: } \max_{\alpha \geq 0} \underbrace{-\frac{\alpha^2}{4} + \alpha b}$$

$$\frac{\partial}{\partial \alpha} = -\frac{\alpha}{2} + b = 0 \Rightarrow \alpha^* = 2b$$

$$\Rightarrow \alpha^* = \max(0, 2b)$$

$$\alpha^*(x^* = b) = 0$$

$$\alpha^* \neq 0 \Rightarrow x^* = b$$

$$\alpha^* = 0$$

$$\alpha^* = 2 \neq 0 \Rightarrow x^* = 1$$

$$\alpha^* = 0 \text{ if } b = -1$$

$$= 2 \text{ if } b = 1$$

# Solving the dual

Solving:

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

Find the dual: Optimization over  $x$  is unconstrained.

$$\begin{aligned} \frac{\partial L}{\partial x} = 2x - \alpha = 0 & \Rightarrow x^* = \frac{\alpha}{2} & L(x^*, \alpha) &= \frac{\alpha^2}{4} - \alpha \left( \frac{\alpha}{2} - b \right) \\ & & &= -\frac{\alpha^2}{4} + b\alpha \end{aligned}$$

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

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

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

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

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

# Dual SVM – linearly separable case

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

- Primal problem: minimize  $\underline{\mathbf{w}}, b$   $\frac{1}{2} \mathbf{w} \cdot \mathbf{w}$   $\leftarrow \frac{\|\mathbf{w}\|^2}{2}$   
 $(\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1, \forall j$   $\leftarrow n$  constraints  
 $(\mathbf{w} \cdot \mathbf{x}_j + b) y_j - 1 \geq 0$   $\alpha_j \geq 0$   
 $\mathbf{w}$  - weights on features (d-dim problem)
- Dual problem (derivation):  $\frac{1}{2} \mathbf{w} \cdot \mathbf{w} - \sum_j \alpha_j ((\mathbf{w} \cdot \mathbf{x}_j + b) y_j - 1) = \mathcal{L}(\mathbf{w}, b, \alpha)$

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

$$\underline{\alpha_j} \geq 0, \forall j$$

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

# Dual SVM – linearly separable case

- Dual problem (derivation):

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

$\alpha_j \geq 0, \forall j$

$d(\alpha)$

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

$$\frac{\partial L}{\partial \mathbf{w}} = 0$$

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

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

$$\frac{\partial L}{\partial b} = 0$$

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



# Dual SVM – linearly separable case

$$\vec{\alpha} = \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_d \end{pmatrix}$$

- Dual problem:

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

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

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

$\in \mathbb{R}^d$        $\in \mathbb{R}^d$

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

$$d(\alpha) = \frac{1}{2} \sum_j \alpha_j y_j \vec{x}_j \cdot \sum_i \alpha_i y_i \vec{x}_i - \sum_j \alpha_j \left( \sum_i \alpha_i y_i \vec{x}_i \right) \cdot \vec{x}_j y_j - b \sum_j \alpha_j y_j + \sum_j \alpha_j$$

$$= -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \vec{x}_i \cdot \vec{x}_j + \sum_j \alpha_j$$

# Dual SVM – linearly separable case

maximize <sub>$\alpha$</sub>   $\sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j$  ←

$\sum_i \alpha_i y_i = 0$  ← {  $\vec{x}_i, y_i$  }<sub>i=1</sub><sup>n</sup>

$\alpha_i \geq 0$  ←

Dual problem is also QP

Solution gives  $\alpha_j$ s



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

What about b?

$$b^* = \frac{1 - \mathbf{w}^* \cdot \mathbf{x}_i y_i}{y_i}$$

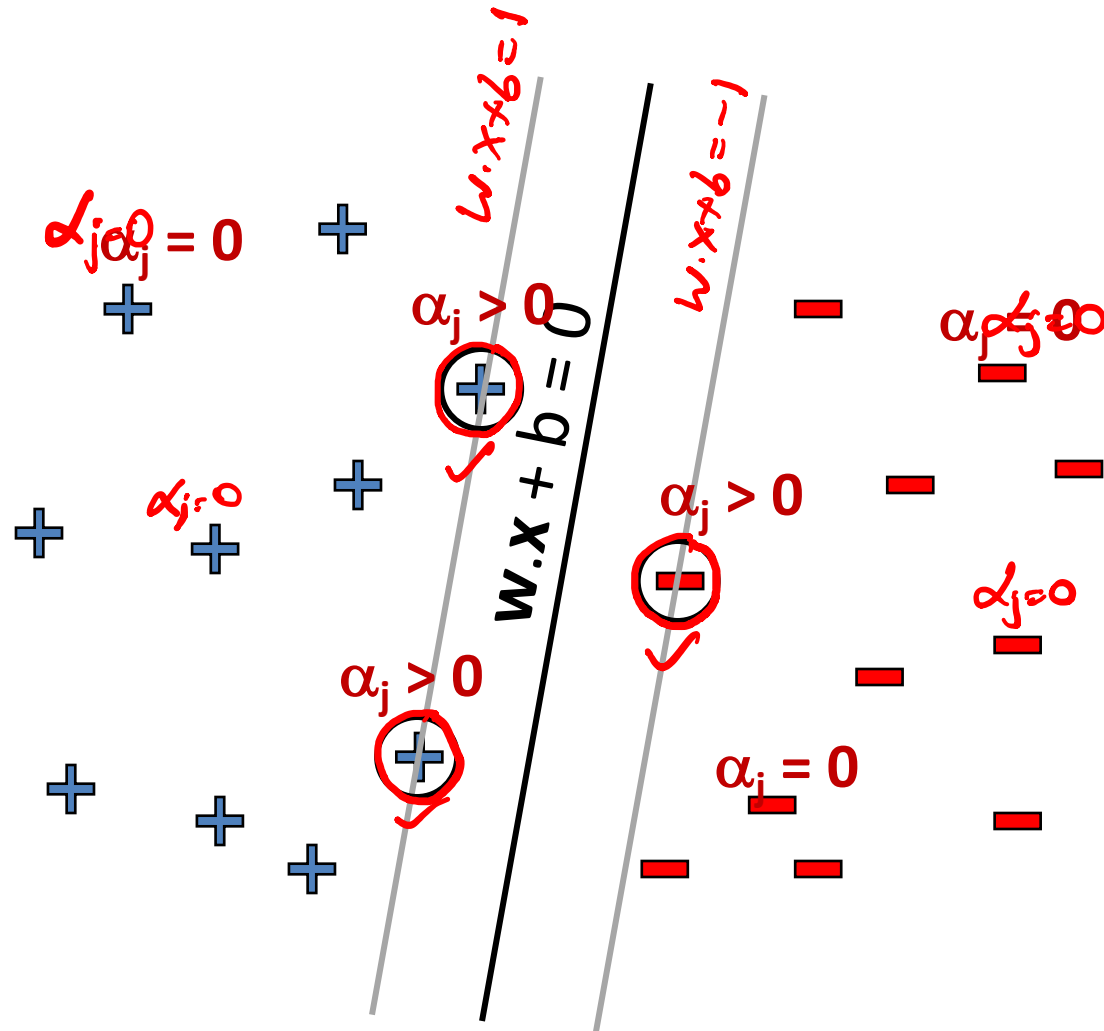
$\alpha_i^* (\underbrace{(\mathbf{w}^* \cdot \mathbf{x}_i + b^*)}_{\text{Comp slackness}} y_i - 1) = 0$

Toy eg:  $\alpha^* (\underbrace{x^* - b^*}_{\geq b}) = 0$

$\mathbf{w}^* \cdot \mathbf{x}_i y_i + b^* y_i = 1$

$\frac{1}{y_i} - \mathbf{w}^* \cdot \mathbf{x}_i$

# Dual SVM: Sparsity of dual solution



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

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

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

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

# Dual SVM – linearly separable case

$$\begin{aligned} \text{maximize}_{\alpha} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j \\ & \sum_i \alpha_i y_i = 0 \\ & \alpha_i \geq 0 \end{aligned}$$

Dual problem is also QP

Solution gives  $\alpha_j$ s  $\longrightarrow$

Use any one of support vectors with  $\alpha_k > 0$  to compute  $b$  since constraint is tight  $(\mathbf{w} \cdot \mathbf{x}_k + b)y_k = 1$

$$\begin{aligned} \mathbf{w} &= \sum_i \alpha_i y_i \mathbf{x}_i \\ b &= y_k - \mathbf{w} \cdot \mathbf{x}_k \\ &\text{for any } k \text{ where } \alpha_k > 0 \end{aligned}$$

# Dual SVM – non-separable case

- Primal problem:

$$\begin{aligned} & \text{minimize}_{\mathbf{w}, b, \{\xi_j\}} \frac{1}{2} \mathbf{w} \cdot \mathbf{w} + C \sum_j \xi_j \\ & (\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1 - \xi_j, \quad \forall j \\ & \xi_j \geq 0, \quad \forall j \end{aligned}$$

$w - d(x)$   
 $b - |x|$   
 $\{\xi_j - |x|\}_n$

$$\begin{aligned} n \quad & \alpha_j \geq 0 \\ n \quad & \mu_j \geq 0 \end{aligned}$$

Lagrange  
Multipliers

- Dual problem:

$$\begin{aligned} & \max_{\alpha, \mu} \min_{\mathbf{w}, b, \{\xi_j\}} L(\mathbf{w}, b, \xi, \alpha, \mu) \leftarrow d(\alpha, \mu) \\ & s.t. \alpha_j \geq 0 \quad \forall j \\ & \mu_j \geq 0 \quad \forall j \end{aligned}$$

HW3!

# Dual SVM – non-separable case

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

$$\sum_i \alpha_i y_i = 0 \quad -$$

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

$$\begin{aligned} \alpha_i &\geq 0 \\ -\alpha_i &\geq -C \end{aligned}$$

comes from  $\frac{\partial L}{\partial \xi} = 0$

Intuition:

If  $C \rightarrow \infty$ , recover hard-margin SVM

Dual problem is also QP

Solution gives  $\alpha_j$



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

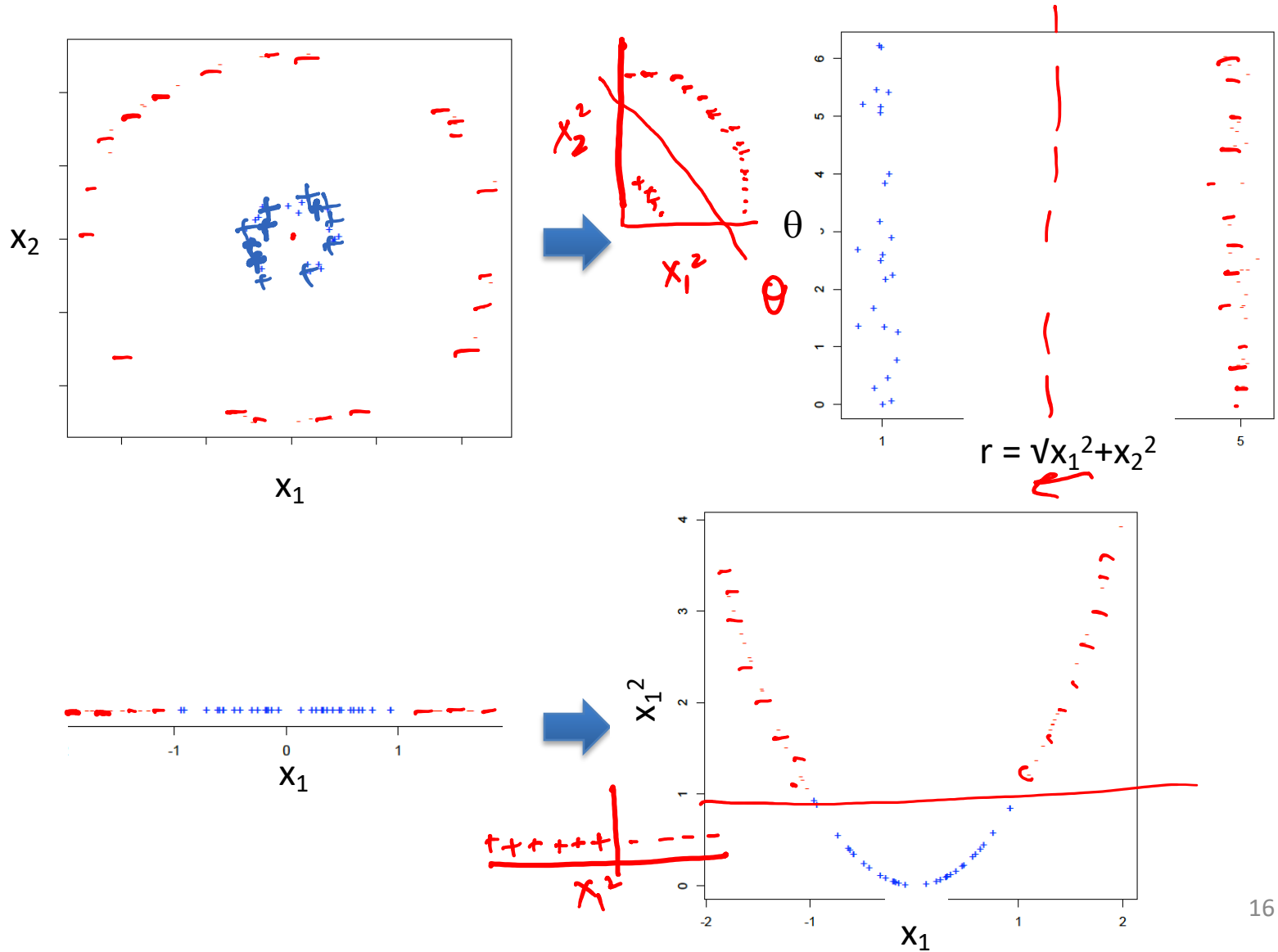
$$b = y_k - \mathbf{w} \cdot \mathbf{x}_k$$

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

# So why solve the dual SVM?

- There are some quadratic programming algorithms that can solve the dual faster than the primal, (specially in high dimensions  $d \gg n$ )
- But, more importantly, the “kernel trick”!!!

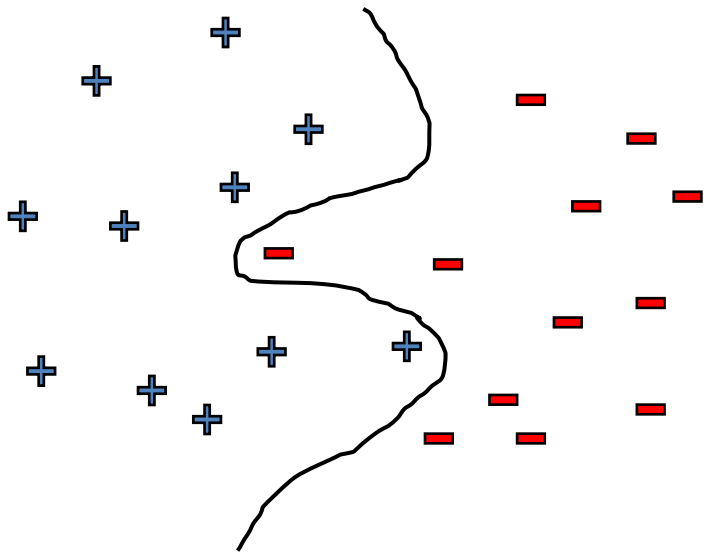
# Separable using higher-order features





# What if data is not linearly separable?

Use features of features  
of features of features....



$$\Phi(\mathbf{x}) = (x_1^2, x_2^2, x_1x_2, \dots, \exp(x_1))$$

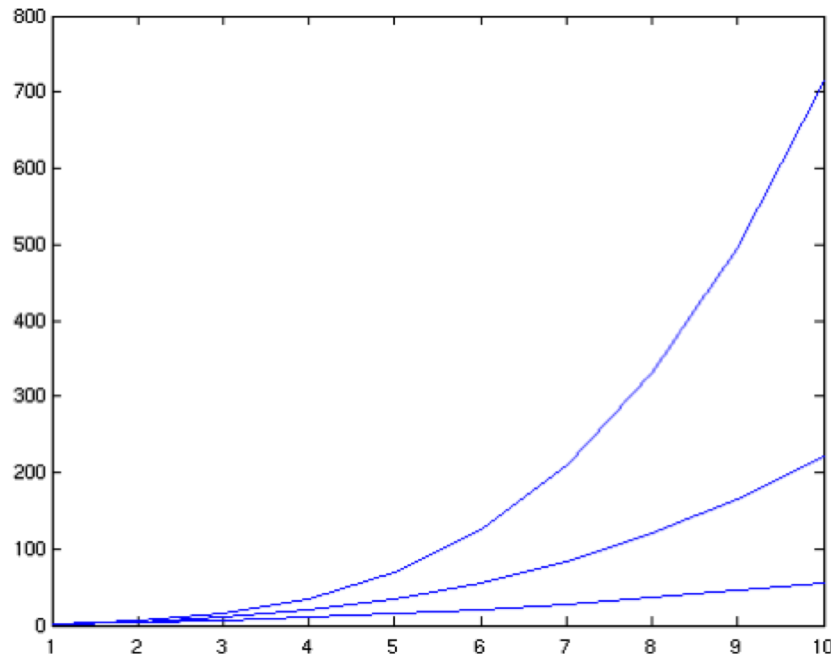
Feature space becomes really large very quickly!

# Higher Order Polynomials

m – input features

d – degree of polynomial

$$\text{num. terms} = \binom{d + m - 1}{d} = \frac{(d + m - 1)!}{d!(m - 1)!} \sim \underline{\underline{m^d}}$$



grows fast!

d = 6, m = 100 ←

about 1.6 billion terms ←

m=3 d=4

$$x_1^{\alpha_1} x_2^{\alpha_2} x_3^{\alpha_3} \quad \alpha_1 + \alpha_2 + \alpha_3 = 4$$
$$x_1^4 + x_2^4 + x_1^2 x_2^2 + x_1^2 x_3^2 + x_1 x_2^2 x_3^2$$

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

$$\begin{aligned} \text{maximize}_{\alpha} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \underbrace{\mathbf{x}_i \cdot \mathbf{x}_j}_{\substack{dx_i \\ dx_j}} \\ & \sum_i \alpha_i y_i = 0 \\ & C \geq \alpha_i \geq 0 \end{aligned}$$



$$\begin{aligned} \text{maximize}_{\alpha} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \underbrace{K(\mathbf{x}_i, \mathbf{x}_j)}_{\substack{\phi(\mathbf{x}_i) = D\mathbf{x}_i \\ D \gg d}} \\ & K(\mathbf{x}_i, \mathbf{x}_j) = \underbrace{\phi(\mathbf{x}_i)} \cdot \underbrace{\phi(\mathbf{x}_j)} \\ & \sum_i \alpha_i y_i = 0 \\ & C \geq \alpha_i \geq 0 \end{aligned}$$

$\Phi(\mathbf{x})$  – 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  $d$

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

$$\phi(\mathbf{x}) = \begin{bmatrix} x_1^d \\ x_1^{d-1} x_2 \\ x_2^{d-1} x_1 \\ x_2^d \\ \vdots \end{bmatrix}$$

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

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

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

$$= \underline{(\mathbf{x} \cdot \mathbf{z})^2} \leftarrow 4 \text{ operations}$$

$K(\mathbf{x}, \mathbf{z})$

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

# Finally: The Kernel Trick!

$$\text{maximize}_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \underbrace{K(\mathbf{x}_i, \mathbf{x}_j)}_{(x_i \cdot x_j)^{d \leftarrow \text{deg}}}$$

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

$$\sum_i \alpha_i y_i = 0$$

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

- Never represent features explicitly
  - Compute dot products in closed form
- Constant-time high-dimensional dot-products for many classes of features

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

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

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

$$\text{sign}(\mathbf{w} \cdot \Phi(\mathbf{x}) + b)$$

$$\mathbf{w} \cdot \Phi(\mathbf{x}) = \sum_i \alpha_i y_i \underbrace{\Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x})}_{K(\mathbf{x}_i, \mathbf{x})}$$

# Common Kernels

- Polynomials of degree  $d$

$$K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v})^d \quad \checkmark$$

- Polynomials of degree up to  $d$

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

$d=2$   
 $x_1 \quad x_2 \quad x_1^2 \quad x_2^2$   
 $x_1 x_2$

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

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

- Sigmoid

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