

**Two SVM tutorials linked in class website
(please, read both):**

- High-level presentation with applications (Hearst 1998)
- Detailed tutorial (Burges 1998)

SVMs, Duality and the Kernel Trick

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

February 27th, 2006

Announcements

■ Third homework

- ☐ is out

- ☐ Due March 1st

start early !!
ü

■ Final assigned by registrar:

- ☐ May 12, 1-4p.m

Friday

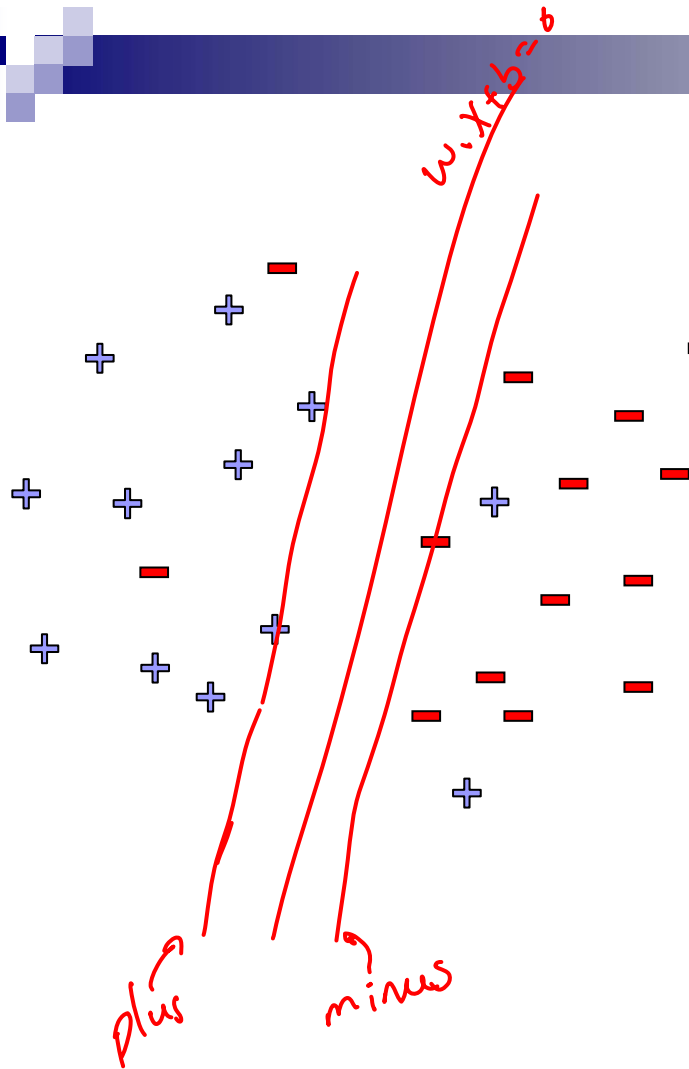
- ☐ Location TBD

■ Midterm

- ☐ March 8th, a week from Wednesday

- ☐ Open book, notes, papers, etc. No computers

SVMs reminder



slack penalty

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

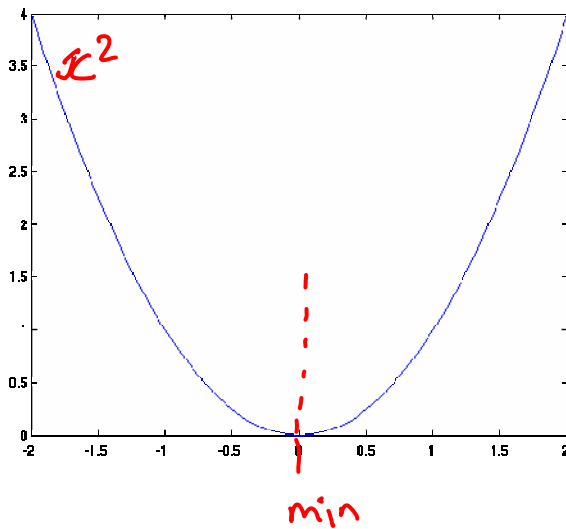
Today's lecture

- Learn one of the most interesting and exciting recent advancements in machine learning
 - The “kernel trick”
 - High dimensional feature spaces at no extra cost!
- But first, a detour
 - Constrained optimization!

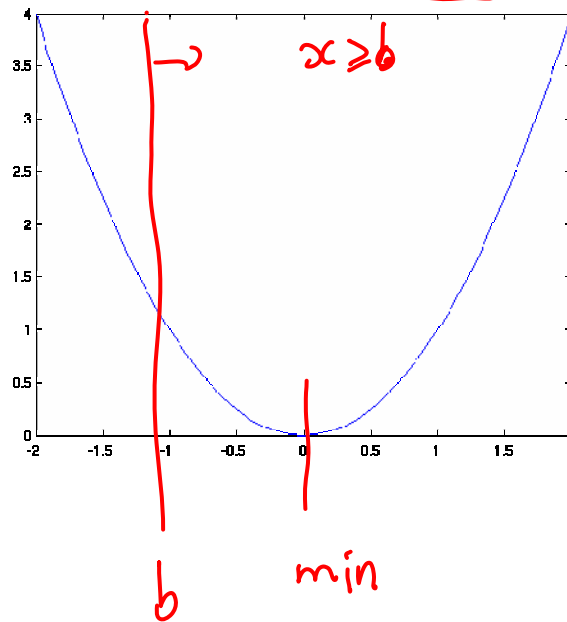
Constrained optimization

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

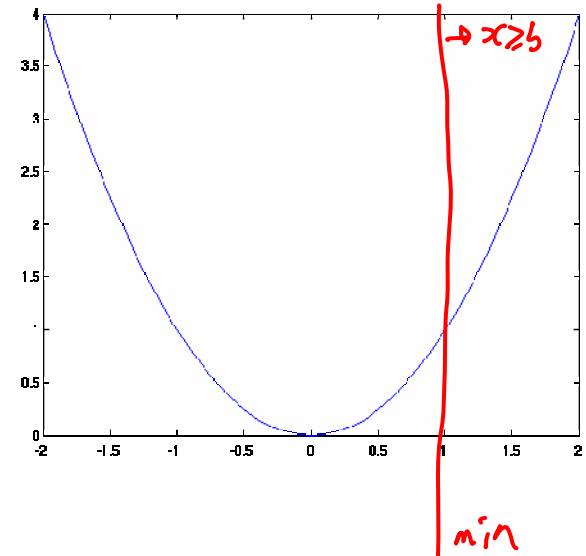
$$\min_x x^2$$



unconstrained

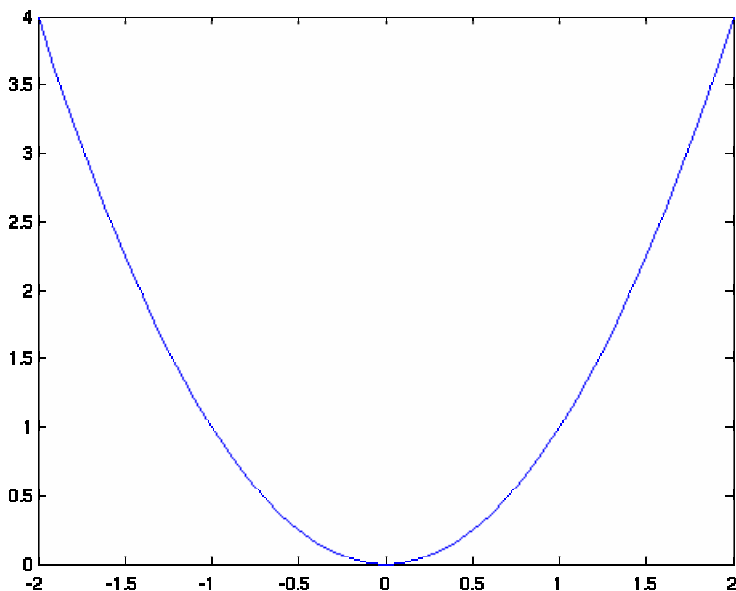


constraint
does not matter



stopped by
constraint

Lagrange multipliers – Dual variables



Solve:

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

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

$$\min_x x^2$$

$$\text{s.t. } x \geq b \equiv x - b \geq 0$$

Moving the constraint to objective function
Lagrangian:

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

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

what if $x > b$: $x - b > 0$
constraint sat.

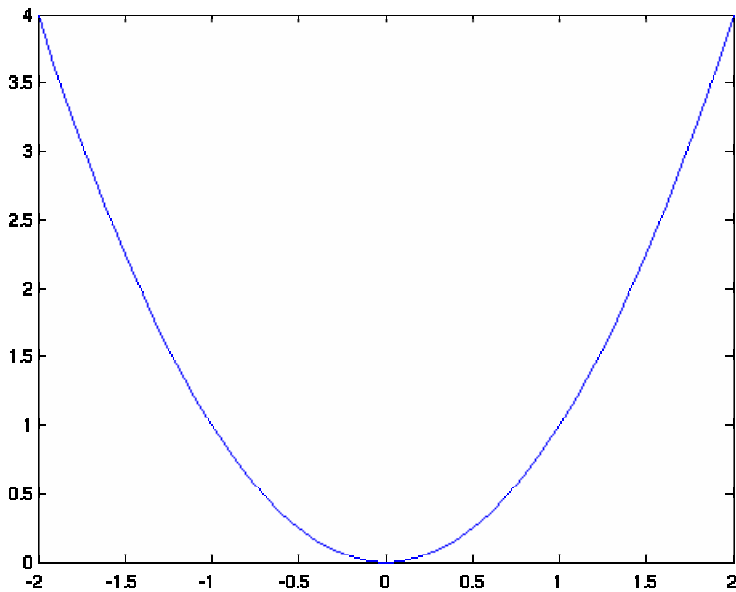
$-\alpha(x-b) \rightarrow$ negative \rightarrow max $\rightarrow \alpha = 0$
 $\min_x \rightarrow$ happy!!

$x < b$: $x - b < 0$ constraint not sat.

$-\alpha(x-b) \rightarrow$ positive \rightarrow max $\rightarrow \alpha = +\infty$
 $\min_x \rightarrow$ v. unhappy!!

$\min_x \rightarrow$ only suggest x that sat. constraint!!

Lagrange multipliers – Dual variables



Solving: $\min_x \max_{\alpha} x^2 - \alpha(x - b)$
 s.t. $\alpha \geq 0$

$$\frac{\partial L}{\partial x} = 2x - \alpha = 0 \Rightarrow x = \frac{\alpha}{2}$$

$$\frac{\partial L}{\partial \alpha} = -(x - b)$$

$$x = 0$$

either $\frac{\partial L}{\partial \alpha} = 0 \Rightarrow x = b$

or if $\alpha = 0$ term $\alpha(x - b)$ is irrelevant

either $\alpha = 0$ and constraint is ignored

or $\alpha \neq 0$ constraint plays a role

Dual SVM derivation (1) – the linearly separable case

✓ simpler eqs.

$$\text{minimize}_w \quad \frac{1}{2} w \cdot w$$
$$(w \cdot x_j + b) y_j \geq 1, \quad \forall j \in \text{training set}$$

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

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

Dual SVM derivation (2) – the linearly separable case

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

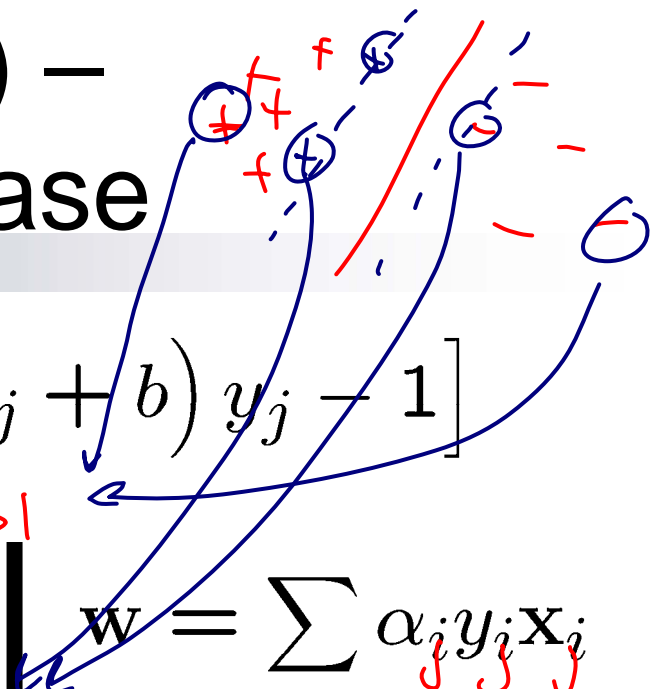
$$\alpha_i \geq 0, \quad \forall j$$

$$(w \cdot x_j + b) y_j \geq 1 \quad \left\{ \begin{array}{l} \text{when } (w \cdot x_j + b) y_j > 1 \\ \Rightarrow \alpha_j = 0 \\ \hline \text{when } (w \cdot x_j + b) y_j = 1 \\ \alpha_j > 0 \end{array} \right.$$

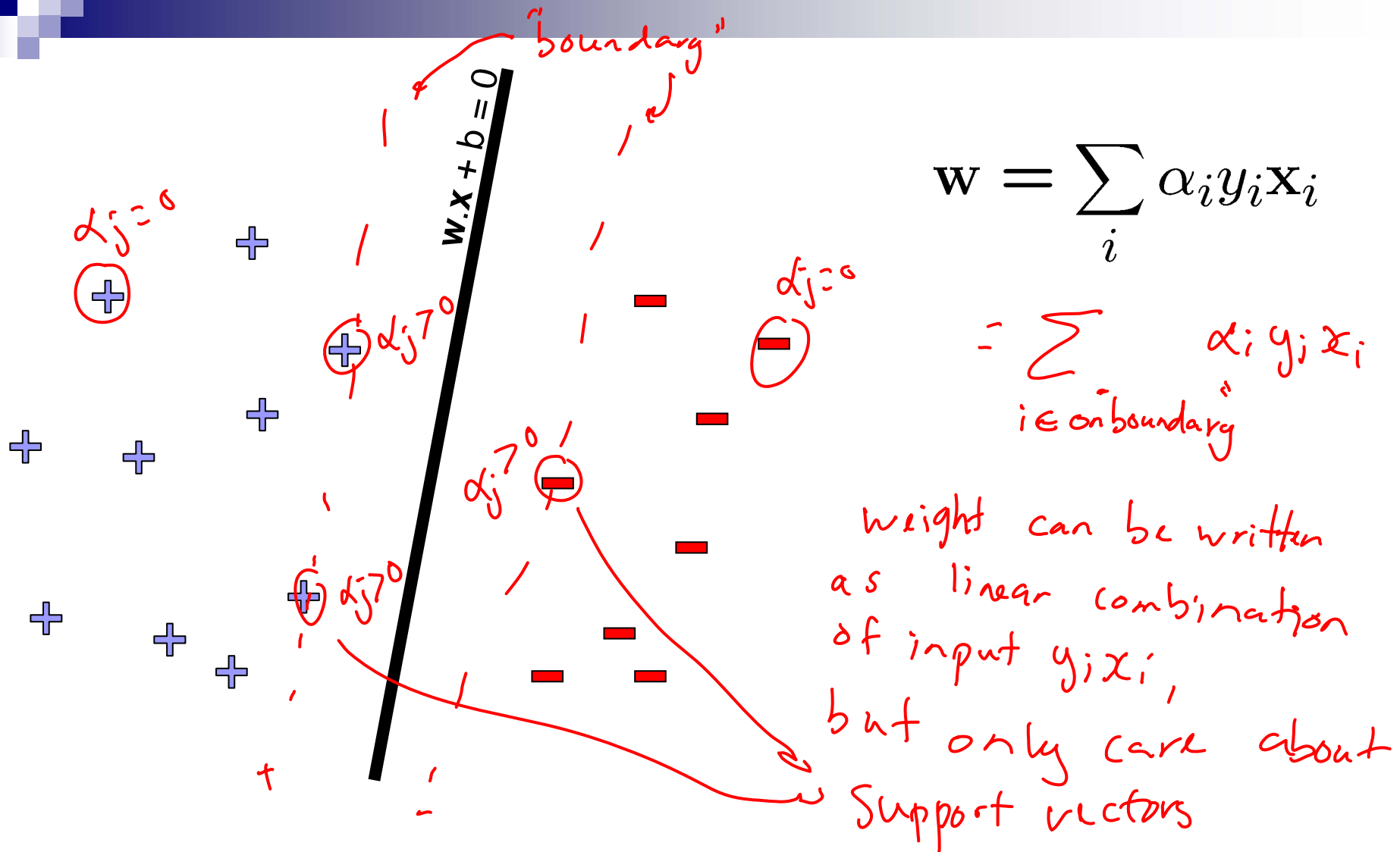
$$\mathbf{w} = \sum_j \alpha_j y_j \mathbf{x}_j$$

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

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



Dual SVM interpretation



Dual SVM formulation – the linearly separable case

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

$$\sum_i \alpha_i y_i = 0$$

$$\alpha_i \geq 0$$

dual program, ^{solve} SVM:
→ solve dual
obtain the α

get w, b

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

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

for any k where $\alpha_k > 0$

obj function dual → quadratic → dual quadratic program

Dual SVM derivation – the non-separable case

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

$$\begin{aligned} L(\mathbf{w}, b, \alpha, \mu) = & \frac{1}{2} \mathbf{w} \cdot \mathbf{w} + C \sum_j \xi_j - \sum_j \alpha_j [(\mathbf{w} \cdot \mathbf{x}_j + b) y_j - 1 + \xi_j] \\ & - \sum_j \mu_j \xi_j \end{aligned}$$

Dual SVM formulation – the non-separable case

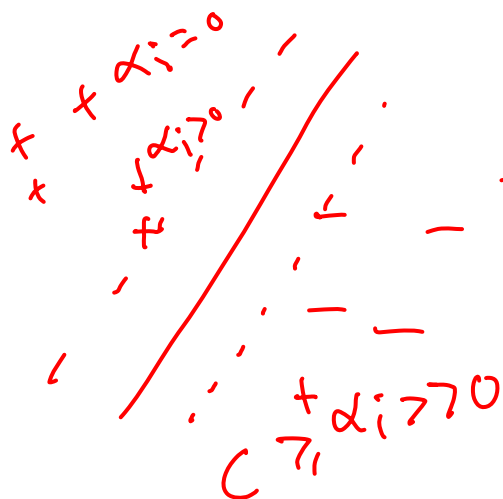
~~minimize~~ α $\sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \mathbf{x}_j$

N training examples $\{\alpha_1, \alpha_2, \dots, \alpha_N\}$

$$\sum_i \alpha_i y_i = 0$$

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

α_i can't be too large



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

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

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

Why did we learn about the dual SVM?

- There are some quadratic programming algorithms that can solve the dual faster than the primal
- But, more importantly, the “kernel trick”!!!
 - Another little detour...

Reminder from last time: What if the data is not linearly separable?

Use features of features of features....

$$\Phi(\mathbf{x}) : \mathbb{R}^m \mapsto F$$

↪ feature mapping $\mathbf{x} = (x_1, x_2)$

$$\phi(\mathbf{x}) = \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ x_1^2 \\ x_2^2 \\ \sin x_1 \\ \vdots \\ x_1^4 \\ \vdots \end{pmatrix}$$

$$\phi(\mathbf{x}) = \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ x_1^2 \\ x_2^2 \\ x_1 x_2 \\ x_1^3 \\ x_2^3 \\ x_1^2 x_2 \\ x_1 x_2^2 \\ x_1^4 \\ \vdots \end{pmatrix}$$

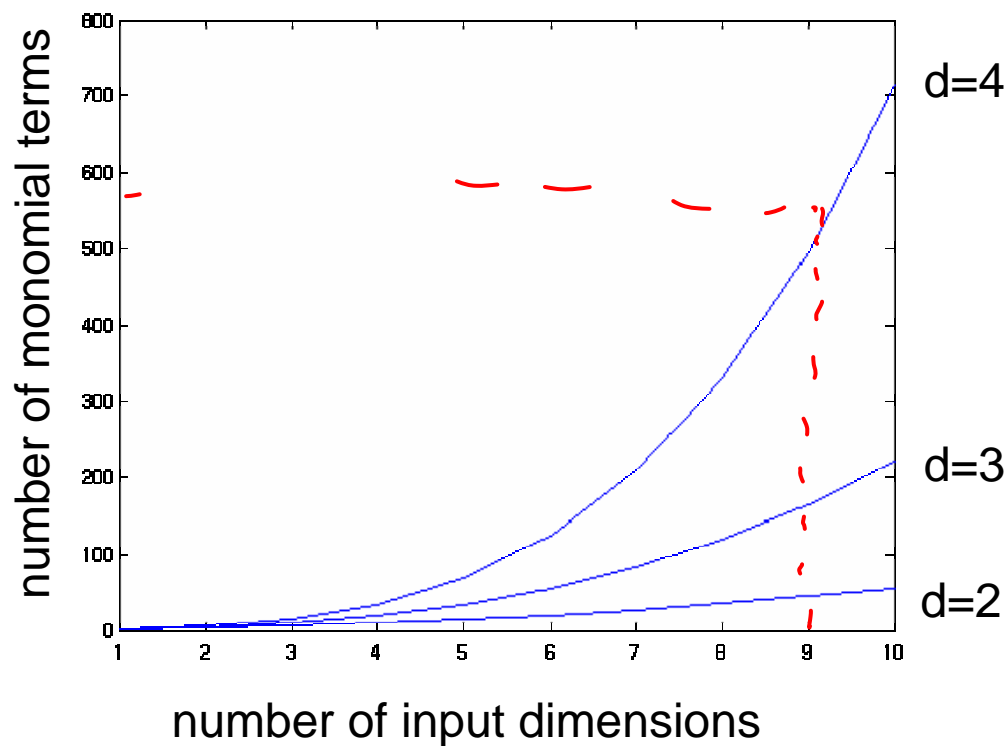
input d dim \rightarrow poly degree p
size of $\phi(\mathbf{x})$

Feature space can get really large really quickly!

Higher order polynomials

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

degree of poly
m – input features
d – degree of polynomial



grows fast!
d = 6, m = 100
about 1.6 billion terms

Dual formulation only depends on dot-products, not on \mathbf{w} !

only thing is \mathbf{x}

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

$$\sum_i \alpha_i y_i = 0$$

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

dot prod.
 $\mathbf{x}_i \mathbf{x}_j = \mathbf{x}_i \cdot \mathbf{x}_j$

no \mathbf{w} !

use features

$\phi(\mathbf{x})$

all I need is $\phi(\mathbf{x}_j) \cdot \phi(\mathbf{x}_i)$

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

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

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

$$\sum_i \alpha_i y_i = 0$$

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

Dot-product of polynomials

$\Phi(\mathbf{u}) \cdot \Phi(\mathbf{v})$ ^{use} ~~=~~ polynomials of degree d ^{exactly} _{u.v}

degree 1 $\Rightarrow \Phi(\mathbf{u}) = \begin{pmatrix} u_1 \\ u_2 \end{pmatrix}$ $\Phi(\mathbf{v}) = \begin{pmatrix} v_1 \\ v_2 \end{pmatrix}$ $\Phi(\mathbf{u}) \cdot \Phi(\mathbf{v}) = u_1 v_1 + u_2 v_2$

degree 2 $\Rightarrow \Phi(\mathbf{u}) = \begin{pmatrix} u_1^2 \\ u_1 u_2 \\ u_2 u_1 \\ u_2^2 \end{pmatrix}$ $\Phi(\mathbf{v}) = \begin{pmatrix} v_1^2 \\ v_1 v_2 \\ v_2 v_1 \\ v_2^2 \end{pmatrix}$ $\Phi(\mathbf{u}) \cdot \Phi(\mathbf{v}) = u_1^2 v_1^2 + u_1 u_2 v_1 v_2 + u_2 u_1 v_2 v_1 + u_2^2 v_2^2$
 \swarrow 12 multiplies

$$\begin{aligned} \Phi(\mathbf{u}) \cdot \Phi(\mathbf{v}) &= u_1^2 v_1^2 + u_1 u_2 v_1 v_2 + u_2 u_1 v_2 v_1 + u_2^2 v_2^2 \\ &= (u_1 v_1 + u_2 v_2)^2 \quad \leftarrow 3 \text{ multiplies} \\ &= (\mathbf{u} \cdot \mathbf{v})^2 \end{aligned}$$

Polynomials of degree exactly d: $\Phi(\mathbf{u}) \cdot \Phi(\mathbf{v}) = (\mathbf{u} \cdot \mathbf{v})^d$

Finally: the “kernel trick”!

i, j all pairs of data points including data points with itself.

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

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

$$\sum_i \alpha_i y_i = 0$$

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

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

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

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

*if using poly. degree exactly d
 $\forall i, j$ compute $\mathbf{x}_i \cdot \mathbf{x}_j$*

Set $K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i \cdot \mathbf{x}_j)^d$

Solve dual Q.P.

get α

Finally: the “kernel trick”!

i, j all pairs of data points including data points with itself.

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

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

$$\sum_i \alpha_i y_i = 0$$

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

- Never represent features explicitly
 - Compute dot products in closed form
- Constant-time high-dimensional dot-products for many classes of features
- Very interesting theory – Reproducing Kernel Hilbert Spaces
 - Not covered in detail in 10701/15781, more in 10702

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

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

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

Polynomial kernels

- All monomials of degree d in $O(d)$ operations:

$$\Phi(\mathbf{u}) \cdot \Phi(\mathbf{v}) = (\mathbf{u} \cdot \mathbf{v})^d = \text{polynomials of degree } d \text{ exactly}$$

- How about all monomials of degree up to d ?

- Solution 0: $\Phi(\mathbf{u}) \cdot \Phi(\mathbf{v}) = \sum_{i=0}^d (\mathbf{u} \cdot \mathbf{v})^i$ $d=2$

- Better solution: $(\mathbf{u} \cdot \mathbf{v} + 1)^2 = (\mathbf{u} \cdot \mathbf{v})^2 + \mathbf{u} \cdot \mathbf{v} + \mathbf{v} \cdot \mathbf{u} + 1$

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

$O(d)$
time

Common kernels

- Polynomials of degree ^{exactly} d

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

- Polynomials of degree up to ^{including} d

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

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

- Sigmoid

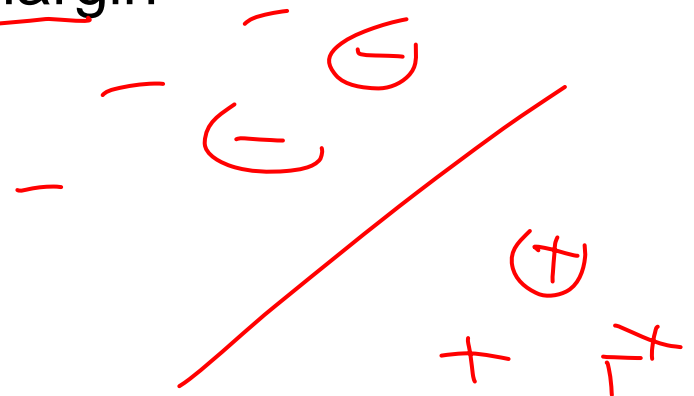
$$K(\mathbf{u}, \mathbf{v}) = \tanh(\underline{\eta \mathbf{u} \cdot \mathbf{v} + \nu})$$

↙ correspond infinite dimensional feature space
 $\dim[\phi(x)] = \infty$

Overfitting?

- Huge feature space with kernels, what about overfitting???
 - Maximizing margin leads to sparse set of support vectors
 - Some interesting theory says that SVMs search for simple hypothesis with large margin
 - Often robust to overfitting

sparse solutions \rightarrow
a few support vectors
 \rightarrow less overfitting



What about at classification time

- For a new input \mathbf{x} , if we need to represent $\Phi(\mathbf{x})$, we are in trouble! *if have to write w, b , too large*
- Recall classifier: $\text{sign}(\mathbf{w} \cdot \Phi(\mathbf{x}) + b)$
- Using kernels we are cool!

$$K(\mathbf{u}, \mathbf{v}) = \Phi(\mathbf{u}) \cdot \Phi(\mathbf{v})$$

$$\mathbf{w} \cdot \Phi(\mathbf{x}) = \sum_i \alpha_i y_i \underbrace{\Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}_i)}_{\text{easy to compute}}$$

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

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

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

SVMs with kernels

- Choose a set of features and kernel function
- Solve dual problem to obtain support vectors α_i
- At classification time, compute:

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

new x old data





$$b = y_k - \sum_i \alpha_i y_i K(\mathbf{x}_k, \mathbf{x}_i)$$

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

Classify as

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

What's the difference between SVMs and Logistic Regression?

	SVMs	Logistic Regression
Loss function		
High dimensional features with kernels		

Kernels in logistic regression

$$P(Y = 1 \mid x, \mathbf{w}) = \frac{1}{1 + e^{-(\mathbf{w} \cdot \Phi(\mathbf{x}) + b)}}$$

- Define weights in terms of support vectors:

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

$$\begin{aligned} P(Y = 1 \mid x, \mathbf{w}) &= \frac{1}{1 + e^{-(\sum_i \alpha_i \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}) + b)}} \\ &= \frac{1}{1 + e^{-(\sum_i \alpha_i K(\mathbf{x}, \mathbf{x}_i) + b)}} \end{aligned}$$

- Derive simple gradient descent rule on α_i

What's the difference between SVMs and Logistic Regression? (Revisited)

	SVMs	Logistic Regression
Loss function	Hinge loss	Log-loss
High dimensional features with kernels	Yes!	Yes!

What you need to know



- Dual SVM formulation
 - How it's derived
- The kernel trick
- Derive polynomial kernel
- Common kernels
- Kernelized logistic regression
- Differences between SVMs and logistic regression

Acknowledgment



- SVM applet:
 - <http://www.site.uottawa.ca/~gcaron/applets.htm>