



SVMs, Duality and the Kernel Trick

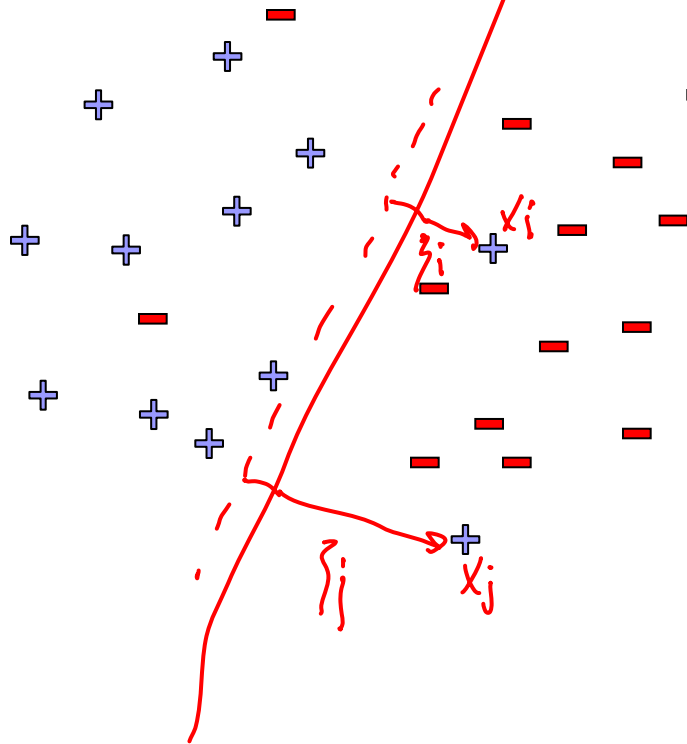
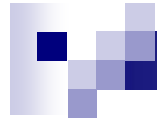
Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

February 26th, 2007

SVMs reminder



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

slack penalty

each point must have margin ≥ 1

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

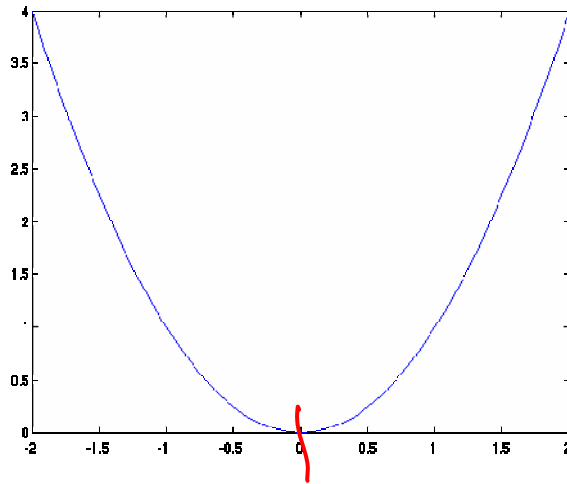
- ☐ Constrained optimization!

*some restrictions apply
see slides for
details...*

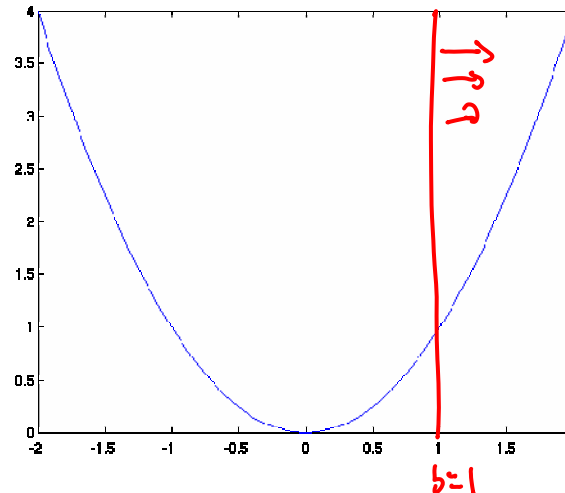
Constrained optimization

$$\min_x x^2$$

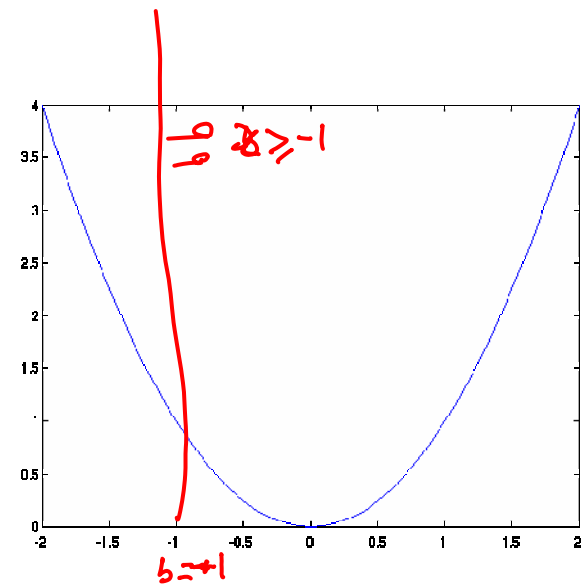
$$\text{s.t. } \underline{x \geq b}$$



Suppose no constraints
 min: $x=0$



$b > 0$ $x \geq b$
 min: $x=1$
 $x=b$



$b < 0$
 min: $x=0$
 constraint irrelevant

Lagrange multipliers – Dual variables



$$\min_x x^2$$

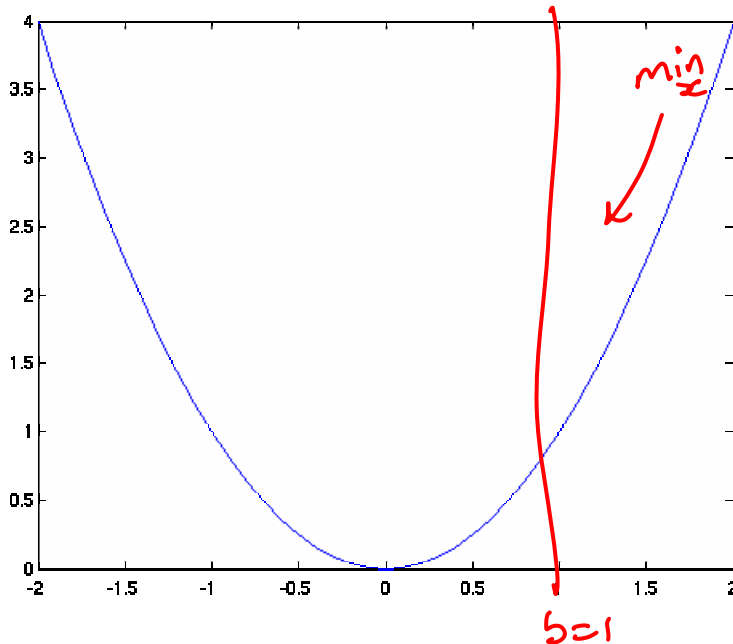
$$\text{s.t. } x \geq b$$

Moving the constraint to objective function

Lagrangian:

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

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



Solve:

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

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

suppose I pick:

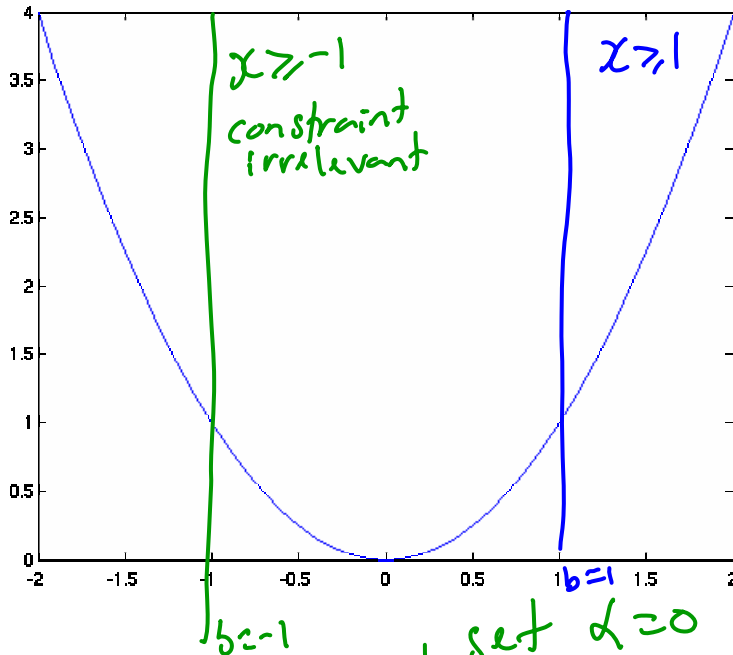
$x < b$ $\alpha \rightarrow +\infty$
violates the constraint $\Rightarrow L \rightarrow +\infty$

constraint not important

$x > b$ $\alpha = 0$
don't violate the constraint $\Rightarrow L(x, \alpha) = x^2$

$x = b$ α anything
constraint important $L(x, \alpha) = x^2$

Lagrange multipliers – Dual variables



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

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

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

$$\alpha = 2x \quad \text{if } x \geq 1$$

$$\frac{\partial L}{\partial \alpha} = 0 \Rightarrow x = 1$$

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

OK!

constraint plays
role $\Rightarrow \alpha > 0$

no way

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

unless $x = -1$
 but $x = -1, \alpha < 0$ unhappy!!

set $\alpha = 0$
 plays no
 role in
 Lagrangian

Dual SVM derivation (1) – the linearly separable case

equations look nicer

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

$$(w \cdot x_j + b) y_j \geq 1, \quad \forall j \in \text{Dataset}$$

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

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

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

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

x_j is d -dim

vector for d features for point j

give me d 's
give you w 's

Dual SVM derivation (2) – the linearly separable case

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

for point j
constraint is irrelevant: $\alpha_j = 0$
 $\Rightarrow (\mathbf{w} \cdot \mathbf{x}_j + b) y_j > 1$

constraint is relevant: $\alpha_j > 0$
 $(\mathbf{w} \cdot \mathbf{x}_j + b) y_j = 1$

Finding b : find any point where $\alpha_j > 0$
and $b \cdot y_j = 1 - \mathbf{w} \cdot \mathbf{x}_j y_j$

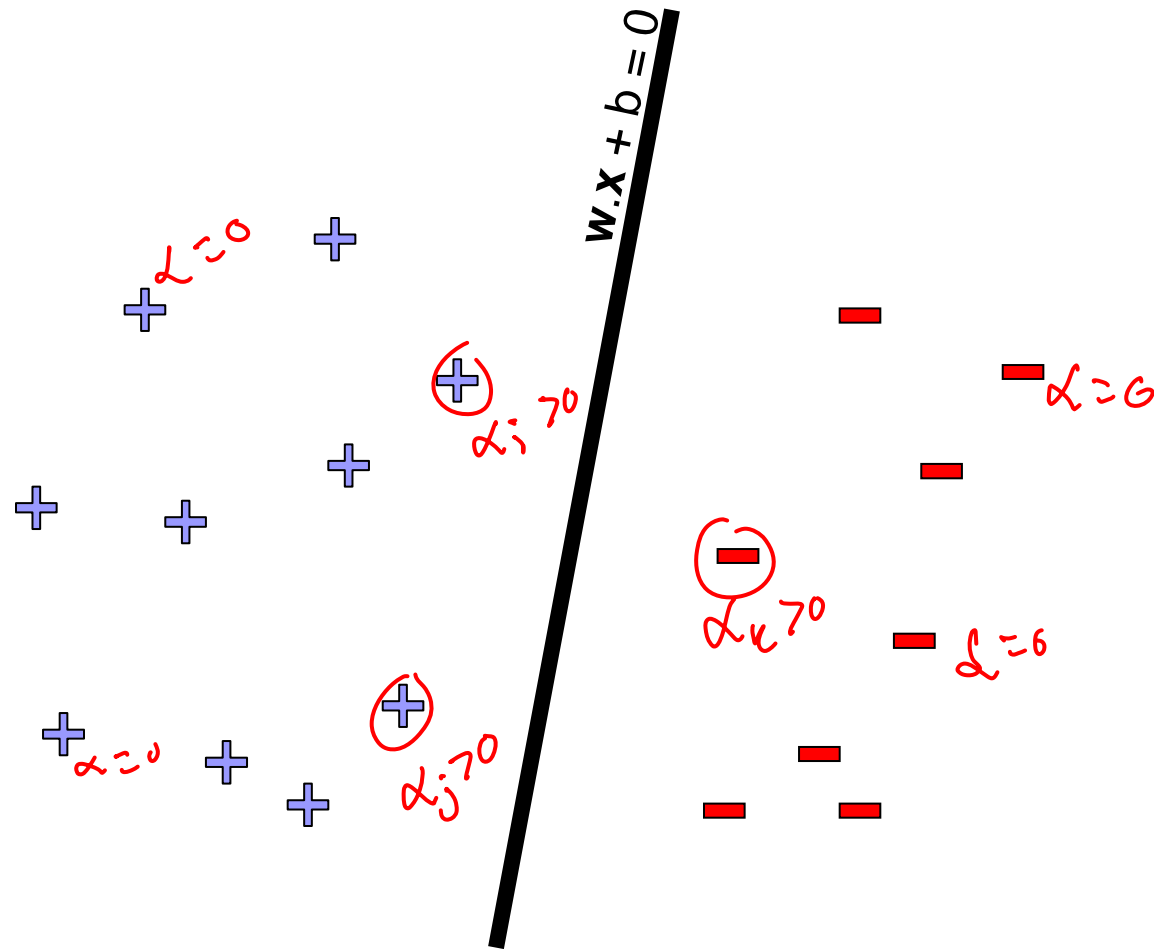
$$\mathbf{w} = \sum_j \alpha_j y_j \mathbf{x}_j$$

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

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

for any k where $\alpha_k > 0$

Dual SVM interpretation



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

$$\mathbf{w} = \sum_{\substack{j: \alpha_j > 0 \\ \text{or} \\ j \in \text{Support} \\ \text{vectors}}} \alpha_j y_j \mathbf{x}_j$$

weights are linear combination of feature values of the support vectors

Dual SVM formulation – the 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
svm

solve this \Rightarrow get α 's

\Rightarrow get w and b

\Rightarrow solved the SVM
problem

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

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

for any k where $\alpha_k > 0$

Solve by Quadratic Programming ...

Dual SVM derivation – the non-separable case

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

\swarrow slack penalty
 \swarrow slack variables
 \swarrow α_j
 \swarrow μ_j

$$\begin{aligned} L(\mathbf{w}, \alpha, \mu) = & \frac{1}{2} \mathbf{w} \cdot \mathbf{w} + C \sum_j \xi_j \\ & - \sum_j \alpha_j [(w \cdot x_j + b) y_j - 1 + \xi_j] \\ & - \sum_j \mu_j [\xi_j - 0] \end{aligned}$$

Dual SVM formulation – the non-separable case

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

$$\sum_i \alpha_i y_i = 0$$

$$\underline{C} \geq \alpha_i \geq 0 \quad \forall_i$$

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

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

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

compared to
separable case: only difference

$$\alpha_i \leq C \quad \forall_i$$

intuitively, don't give me alphas that
are too large

Announcements

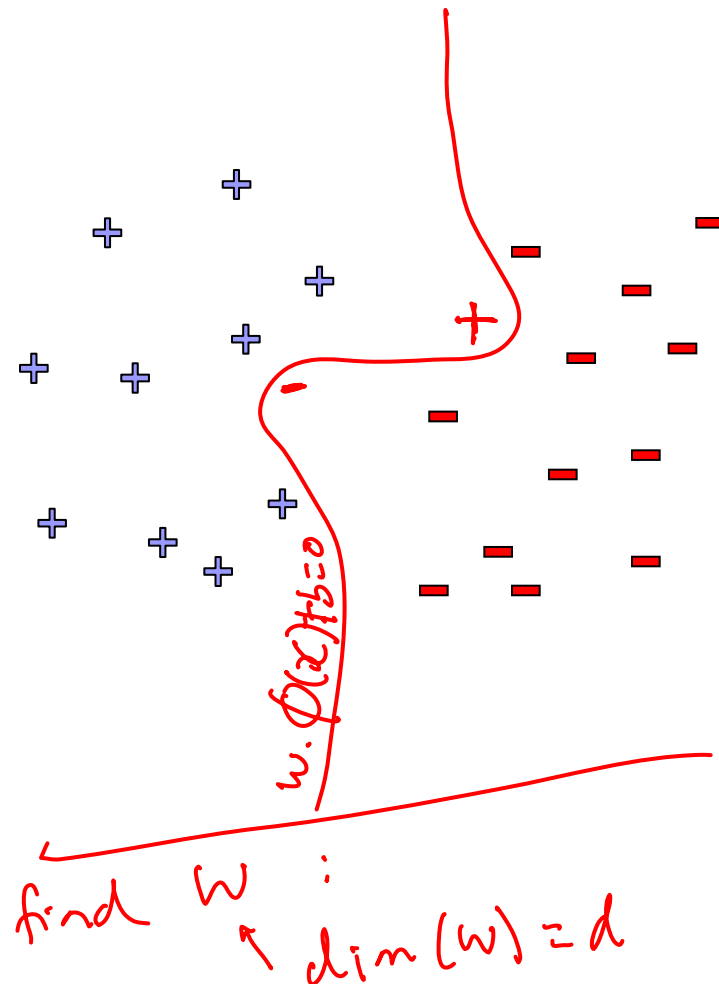


- Class projects out later this week

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....

$$\Phi(x) : R^m \mapsto F$$

$\phi(x) =$

mapping to high dim features

x_1
 x_2
 $x_1 x_2$
 x_1^2
 $\sin x_1$
 e^{-x_1/x_2}

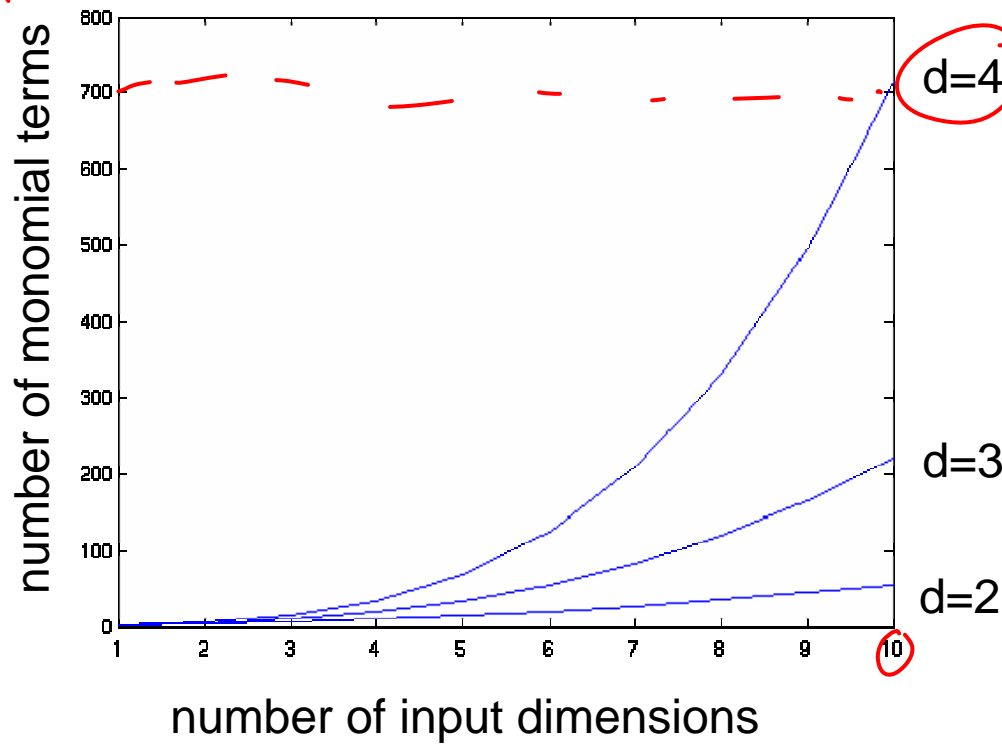
d

Feature space can get really large really quickly!

Higher order polynomials

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

$\phi(x)$



m – input features
d – degree of polynomial

$d=2$
 $m=2$

$x = \begin{pmatrix} v_1 \\ v_2 \end{pmatrix}$

$\phi(x) = \left\{ \begin{array}{l} v_1 \\ v_2 \\ v_1 v_2 \\ v_1^2 \\ v_2^2 \end{array} \right\}$

grows fast!

d = 6, m = 100

about 1.6 billion terms

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

$$\mathbf{x}_i = \begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_d \end{pmatrix} \quad \mathbf{x}_j = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_d \end{pmatrix}$$

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

$$\sum_i \alpha_i y_i = 0$$

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

$$\mathbf{x}_i \cdot \mathbf{x}_j = v_1 \mu_1 + v_2 \mu_2 + \dots + v_d \mu_d$$

$$\phi(\mathbf{x}) = \begin{pmatrix} v_1 \\ v_2 \\ v_1 v_2 \\ v_1^2 \\ v_2^2 \\ \vdots \end{pmatrix}$$

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

how many terms?

$$n^2$$

$n \in \# \text{ data points}$

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

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

$$\sum_i \alpha_i y_i = 0$$

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

Dot-product of polynomials

$$\mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$$
$$v = \begin{pmatrix} v_1 \\ v_2 \end{pmatrix}$$

$\Phi(u) \cdot \Phi(v) =$ polynomials of degree d

$$d=1 \quad \phi(\mu) = \mu$$

$$\phi(\mu) \cdot \phi(v) = \mu \cdot v = \mu_1 v_1 + \mu_2 v_2$$

$$d \geq 2 \quad \phi(\mu) = \begin{pmatrix} \mu_1^2 \\ \mu_1 \mu_2 \\ \mu_2 \mu_1 \\ \mu_2^2 \end{pmatrix}$$

$$\phi(\mu) \cdot \phi(v) =$$

$$\mu_1^2 \cdot v_1^2 + 2\mu_1 \mu_2 \cdot v_1 v_2 + \mu_2^2 \cdot v_2^2$$

$$= (\mu_1 v_1 + \mu_2 v_2)^2 = (\mu \cdot v)^2$$

degree $= d$

$$K(\mu, v) = \phi(\mu) \cdot \phi(v) = (\mu \cdot v)^d$$

$\equiv O(d)$ multiplications

Finally: the “kernel trick”!

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

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

$$\sum_i \alpha_i y_i = 0$$

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

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

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

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

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

- Very interesting theory – Reproducing Kernel Hilbert Spaces
 - Not covered in detail in 10701/15781, more in 10702

Polynomial kernels

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

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

- How about all monomials of degree up to d?

- Solution 0: $\sum_{i=1}^d (\mathbf{u} \cdot \mathbf{v})^i \equiv O(d^2)$

- Better solution: $(\mathbf{u} \cdot \mathbf{v} + 1)^2 = \underbrace{(\mathbf{u} \cdot \mathbf{v})^2 + 2\mathbf{u} \cdot \mathbf{v} + 1}_{\text{got poly monials upto } d=2}$

$$\text{poly degree upto } d \equiv \Phi(\mathbf{u}) \cdot \Phi(\mathbf{v}) = (\mathbf{u} \cdot \mathbf{v} + 1)^d \equiv O(d)$$

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 kernels

Gaussian kernels

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

equivalent to
 $\phi(\mathbf{u})$ has
infinite
dimensionality

- Sigmoid

Sigmoid:

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

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

What about at classification time

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

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

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

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

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:

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

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

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


Classify as

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

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