## Support Vector Machines: Kernels



Covering chapters (HTF): 4.5, 12

# R Greiner Department of Computing Science University of Alberta

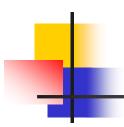
Much of this is taken from Andrew W. Moore; CMU http://www.cs.cmu.edu/~awm/tutorials

+ Nello Cristianini, Ron Meir, Ron Parr, Barnabas Poczos



#### **Outline**

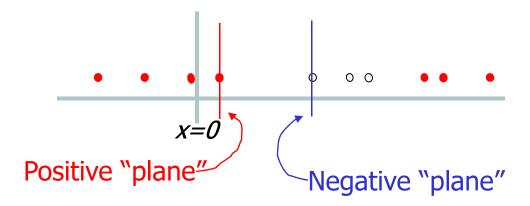
- Foundations
  - Linear Programming;Primal/Dual
  - Constrained Optimization
    - Lagrange Multipliers
    - KKT
  - Perceptron Factoids
    - Dual Representation
- <u>"Best"</u> Linear Separator: Max Margin!
- Coping with Non-Linearly Separated Data
  - Slack Variables
- Kernel Trick
- Regression



#### Hard 1-dimensional Dataset

What would SVMs do with this data?

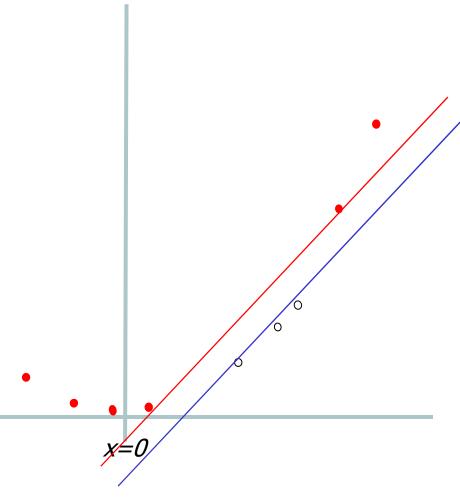
Not a big surprise



Doesn't look like slack variables will save us this time...



#### Hard 1-dimensional Dataset



Make up a new feature!

$$z_k = (x_k, x_k^2)$$

... computed from original feature(s)

Separable! ... MAGIC!

Now run linear SVM on this "augmented" data

# 4

#### ... New Features from Old ...

- Here: mapped  $\Re \to \Re^2$  by  $\Phi: x \to [x, x^2]$ 
  - Using "extra dimensions" ⇒ linearly separable!
- In general,
  - Start with vector  $\mathbf{x} \in \Re^k$
  - Want to add in  $x_1^2$ ,  $x_2^2$ , ...
  - Probably want other terms eg  $x_2 \cdot x_7$ , ...
  - Which ones to include?
- Why not ALL OF THEM? (If some features are linearly-separable, then any SUPERSET is)
- - $\mathbf{R}^3 o \mathfrak{R}^{10}$
  - In general, #features grows:

$$1+m \rightarrow 1+m+m+{m \choose 2} = \frac{(m+2)(m+1)}{2} \approx \frac{m^2}{2}$$

## Implied Algorithm

- Training: Given R training instances, each in  $\mathfrak{R}^m$ 
  - 1. Map each  $\Re^m$  –tuple  $\mathbf{x}_i$  to  $\Re^{\frac{m}{2}}$  –tuple  $\Phi(\mathbf{x}_i)$
  - Learn SVM classifier wrt these  $\Phi(\mathbf{x}_i)$  tuples
- Performance: Given new ℜ<sup>m</sup> –tuple x
  - 1. Map this **x** to  $\Re^{\frac{m}{2}}$  tuple  $\Phi(\mathbf{x})$
  - 2. Apply learned SVM classifier to  $\Phi(\mathbf{x})$

## Original "Linear" QP

$$\max_{\lambda} \sum_{k=1}^{R} \lambda_k - \frac{1}{2} \sum_{k=1}^{R} \sum_{r=1}^{R} \lambda_k \lambda_r Q_{kr} \qquad Q_{kr} = y_k y_r (\mathbf{x}_k \cdot \mathbf{x}_r)$$

$$0 \le \lambda_k \le C \quad \forall k \qquad \sum \lambda_k y_k = 0$$

$$\sum_{k=1}^{R} \lambda_k y_k = 0$$

#### Then define:

$$\mathbf{w} = \sum_{k=1}^{K} \lambda_k y_k \, \mathbf{x}_k$$

$$b = y_K(1 - \varepsilon_K) - \mathbf{w} \cdot \mathbf{x}_K$$

where 
$$K = arg \max_{k} \lambda_{k}$$

Then classify with:

$$f(\mathbf{x}, \mathbf{w}, \mathbf{b}) = sign(\mathbf{w} \cdot \mathbf{x} + \mathbf{b})$$

### **QP** using Basis Functions

$$\max_{\lambda} \sum_{k=1}^{R} \lambda_k - \frac{1}{2} \sum_{k=1}^{R} \sum_{r=1}^{R} \lambda_k \lambda_r Q_{kr} \qquad Q_{kr} = y_k y_r (\Phi(\mathbf{x}_k) \cdot \Phi(\mathbf{x}_r))$$

Subject to the constraints:

$$0 \le \lambda_k \le C \quad \forall k \qquad \sum \lambda_k y_k = 0$$

$$\sum_{k=1}^{R} \lambda_k y_k = 0$$

Then define:

$$\mathbf{w} = \sum_{k=1}^{K} \lambda_k y_k \, \Phi(\mathbf{x}_k)$$

$$b = y_K(1 - \varepsilon_K) - \mathbf{w} \cdot \Phi(\mathbf{x}_K)$$

where 
$$K = \underset{k}{\text{arg max }} \lambda_k$$

Then classify with:

$$f(\mathbf{x}, \mathbf{w}, \mathbf{b}) = sign(\mathbf{w} \cdot \Phi(\mathbf{x}) + \mathbf{b})$$



## **QP using Basis Functions**

$$\max_{\lambda} \sum_{k=1}^{R} \lambda_k - \frac{1}{2} \sum_{k=1}^{R} \sum_{r=1}^{R} \lambda_k \lambda_r Q_{kr}$$

$$Q_{kr} = y_k y_r (\Phi(\mathbf{x}_k) \cdot \Phi(\mathbf{x}_r))$$

Subject to the constraints:

$$0 \le \lambda_k \le C \quad \forall k$$

$$\sum_{k=1}^{R} \lambda_k y_k$$

Then define:

$$\mathbf{w} = \sum_{k=1}^{K} \lambda_k y_k \, \Phi(\mathbf{x}_k)$$

$$b = y_K(1 - \varepsilon_K) - \mathbf{w} \cdot \Phi(\mathbf{x}_K)$$

where 
$$K = \arg \max_{k} \lambda_k$$

 $\Phi(\mathbf{x}_k)$  only appears within dot product!

Then classify with:
$$f(\mathbf{x}, \mathbf{w}, \mathbf{b}) = \text{sign}(\mathbf{w} \cdot \mathbf{b})$$

$$= \text{sign}\left(\sum_{k} \lambda_{k} y_{k} \phi(\mathbf{x}_{k}) \cdot \phi(\mathbf{x}) + \mathbf{b}\right)$$

$$= \text{sign}\left(\sum_{k} \lambda_{k} y_{k} \left[\phi(\mathbf{x}_{k}) \cdot \phi(\mathbf{x})\right] + \mathbf{b}\right)$$

## **Implied Algorithm**

- Training: Given R training instances, each in  $\mathfrak{R}^{\mathsf{m}}$ 
  - 1. Map each  $\Re^m$  –tuple  $\mathbf{x}_i$  to  $\Re^{m^*m/2}$  –tuple  $\Phi(\mathbf{x}_i)$
  - Learn SVM classifier wrt these  $\Phi(\mathbf{x}_i)$  tuples
- Performance: Given new \( \partial^m \text{tuple } \x \)
  - 1. Map this **x** to  $\Re^{m*m/2}$  —tuple  $\Phi(\mathbf{x})$
  - 2. Apply learned SVM classifier to  $\Phi(\mathbf{x})$

#### Issue:

One more trick!

- This  $\Phi(.)$  operation is expensive!!  $O(m^2)$
- What if want  $\Phi'(.)$  that deals with "x<sup>3</sup>", or "x<sup>4</sup>", or ...

# 4

## Higher Dimensional Space

#### Given

- $\mathbf{x} = (x_1, x_2, x_3)^T$
- $\mathbf{y} = (y_1, y_2, y_3)^T$

#### Standard model:

$$\mathbf{x} \cdot \mathbf{y} = \sum_{i} x_{i} y_{i} = x_{1}y_{1} + x_{2}y_{2} + x_{3}y_{3}$$

What if some squared terms might matter:

$$\Phi(\mathbf{x}) = (x_1 x_2, x_1 x_3, x_2 x_3, x_1^2, x_2^2, x_3^2)^T ?$$

## Quadratic Basis Function

$$(\mathbf{a}^{T}\mathbf{b})^{2} = (\mathbf{a}_{1} \ \mathbf{b}_{1} + \mathbf{a}_{2} \ \mathbf{b}_{2})^{2}$$

$$= \mathbf{a}_{1}^{2} \mathbf{b}_{1}^{2} + 2 \mathbf{a}_{1} \mathbf{b}_{1} \mathbf{a}_{2} \mathbf{b}_{2} + \mathbf{a}_{2}^{2} \mathbf{b}_{2}^{2}$$

$$= (\mathbf{a}_{1}^{2}, \sqrt{2} \mathbf{a}_{1} \mathbf{a}_{2}, \mathbf{a}_{2}^{2}) \begin{pmatrix} \mathbf{b}_{1}^{2} \\ \sqrt{2} \mathbf{b}_{1} \mathbf{b}_{2} \\ \mathbf{b}_{2}^{2} \end{pmatrix}$$

$$= \phi(\mathbf{a}) \bullet \phi(\mathbf{b})$$

where

$$\phi(z) = \phi([z_1, z_2]) = (z_1^2, \sqrt{2} z_1 z_2, z_2^2)$$

# 4

## Higher Dimensional Space

#### Want to consider

$$\Phi(\mathbf{x}) = \left( \sqrt{2} \mathbf{X}_1 \mathbf{X}_2 \sqrt{2} \mathbf{X}_1 \mathbf{X}_3 \sqrt{2} \mathbf{X}_2 \mathbf{X}_3, \quad \mathbf{X}_1^2, \quad \mathbf{X}_2^2, \mathbf{X}_3^2 \right)^{\mathrm{T}}$$

$$\Phi(\mathbf{x})^{\mathrm{T}} \Phi(\mathbf{y}) = \mathbf{2} \mathbf{x}_1 \mathbf{x}_2 \mathbf{y}_1 \mathbf{y}_2 + \mathbf{2} \mathbf{x}_1 \mathbf{x}_3 \mathbf{y}_1 \mathbf{y}_3 + \mathbf{2} \mathbf{x}_2 \mathbf{x}_3 \mathbf{y}_2 \mathbf{y}_3 + \mathbf{x}_1^2 \mathbf{y}_1^2 + \mathbf{x}_2^2 \mathbf{y}_2^2 + \mathbf{x}_3^2 \mathbf{y}_3^2$$

#### Computation complexity:

- Compute  $\Phi(\mathbf{x}) O(m^2)$  multiplications
- Compute  $\Phi(y) O(m^2)$  multiplications
- Compute  $\Phi(\mathbf{x})^T \Phi(\mathbf{y}) O(m^2)$  multiplications & additions
- $\rightarrow O(m^2)$  time

#### Small tweak

- Note  $\Phi(\mathbf{x})^{\mathrm{T}}\Phi(\mathbf{y}) = (\mathbf{x}^{\mathrm{T}}\mathbf{y})^2$
- Computational complexity: O( m )

# 4

## Higher Dimensional Space

Want to consider

$$\Phi(\mathbf{x}) = \left( \mathbf{\overline{2}} \mathbf{X}_{1} \mathbf{X}_{2} \mathbf{\overline{2}} \mathbf{X}_{1} \mathbf{X}_{3}, \mathbf{\overline{2}} \mathbf{X}_{2} \mathbf{X}_{3}, \mathbf{X}_{1}^{2}, \mathbf{X}_{2}^{2}, \mathbf{X}_{3}^{2} \right)^{T}$$

$$\Phi(\mathbf{x})^{T} \Phi(\mathbf{y}) = \mathbf{2} \mathbf{x}_{1} \mathbf{x}_{2} \mathbf{y}_{1} \mathbf{y}_{2} + \mathbf{2} \mathbf{x}_{1} \mathbf{x}_{3} \mathbf{y}_{1} \mathbf{y}_{3} + \mathbf{2} \mathbf{x}_{2} \mathbf{x}_{3} \mathbf{y}_{2} \mathbf{y}_{3} + \mathbf{x}_{1}^{2} \mathbf{y}_{1}^{2} + \mathbf{x}_{2}^{2} \mathbf{y}_{2}^{2} + \mathbf{x}_{3}^{2} \mathbf{y}_{3}^{2}$$

- Note  $\Phi(\mathbf{x})^{\mathrm{T}}\Phi(\mathbf{y}) = (\mathbf{x}^{\mathrm{T}}\mathbf{y})^2$
- Advantage:
  - Computational complexity: O(m) [not O(m²)]
- But...
  - Only considering THIS COMBINATION of terms
  - Why not  $(17 x_1 x_2, 5 x_1 x_3, 2 x_2 x_3, 4 x_1^2, x_2^2, 0 x_3^2)^T$



#### **Constant Term**

Linear Terms

**Pure** Quadratic Terms

## Another **Quadratic Basis Function**

 $\Phi(\mathbf{x}) =$ 

 $\sqrt{2}x_m$ 

Quadratic **Cross-Terms**  What about those  $\sqrt{2}$  s ?? ... stay tuned

# Quadratic Dot Products

$$\Phi(\mathbf{a}) \bullet \Phi(\mathbf{b}) =$$

$$\begin{pmatrix} 1 \\ \sqrt{2}a_1 \\ \sqrt{2}a_2 \\ \vdots \\ \sqrt{2}a_m \\ a_1^2 \\ a_2^2 \\ \vdots \\ \sqrt{2}b_m \\ b_1^2 \\ b_2^2 \\ \vdots \\ \sqrt{2}b_1b_2 \\ \sqrt{2}b_1b_2 \\ \sqrt{2}b_1b_3 \\ \vdots \\ \sqrt{2}a_1a_m \\ \sqrt{2}a_2a_3 \\ \vdots \\ \sqrt{2}a_2a_m \\ \vdots \\ \sqrt{2}b_2b_m \\ \vdots \\ \sqrt{2}b_mb_m \\ \sqrt{2}b_2b_3 \\ \vdots \\ \sqrt{2}b_2b_m \\ \vdots \\ \sqrt{2}b_mb_m \\ \sqrt{2}b_mb_m \\ \sqrt{2}b_mb_m \\ \end{bmatrix}$$

$$\begin{array}{c}
1 \\
+ \\
\sum_{i=1}^{m} 2a_{i}b_{i} \\
+ \\
+ \\
\sum_{i=1}^{m} \sum_{j=i+1}^{m} 2a_{i}a_{j}b_{i}b_{j}
\end{array}$$

# Quadratic Dot Products

$$\Phi(\mathbf{a}) \bullet \Phi(\mathbf{b}) =$$

$$1 + 2\sum_{i=1}^{m} a_i b_i + \sum_{i=1}^{m} (a_i b_i)^2 + \sum_{i=1}^{m} \sum_{j=i+1}^{m} 2a_i a_j b_i b_j$$

#### Now consider another fn of **a** and **b**:

$$(\mathbf{a} \cdot \mathbf{b} + 1)^2$$

$$= (\mathbf{a} \cdot \mathbf{b})^2 + 2\mathbf{a} \cdot \mathbf{b} + 1$$

$$= \left(\sum_{i=1}^{m} a_i b_i\right)^2 + 2\sum_{i=1}^{m} a_i b_i + 1$$

$$= \sum_{i=1}^{m} \sum_{j=1}^{m} a_i b_i a_j b_j + 2 \sum_{i=1}^{m} a_i b_i + 1$$

$$= \sum_{i=1}^{m} (a_i b_i)^2 + 2 \sum_{i=1}^{m} \sum_{j=i+1}^{m} a_i b_i a_j b_j + 2 \sum_{i=1}^{m} a_i b_i + 1$$

#### They're the same!

But this is only O(m) to compute... not  $O(m^2)$ 



## Higher Order Polynomials

$$Q_{kr} = y_k y_r \phi(x_k) \cdot \phi(x_r)$$

Poly- nomial	φ( <b>x</b> )	Cost to build Q <sub>kr</sub> matrix: traditional	Cost if m=100 inputs	φ(a)•φ(b)	Cost to build Q <sub>kr</sub> matrix: sneaky	Cost if m=100 inputs
Quadratic	All <i>m²/2</i> terms ≤ degree 2	$m^2R^2/4$	2 500 <i>R</i> <sup>2</sup>	( <b>a·b</b> +1) <sup>2</sup>	m I <sup>R2</sup> / 2	50 <i>R</i> <sup>2</sup>
Cubic	All <i>m³/6</i> terms ≤ degree 3	$m^3R^2/12$	83 000 <i>R</i> <sup>2</sup>	( <b>a·b</b> +1) <sup>3</sup>	$mR^2/2$	50 <i>R</i> <sup>2</sup>
Quartic	All <i>m⁴/24</i> terms ≤ degree 4	m <sup>4</sup> R <sup>2</sup> /48	1 960 000 <i>R</i> <sup>2</sup>	( <b>a·b</b> +1) <sup>4</sup>	m l <sup>R2</sup> / 2	50 <i>R</i> <sup>2</sup>

Kinda ... coefficients for terms are not independent ...



## QP using Quintic Basis Functions

This matrix requires R<sup>2</sup>/2 dot products.

not 75 million

But still worries...

In 100-d, each dot product requires 103 ops, ...  $Q_{kr} = y_k y_r (\Phi(x_k) \cdot \Phi(x_r))$ 

The use of Maximum Margin magically reduces this problem

#### Subject to the constraints:

- Overfitting due to enormous number of terms
  - The evaluation phase (doing a predictions on a test instance x) seems expensive...

as  $\mathbf{w} \cdot \phi(\mathbf{x})$  needs 75 million operations

#### Then define:

$$w = \sum_{k=1}^{K} \lambda_k y_k \, \Phi(\mathbf{x}_k)$$

$$\mathbf{w} \cdot \Phi(\mathbf{x}) = \sum_{k \text{ s.t. } \lambda_k > 0} \lambda_k y_k \Phi(\mathbf{x}_k) \Phi(\mathbf{x})$$

Only *S m* operations (*S*=#support vectors)

$$f(\mathbf{x}, \mathbf{w}, \mathbf{b}) = sign(\mathbf{w} \cdot \Phi(\mathbf{x}) + \mathbf{b})$$

$$= \operatorname{sign}\left(\sum_{k} \lambda_{k} y_{k} \, \boldsymbol{\phi}(\, \mathbf{x}_{k}) \cdot \boldsymbol{\phi}(\mathbf{x}) \, + \mathbf{b}\right)$$

$$= \operatorname{sign}\left(\sum_{k} \lambda_{k} y_{k} \, \left[\, \boldsymbol{\phi}(\, \mathbf{x}_{k}) \cdot \boldsymbol{\phi}(\mathbf{x}) \right] + \mathbf{b}\right)$$



#### The "Kernel Trick"!

$$\max_{\lambda} \sum_{k=1}^{R} \lambda_{k} - \frac{1}{2} \sum_{k=1}^{R} \sum_{r=1}^{R} \lambda_{k} \lambda_{r} y_{k} y_{r} K(\mathbf{x}_{k}, \mathbf{x}_{r})$$

$$K(\mathbf{x}_{k}, \mathbf{x}_{r}) = \boldsymbol{\phi}(\mathbf{x}_{i}) \cdot \boldsymbol{\phi}(\mathbf{x}_{j})$$

$$\sum_{\lambda_{i}} \lambda_{i} y_{i} = 0 \qquad 0 \leq \lambda_{k} \leq C \qquad \boldsymbol{w} = \sum_{i} \lambda_{i}$$

- $\sum_{i} \lambda_{i} y_{i} = 0 \qquad 0 \le \lambda_{k} \le C$  Never represent features explicitly
  - Compute dot products in closed form

$$m{w} = \sum_i \lambda_i y_i \ m{\phi}(m{x_i})$$
 $b = y_k - m{w} \cdot m{\phi}(m{x_k})$ 
for any k s.t.  $0 < \lambda_k < C$ 

■ ∃ constant-time high-dimensional dot-products for many classes of features

#### ... at classification time

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

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

$$w \cdot \Phi(x) = k \sum_{i} \lambda_{k} y_{k} K(x, x_{k})$$

$$b = y_{k} - \sum_{i} \lambda_{i} y_{i} K(x_{k}, x_{i})$$
for any  $k$  s.  $t$ .  $\lambda_{k} > 0$ 

$$\mathbf{w} = \sum_{k} \lambda_{k} y_{k} \, \boldsymbol{\phi}(\mathbf{x}_{k})$$

$$b = y_k - \boldsymbol{w} \cdot \boldsymbol{\phi}(x_k)$$

for any k s.t.  $0 < \lambda_k < C$ 

classify as 
$$sign(\sum_{k} \lambda_{k} y_{k} [\phi(\mathbf{x}_{k}) \cdot \phi(\mathbf{x})] + b)$$

$$= sign(\sum_{k} \lambda_{k} y_{k} K(\mathbf{x}_{k}, \mathbf{x}) + b)$$

Nb: never have to deal with  $\phi(x)$  ... only  $K(\mathbf{x}_k, \mathbf{x}_i)$ 

# 4

#### Classifying using SVMs with Kernels

- Training time: Given { [x<sub>i</sub>, y<sub>i</sub>] }
  - Choose kernel function  $\phi(\cdot)$  of kernel  $K(\cdot, \cdot)$
  - Solve dual problem to obtain parameter values  $\lambda_i$  (of support vectors)
  - Compute  $\mathbf{b} = \mathbf{y_k} \sum_i \lambda_i \mathbf{y_i} \mathbf{K}(\mathbf{x_i}, \mathbf{x_k})$  for any k s.t.  $0 < \lambda_k < C$
- Performance time, given x, compute (and return):

$$sign(\sum_{k} \lambda_{k} y_{k} K(\mathbf{x}_{k}, \mathbf{x}) + \mathbf{b})$$

### What makes a valid kernel?

- A *sufficient* (but not necessary) condition is for K(·, ·) to behave like a distance metric
  - Non-negative: K(x,y) ≥ 0
     K(x,x) = 0
  - Symmetric K(x,y) = K(y,x)
  - Obeys triangle inequality  $K(x,y) + K(y,z) \ge K(x,z)$
- In general,  $K(\cdot, \cdot)$  matrix must be symmetric, positive semi-definite
  - Given  $X = \{x_1, ..., x_n\}, K(X) = [K(x_i, x_j)]$
  - Matrix K(X) is positive semi-definite iff  $\forall y \ y^T K(X) y \ge 0$

#### Mercer Theorem

- Any symmetric positive definite matrix K
   can be a kernel matrix
  - ie,  $\exists \phi(\cdot)$  s.t.  $K(x,y) = \phi(x) \cdot \phi(y)$
  - as an inner product matrix in some space

- Matrix M is:
  - symmetric iff  $\mathbf{M} = \mathbf{M}^{\mathrm{T}}$
  - positive definite iff  $\forall y \ y^T \mathbf{M} \ y \ge 0$

#### **Common Kernels**

Polynomials of degree d

$$K_d^{poly}(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v})^d$$

Polynomials of degree up to d

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

Sigmoid

$$K_{\alpha,\beta}^{sigmoid}(\mathbf{u}, \mathbf{v}) = \tanh(\alpha \mathbf{u} \cdot \mathbf{v} + \beta)$$

Gaussian kernels

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

Corresponding  $\phi(x)$  is of infinite dimensionality!

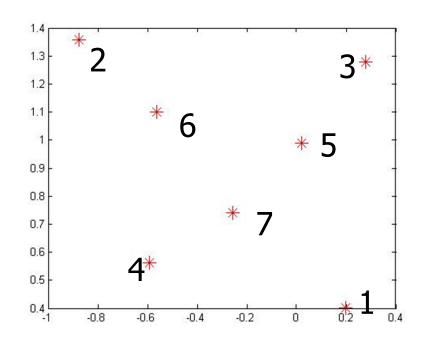
## -

## Finite example

- Given
  - a kernel function  $k: D \times D \to \Re$
  - FINITE set  $X = \{x_1, ..., x_r\} \subset D$
- construct associated  $\phi(\cdot)$
- Build Gram matrix  $G \in \Re^{r \times r}$  where

$$G_{i,j} = k(x_i, x_j)$$

Note G is symmetric, and positive semi-definite



### Finite example

#### Given

- 7 2D-points
- Choose a kernel k

Calculate 
$$G_{i,j} = \exp(-\frac{|x_i - x_j|^2}{10})$$

# -

#### Source of Kernels

- $K(\mathbf{x}, \mathbf{x}') = \mathbf{x}^{\mathsf{T}} \mathsf{A} \mathbf{x}'$ 
  - for any positive-semidefinite matrix A
- The set of kernels is closed under certain operations
  - ⇒ can make complex kernels from simple ones: modularity!
- If K, K' are kernels, then so are:
  - K + K'
  - cK for c>0
    - aK+bK′ for a,b >0
  - ...

#### Source of Kernels

- Can generate new kernels from old:
- If  $k_1(\mathbf{x},\mathbf{x}')$ ,  $k_2(\mathbf{x},\mathbf{x}')$  are kernels, then so is:
  - **x**<sup>T</sup> A **x**′ A any positive semidefinite matrix
  - $c \in \mathfrak{R}^+$
  - $f(x) k_1(x,x') f(x') f(.)$  any function
  - $q(k_1(x,x'))$  q(.) any poly function w/c coeff's  $\geq 0$
  - $\bullet$  exp(  $k_1(\mathbf{x},\mathbf{x}')$  )
  - $\mathbf{k}_1(\mathbf{x},\mathbf{x}') + \mathbf{k}_2(\mathbf{x},\mathbf{x}')$
  - $k_a(\mathbf{x}_a, \mathbf{x}'_a) + k_b(\mathbf{x}_b, \mathbf{x}'_b)$  - $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$  and  $k_a(.,.)$  kernel over "a" space,  $k_b(.,.)$  kernel over "b" space
  - **...**

# 4

#### Must Domain be $\Re^n$ ?

- So far, considering  $x \in D = \Re^n$
- But note: no computation uses x, just K(x, y)
  - $\Rightarrow$  so can have domain = strings, graphs, ...
- So Kernels for ...
  - Text
  - DNA strings
  - Graphs
  - **...**



## 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



#### **SVM** Performance

- Anecdotally SVMs do work very very well indeed.
  - Eg1: The best-known classifier on a well-studied hand-written-character recognition benchmark
  - Eg2: Many people doing practical real-world work claim that SVMs have saved them... when their other favorite classifiers did poorly.
- Lots of excitement and religious fervor about SVMs as of 2001...
- Still... some practitioners are a little skeptical...

# 4

## 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:  $\mathbf{w} = \sum_{i} \alpha_{i} \Phi(\mathbf{x}_{i})$ 

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

lacksquare Derive simple gradient descent rule on  $lpha_i$ 



# Differences between SVM and Logistic Regression

	SVMs	Logistic Regression	
Loss function	Hinge loss	Log-loss	
High dimensional features with kernels	Yes!	No	
Solution sparse	Often yes!	Almost always no!	
Semantics of output	"Margin"	Real probabilities	



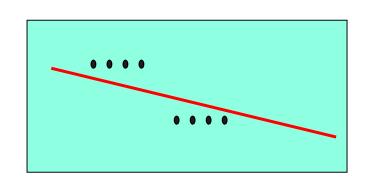
#### **Outline**

- Foundations
  - Linear Programming;Primal/Dual
  - Constrained Optimization
    - Lagrange Multipliers
    - KKT
  - Perceptron Factoids
    - Dual Representation
- <u>"Best"</u> Linear Separator: Max Margin!
- Coping with Non-Linearly Separated Data
  - Slack Variables
- Kernel Trick
- Regression

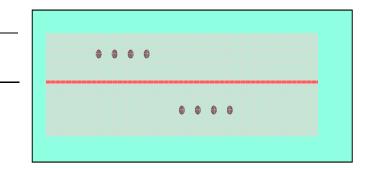


## Simple Example

- What is best linear fit:
  - Minimizes the sum of the squared error:  $\sum_{i} (y_i t_i)^2$

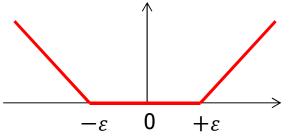


• What if no penalty if within  $\varepsilon$ 



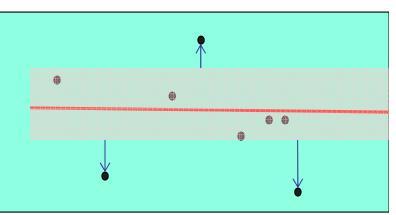
For each point,

$$E_{\varepsilon}(y,t) = \begin{cases} 0 & \text{if } |y-t| < \varepsilon \\ |y-t| - \varepsilon & \text{otherwise} \end{cases}$$





# **Error Function**

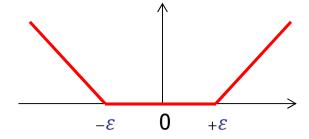


• If point is within  $\varepsilon$  of line:  $\overline{0}$  error

• If point is NOT within  $\varepsilon$  of line: Just that additional distance:  $|y - t| - \varepsilon$ 

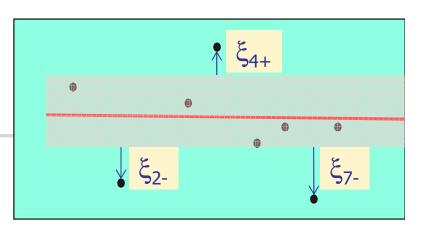
$$E_{\varepsilon}(y,t) = \begin{cases} 0 & \text{if } |y-t| < \varepsilon \\ |y-t| - \varepsilon \text{ otherwise} \end{cases}$$

• Error is  $\sum_{i} E_{\varepsilon}(y_{i}, t_{i})$ 





## Loss Function 🕏

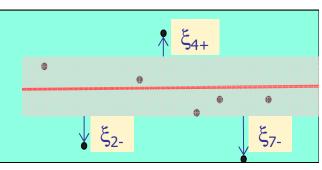


- Want "best"  $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$
- Typically:  $L_{LMS}(w) = C \sum_{i} (y_i(\mathbf{x}) t_i)^2 + \frac{1}{2} |w|^2$
- Instead:  $L_{SVR}(\mathbf{w}) = C \sum_i E_{\varepsilon}(y_i(\mathbf{x}), t_i) + \frac{1}{2} |\mathbf{w}|^2$
- wrt  $E_{\varepsilon}(y_i(x), t_i)$ :
  - No penalty if  $t_i \in [y_i \varepsilon, y_i + \varepsilon]$
  - Slack variables:  $\{\xi_{i+}, \xi_{i-}\}$ 
    - $t_i \leq y_i + \varepsilon + \xi_{i+1}$
    - $t_i \ge y_i \epsilon \xi_{i-}$

$$\sum_{i} \xi_{i+} + \xi_{i-}$$



#### **SVR Formulation**



- $\min_{\mathbf{w}} L_{SVR}(\mathbf{w}) = C \sum_{i} (\xi_{i+} + \xi_{i-}) + \frac{1}{2} |\mathbf{w}|^{2}$ s.t.
  - $t_i \leq y_i + \varepsilon + \xi_{i+1}$
  - $t_i \ge y_i \varepsilon \xi_{i-}$
  - $\xi_{i+}$ ,  $\xi_{i-} \geq 0$
- Use Lagrange Multipliers  $\{a_{n+}, a_{n-}, \mu_{n+}, \mu_{n-}\} \ge 0$

$$L_{SVR}(\mathbf{w}, ...) = C \sum_{i} (\xi_{i+} + \xi_{i-}) + \frac{1}{2} |\mathbf{w}|^{2}$$

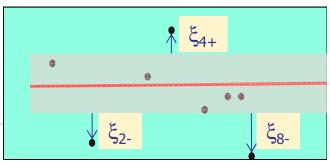
$$- \sum_{i} (\mu_{i+} \xi_{i+} + \mu_{i-} \xi_{i-})$$

$$- \sum_{i} a_{i+} (\epsilon + \xi_{i+} + y_{i} - t_{i})$$

$$- \sum_{i} a_{i-} (\epsilon + \xi_{i-} - y_{i} + t_{i})$$



#### SVR Solution



$$L_{SVR}(\mathbf{w}, \dots) = C \sum_{i} (\xi_{i+} + \xi_{i-}) + \frac{1}{2} |\mathbf{w}|^2 - \sum_{i} (\mu_{i+} \xi_{i+} + \mu_{i-} \xi_{i-})$$
$$- \sum_{i} a_{i+} (\epsilon + \xi_{i+} + y_i - t_i) - \sum_{i} a_{i-} (\epsilon + \xi_{i-} - y_i + t_i)$$

- Set derivatives to 0, solve for  $\{a_{n+}, a_{n-}, \mu_{n+}, \mu_{n-}\}$
- $L'(a_+, a_-) = \frac{1}{2} \sum_i \sum_j (a_{i+} a_{i-}) (a_j a_j) \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)$ Quadratic Program!!  $-\epsilon \sum_{i} (a_{i+} - a_{i-}) + \sum_{i} (a_{i+} - a_{i-}) t_{i}$ • s.t.  $0 \le a_{i+} \le C$ ,  $0 \le a_{i-} \le C$
- Predictions for new x:

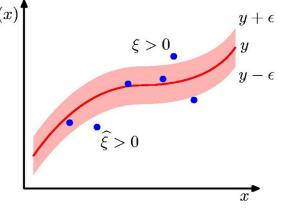
$$y(\mathbf{x}) = \sum_{i} (a_{i+} - a_{i-}) \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x})$$

<u>↑</u>



# SVM Regression, con't y(x)

$$y(\mathbf{x}) = \sum_{i} (a_{i+} - a_{i-}) K(\mathbf{x}_i, \mathbf{x})$$



- Can ignore  $x_n$  unless either  $a_{n+}>0$  or  $a_{n-}>0$ 
  - $a_{n+}>0$  only if  $t_n = y_n + \varepsilon + \xi_{n+}$ ie, if on upper boundary of  $\varepsilon$ -tube ( $\xi_{n+}=0$ ) or above ( $\xi_{n+}>0$ )
  - $a_{n-}>0$  only if  $t_n=y_n-\epsilon-\xi_{n-}$  ie, if on lower boundary of  $\epsilon$ -tube ( $\xi_{n-}=0$ ) or below ( $\xi_{n-}>0$ )



#### **SVM Implementations**

- Sequential Minimal Optimization, SMO [Platt]
  - efficient implementation of SVMs
  - in Weka

- SVMlight
  - http://svmlight.joachims.org/



### Summary I: Advantages

- Systematic implementation through quadratic programming
  - ∃ *very efficient* implementations
- Excellent data-dependent generalization bounds
- Regularization built into cost function
- Statistical performance is independent of dimensions of feature space
  - Runtime typically quadratic in the # of data points
    - ... not #dimensions
  - ... less if # support vectors is small
- Theoretically related to widely studied fields of regularization theory and sparse approximation
- Fully adaptive procedures available for determining hyper-parameters



#### Summary II: **Drawbacks**

- Treatment of non-separable case somewhat heuristic
- Number of support vectors may depend strongly on
  - the kernel type and
  - the hyper-parameters (σ for RBF, k for poly)
- Systematic choice of kernels is difficult (prior information)
  - ... some ideas exist ...
- Optimization may require clever heuristics for large problems



# Summary III: Extensions

- Online algorithms
- Systematic choice of kernels using generative statistical models
- Applications to
  - Clustering
  - Non-linear principal component analysis
  - Independent component analysis
- Generalization bounds constantly improving
  - (some even practically useful!)

# 4

#### Key SVM Ideas

- Maximize the margin between + and examples
  - connects to PAC theory
- Sparse:
   Only the support vectors contribute to solution
- Handles non-separable case ... with penalties
- Kernels map examples into a new space ... usually nonlinear
  - Implicitly: dot products in this new space (in the "dual" form of the SVM program)
  - Kernels ≠ SVMs
    - ... but they combine very nicely with SVMs

#### What You Should Know

- Definition of a maximum margin classifier
- Sparse version: (Linear) SVMs
- How Maximum Margin = a QP problem
- What QP can do for you
  - Better if you know how it works!
- How to deal with noisy (non-separable) data
  - Slack variables
- How to permit "non-linear boundaries"
  - Kernel trick
- How SVM Kernel functions permit us to pretend we're working with a zillion features



#### An excellent tutorial on VC-dimension and Support Vector Machines:

C.J.C. Burges. A tutorial on support vector machines for pattern recognition. Data Mining and Knowledge Discovery, 2(2):955-974, 1998. http://citeseer.nj.nec.com/burges98tutorial.html

#### The VC/SRM/SVM Bible:

Statistical Learning Theory by Vladimir Vapnik, Wiley-Interscience; 1998.

But not easy reading ...