

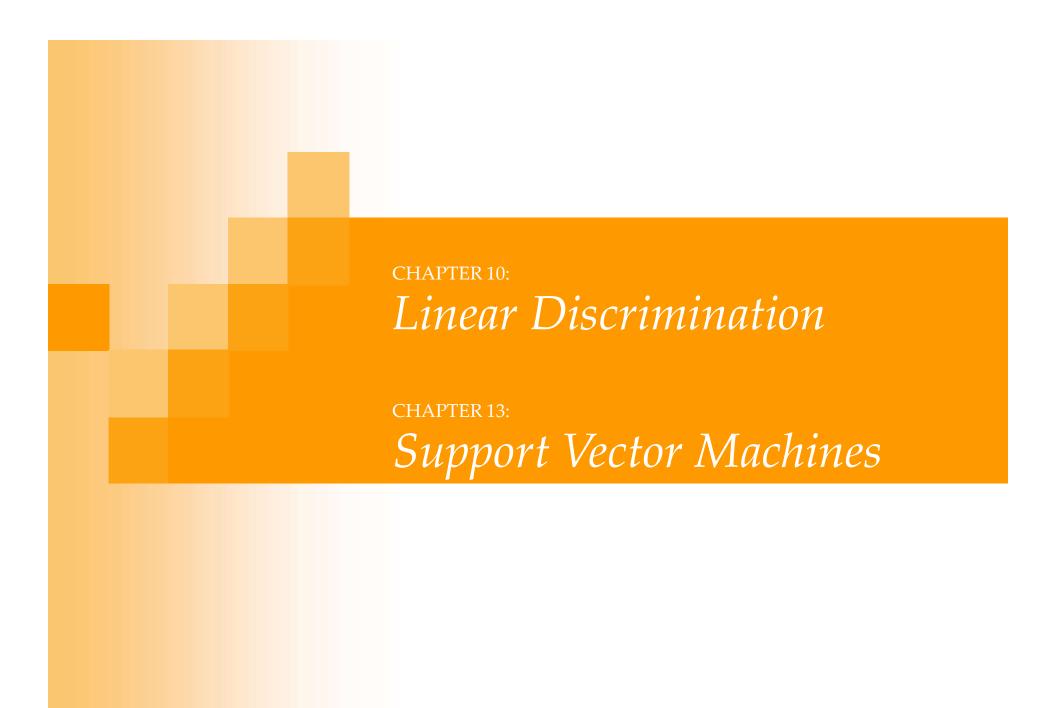
Lecture Slides for

INTRODUCTION TO

Machine Learning

ETHEM ALPAYDIN
© The MIT Press, 2004

alpaydin@boun.edu.tr http://www.cmpe.boun.edu.tr/~ethem/i2ml



Likelihood- vs. Discriminant-based Classification

- Likelihood-based: Assume a model for $p(\mathbf{x}|C_i)$, use Bayes' rule to calculate $P(C_i|\mathbf{x})$ Choose C_i if $g_i(\mathbf{x}) = \log P(C_i|\mathbf{x})$ is maximum
- Discriminant-based: Assume a model for the discriminant $g_i(\mathbf{x}|\Phi_i)$; no density estimation
 - Estimating the boundaries is enough; no need to accurately estimate the densities inside the boundaries

Linear Discriminant

Linear discriminant:

$$g_i(\mathbf{x} \mid \mathbf{w}_i, \mathbf{w}_{i0}) = \mathbf{w}_i^T \mathbf{x} + \mathbf{w}_{i0} = \sum_{j=1}^d \mathbf{w}_{ij} \mathbf{x}_j + \mathbf{w}_{i0}$$

- Advantages:
 - \square Simple: O(*d*) space/computation
 - Knowledge extraction: Weighted sum of attributes; positive/negative weights, magnitudes (credit scoring)
 - □ Optimal when $p(x|C_i)$ are Gaussian with shared cov matrix; useful when classes are (almost) linearly separable

Generalized Linear Model

Quadratic discriminant (O(d²))

$$g_i(\mathbf{x} \mid \mathbf{W}_i, \mathbf{w}_i, \mathbf{w}_{i0}) = \mathbf{x}^T \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^T \mathbf{x} + \mathbf{w}_{i0}$$

- Instead of higher complexity, we can still use a linear classifier if we use higher-order (product) terms.
- Map from x to z using nonlinear basis functions and use a linear discriminant in z-space

$$Z_1 = X_1$$
, $Z_2 = X_2$, $Z_3 = X_1^2$, $Z_4 = X_2^2$, $Z_5 = X_1X_2$

The linear function defined in the z space corresponds to a non-linear function in the x space.

$$g_i(\mathbf{x}) = \sum_{j=1}^k w_{ij} \phi_j(\mathbf{x})$$

Two Classes

choose
$$C_1$$
 if $g_1(x) > g_2(x)$
 C_2 if $g_2(x) > g_1(x)$

x_1

Define:

$$g(\mathbf{x}) = g_1(\mathbf{x}) - g_2(\mathbf{x})$$

$$= (\mathbf{w}_1^T \mathbf{x} + \mathbf{w}_{10}) - (\mathbf{w}_2^T \mathbf{x} + \mathbf{w}_{20})$$

$$= (\mathbf{w}_1 - \mathbf{w}_2)^T \mathbf{x} + (\mathbf{w}_{10} - \mathbf{w}_{20})$$

$$= \mathbf{w}^T \mathbf{x} + \mathbf{w}_0$$

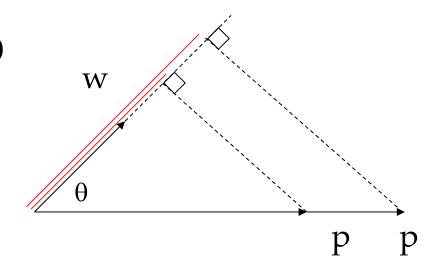
choose
$$\begin{cases} C_1 & \text{if } g(\mathbf{x}) > 0 \\ C_2 & \text{otherwise} \end{cases}$$

A Bit of Geometry



Dot Product and Projection

- $< w, p > = w^T p = ||w||||p|| \cos\theta$
- proj. of p onto w $= ||\mathbf{p}|| \mathbf{Cos}\theta$ $= \mathbf{w}^{\mathsf{T}} \cdot \mathbf{p} / ||\mathbf{w}||$



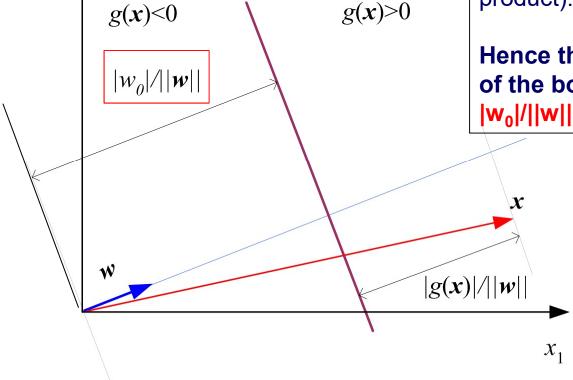
Geometry

Lectu

The points \mathbf{x} on the separating hyperplane have $g(\mathbf{x}) = \mathbf{w}^{\mathsf{T}}\mathbf{x} + \mathbf{w}_0 = 0$. Hence for the points on the boundary $\mathbf{w}^{\mathsf{T}}\mathbf{x} = -\mathbf{w}_0$.

Thus, these points also have the same projection onto the weight vector \mathbf{w} , namely $\mathbf{w}^{\mathsf{T}}\mathbf{x}/\|\mathbf{w}\|$ (by definition of projection and dot product). But this is equal to $-\mathbf{w}_0/\|\mathbf{w}\|$.

Hence the perpendicular distance of the boundary to the origin is $|\mathbf{w}_0|/||\mathbf{w}||$.



g(x)=0

Geometry

Lectu

In general, any point \mathbf{x} have $g(\mathbf{x}) = \mathbf{w}^{\mathsf{T}}\mathbf{x} + \mathbf{w}_0$.

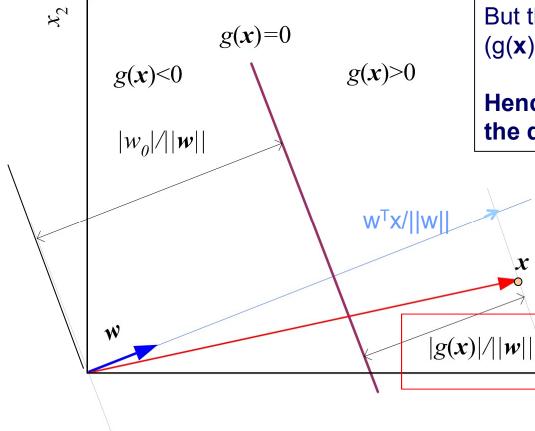
When we consider \mathbf{x} s projection onto \mathbf{w} , we have $\mathbf{w}^T\mathbf{x}/||\mathbf{w}||$ by definition of projection and dot product.

But this projection is equal to $(g(\mathbf{x})-w_0)/||\mathbf{w}||$ since $\mathbf{w}^\mathsf{T}\mathbf{x}=g(\mathbf{x})-w_0$.

 x_1

1)

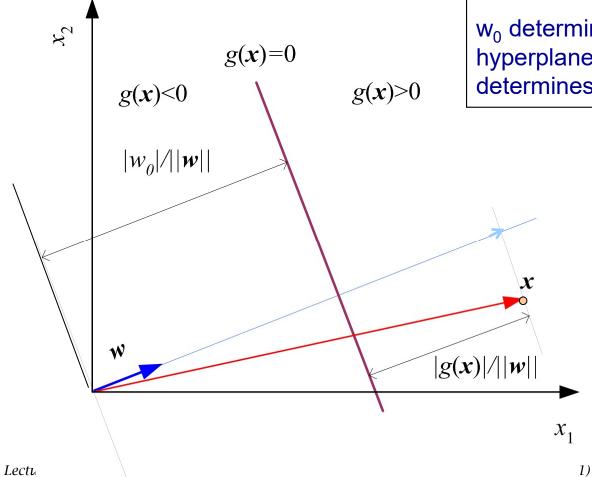
Hence the distance of any point x to the decision boundary is |g(x)|/||w||.



The distance of any point x to the decision boundary is g(x)/||w||

Positive if x is on the positive side Negative if x is on the negative side

w₀ determines the location of the hyperplane w.r.t the origin and **w** determines its orientation



How to Deal with Multi-class Problems



■ <u>When K > 2</u>

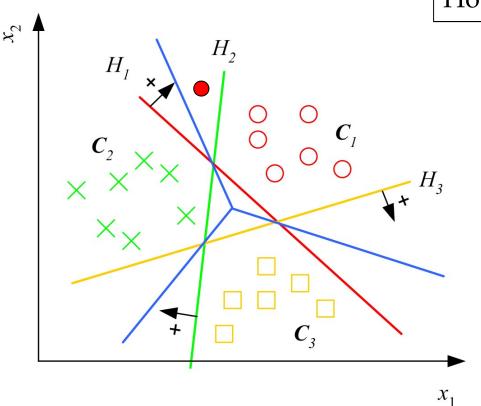
□ Combine K two-class problems, each one separating one class from all other classes

Multiple Classes

$$g_i(\mathbf{x} \mid \mathbf{w}_i, \mathbf{w}_{i0}) = \mathbf{w}_i^T \mathbf{x} + \mathbf{w}_{i0}$$

How to train?

How to decide on a test?



Choose C_i if

$$g_i(\mathbf{x}) = \max_{j=1}^K g_j(\mathbf{x})$$

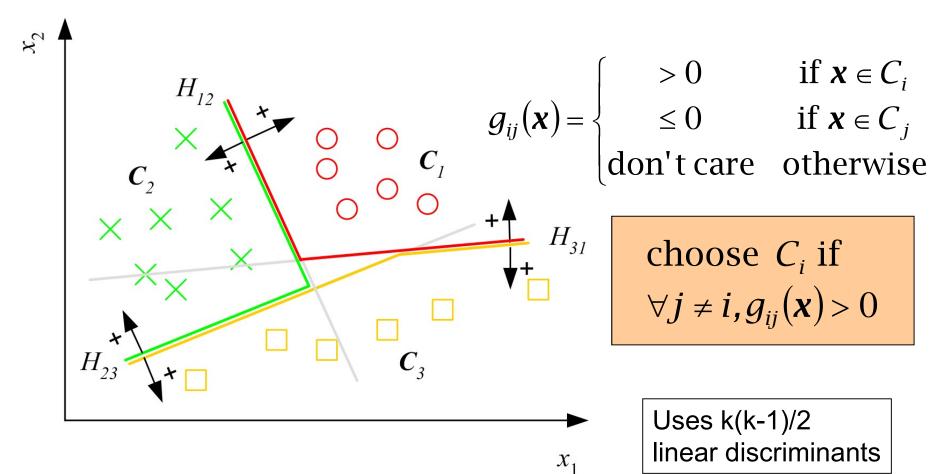
Convex decision regions based on g_i s (indicated with blue) dist is $|g_i(x)|/||wi||$

Assumes that classes are linearly separable: reject may be used

Pairwise Separation

If the classes are not linearly separable:

$$g_{ij}(\mathbf{x} \mid \mathbf{w}_{ij}, \mathbf{w}_{ij0}) = \mathbf{w}_{ij}^T \mathbf{x} + \mathbf{w}_{ij0}$$





- None of the classes may satisfy that condition
- Use instead:

choose
$$C_i$$
 if choose C_i maximizing
$$\forall j \neq i, g_{ij}(\mathbf{x}) > 0 \longrightarrow g_i(\mathbf{x}) = \sum_{j \neq i} g_{ij}(\mathbf{x})$$

Learning the Discriminants

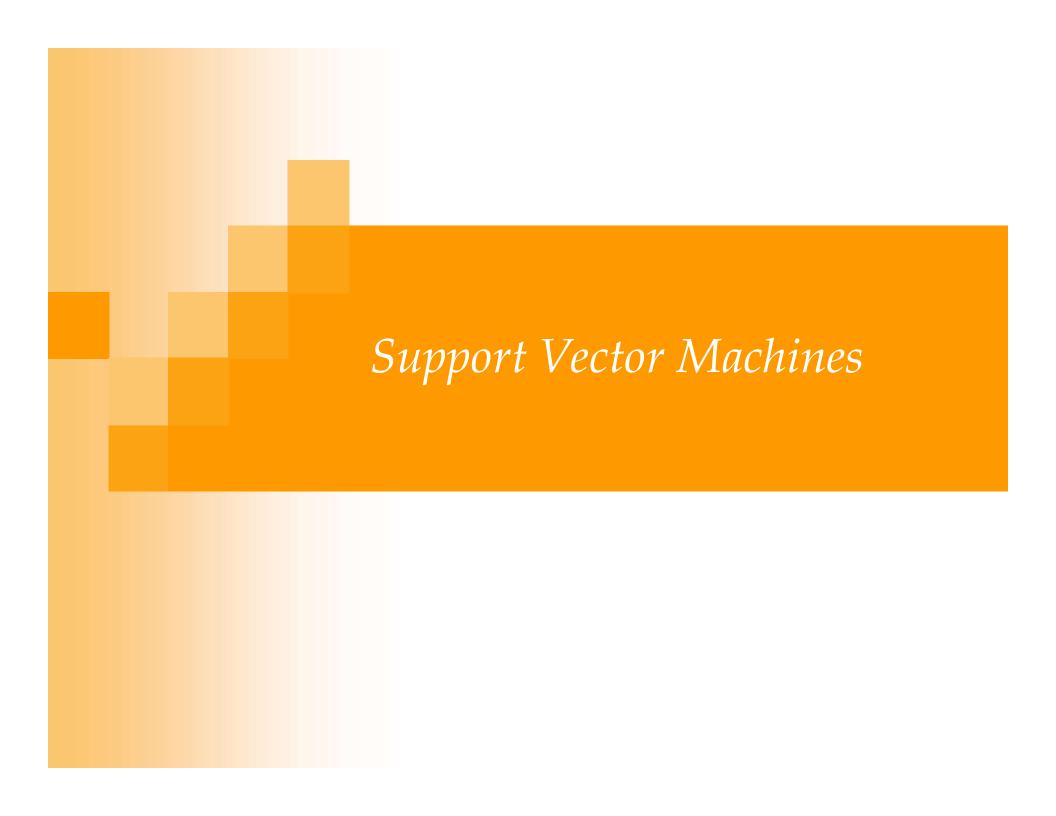
As we have seen before, when $p(x \mid C_i) \sim \mathcal{N}(\mu_i, \Sigma)$, the optimal discriminant is a linear one:

$$g_i(\mathbf{x} \mid \mathbf{w}_i, \mathbf{w}_{i0}) = \mathbf{w}_i^T \mathbf{x} + \mathbf{w}_{i0}$$

$$\mathbf{w}_i = \mathbf{\Sigma}^{-1} \mathbf{\mu}_i \quad \mathbf{w}_{i0} = -\frac{1}{2} \mathbf{\mu}_i^T \mathbf{\Sigma}^{-1} \mathbf{\mu}_i + \log P(C_i)$$

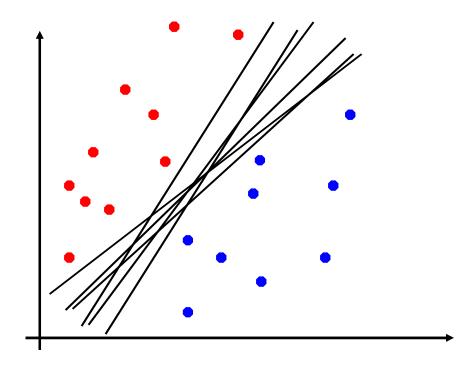
So, estimate μi and Σ from data, and plug into the gi's to find the linear discriminant functions.

Of course any way of learning can be used.



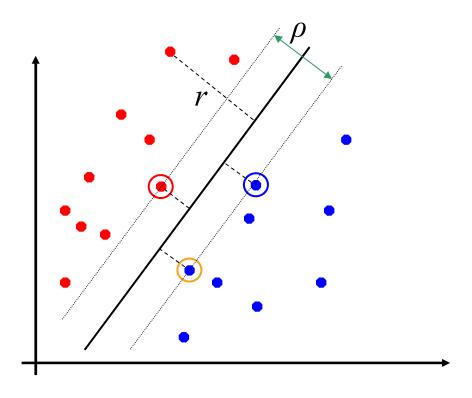
Linear Separators

Which of the linear separators is optimal?



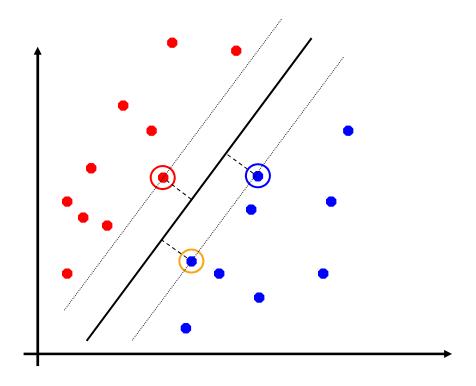
- Classification Margin

- Distance from example \mathbf{x}_i to the separator is $r = \frac{\mathbf{w} \cdot \mathbf{x}_i + \mathbf{v}}{\|\mathbf{w}\|}$
- Examples closest to the hyperplane are support vectors.
- $Margin \rho$ of the separator is the distance between support vectors from two classes.



Maximum Margin Classification

- Maximizing the margin is good according to intuition.
- Implies that only support vectors matter; other training examples are ignorable.



SVM as 2-class Linear Classifier

(Cortes and Vapnik, 1995; Vapnik, 1995)

$$\mathbf{X} = \left\{ \mathbf{x}^{t}, r^{t} \right\}_{t} \text{ where } r^{t} = \begin{cases} +1 & \text{if } \mathbf{x}^{t} \in C_{1} \\ -1 & \text{if } \mathbf{x}^{t} \in C_{2} \end{cases}$$

find w and w_0 such that

$$\mathbf{w}^T \mathbf{x}^t + w_0 \ge +1 \text{ for } r^t = +1$$

$$\mathbf{w}^T \mathbf{x}^t + w_0 \le -1 \quad \text{for } r^t = -1$$

Note the condition >= 1 (not just 0). We can always do this if the classes are linearly separable by rescaling w and w0, without affecting the separating hyperplane: $\mathbf{w}^T \mathbf{x}^t + w_0 = 0$

r^t is 1 for x^t if x belongs to class +1 r^t is -1 for x^t if x belongs to class -1

Optimal separating hyperplane:

Separating hyperplane maximizing the margin

Optimal Separating Hyperplane

Must satisfy:

$$\mathbf{w}^{T}\mathbf{x}^{t} + w_{0} \ge +1 \text{ for } r^{t} = +1$$

$$\mathbf{w}^T \mathbf{x}^t + w_0 \le -1 \text{ for } r^t = -1$$

which can be rewritten as

$$r^t \left(\mathbf{w}^T \mathbf{x}^t + w_0 \right) \ge +1$$

(Cortes and Vapnik, 1995; Vapnik, 1995)

Maximizing the Margin

Distance from the discriminant to the closest instances on either side is called the margin

In general this relationship holds (geometry): $d = \frac{|g(x)|}{\|\mathbf{w}\|}$

So, for the support vectors, we have:

$$d = \begin{cases} \frac{1}{\|\mathbf{w}\|} \\ \frac{|-1|}{\|\mathbf{w}\|} \end{cases}$$

$$\rho = 2d = \frac{2}{\|\mathbf{w}\|}$$

To maximize margin, minimize the Euclidian norm of the weight vector w



- The SVM formulation starts as an optimization function subject to some constraints (red box)
- Then the constraints are added into the optimization by adding extra variables (called Lagrange multipliers)

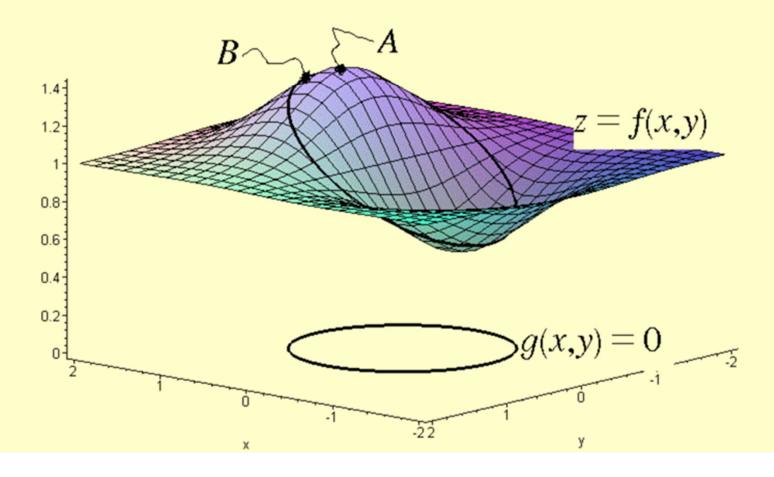
$$\min \frac{1}{2} \|\mathbf{w}\|^2 \text{ subject to } r^t (\mathbf{w}^T \mathbf{x}^t + w_0) \ge +1, \forall t$$

$$L_p = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{t=1}^{N} \alpha^t \left[r^t \left(\mathbf{w}^T \mathbf{x}^t + w_0 \right) - 1 \right]$$
$$= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{t=1}^{N} \alpha^t r^t \left(\mathbf{w}^T \mathbf{x}^t + w_0 \right) + \sum_{t=1}^{N} \alpha^t$$

Unconstrained problem using Lagrange multipliers (+)



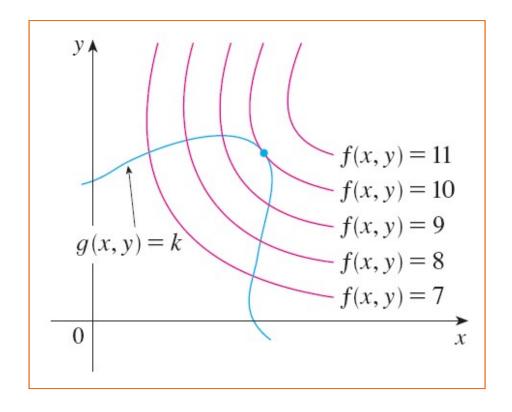
In the figure below we have illustrated an extreme value problem with constraints. The point A is the largest value of the function z=f(x,y) while the point B is the largest value of the function under the *constraint* g(x,y)=0.



W

LAGRANGE MULTIPLIERS—TWO VARIABLES

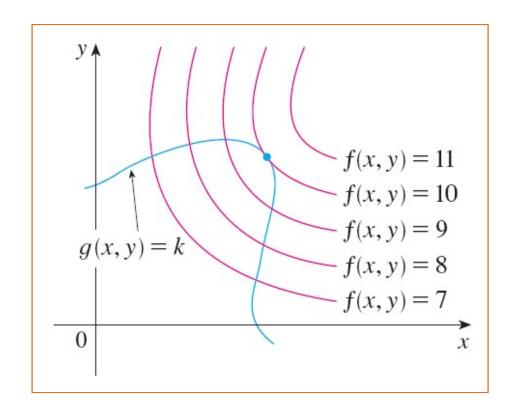
- The figure shows this curve together with several level curves of f.
- These have the equations f(x, y) = c, where c = 7, 8, 9, 10, 11



м

LAGRANGE MULTIPLIERS—TWO VARIABLES

- To maximize f(x, y) subject to g(x, y) = k is to find:
- The largest value of c such that the level curve f(x, y) = c intersects g(x, y) = k.





$$\min \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{t=1}^{N} \alpha^t r^t (\mathbf{w}^T \mathbf{x}^t + w_0) + \sum_{t=1}^{N} \alpha^t$$

At the minimum, the following minimum requirements have to be satisfied:

$$\frac{\partial L_p}{\partial \mathbf{w}} = 0 \Longrightarrow \mathbf{w} = \sum_{t=1}^N \alpha^t r^t \mathbf{x}^t$$

$$\frac{\partial L_p}{\partial w_0} = 0 \Rightarrow \sum_{t=1}^{N} \alpha^t r^t = 0$$

Going from primal to dual form

$$\begin{aligned} L_p &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{t=1}^N \alpha^t \Big[r^t \Big(\mathbf{w}^T \mathbf{x}^t + w_0 \Big) - 1 \Big] & \frac{\partial L_p}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_{t=1}^N \alpha^t r^t \mathbf{x}^t \\ &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{t=1}^N \alpha^t r^t \Big(\mathbf{w}^T \mathbf{x}^t + w_0 \Big) + \sum_{t=1}^N \alpha^t & \frac{\partial L_p}{\partial w_0} = 0 \Rightarrow \sum_{t=1}^N \alpha^t r^t = 0 \end{aligned}$$

$$L_{d} = \frac{1}{2} (\mathbf{w}^{T} \mathbf{w}) - \mathbf{w}^{T} \sum_{t} \alpha^{t} r^{t} \mathbf{x}^{t} - w_{0} \sum_{t} \alpha^{t} r^{t} + \sum_{t} \alpha^{t}$$

$$= -\frac{1}{2} (\mathbf{w}^{T} \mathbf{w}) + \sum_{t} \alpha^{t}$$

$$= -\frac{1}{2} \sum_{t} \sum_{s} \alpha^{t} \alpha^{s} r^{t} r^{s} (\mathbf{x}^{t})^{T} \mathbf{x}^{s} + \sum_{t} \alpha^{t}$$
subject to $\sum_{t} \alpha^{t} r^{t} = 0$ and $\alpha^{t} \geq 0$, $\forall t$

- •Maximize L_d with respect to α^t only
- Quadratic programming problem
- Thanks to the convexity of the problem, optimal value of L_p = L_d

Calculating the parameters w and w_0

- Once we solve for α^t , we see that
 - most of them are 0
 - \Box only a small number have $\alpha^t > 0$ and the corresponding x^t s are called the **support vectors**

$$L_{p} = \frac{1}{2} \|\mathbf{w}\|^{2} - \sum_{t=1}^{N} \alpha^{t} [r^{t} (\mathbf{w}^{T} \mathbf{x}^{t} + w_{0}) - 1]$$

Calculating the parameters w and w_0

Once we have the α^t s, we can reconstruct the weight vectors:

$$\mathbf{w} = \sum_{t=1}^{N} \alpha^t r^t \mathbf{x}^t = \sum_{t \in SV} \alpha^t r^t \mathbf{x}^t$$

where SV is the set of the Support Vectors.

After we compute **w**, from any x^t (or better averaging over all of them) we can compute w_{0:}

$$w_0 = r^t - \mathbf{w}^T \mathbf{x}^t$$



 We make decisions by comparing each new example x with only the support vectors {x_i}_{i∈SV}:

$$\hat{y} = \operatorname{sign}\left(\sum_{i \in SV} \hat{\alpha}_i y_i \left(\mathbf{x}_i^T \mathbf{x}\right) + \hat{w}_0\right)$$

Choose C1 if > 0
C2 otherwise

41

Not-Linearly Separable Case



- The non-separable case cannot find a feasible solution using the previous approach
 - □ The objective function (L_D) grows arbitrarily large.
- Relax the constraints, but only when necessary
 - Introduce a further cost for this

Soft Margin Hyperplane

Not linearly separable

$$r^{t}\left(\mathbf{w}^{T}x^{t}+w_{0}\right)\geq1-\xi^{t}$$

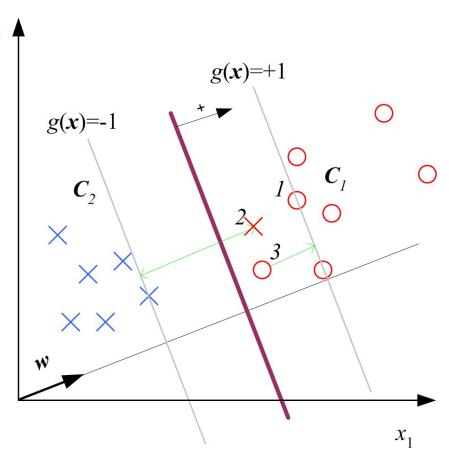
$$\xi \geq 0$$

Three cases (shown in fig):

Case 1: $\xi = 0$

Case 2: $\xi \ge 1$

Case 3: $0 \le \xi^{\epsilon} < 1$



Soft Margin Hyperplane

Define Soft error

$$\sum_t \xi^t$$

New primal to minimize is

Upper bound on the number of training errors

Lagrange multipliers to enforce positivity of ξ

$$L_{p} = \frac{1}{2} \|\mathbf{w}\|^{2} + C \sum_{t} \xi^{t} - \sum_{t} \alpha^{t} [r^{t} (\mathbf{w}^{T} \mathbf{x}^{t} + \mathbf{w}_{0}) - 1 + \xi^{t}] - \sum_{t} \mu^{t} \xi^{t}$$

- Parameter C can be viewed as a way to control overfitting: it "trades off" the relative importance of maximizing the margin and fitting the training data.
- μ^{t} are the new Lagrange parameters to guarantee positivity of ξ^{t}

Soft Margin Hyperplane

New dual is the same as the old one

$$L_d = -\frac{1}{2} \sum_{t} \sum_{s} \alpha^t \alpha^s r^t r^s (\mathbf{x}^t)^T \mathbf{x}^s + \sum_{t} \alpha^t$$

subject to

$$\sum_{t} \alpha^{t} r^{t} = 0 \text{ and } 0 \le \alpha^{t} \le C, \forall t$$

As in the separable case, instances that are not support vectors vanish with their at=0 and the remaining define the boundary.



When we use the soft margin hyperplane

- If C is too large, too high a penalty for non-separable points (too many support vectors)
- ☐ If C is too small, we may have underfitting

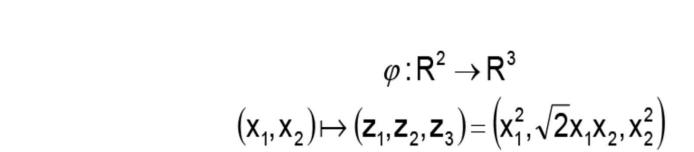
Kernel Functions in SVM

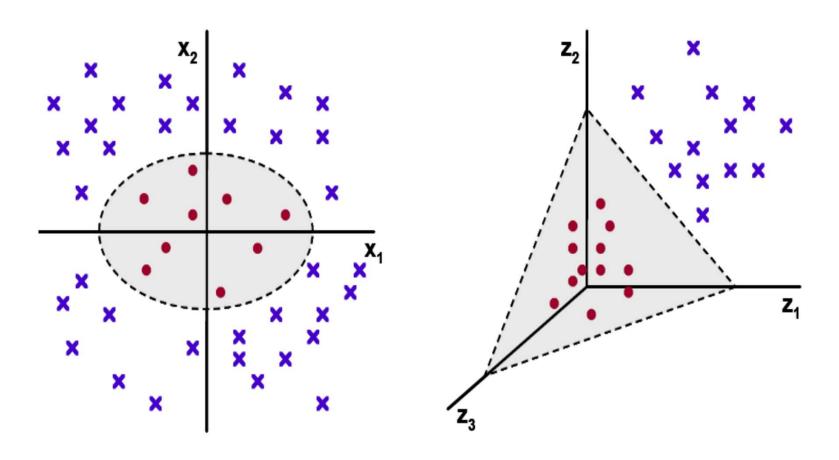
Kernel Functions

- Instead of trying to fit a non-linear model, we can
 - map the problem to a new space through a non-linear transformation and
 - use a linear model in the new space
- Say we have the new space calculated by the basis functions $\mathbf{z} = \boldsymbol{\varphi}(\mathbf{x})$ where $z_i = \phi_i(\mathbf{x})$, j=1,...,k

d-dimensional **x** space \longrightarrow k-dimensional **z** space

$$\boldsymbol{\varphi}(\mathbf{x}) = [\varphi_1(\mathbf{x}) \ \varphi_2(\mathbf{x}) \ \dots \ \varphi_k(\mathbf{x})]$$





Н

Kernel Functions

$$g(\mathbf{x}) = \sum_{k=1} w_k \varphi_k(\mathbf{x}) + b$$

$$g(\mathbf{x}) = \sum_{k=0} w_k \varphi_k(\mathbf{x})$$

if we assume
$$\varphi_0(\mathbf{x}) = 1$$
 for $\forall \mathbf{x}$

Kernel Machines

Preprocess input x by basis functions

$$z = \varphi(x)$$
 $g(z) = w^T z$
 $g(x) = w^T \varphi(x)$

SVM solution: Find Kernel functions K(x,y) such that the inner product of basis functions are replaced by a Kernel function in the original input space

$$\mathbf{w} = \sum_{t} \alpha^{t} \mathbf{r}^{t} \mathbf{z}^{t} = \sum_{t} \alpha^{t} \mathbf{r}^{t} \boldsymbol{\varphi}(\mathbf{x}^{t})$$

$$g(\mathbf{x}) = \mathbf{w}^{T} \boldsymbol{\varphi}(\mathbf{x}) = \sum_{t} \alpha^{t} \mathbf{r}^{t} \boldsymbol{\varphi}(\mathbf{x}^{t})^{T} \boldsymbol{\varphi}(\mathbf{x})$$

$$g(\mathbf{x}) = \sum_{t} \alpha^{t} \mathbf{r}^{t} K(\mathbf{x}^{t}, \mathbf{x})$$

Kernel Functions

Consider polynomials of degree q:

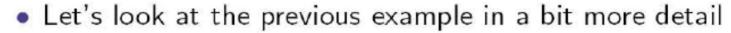
$$K(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^{T} \mathbf{y} + 1)^{2}$$

$$= (x_{1}y_{1} + x_{2}y_{2} + 1)^{2}$$

$$= 1 + 2x_{1}y_{1} + 2x_{2}y_{2} + 2x_{1}x_{2}y_{1}y_{2} + x_{1}^{2}y_{1}^{2} + x_{2}^{2}y_{2}^{2}$$

$$\phi(\mathbf{x}) = \begin{bmatrix} 1, \sqrt{2}x_{1}, \sqrt{2}x_{2}, \sqrt{2}x_{1}x_{2}, x_{1}^{2}, x_{2}^{2} \end{bmatrix}^{T}$$

(Cherkassky and Mulier, 1998)



$$\mathbf{x} \to \phi(\mathbf{x}) = [x_1^2 \ x_2^2 \ \sqrt{2}x_1x_2 \ \sqrt{2}x_1 \ \sqrt{2}x_2 \ 1]$$

 The SVM classifier deals only with inner products of examples (or feature vectors). In this example,

$$\phi(\mathbf{x})^T \phi(\mathbf{x}') = x_1^2 x_1'^2 + x_2^2 x_2'^2 + 2x_1 x_2 x_1' x_2' + 2x_1 x_1' + 2x_2 x_2' + 1$$

$$= (1 + x_1 x_1' + x_2 x_2')^2$$

$$= (1 + (\mathbf{x}^T \mathbf{x}'))^2$$

so the inner products can be evaluated without ever explicitly constructing the feature vectors $\phi(\mathbf{x})$!

• $K(\mathbf{x}, \mathbf{x}') = (1 + (\mathbf{x}^T \mathbf{x}'))^2$ is a kernel function (inner product in the feature space)

Examples of Kernel Functions

Linear: $K(x_i,x_i)=x_i^Tx_i$

Mapping Φ : $x \to \phi(x)$, where $\phi(x)$ is x itself

Polynomial of power p: $K(x_i,x_j) = (1 + x_i^T x_j)^p$ Mapping Φ : $x \to \phi(x)$, where $\phi(x)$ has $\binom{d+p}{p}$ dimensions

Gaussian (radial-basis function):
$$K(x_i, x_j) = e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$$

Mapping $\Phi: x \to \varphi(x)$, where $\varphi(x)$ is infinite-dimensional: every example of the property of t

Mapping Φ : $x \to \phi(x)$, where $\phi(x)$ is infinite-dimensional: every point is mapped to a function (a Gaussian); combination of functions for support vectors is the separator.

Higher-dimensional space still has intrinsic dimensionality d, but linear separators in it correspond to non-linear separators in original space.



$$x=(x_1,x_2);$$

$$z=(z_1,z_2);$$

$$\langle x, z \rangle^2 = (x_1 z_1 + x_2 z_2)^2 =$$

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

$$= \langle (x_1^2, x_2^2, \sqrt{2} x_1 x_2), (z_1^2, z_2^2, \sqrt{2} z_1 z_2) \rangle =$$

$$= \langle \phi(x), \phi(z) \rangle$$
www.support-vector.net



Other Kernel Functions

Polynomials of degree q:

$$K(\mathbf{x}^t, \mathbf{x}) = (\mathbf{x}^T \mathbf{x}^t)^q$$

$$K(\mathbf{x}^t, \mathbf{x}) = (\mathbf{x}^T \mathbf{x}^t + 1)^q$$

Radial-basis functions:

$$K(\mathbf{x}^t, \mathbf{x}) = \exp \left[-\frac{\|\mathbf{x}^t - \mathbf{x}\|^2}{\sigma^2} \right]$$

Sigmoidal functions such as:

$$K(\mathbf{x}^t, \mathbf{x}) = \tanh(2\mathbf{x}^T\mathbf{x}^t + 1)$$

Sample solutions

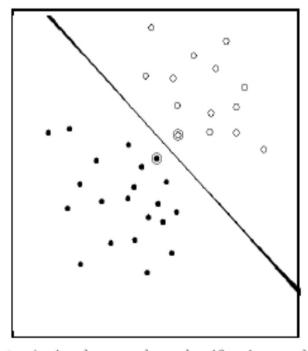


Figure 3: A simple two-class classification problem as solved by the Support Vector algorithm ($c_i = 1$ for all i; cf. Eq. 1). Note that the RBF centers (indicated by extra circles) are closest to the decision boundary.

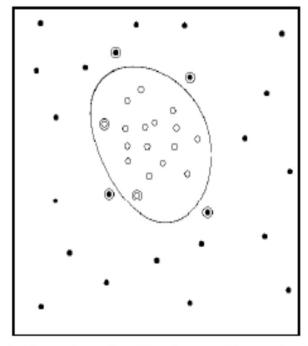


Figure 4: Two-class classification problem solved by the Support Vector algorithm ($c_i = 1$ for all i; cf. Eq. 1).



- Informally, kernel methods implicitly define the class of possible patterns by introducing a notion of similarity between data
 - Choice of similarity -> Choice of relevant features
- More formally, kernel methods exploit information about the inner products between data items
 - Many standard algorithms can be rewritten so that they only require inner products between data (inputs)
 - Kernel functions = inner products in some feature space (potentially very complex)
 - □ If kernel given, no need to specify what features of the data are being used
 - □ Kernel functions make it possible to use infinite dimensions
 - efficiently in time / space

SVM summary

- SVMs were originally proposed by Boser, Guyon and Vapnik in 1992 and gained increasing popularity in late 1990s.
- SVMs represent a general methodology for many PR problems: classification,regression, feature extraction, clustering, novelty detection, etc.
- SVMs can be applied to complex data types beyond feature vectors (e.g. graphs, sequences, relational data) by designing kernel functions for such data.
- SVM techniques have been extended to a number of tasks such as regression [Vapnik *et al.* '97], principal component analysis [Schölkopf *et al.* '99], etc.
- Most popular optimization algorithms for SVMs use *decomposition* to hill-climb over a subset of α_i 's at a time, e.g. SMO [Platt '99] and [Joachims '99]

Advantages of SVMs

- □ There are no problems with local minima, because the solution is a Qaudratic Programming problem
- ☐ The optimal solution can be found in polynomial time
- There are few model parameters to select: the penalty term C, the kernel function and parameters (e.g., spread σ in the case of RBF kernels)
- □ The final results are stable and repeatable (e.g., no random initial weights)
- The SVM solution is sparse; it only involves the support vectors
- □ SVMs rely on elegant and principled learning methods
- SVMs provide a method to control complexity independently of dimensionality
- SVMs have been shown (theoretically and empirically) to have excellent generalization capabilities

Challenges

- Can the kernel functions be selected in a principled manner?
- SVMs still require selection of a few parameters, typically through cross-validation
- How does one incorporate domain knowledge?
 - Currently this is performed through the selection of the kernel and the introduction of "artificial" examples
- How interpretable are the results provided by an SVM?
- What is the optimal data representation for SVM? What is the effect of feature weighting? How does an SVM handle categorical or missing features?
- Do SVMs always perform best? Can they beat a hand-crafted solution for a particular problem?
- Do SVMs eliminate the model selection problem?

- More explanations or demonstrations can be found at:
 - http://www.support-vector-machines.org/index.html
 - □ Haykin Chp. 6 pp. 318-339
 - Burges tutorial (under/reading/)
 - Burges, CJC "<u>A Tutorial on Support Vector Machines for Pattern Recognition</u>" Data Mining and Knowledge Discovery, Vol 2 No 2, 1998.
 - □ http://www.dtreg.com/svm.htm

Software

- □ **SVM***light*, by Joachims, is one of the most widely used SVM classification and regression package. Distributed as C++ source and binaries for Linux, Windows, Cygwin, and Solaris. Kernels: polynomial, radial basis function, and neural (tanh).
- LibSVM http://www.csie.ntu.edu.tw/~cjlin/libsvm/ LIBSVM (Library for Support Vector Machines), is developed by Chang and Lin; also widely used. Developed in C++ and Java, it supports also multi-class classification, weighted SVM for unbalanced data, cross-validation and automatic model selection. It has interfaces for Python, R, Splus, MATLAB, Perl, Ruby, and LabVIEW. Kernels: linear, polynomial, radial basis function, and neural (tanh).
- Applet to play with:
 - http://lcn.epfl.ch/tutorial/english/svm/html/index.html
- Some non-SVM good applets:

SVM Applications

- Cortes and Vapnik 1995:
 - □ Handwritten digit classification
 - □ 16x16 bitmaps -> 256 dimensions
 - □ Polynomial kernel where q=3 -> feature space with 10⁶ dimensions
 - □ No overfitting on a training set of 7300 instances
 - □ Average of 148 support vectors over different training sets

Expected test error rate:

 $Exp_N[P(error)] = Exp_N[#support vectors] / N$ (= 0.02 for the above example)



- The main things to know are:
 - □ The aim of SVMs (maximizing the margins)
 - How it is specified as an optimization problem
 - $\hfill \square$ Once solved, the optimization finds the αs , that in turn specifies the weights
 - α s also indicate which input data are the support vectors (those for which the corresponding α is non-zero)
 - The non-linearly separable case can be handled in a similar manner, with a slight extension that involves a parameter C
 - It is important to understand <u>what C does</u> in order to properly use SVM software
 - The <u>use of kernels</u> makes it possible to go to higher dimensional space, without actually computing those features explicitly; nor being badly affected by the curse of dimensionality