

Support Vector Machines

srihari@buffalo.edu

SVM Discussion Overview

1. Overview of SVMs
2. Margin Geometry
3. SVM Optimization
4. Overlapping Distributions
5. Relationship to Logistic Regression
6. Dealing with Multiple Classes
7. SVM for Regression
8. SVM and Computational Learning Theory
9. Relevance Vector Machines

Kernel Methods

- In the kernel method for regression

- Given data set $\{\mathbf{x}_i, t_i\}$, $i = 1..N$

- Gram matrix \mathbf{K} has elements

$$K_{nm} = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m)$$

- Output y predicted for \mathbf{x} is

$$y(\mathbf{x}) = k(\mathbf{x})^T (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}, \text{ where } k(\mathbf{x}) \text{ has elements } k_n(\mathbf{x}) = (\mathbf{x}_n, \mathbf{x})$$

- Many choices of kernel

- Gaussian: $k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2 / 2\sigma^2)$

- Cosine similarity (BOW, TF-IDF):

$$k(\mathbf{x}_i, \mathbf{x}_{i'}) = \frac{\mathbf{x}_i^T \mathbf{x}_{i'}}{\|\mathbf{x}_i\|_2 \|\mathbf{x}_{i'}\|_2}$$

$$\mathbf{x}_i = [x_{i1}, \dots, x_{iD}]$$

- Mercer kernel:
substrings shared

$$k(x, x') = \sum_{s \in A^*} w_s \phi_s(x) \phi_s(x')$$

$\phi_s(x)$ is no of times
substring s appears in x

Primal optimization

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{ \mathbf{w}^T \phi(\mathbf{x}_n) - t_n \}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

$$\mathbf{w} = \Phi^T \mathbf{a} \quad a_n = -\frac{1}{\lambda} \{ \mathbf{w}^T \phi(\mathbf{x}_n) - t_n \}$$

\mathbf{w} appearing in \mathbf{a} eliminated using kernel

$$J(\mathbf{w}) = \frac{1}{2} \mathbf{a}^T \Phi \Phi^T \Phi \Phi^T \mathbf{a} - \mathbf{a}^T \Phi \Phi^T \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \Phi \Phi^T \mathbf{a}$$

Dual optimization

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \mathbf{K} \mathbf{K} \mathbf{a} - \mathbf{a}^T \mathbf{K} \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K} \mathbf{a}$$

$$\mathbf{a} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}$$

Sparse Kernel Methods

- Kernel methods have a limitation
 - that $k(\mathbf{x}_n, \mathbf{x}_m)$ must be evaluated for all possible pairs $(\mathbf{x}_n, \mathbf{x}_m)$ of training points
 - Which is computationally infeasible
- In sparse solutions, predictions for new inputs depend only on kernels evaluated at a subset of training data points

Alternative Names for SVM

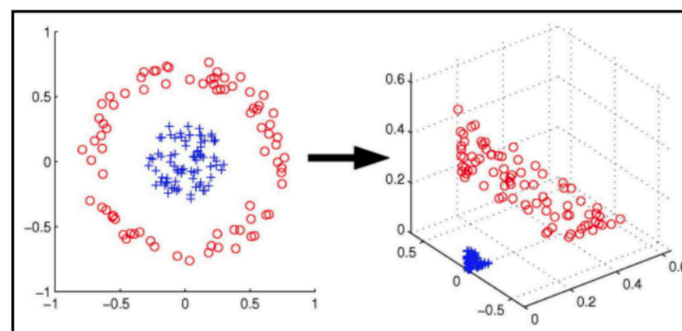
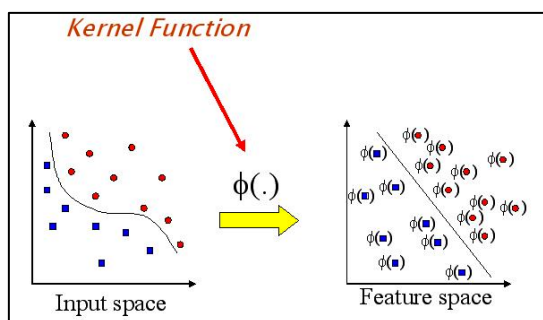
- Also called Sparse kernel machines
 - Prediction is based on a linear combination of a kernel evaluated at the training points
 - Sparse because not all pairs of training points need be used
- Also called Maximum margin classifiers

Importance of SVMs

- SVM is a discriminative method
- SVM brings together:
 1. Computational learning theory
 2. Previous methods in linear discriminants
 3. Optimization theory
- Widely used for ML problems:
 - Classification, regression, novelty detection

Ideas Used in SVM

- Linearly separable case considered
 - By using higher-dimensions than original feature space

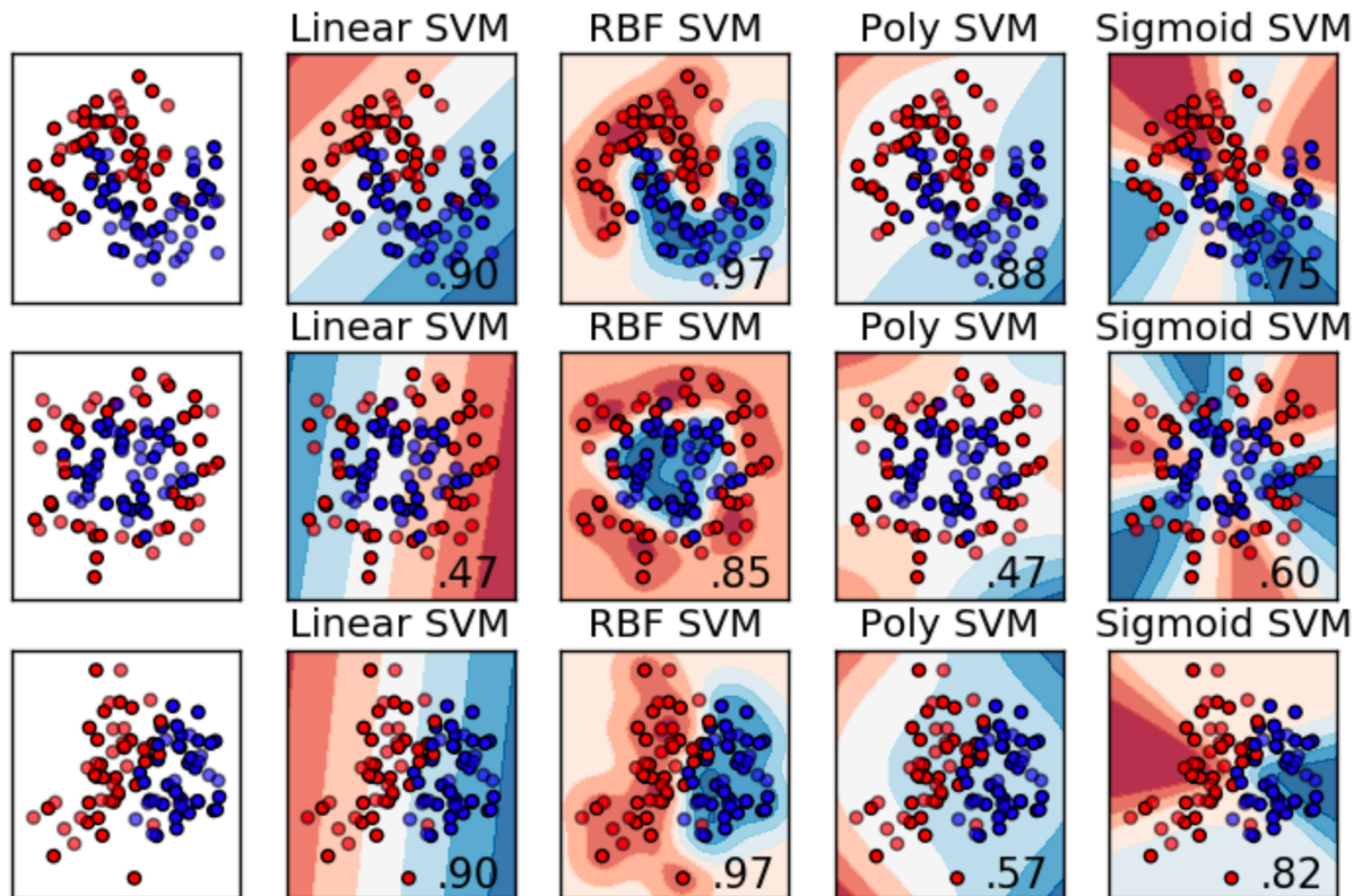


- Kernel trick reduces computational overhead

$$k(\mathbf{y}_j, \mathbf{y}_k) = \mathbf{y}_j^t \cdot \mathbf{y}_k = \phi(\mathbf{x}_j)^t \cdot \phi(\mathbf{x}_k)$$

- Determination of model parameters is a convex optimization problem
 - Extensive use of Lagrange multipliers

Boundaries induced by kernels



Source: <https://gist.github.com/WittmannF/60680723ed8dd0cb993051a7448f7805>

2. Maximum Margin Classifiers

- Begin with 2-class linear classifier

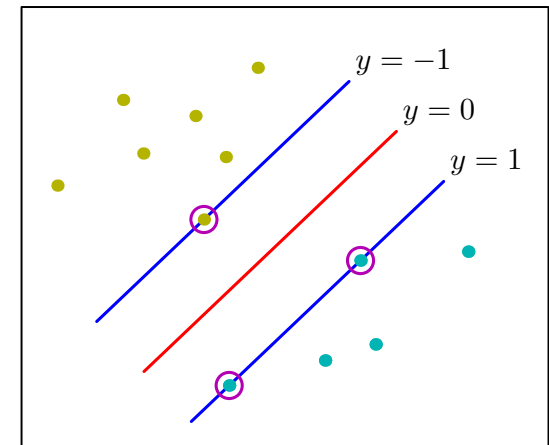
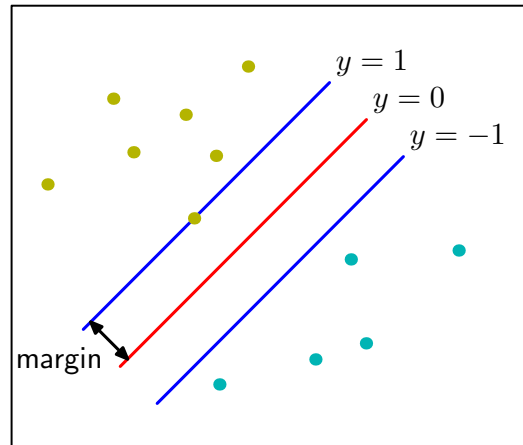
$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$$

- where $\phi(\mathbf{x})$ is a feature space transformation
- We will introduce a dual representation which avoids working with feature space
 - Training set: $\mathbf{x}_1, \dots, \mathbf{x}_N$ with targets t_1, \dots, t_N , $t_n \in \{-1, 1\}$
 - New \mathbf{x} classified according to sign of $y(\mathbf{x})$
- Assume training data is linearly separable in feature space,
 - i.e., there exists at least one choice of \mathbf{w} and b such that $t_n y(\mathbf{x}_n) > 0$ for all points

Definition of margin

- Perceptron guarantees a solution in finite no. of steps, which depends on
 1. Initial values for w, b
 2. Order in which data are presented
- If there are multiple solutions we need one with smallest generalization error
 - SVM approaches this with margin concept
- Defined as smallest distance between decision boundary and any of the samples

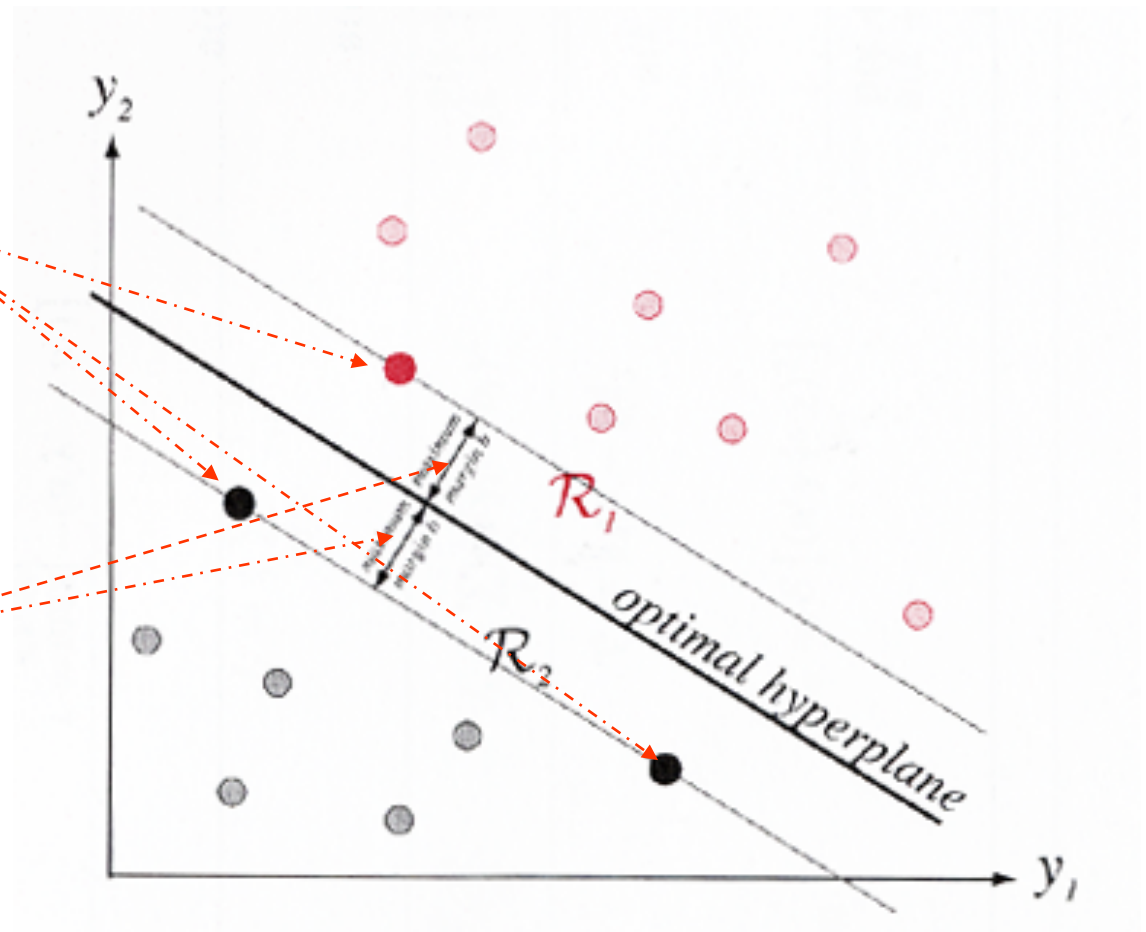
Support Vector Definition



- Margin:
 - Perpendicular distance between boundary and the closest data point (left)
- Maximizing margin leads to a particular choice of decision boundary
 - Determined by a subset of the data points
 - Which are known as *support vectors* (circles)

Support Vectors and Margin

- Support vectors are those nearest patterns at distance b from hyperplane
- SVM finds hyperplane with maximum distance
(margin distance b)
from nearest training patterns



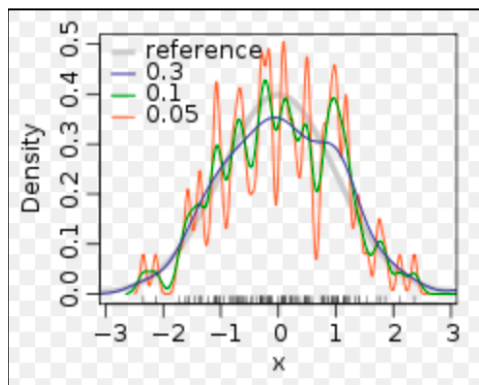
Three support vectors are shown as solid dots

Motivation for maximum margin

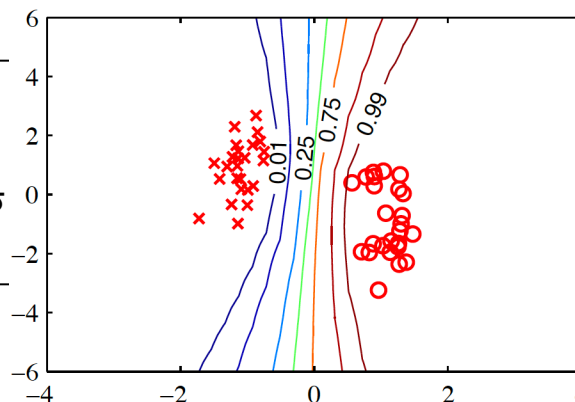
- Maximum margin solution is motivated by computational learning theory
- Simple insight
 - Based on a hybrid of generative and discriminative approaches
 - Seen next (Tong, Koller 2000)

Max Margin Insight

- Model class-conditionals $p(\mathbf{x} | C_i)$
 - Using Parzen estimator with Gaussian kernels
 - With common parameter σ^2



100 normally distributed nos with decreasing σ
 $p(\mathbf{x}|C_i)$



Different priors $p(C_i)$ for linearly separable data leads to a decision boundary that lies in the middle

- $P(C_1 | \mathbf{x}) = \frac{P(\mathbf{x} | C_1)P(C_1)}{P(\mathbf{x} | C_1)P(C_1) + P(\mathbf{x} | C_2)P(C_2)}$ defines the Bayes classifier
 - Determines best hyperplane for learned density model
- As $\sigma^2 \rightarrow 0$ hyperplane has maximum margin
 - In the limit hyperplane becomes independent of points that are not support vectors

Distance of a point \mathbf{x} from plane

Theorem: Given a plane $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = 0$, perpendicular distance from \mathbf{x} to the plane is

$$r = \frac{y(\mathbf{x})}{\|\mathbf{w}\|}$$

Proof:

$$\mathbf{x} = \mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

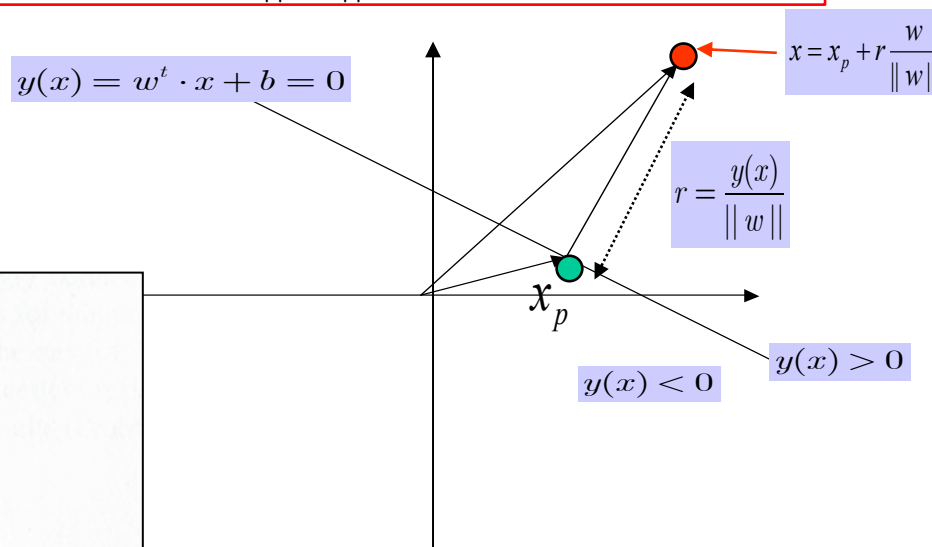
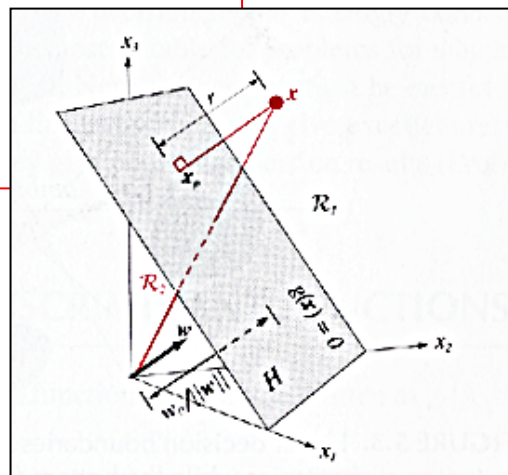
$$= 0$$

$$y(\mathbf{x}) = \mathbf{w}^T \left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|} \right) + b = \mathbf{w}^T \mathbf{x}_p + w_0 + r \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|} = y(\mathbf{x}_p) + r \frac{\|\mathbf{w}\|^2}{\|\mathbf{w}\|} = r \|\mathbf{w}\| \quad \text{QED}$$

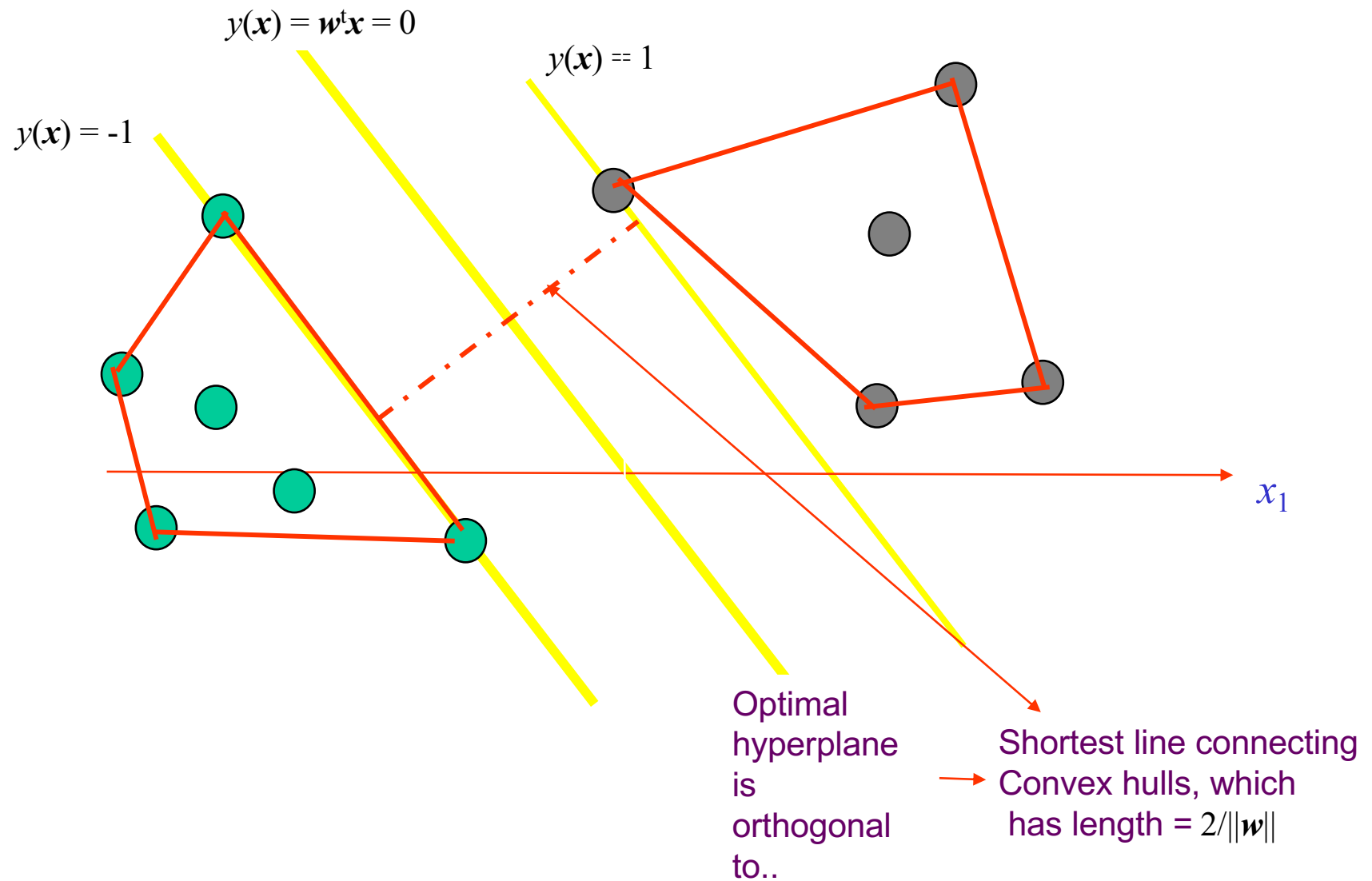
Corollary: Distance of origin to plane is

$$r = y(\mathbf{0}) / \|\mathbf{w}\| = b / \|\mathbf{w}\|$$

Since $y(\mathbf{0}) = \mathbf{w}^T \mathbf{0} + b = b$
then $b=0$ implies that
plane passes through origin



SVM Margin geometry



Maximum Margin Solution

- We interested in solutions for which

$$t_n y(\mathbf{x}_n) > 0 \text{ for all } n$$

- Distance of point \mathbf{x}_n to decision surface is

$$\frac{t_n y(\mathbf{x}_n)}{\|\mathbf{w}\|} = \frac{t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|}$$

- Margin is closest point \mathbf{x}_n of data set
 - We wish to optimize \mathbf{w}, b to maximize distance
- Maximum margin solution is found from
$$\arg \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n \left[t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \right] \right\}$$
- Direct solution is very complex
 - So convert into equivalent problem

Equivalent problem

- If we rescale $\mathbf{w} \rightarrow \kappa \mathbf{w}$ and $b \rightarrow \kappa b$ then the distance from any point \mathbf{x}_n to the decision surface given by $t_n y(\mathbf{x}_n) / \|\mathbf{w}\|$ is unchanged
 - We can use this freedom to set, for the point closest to the surface, $t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) = 1$
 - In this case all points will satisfy
$$t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1 \quad n=1, \dots, N$$
- This is the *canonical representation* of the hyperplane

Canonical Representation

$$t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1 \quad n=1, \dots, N$$

- When equality holds, constraints are *active*
 - For the remainder they are *inactive*
 - There will always be at least one active constraint, because there will be a closest point
 - Once margin is maximized there will be two active constraints

Optimization Problem (OP) Restated

- Augmented space:

$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ by choosing $w_0 = b$ and $x_0 = 1$, i.e, plane passes through origin

- Canonical representation of decision hyperplane is then

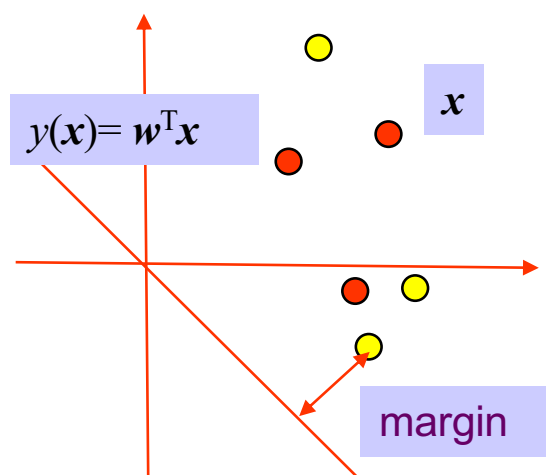
$$t_n y(\mathbf{x}_n) \geq 0, \quad n = 1, \dots, N$$

- Since distance of a point y to hyperplane $y(\mathbf{x})=0$ is

Maximizing the margin is equivalent to maximizing $\|\mathbf{w}\|^{-1}$

Which is equivalent to minimizing $\|\mathbf{w}\|^2$

$$\frac{y(\mathbf{x})}{\|\mathbf{w}\|}$$



So the optimization problem

$$\arg \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n [t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b)] \right\}$$

is equivalent to

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

The factor $\frac{1}{2}$ is for later convenience

subject to the constraints

$$t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1 \quad n=1, \dots, N$$

This is an example of quadratic programming

We are trying to maximize a quadratic function
subject to a set of linear inequality constraints

SVM Constrained Optimization

Optimize

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

subject to constraints

$$t_n(\mathbf{w}^t \phi(\mathbf{x}_n) + b) \geq 1, \quad n = 1, \dots, N$$

- Can be cast as an unconstrained problem
- by introducing Lagrange undetermined multipliers with one multiplier $a_n \geq 0$ for each constraint
- The Lagrange function we wish to minimize is

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n [t_n(\mathbf{w}^t \phi(\mathbf{x}_n) + b) - 1]$$

where $\mathbf{a} = (a_1, \dots, a_N)^T$

Note the minus sign in front of the Lagrange multiplier term

Because we are minimizing wrt \mathbf{w} and b , and maximizing wrt \mathbf{a}

Optimization of Lagrange function

- We wish to minimize

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n \left[t_n \mathbf{w}^t (\phi(\mathbf{x}_n) + b) - 1 \right]$$

- We seek to minimize $L ()$
 - with respect to \mathbf{w} and b
 - maximize it w.r.t . the undetermined multipliers $a_n \geq 0$
- Last term represents the goal of classifying the points correctly
- Karush-Kuhn-Tucker construction shows that this can be recast as a maximization problem which is computationally better

Eliminating \mathbf{w} and b

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n [t_n \mathbf{w}^t (\phi(\mathbf{x}_n) + b) - 1]$$

- Setting derivatives of $L(\mathbf{w}, b, \mathbf{a})$ wrt \mathbf{w} and b equal to zero
 - We obtain the following conditions

$$\mathbf{w} = \sum_n a_n t_n \mathbf{x}_n$$

$$0 = \sum_n a_n t_n$$

- Eliminating \mathbf{w} and b from $L(\mathbf{w}, b, \mathbf{a})$ using these conditions then gives the dual representation of the maximum margin problem

Dual O.P.

- Problem is one of maximizing

$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

- Subject to the constraints

$$\sum_{n=1}^N a_n t_n = 0 \quad a_n \geq 0, \quad n = 1, \dots, N$$

given the training data

- Here the kernel function is defined by

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^t \cdot \phi(\mathbf{x}')$$

- Again a quadratic optimization problem
 - Optimize a quadratic function subject to a set of inequality constraints

Complexity of Primal and Dual

- Quadratic programming complexity is $O(M^3)$
 - Primal OP of minimizing
 - $\arg \min_{w,b} \frac{1}{2} \|w\|^2$ over M basis functions
 - Dual OP of maximizing
 - $\tilde{L}(a) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_m a_n t_n t_m k(x_n, x_m)$ over N data points
- For $M \ll N$ it appears disadvantageous
 - But it allows model to be reformulated using kernels
 - So it can be applied to feature spaces whose dimensionality M exceeds no of data points N
 - Kernel formulation also makes clear that $k(x, x')$ be positive definite, which bounds $L(a)$ is lower-bounded

Classifying new points

- To classify points using trained model:
 - Evaluate sign of $y(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) + b$
 - This can be expressed in terms of the parameters $\{a_n\}$ and the kernel function by substituting for \mathbf{w} using $\mathbf{w} = \sum_n a_n t_n \mathbf{x}_n$ to give

$$y(\mathbf{x}) = \sum_{n=1}^N a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b$$

Definition of Support Vectors

- This OP satisfies KKT conditions

- which requires the following hold:

$$\begin{aligned} a_n &\geq 0 \\ t_n y(\mathbf{x}_n) - 1 &\geq 0 \\ a_n \{t_n y(\mathbf{x}_n) - 1\} &= 0 \end{aligned}$$

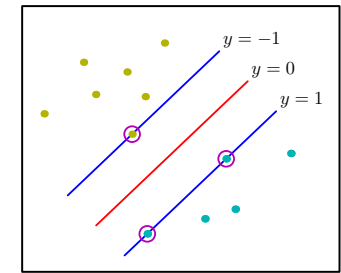
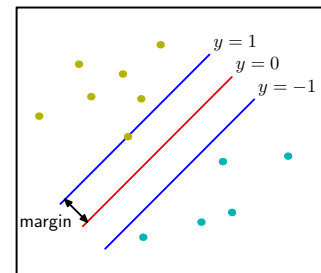
- Thus for every data point either $a_n=0$ or $t_n y(\mathbf{x}_n)=1$

- A point with $a_n=0$ will not appear in sum

$$y(\mathbf{x}) = \sum_{n=1}^N a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b$$

- Remaining data points are called support vectors

- Because they satisfy $t_n y(\mathbf{x}_n)=1$ they correspond to points that lie on the maximum margin hyperplanes in feature space



This property is central to practical applicability of SVMs. Once it is trained a significant proportion of Data points can be discarded and only the support vectors retained

Solution for b

- Having solved the quadratic programming problem and found a we can determine value of threshold parameter b
- Noting that for any support vector $t_n y(\mathbf{x}_n) = 1$
- Using $y(\mathbf{x}) = \sum_{n=1}^N a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b$ gives $t_n \left(\sum_{m \in S} a_m t_m k(\mathbf{x}_n, \mathbf{x}_m) + b \right) = 1$
 - Where S denotes the indices of support vectors
- Solving for b gives $b = \frac{1}{N_S} \left(t_n - \sum_{m \in S} a_m t_m k(\mathbf{x}_n, \mathbf{x}_m) \right)$
 - Where N_S is the total no of support vectors

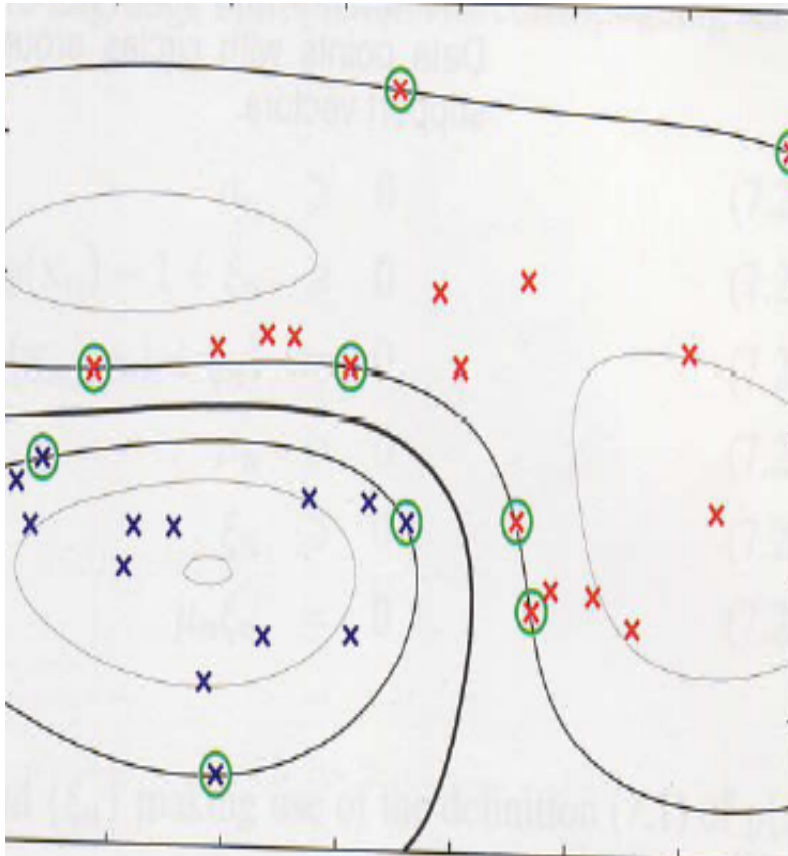
Equivalent Error Function of SVM

- For comparison with alternative models
 - Express the maximum margin classifier in terms of minimization of an error function:

$$\sum_{n=1}^N E_{\infty}(y(\mathbf{x}_n)t_n - 1) + \lambda \|\mathbf{w}\|^2$$

- Where $E_{\infty}(z)$ is zero if $z \geq 0$ and ∞ otherwise
 - Ensures that $t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1 \quad n=1, \dots, N$ are satisfied
- Gives infinite error if point is misclassified
- Zero error if it was classified correctly
- We optimize parameters to maximize margin

Example of SVM Classifier



- Two classes in two dimensions
- Synthetic Data
- Shows contours of constant $y(\mathbf{x})$
- Obtained from SVM with Gaussian kernel function
- Decision boundary is shown
- Margin boundaries are also shown
- Support vectors are shown
- Shows sparsity of SVM

Summary of SVM Optimization Problems

Dual Optimization Problem

$$\begin{aligned} \text{Primal OP: minimize } P(\vec{w}, b, \xi) &= \frac{1}{2} \vec{w} \cdot \vec{w} + C \sum_{i=1}^n \xi_i \\ \text{s. t. } y_i [\vec{w} \cdot \vec{x}_i + b] &\geq 1 - \xi_i \text{ and } \xi_i \geq 0 \end{aligned}$$

Lemma: The solution w° can always be written as a linear combination

$$\vec{w}^\circ = \sum_{i=1}^n \alpha_i y_i \vec{x}_i \quad \alpha_i \geq 0$$

of the training data.

$$\begin{aligned} \text{Dual OP: maximize } D(\vec{\alpha}) &= \left(\sum_{i=1}^n \alpha_i \right) - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) \\ \text{s.t. } \sum_{i=1}^n \alpha_i y_i &= 0 \quad \text{and} \quad 0 \leq \alpha_i \leq C \end{aligned}$$

==> positive semi-definite quadratic program

Different
Notation here!

Quadratic term

SVM with Kernels

Training: maximize $D(\vec{\alpha}) = \left(\sum_{i=1}^n \alpha_i \right) - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(\vec{x}_i, \vec{x}_j)$

s. t. $\sum_{i=1}^n \alpha_i y_i = 0 \quad \text{und} \quad 0 \leq \alpha_i \leq C$

Classification: For new example x $h(\vec{x}) = \text{sign} \left(\sum_{x_i \in SV} \alpha_i y_i K(\vec{x}_i, \vec{x}) + b \right)$

New hypotheses spaces through new Kernels:

Linear: $K(\vec{x}_i, \vec{x}_j) = \vec{x}_i \cdot \vec{x}_j$

Polynomial: $K(\vec{x}_i, \vec{x}_j) = [\vec{x}_i \cdot \vec{x}_j + 1]^d$

Radial Basis Functions: $K(\vec{x}_i, \vec{x}_j) = \exp(-|\vec{x}_i - \vec{x}_j|^2 / \sigma^2)$

Sigmoid: $K(\vec{x}_i, \vec{x}_j) = \tanh(\gamma(\vec{x}_i - \vec{x}_j) + c)$

Kernel Function: key property

- If kernel function is chosen with property

$$k(\mathbf{x}, \mathbf{x}') = (\phi(\mathbf{x})^T \cdot \phi(\mathbf{x}'))$$

then computational expense of increased dimensionality is avoided.

- **Polynomial kernel** $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \cdot \mathbf{x}')^p$
can be shown (next slide) to correspond to a map ϕ into the space spanned by **all** products of exactly d dimensions.

A Polynomial Kernel Function

Suppose $\mathbf{x} = (x_1, x_2)$ is a 2 - dimensional input vector

The feature space is 3 - dimensional :

$$\phi(\mathbf{x}) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)^t$$

Then inner product is

$$\phi(\mathbf{x})^T \phi(\mathbf{y}) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)^t (y_1^2, \sqrt{2}y_1y_2, y_2^2)^t = (x_1y_1 + x_2y_2)^2$$

Polynomial kernel function to compute the same value is

$$k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T \cdot \mathbf{y})^2 = \left((x_1, x_2)^T (y_1, y_2) \right)^2 = (x_1y_1 + x_2y_2)^2$$

or $k(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x})^T \phi(\mathbf{y})$

same

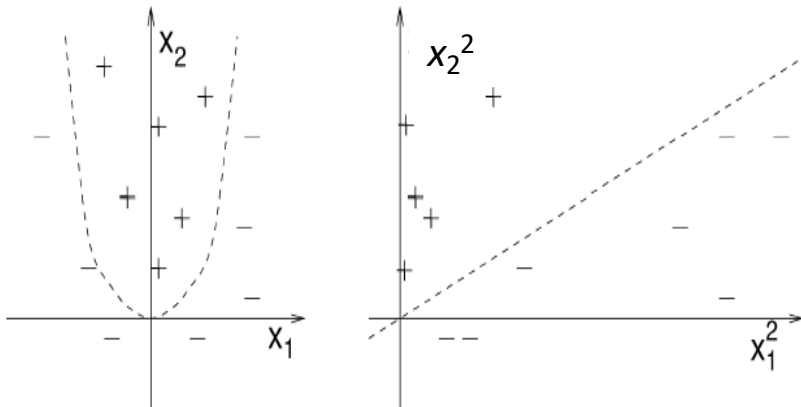
- Inner product $\phi(\mathbf{x})^T \phi(\mathbf{y})$ needs computing six feature values and $3 \times 3 = 9$ multiplications
- Kernel function $k(\mathbf{x}, \mathbf{y})$ has 2 multiplications, one add and a squaring

Another Polynomial kernel function

Example

Input Space: $\mathbf{x} = (x_1, x_2)$ (2 Attributes)

Feature Space: $\Phi(\mathbf{x}) = (x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1x_2, 1)$ (6 Attributes)



- $k(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y} + 1)^2$
- This one maps $d = 2, p = 2$ into a six-dimensional space
- Contains all the powers of \mathbf{x}

$$k(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x})\phi(\mathbf{y})$$

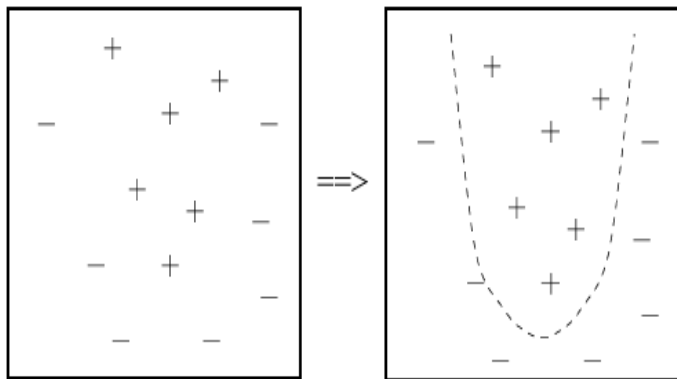
where

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

- Inner product needs 36 multiplications
- Kernel function needs 4 multiplications

Non-Linear Case

Non-Linear Problems



Problem:

- some tasks have non-linear structure
- no hyperplane is sufficiently accurate

How can SVMs learn non-linear classification rules?

- Mapping function $\phi(.)$ to a sufficiently high dimension
- So that data from two categories can always be separated by a hyperplane
- Assume each pattern x_k has been transformed to $\phi(x_k)$, for $k=1, \dots, n$
- First choose the non-linear ϕ functions
 - To map the input vector to a higher dimensional feature space
- Dimensionality of space can be arbitrarily high only limited by computational resources

Mapping into Higher Dimensional Feature Space

- Mapping each input point \mathbf{x} by map

$$\Phi(\mathbf{x}) = \begin{pmatrix} 1 \\ x \\ x^2 \end{pmatrix}$$

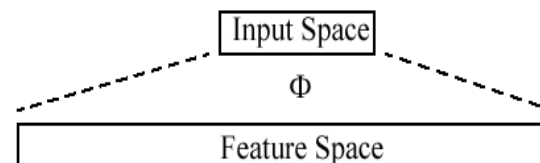
Points on 1-d line are mapped onto curve in 3-d.

- Linear separation in 3-d space is possible. Linear discriminant function in 3-d is in the form

$$y(\mathbf{x}) = a_1x_1 + a_2x_2 + a_3x_3$$

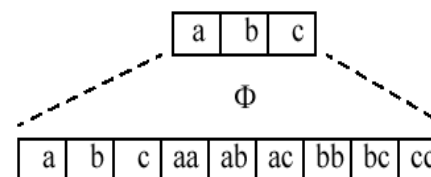
Extending the Hypothesis Space

Idea:



==> Find hyperplane in feature space!

Example:



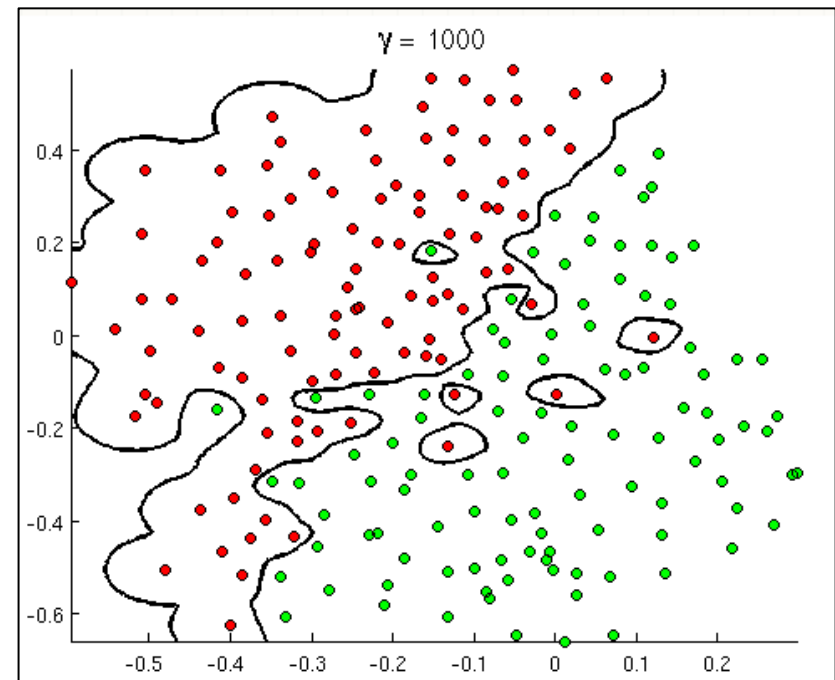
==> The separating hyperplane in features space is a degree two polynomial in input space.

Gaussian (RBF) Kernel

- $\phi(\mathbf{x})$ maps \mathbf{x} to an infinite space

$$K(x, y) = \exp\left(-\frac{\|x - y\|_2^2}{2\sigma^2}\right) = \phi(x)^T \phi(y)$$

$$\phi(x) = e^{-x^2/2\sigma^2} \left[1, \sqrt{\frac{1}{1!\sigma^2}}x, \sqrt{\frac{1}{2!\sigma^4}}x^2, \sqrt{\frac{1}{3!\sigma^6}}x^3, \dots \right]^T$$



- Achieves perfect separation

For sufficiently small kernel bandwidth, decision boundary will look like you just drew little circles around the points whenever they are needed to separate positive and negative examples.

RBF Features from Taylor's Series

Can be inner product in **infinite** dimensional space

Assume $x \in R^1$ and $\gamma > 0$.

$$\begin{aligned}
 e^{-\gamma\|x_i-x_j\|^2} &= e^{-\gamma(x_i-x_j)^2} = e^{-\gamma x_i^2 + 2\gamma x_i x_j - \gamma x_j^2} \\
 &= e^{-\gamma x_i^2 - \gamma x_j^2} \left(1 + \frac{2\gamma x_i x_j}{1!} + \frac{(2\gamma x_i x_j)^2}{2!} + \frac{(2\gamma x_i x_j)^3}{3!} + \dots \right) \\
 &= e^{-\gamma x_i^2 - \gamma x_j^2} \left(1 \cdot 1 + \sqrt{\frac{2\gamma}{1!}} x_i \cdot \sqrt{\frac{2\gamma}{1!}} x_j + \sqrt{\frac{(2\gamma)^2}{2!}} x_i^2 \cdot \sqrt{\frac{(2\gamma)^2}{2!}} x_j^2 \right. \\
 &\quad \left. + \sqrt{\frac{(2\gamma)^3}{3!}} x_i^3 \cdot \sqrt{\frac{(2\gamma)^3}{3!}} x_j^3 + \dots \right) = \phi(x_i)^T \phi(x_j),
 \end{aligned}$$

where

$$\phi(x) = e^{-\gamma x^2} \left[1, \sqrt{\frac{2\gamma}{1!}} x, \sqrt{\frac{(2\gamma)^2}{2!}} x^2, \sqrt{\frac{(2\gamma)^3}{3!}} x^3, \dots \right]^T.$$



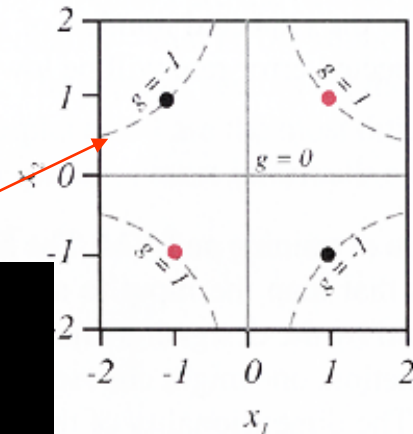
Pattern Transformation using Kernels

- Problem with high-dimensional mapping
 - Very many parameters
 - Polynomial of degree p over d variables leads to $O(d^p)$ variables in feature space
 - Example: if $d = 50$ and $p = 2$ we need a feature space of size 2500
- Solution:
 - Dual Optimization problem needs only inner products
 - Each pattern \mathbf{x}_k transformed into $\phi(\mathbf{x}_k)$
 - Dimensionality of mapped space can be arbitrarily high

SVM for the XOR problem

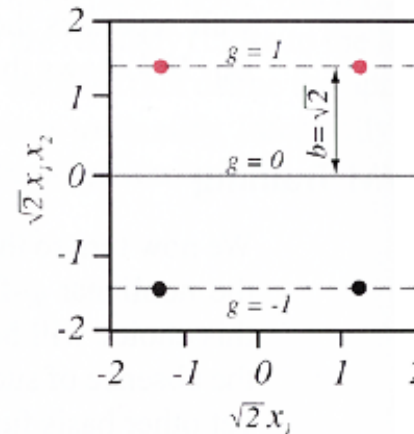
- XOR: binary valued features x_1, x_2
- Not solved by linear discriminant
- Function ϕ maps input $\mathbf{x} = [x_1, x_2]$ into six-D feature space $\mathbf{y} = [1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1x_2, x_1^2, x_2^2]^T$

input space



Hyperbolas
Corresponding
to
 $x_1 x_2 = \pm 1$

feature sub-space



Hyperplanes
Corresponding
to
 $\sqrt{2}x_1 x_2 = \pm 1$

SVM for XOR: maximization problem

- We seek to maximize

$$\sum_{k=1}^4 a_k - \frac{1}{2} \sum_{k=1}^4 \sum_{j=1}^4 a_k a_j t_j t_k \phi(\mathbf{x}_j) \phi(\mathbf{x}_k)$$

- Subject to the constraints $0 = \sum_n a_n t_n$, i.e.,

$$\begin{aligned} a_1 - a_2 + a_3 - a_4 &= 0 \\ 0 \leq a_k \quad k &= 1, 2, 3, 4 \end{aligned}$$

- From problem symmetry, at the solution

$$a_1 = a_3 \text{ and } a_2 = a_4$$

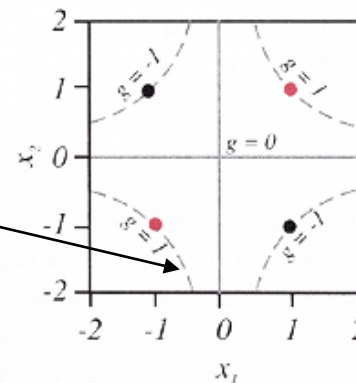
SVM for XOR: maximization problem

- Can use iterative gradient descent
- Or use analytical techniques for small problem
- The solution is

$$\mathbf{a}^* = (1/8, 1/8, 1/8, 1/8)$$

- Last term of Optimizn Problem implies that all four points are support vectors (unusual and due to symmetric nature of XOR)
- The final discriminant function is $y(x_1, x_2) = x_1 \times x_2$
- Decision hyperplane is defined by $y(x_1, x_2) = 0$
- Margin is given by $b = 1/\|\mathbf{a}\| = \sqrt{2}$

Hyperbolas
Corresponding
to
 $x_1 x_2 = \pm 1$



Hyperplanes
Corresponding
to

$$\sqrt{2}x_1 x_2 = \pm 1$$

