

# SPARSE KERNEL MACHINES

J. Elder

CSE 6390/PSYC 6225 Computational Modeling of Visual Perception

# Motivation

Inference can be slow for kernel methods, as the kernel  $k(\mathbf{x}, \mathbf{x}_n)$  must be evaluated for the new data point  $\mathbf{x}$  against **all** training data points  $\mathbf{x}_n$ .

In a sparse kernel machine, the kernel  $k(\mathbf{x}, \mathbf{x}_n)$  need only be evaluated for a subset of the training data.

We will focus in particular on the **Support Vector Machine** (SVM), applied to **classification** problems.

SVMs are **discriminative decision machines**: they do not provide posterior probabilities.

# Support Vector Machines

3

Sparse Kernel Machines

SVMs are based on the linear model  $y(\mathbf{x}) = \mathbf{w}^t \phi(\mathbf{x}) + b$

Assume training data  $\mathbf{x}_1, \dots, \mathbf{x}_N$  with corresponding target values

$t_1, \dots, t_N, t_n \in \{-1, 1\}$ .

$\mathbf{x}$  classified according to sign of  $y(\mathbf{x})$ .

Assume for the moment that the training data are linearly separable in feature space.

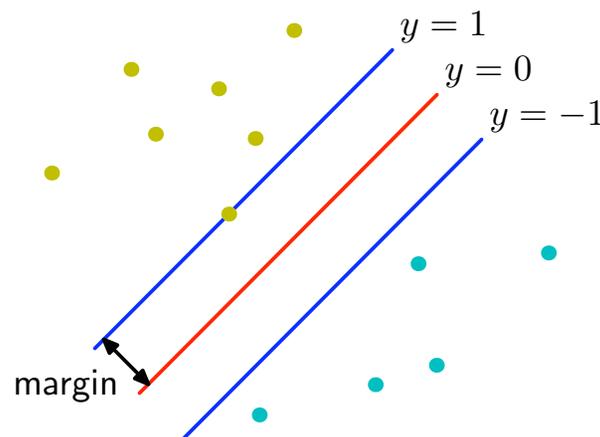
Then  $\exists \mathbf{w}, b : t_n y(\mathbf{x}_n) > 0 \quad \forall n \in [1, \dots, N]$

# Maximum Margin Classifiers

4

Sparse Kernel Machines

- When the training data are linearly separable, there are generally an infinite number of solutions for  $(\mathbf{w}, b)$  that separate the classes exactly.
- The **margin** of such a classifier is defined as the orthogonal distance in feature space between the decision boundary and the closest training vector.
- SVMs are an example of a **maximum margin classifier**, which finds the linear classifier that maximizes the margin.



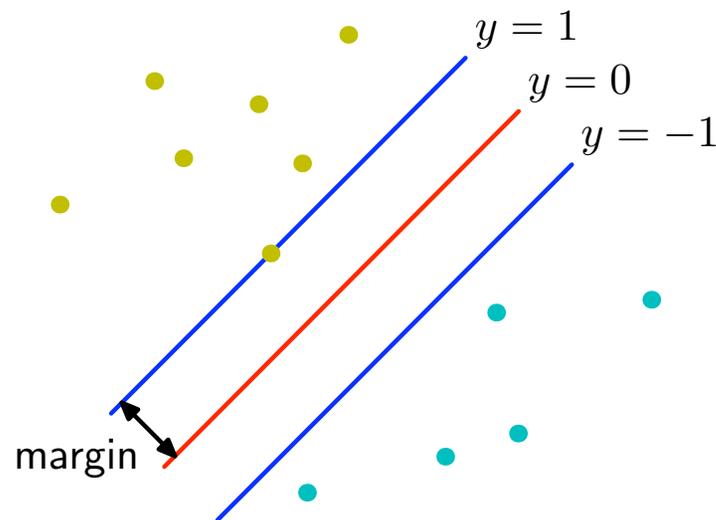
# Probabilistic Motivation

5

Sparse Kernel Machines

- The maximum margin classifier has a probabilistic motivation.

If we model the class-conditional densities with a KDE using Gaussian kernels with variance  $\sigma^2$ , then in the limit as  $\sigma \rightarrow 0$ , the optimal linear decision boundary  $\rightarrow$  maximum margin linear classifier.



# Two Class Discriminant Function

6

Sparse Kernel Machines

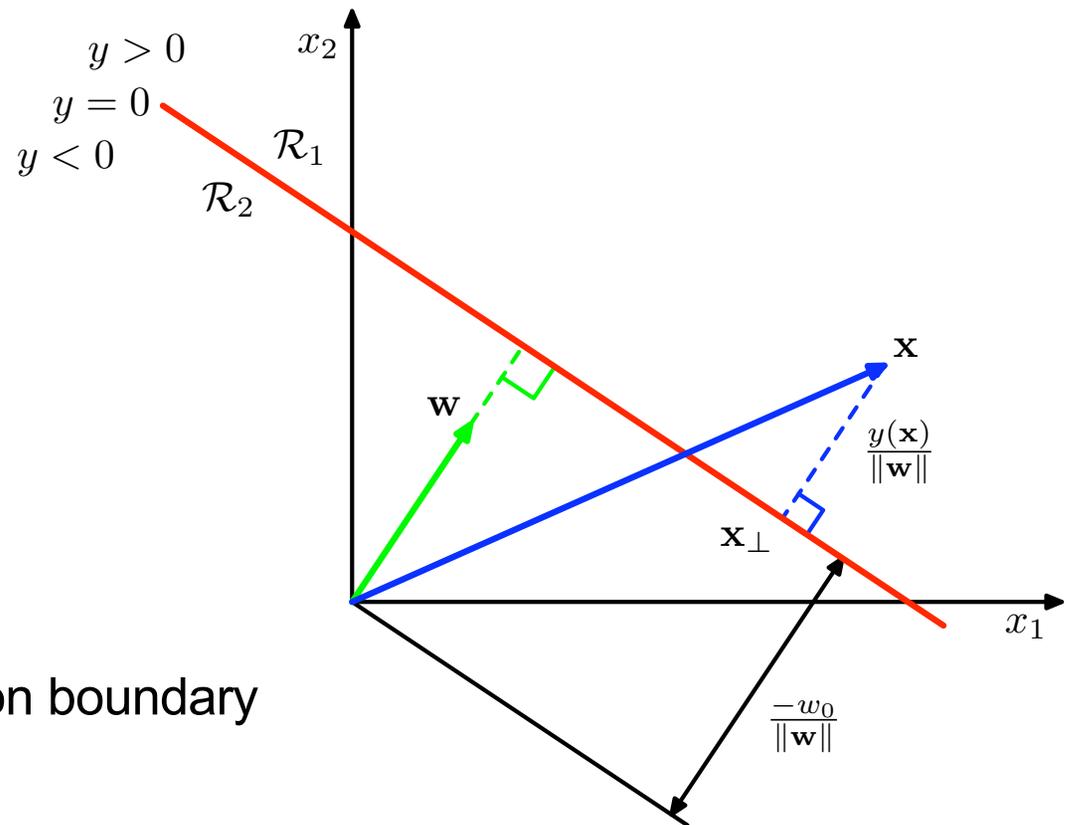
Let  $f(\cdot)$  be the identity:

$$y(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + w_0$$

$y(\mathbf{x}) \geq 0 \rightarrow \mathbf{x}$  assigned to  $C_1$

$y(\mathbf{x}) < 0 \rightarrow \mathbf{x}$  assigned to  $C_2$

Thus  $y(\mathbf{x}) = 0$  defines the decision boundary



# Maximum Margin Classifiers

7

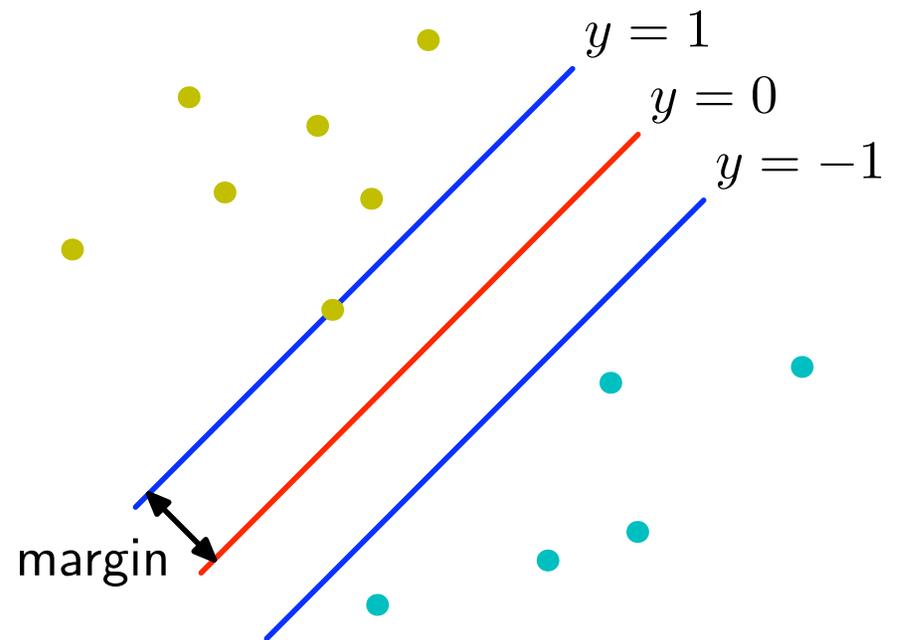
Sparse Kernel Machines

Distance of point  $\mathbf{x}_n$  from decision surface is given by:

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

Thus we seek:

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



# Maximum Margin Classifiers

Distance of point  $\mathbf{x}_n$  from decision surface is given by:

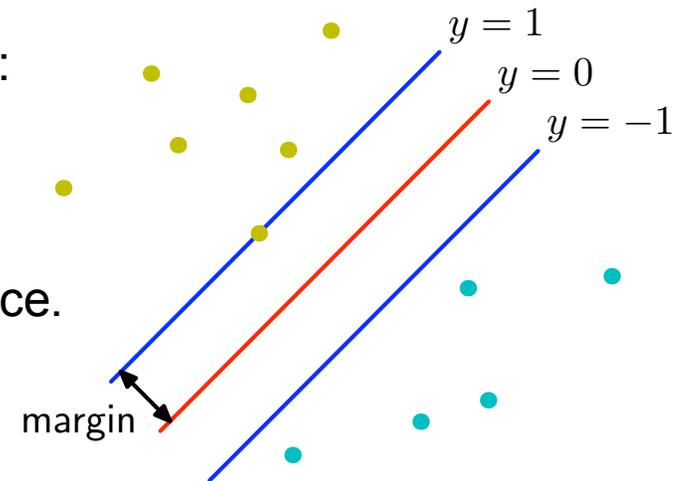
$$\frac{t_n y(\mathbf{x}_n)}{\|\mathbf{w}\|} = \frac{t_n (\mathbf{w}^t \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|}$$

Note that rescaling  $\mathbf{w}$  and  $b$  by the same factor leaves the distance to the decision surface unchanged.

Thus, wlog, we consider only solutions that satisfy:

$$t_n (\mathbf{w}^t \phi(\mathbf{x}_n) + b) = 1.$$

for the point  $\mathbf{x}_n$  that is closest to the decision surface.



# Quadratic Programming Problem

9

Sparse Kernel Machines

Then all points  $\mathbf{x}_n$  satisfy  $t_n(\mathbf{w}^t \phi(\mathbf{x}_n) + b) \geq 1$

Points for which equality holds are said to be **active**.

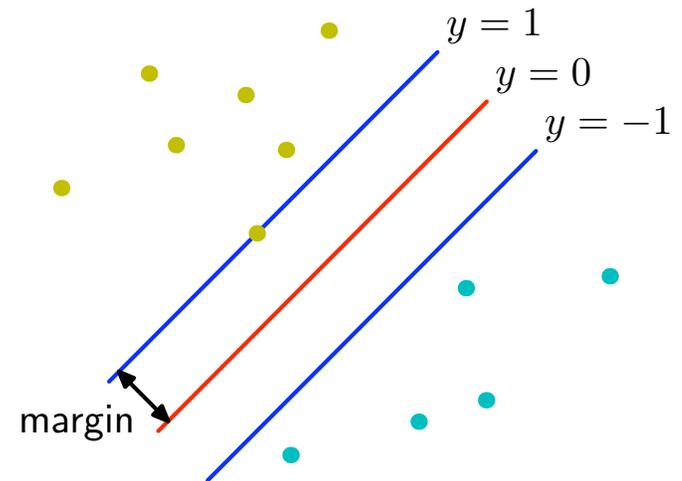
All other points are **inactive**.

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

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

$$\text{Subject to } t_n(\mathbf{w}^t \phi(\mathbf{x}_n) + b) \geq 1 \quad \forall \mathbf{x}_n$$

This is a **quadratic programming** problem.



# Quadratic Programming Problem

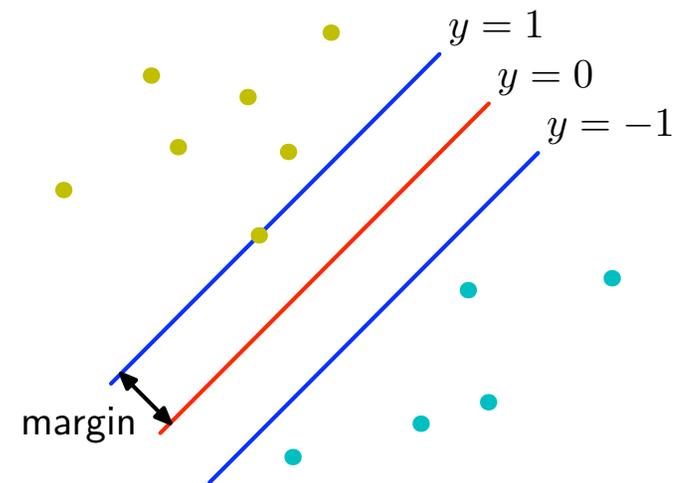
10

Sparse Kernel Machines

$$\frac{1}{2} \arg \min_{\mathbf{w}} \|\mathbf{w}\|^2, \text{ subject to } t_n (\mathbf{w}^t \phi(\mathbf{x}_n) + b) \geq 1 \quad \forall \mathbf{x}_n$$

Solve using Lagrange multipliers  $a_n$  :

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



END OF LECTURE  
NOV 8, 2010

J. Elder

CSE 6390/PSYC 6225 Computational Modeling of Visual Perception

# Dual Representation

12

Sparse Kernel Machines

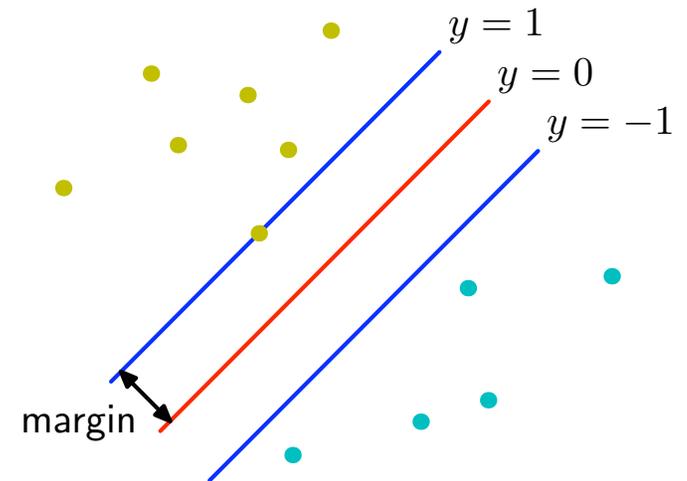
Solve using Lagrange multipliers  $a_n$  :

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

Setting derivatives with respect to  $\mathbf{w}$  and  $b$ , we get:

$$\mathbf{w} = \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n)$$

$$\sum_{n=1}^N a_n t_n = 0$$



# Dual Representation

13

Sparse Kernel Machines

Substituting for  $\mathbf{w}$  and  $b$  leads to the dual representation of the maximum margin problem, in which we maximize:

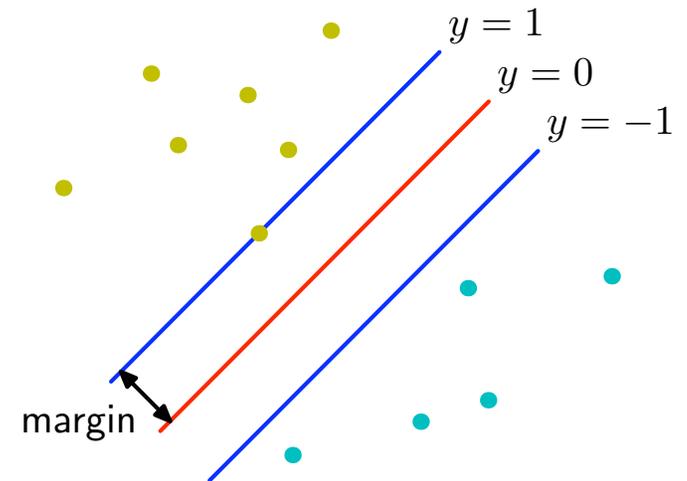
$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

with respect to  $\mathbf{a}$ , subject to:

$$a_n \geq 0 \quad \forall n$$

$$\sum_{n=1}^N a_n t_n = 0$$

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



# Dual Representation

14

Sparse Kernel Machines

Using  $\mathbf{w} = \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n)$ , a new point  $\mathbf{x}$  is classified by computing

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

The Karush-Kuhn-Tucker (KKT) conditions for this constrained optimization problem are:

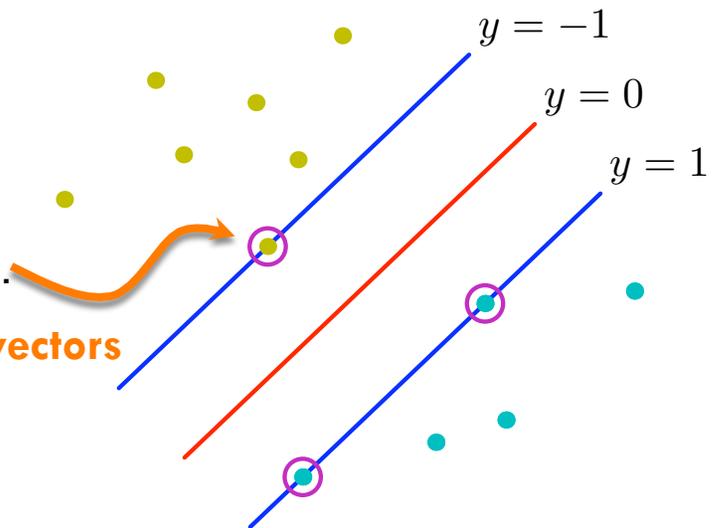
$$a_n \geq 0$$

$$t_n y(\mathbf{x}_n) - 1 \geq 0$$

$$a_n \{t_n y(\mathbf{x}_n) - 1\} = 0$$

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

support vectors



# Solving for the Bias

Once the optimal  $\mathbf{a}$  is determined, the bias  $b$  can be computed from

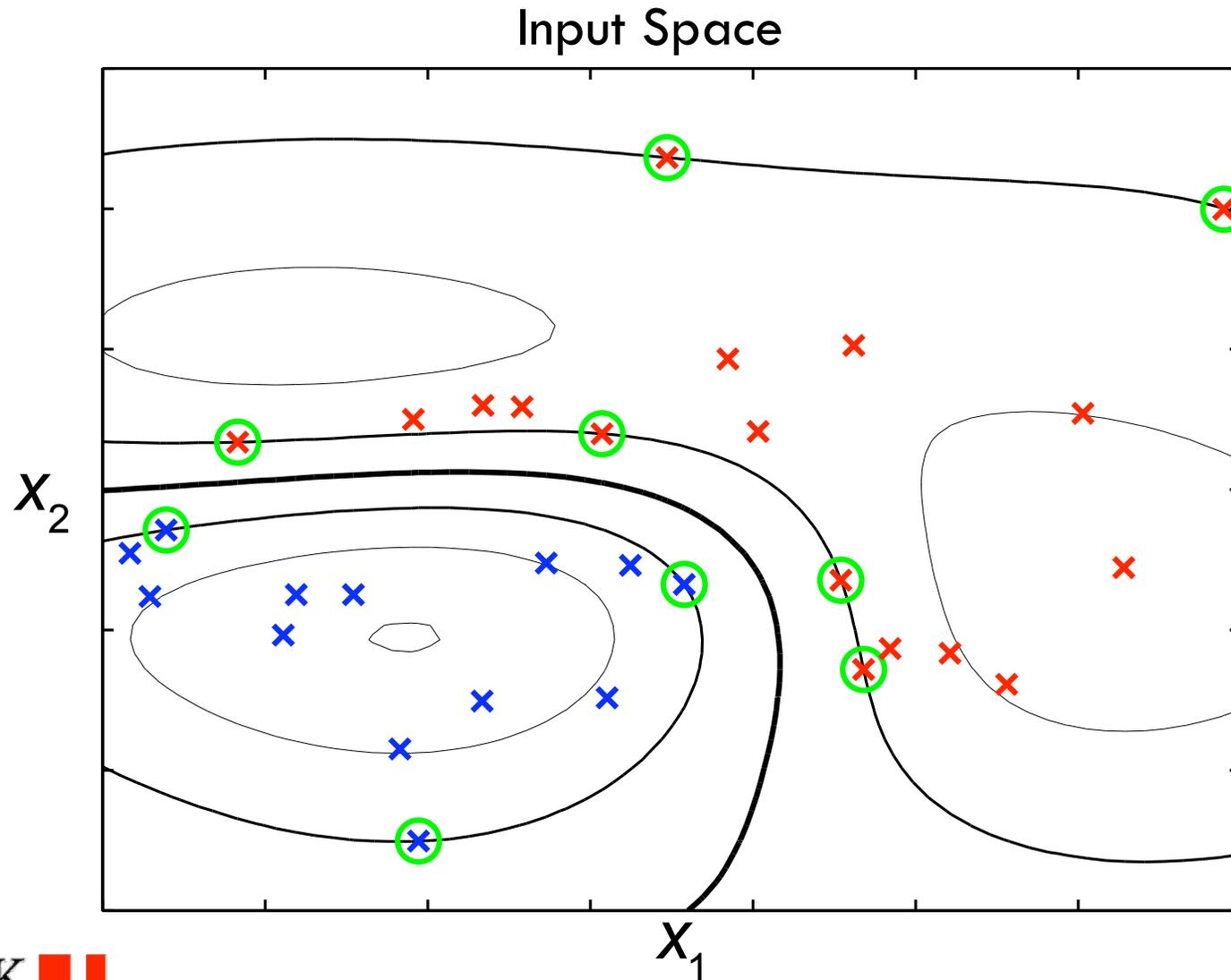
$$b = \frac{1}{N_S} \sum_{n \in S} \left( t_n - \sum_{m \in S} a_m t_m k(\mathbf{x}_n, \mathbf{x}_m) \right)$$

where  $S$  is the index set of support vectors and  $N_S$  is the number of support vectors.

# Example (Gaussian Kernel)

16

Sparse Kernel Machines



# Overlapping Class Distributions

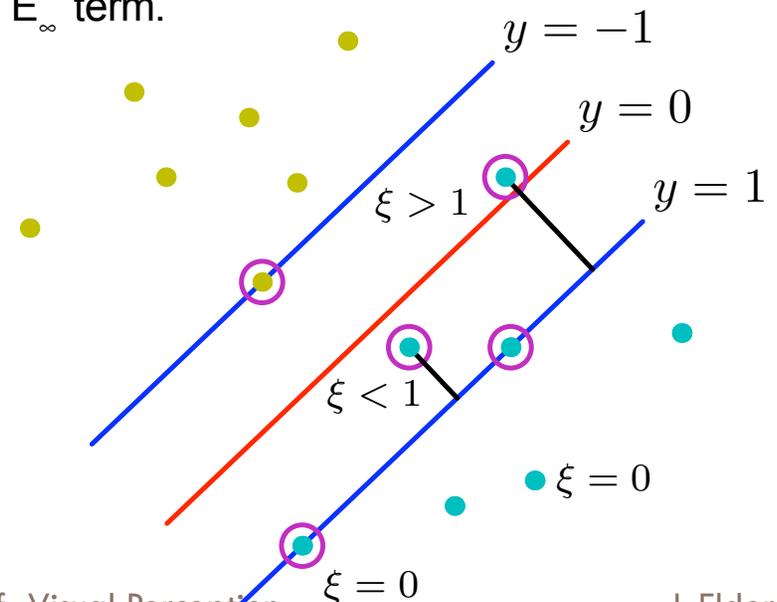
The SVM for non-overlapping class distributions can be expressed as the minimization of

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

where  $E_{\infty}(z)$  is 0 if  $z \geq 0$ , and  $\infty$  otherwise.

This forces all points to lie on or outside the margins, on the correct side for their class.

To allow for misclassified points, we have to relax this  $E_{\infty}$  term.



# Slack Variables

18

Sparse Kernel Machines

To this end, we introduce  $N$  **slack variables**  $\xi_n \geq 0$ ,  $n = 1, \dots, N$ .

$\xi_n = 0$  for points on or on the correct side of the margin boundary for their class

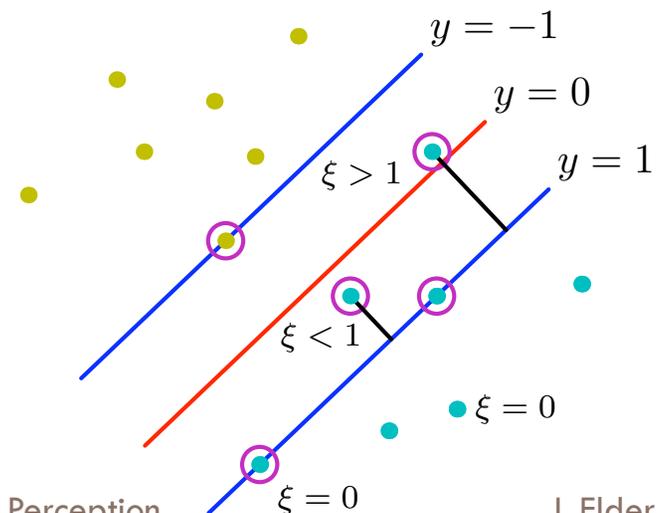
$\xi_n = |t_n - y(\mathbf{x}_n)|$  for all other points.

Thus  $\xi_n < 1$  for points that are correctly classified

$\xi_n > 1$  for points that are incorrectly classified

We now minimize  $C \sum_{n=1}^N \xi_n + \frac{1}{2} \|\mathbf{w}\|^2$ , where  $C > 0$ .

subject to  $t_n y(\mathbf{x}_n) \geq 1 - \xi_n$ , and  $\xi_n \geq 0$ ,  $n = 1, \dots, N$



# Dual Representation

19

Sparse Kernel Machines

This leads to a dual representation, where we maximize

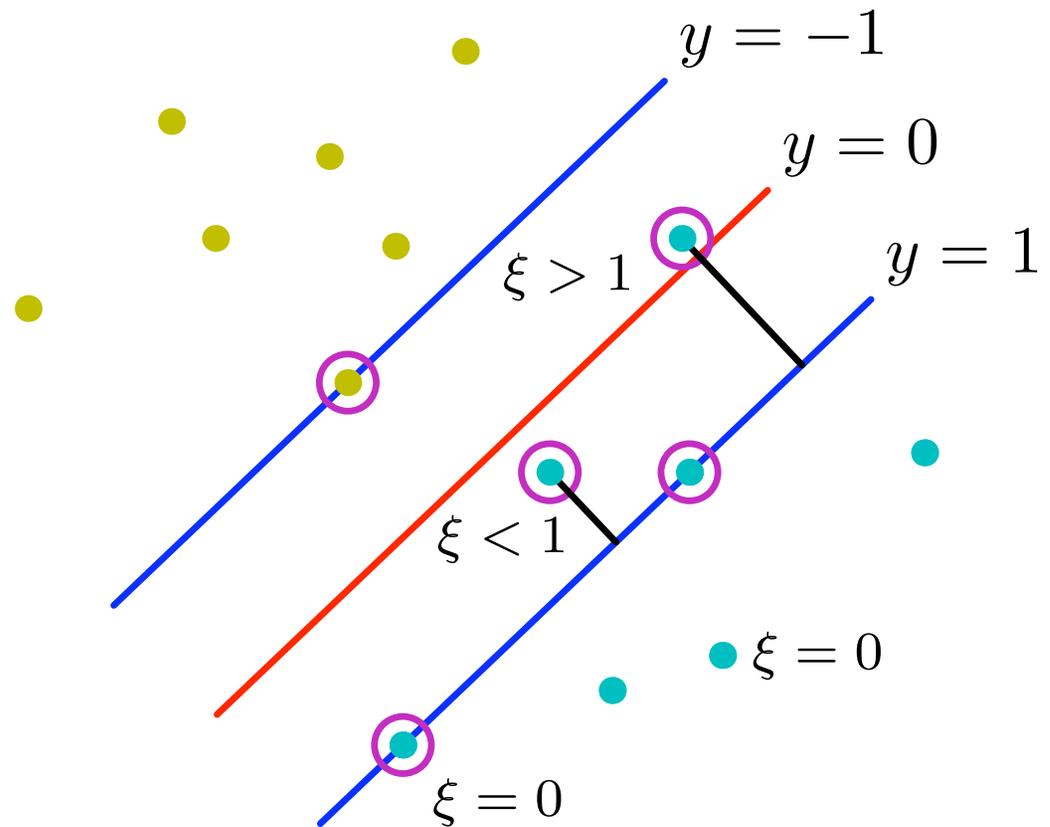
$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

with constraints

$$0 \leq a_n \leq C$$

and

$$\sum_{n=1}^N a_n t_n = 0$$



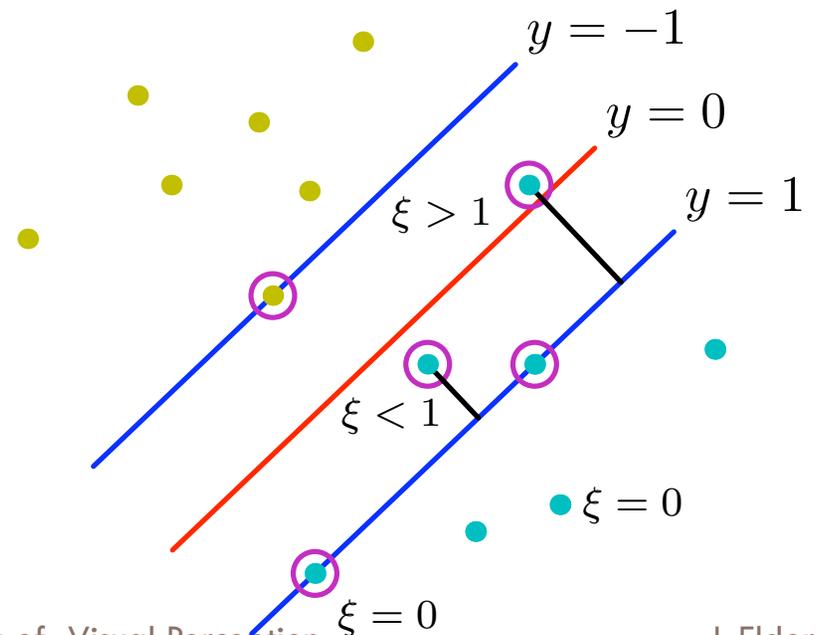
# Support Vectors

Again, a new point  $\mathbf{x}$  is classified by computing

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

For points that are on the correct side of the margin,  $a_n = 0$ .

Thus support vectors consist of points between their margin and the decision boundary, as well as misclassified points.



# Bias

Again, a new point  $\mathbf{x}$  is classified by computing

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

Once the optimal  $\mathbf{a}$  is determined, the bias  $b$  can be computed from

$$b = \frac{1}{N_M} \sum_{n \in M} \left( t_n - \sum_{m \in S} a_m t_m k(\mathbf{x}_n, \mathbf{x}_m) \right)$$

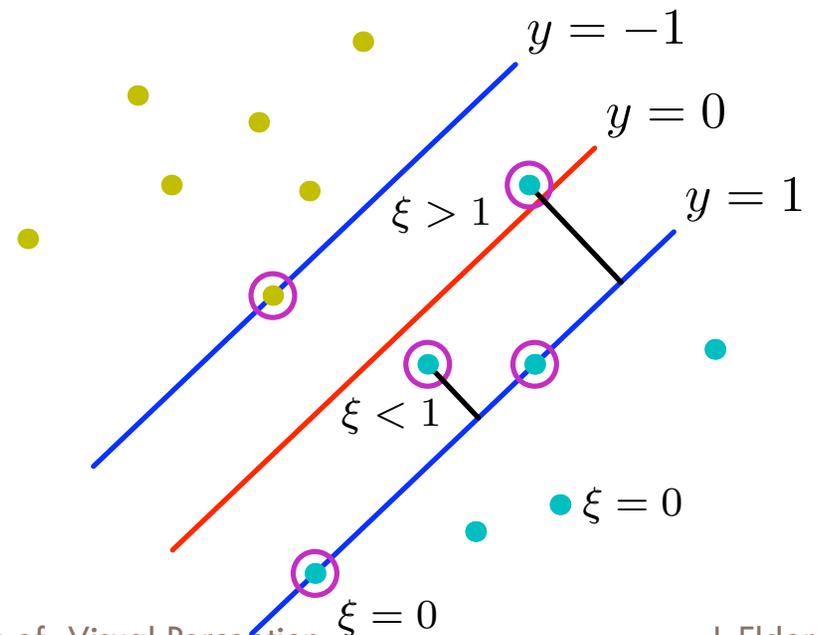
where

$S$  is the index set of support vectors

$N_S$  is the number of support vectors

$M$  is the index set of points on the margins

$N_M$  is the number of points on the margins



# Solving the Quadratic Programming Problem

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

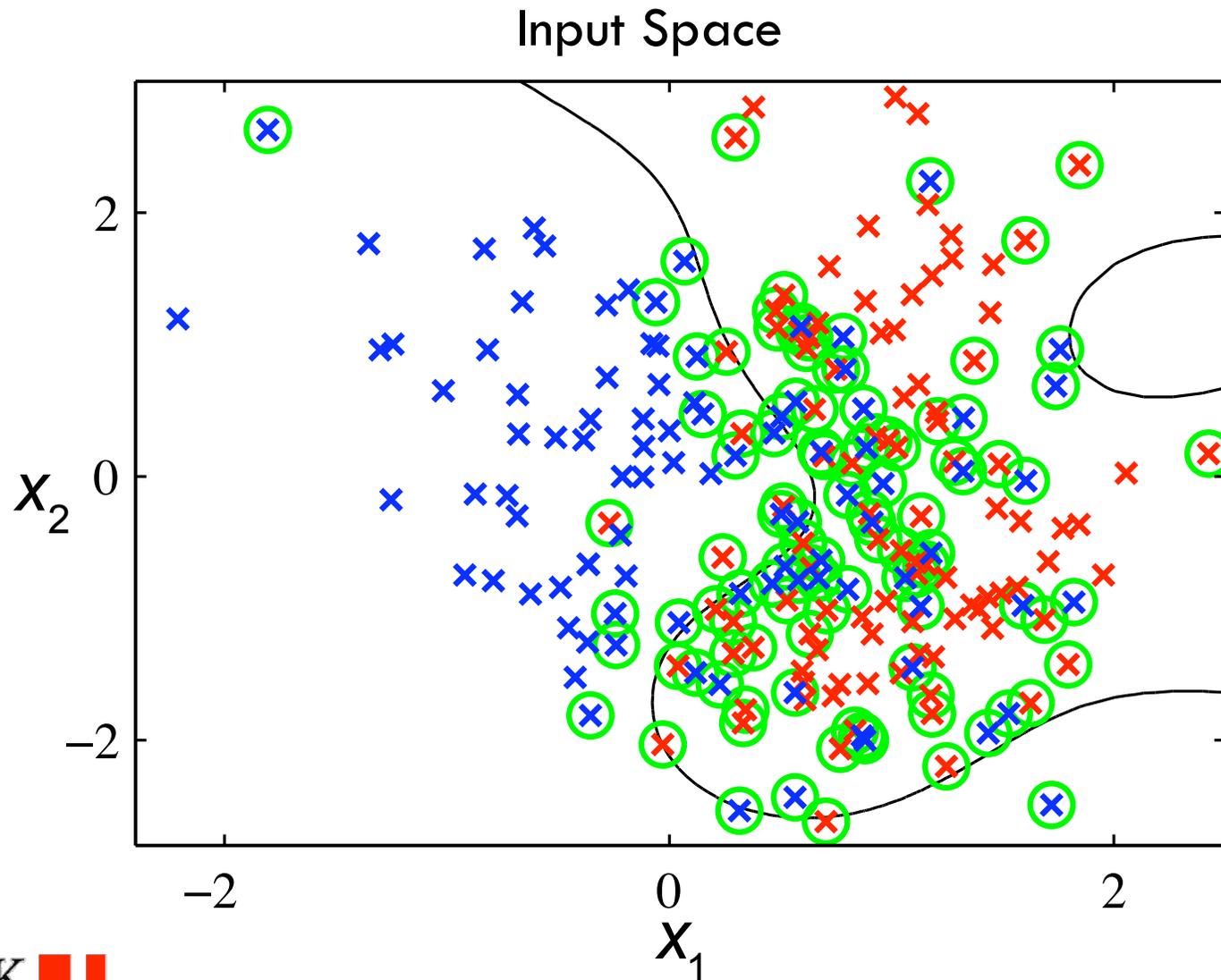
subject to  $0 \leq a_n \leq C$  and  $\sum_{n=1}^N a_n t_n = 0$

- Problem is convex.
- Solutions are generally  $O(N^3)$ .
- Traditional quadratic programming techniques often infeasible due to computation and memory requirements.
- Instead, heuristic methods such as sequential minimal optimization can be used, that in practice are found to scale as  $O(N) - O(N^2)$ .

# Example

23

Sparse Kernel Machines



# Relation to Logistic Regression

24

Sparse Kernel Machines

The objective function for the soft-margin SVM can be written as:

$$\sum_{n=1}^N E_{SV}(y_n t_n) + \lambda \|\mathbf{w}\|^2$$

where  $E_{SV}(z) = [1 - z]_+$  is the **hinge error function**,

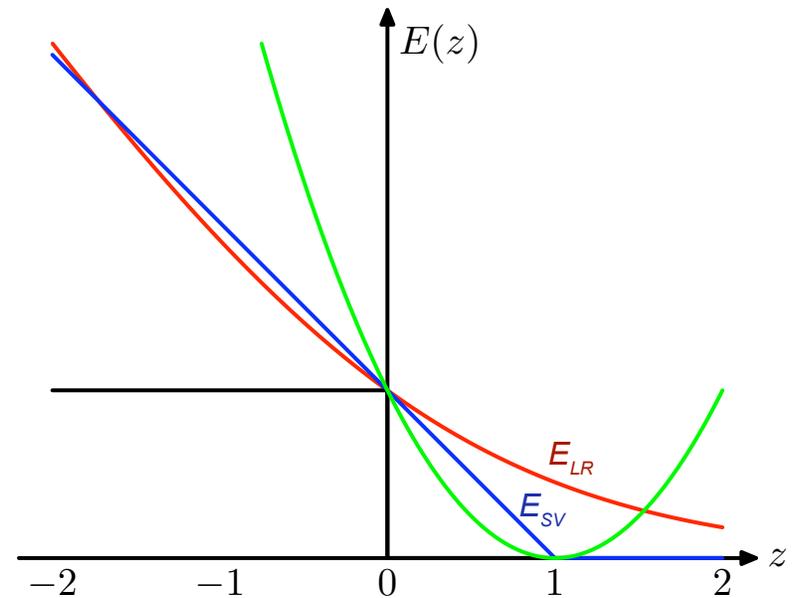
and  $[z]_+ = z$  if  $z \geq 0$

= 0 otherwise.

For  $t \in \{-1, 1\}$ , the objective function for a regularized version of logistic regression can be written as:

$$\sum_{n=1}^N E_{LR}(y_n t_n) + \lambda \|\mathbf{w}\|^2$$

where  $E_{LR}(z) = \log(1 + \exp(-z))$ .



# Multiclass SVMs

25

Sparse Kernel Machines

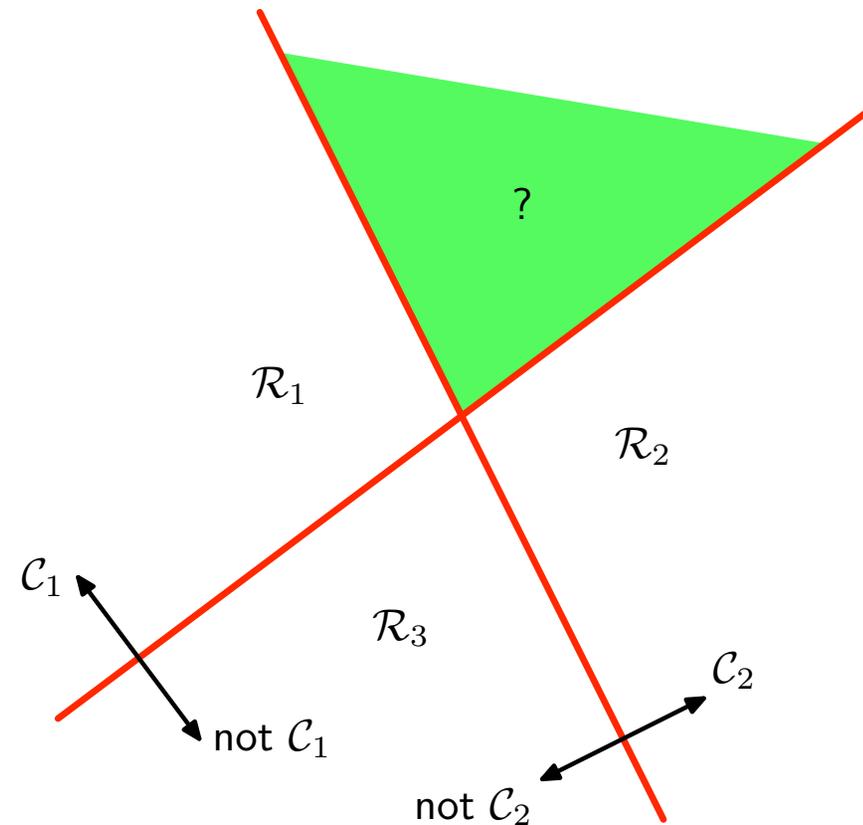
- We encounter the same problems we experienced with least-squares.

# One-Versus-The-Rest

26

Sparse Kernel Machines

- Idea #1: Just use  $K-1$  discriminant functions, each of which separates one class  $C_k$  from the rest.
- Problem: Ambiguous regions

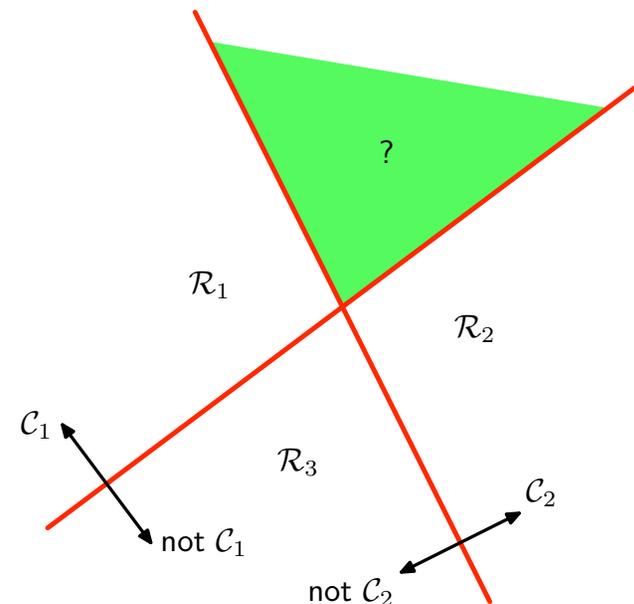


# One-Versus-The-Rest

27

Sparse Kernel Machines

- Possible Solution: select class according to:  $\operatorname{argmax}_k y_k(\mathbf{x})$
- Problems:
  - ▣ Classifiers were all trained separately.
    - Methods for joint training have been proposed – slows training.
  - ▣ Training is imbalanced (e.g., for  $K=10$  classes, 10% in-class, 90% out-of-class)
    - Can be solved by using  $t_n \in \left\{ -\frac{1}{K-1}, 1 \right\}$ .

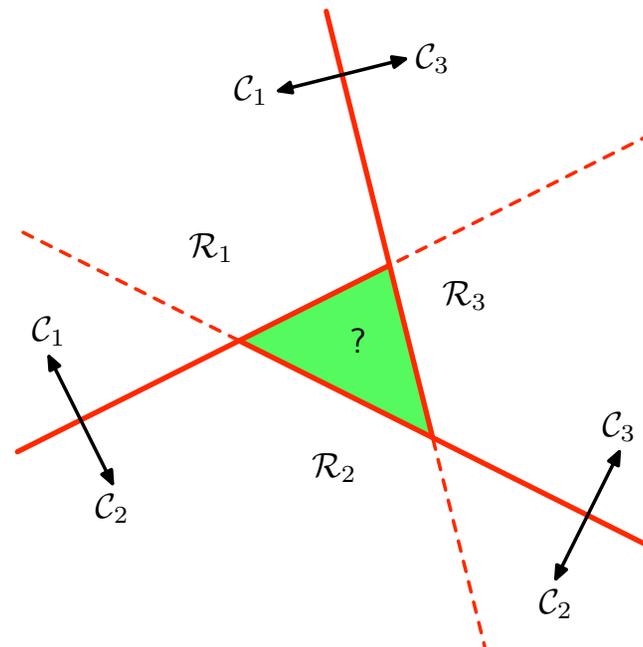


# One-Versus-One

28

Sparse Kernel Machines

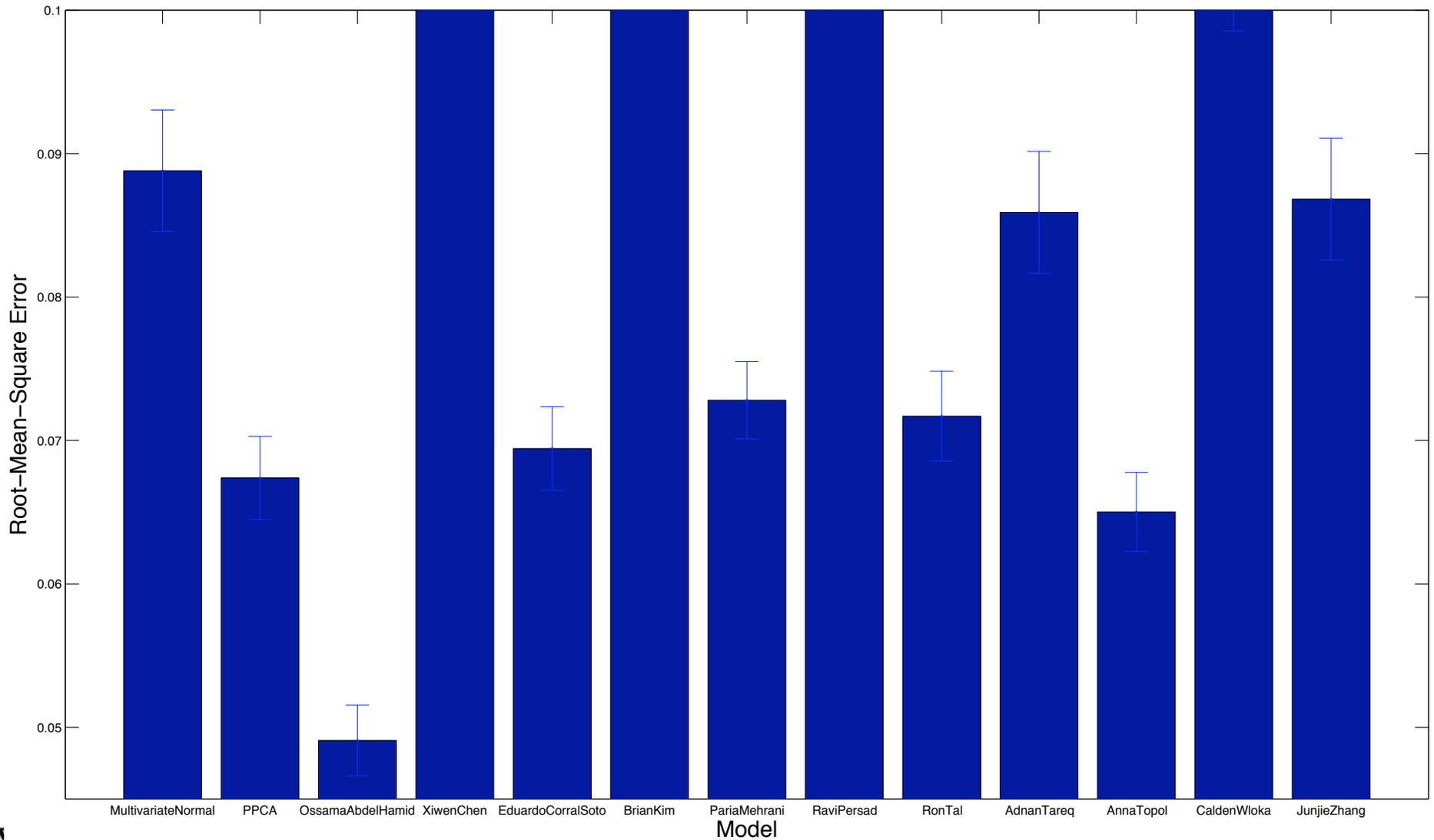
- Idea #2: Use  $K(K-1)/2$  discriminant functions, each of which separates two classes  $C_j, C_k$  from each other.
- Each point classified by majority vote.
- Problems:
  - ▣ Ambiguous regions
  - ▣ Expensive



# Assignment 1 Results

29

Sparse Kernel Machines



# Methods Submitted

- Hierarchy of Gaussian models
- Treat  $x$  and  $y$  coordinates as independent
- Probabilistic PCA
- Gaussian mixtures
- Mean shift
- Use sample mean rather than theoretical mean
- Approximate mean as an ellipse
- Local Gaussian model
- Bi-arc interpolation

# Some Things We've Learned

31

Sparse Kernel Machines

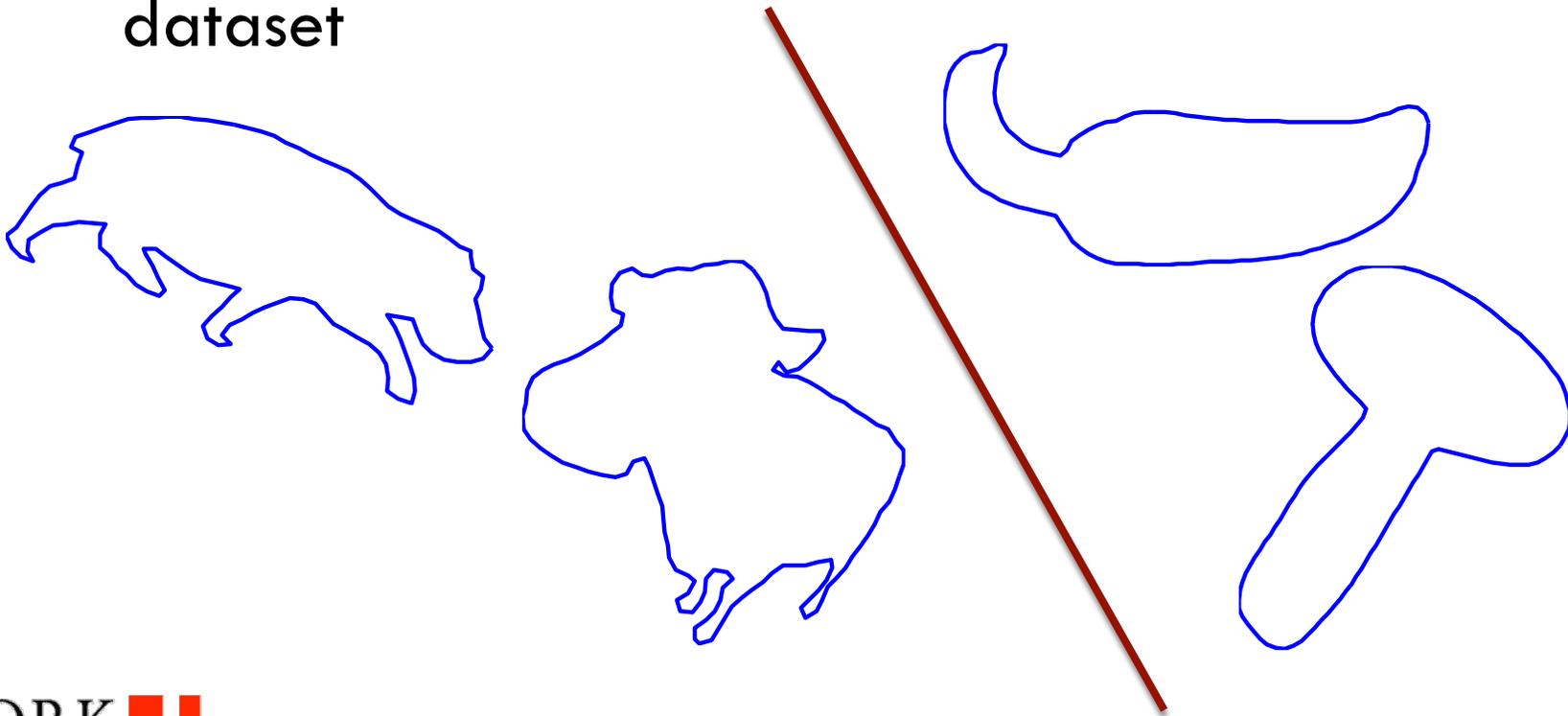
- Use the book!
- The curse of dimensionality
- Probabilistic PCA
- The importance of coding correctly!

# Assignment 2

32

Sparse Kernel Machines

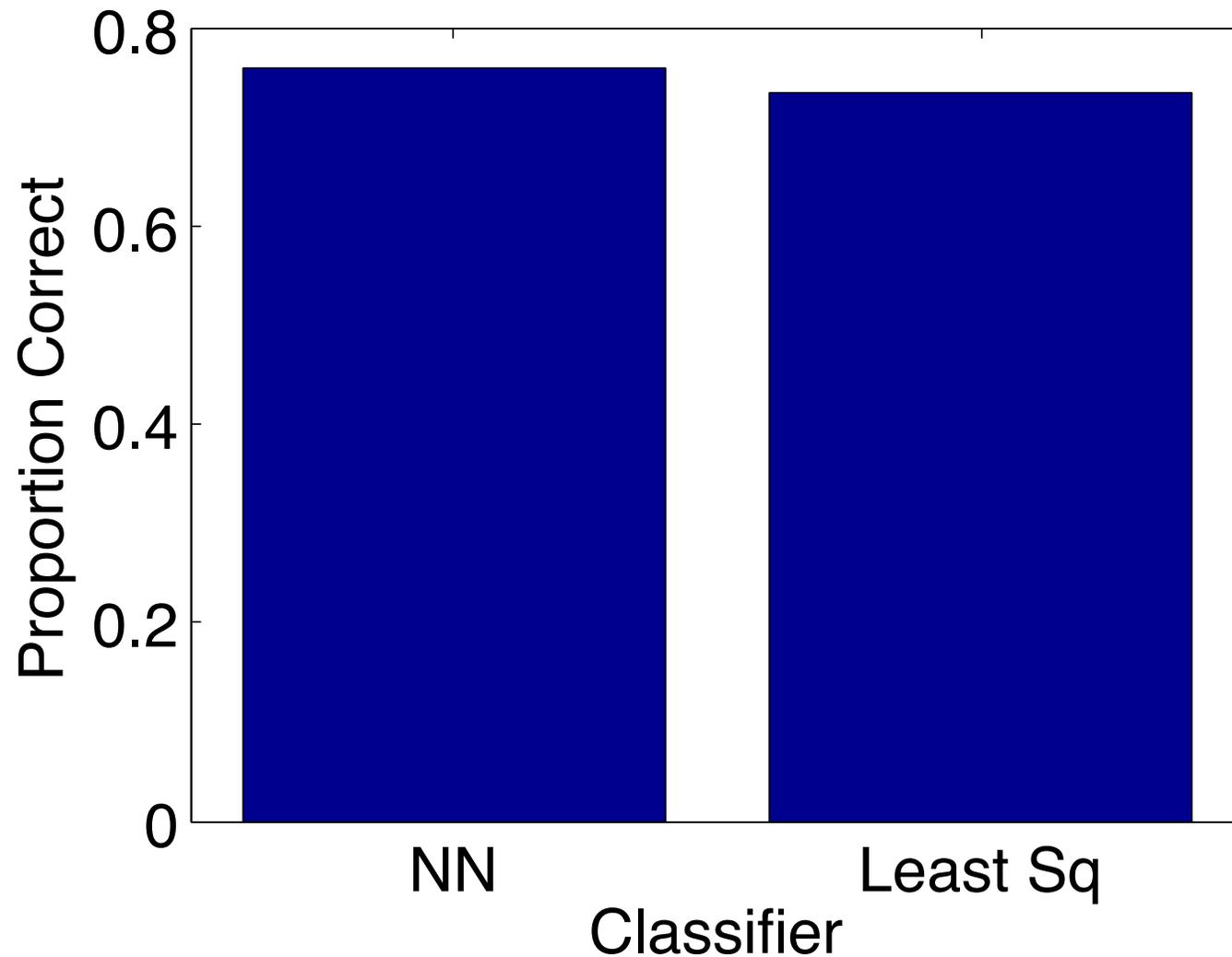
- Classify shapes as ‘animal’ or ‘vegetable’
- Winner has the highest proportion correct
- May be tough to beat nearest-neighbour for this dataset



# Classifiers Provided

33

Sparse Kernel Machines



# SVMs for Regression

34

Sparse Kernel Machines

In standard linear regression, we minimize

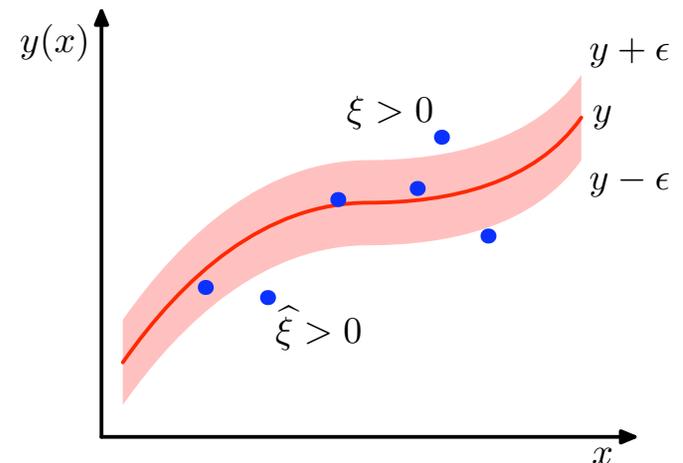
$$\frac{1}{2} \sum_{n=1}^N (y_n - t_n)^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

This penalizes all deviations from the model.

To obtain sparse solutions, we replace the quadratic error function by an  $\epsilon$ -insensitive error function, e.g.,

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

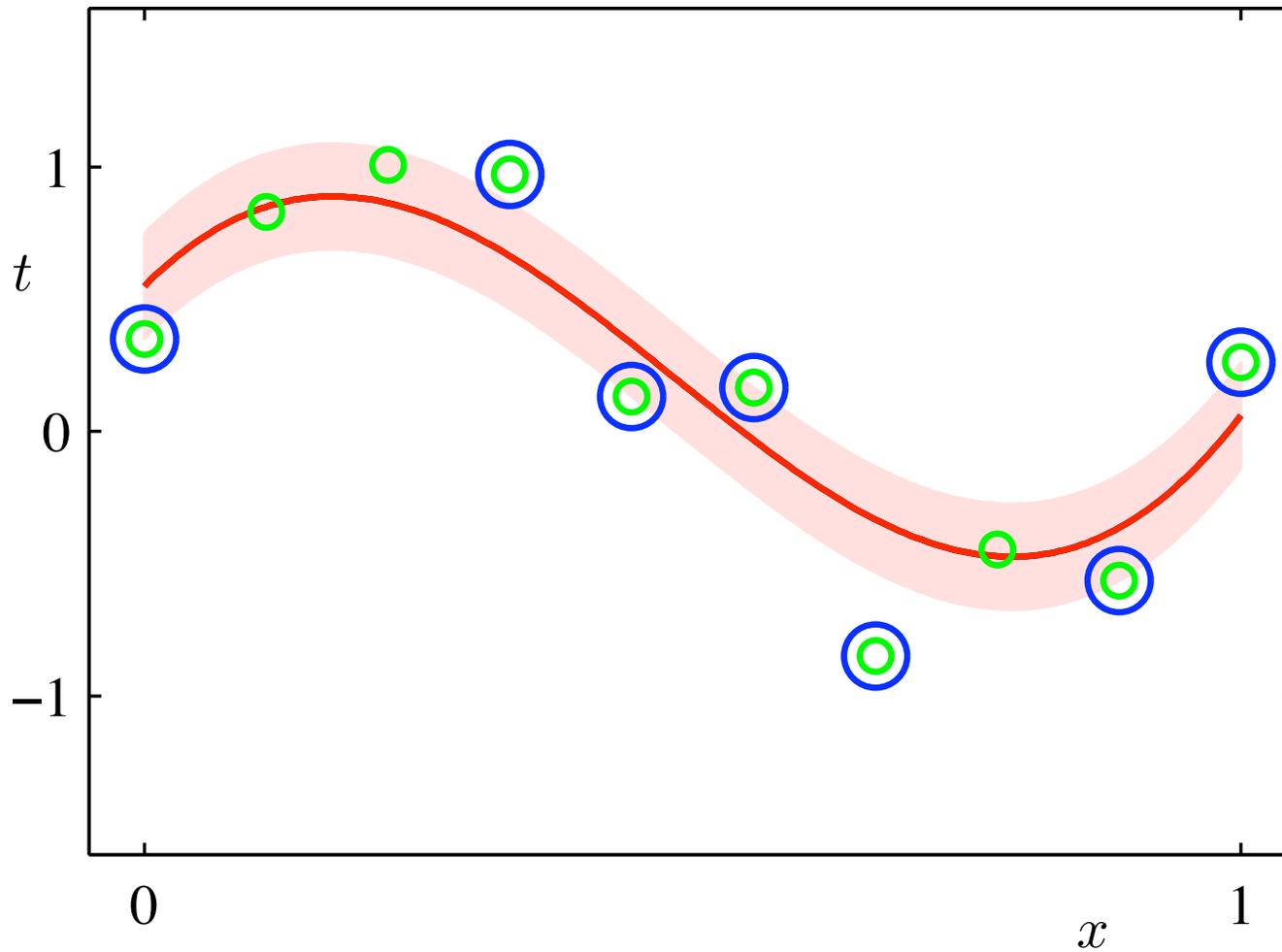
See text for details of solution.



# Example

35

Sparse Kernel Machines



# Relevance Vector Machines

- Some drawbacks of SVMs:
  - ▣ Do not provide posterior probabilities.
  - ▣ Not easily generalized to  $K > 2$  classes.
  - ▣ Parameters  $(C, \epsilon)$  must be learned by cross-validation.
- The **Relevance Vector Machine** is a sparse Bayesian kernel technique that avoids these drawbacks.
- RVMs also typically lead to sparser models.

# RVMs for Regression

$$p(t | \mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t | y(\mathbf{x}), \beta^{-1})$$

where  $y(\mathbf{x}) = \mathbf{w}^t \phi(\mathbf{x})$

In an RVM, the basis functions  $\phi(\mathbf{x})$  are kernels  $k(\mathbf{x}, \mathbf{x}_n)$ :

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

However, unlike in SVMs, the kernels need not be positive definite, and the  $\mathbf{x}_n$  need not be the training data points.

# RVMs for Regression

Likelihood:

$$p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N p(t_n | \mathbf{x}_n, \mathbf{w}, \beta)$$

where the  $n^{\text{th}}$  row of  $\mathbf{X}$  is  $\mathbf{x}_n^t$ .

Prior:

$$p(\mathbf{w} | \alpha) = \prod_{i=1}^M N(w_i | 0, \alpha_i^{-1})$$

- Note that each weight parameter has its own precision hyperparameter.

# RVMs for Regression

39

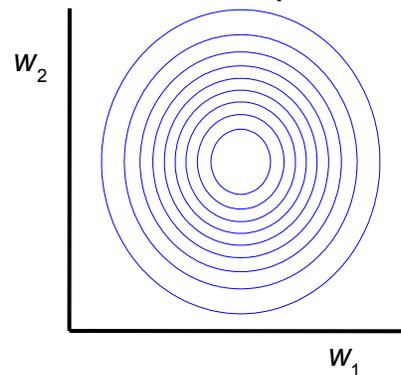
Sparse Kernel Machines

$$p(w_i | \alpha_i) = N(w_i | 0, \alpha_i^{-1})$$

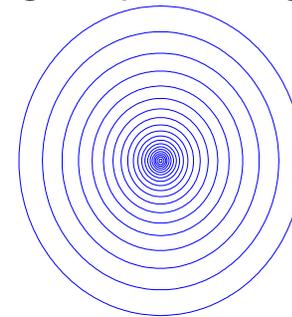
$$p(\alpha_i) = \text{Gam}(\alpha_i | a, b)$$

$$p(w_i) = \text{St}(w_i | 0, a / b, 2a)$$

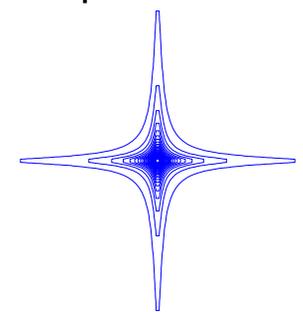
Gaussian prior



Marginal prior: single  $\alpha$



Independent  $\alpha$



- The conjugate prior for the precision of a Gaussian is a gamma distribution.
- Integrating out the precision parameter thus leads to a Student's  $t$  distribution over  $w_i$ .
- Thus the distribution over  $\mathbf{w}$  is a product of Student's  $t$  distributions.
- As a result, maximizing the evidence will yield a sparse  $\mathbf{w}$ .
- Note that to achieve sparsity it is critical that each parameter  $w_i$  has a separate precision  $\alpha_i$ .

# RVMs for Regression

40

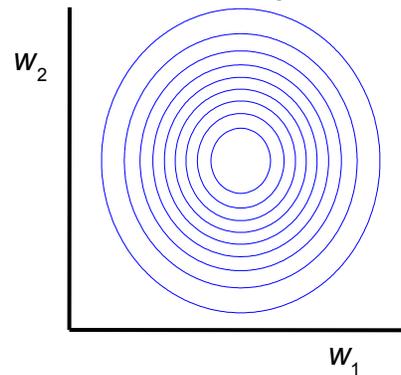
Sparse Kernel Machines

$$p(w_i | \alpha_i) = N(w_i | 0, \alpha_i^{-1})$$

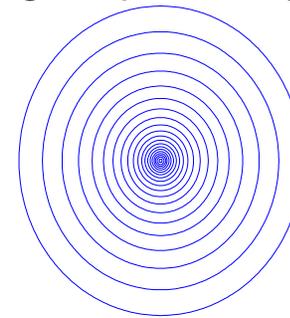
$$p(\alpha_i) = \text{Gam}(\alpha_i | a, b)$$

$$p(w_i) = \text{St}(w_i | 0, a / b, 2a)$$

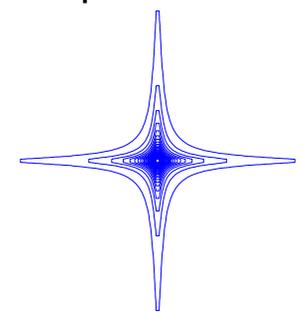
Gaussian prior



Marginal prior: single  $\alpha$



Independent  $\alpha$



If we let  $a \rightarrow 0, b \rightarrow 0$ , then  $p(\log \alpha_i) \rightarrow \text{uniform}$  and  $p(w_i) \rightarrow |w_i|^{-1}$ .

# RVMs for Regression

Likelihood:

$$p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N p(t_n | \mathbf{x}_n, \mathbf{w}, \beta)$$

where the  $n^{\text{th}}$  row of  $\mathbf{X}$  is  $\mathbf{x}_n^t$ .

Prior:

$$p(\mathbf{w} | \alpha) = \prod_{i=1}^M N(w_i | 0, \alpha_i^{-1})$$

- In practice, it is difficult to integrate  $\alpha$  out exactly.
- Instead, we use Type II Maximum Likelihood, finding ML values for each  $\alpha_i$ .
- When we maximize the evidence with respect to these hyperparameters, many will  $\rightarrow \infty$ .
- As a result, the corresponding weights will  $\rightarrow 0$ , yielding a sparse solution.

# RVMs for Regression

- Since both the likelihood and prior are normal, the posterior over  $\mathbf{w}$  will also be normal:

Posterior:

$$p(\mathbf{w} | \mathbf{t}, \mathbf{X}, \alpha, \beta) = \mathcal{N}(\mathbf{w} | \mathbf{m}, \Sigma)$$

where

$$\mathbf{m} = \beta \Sigma \Phi^t \mathbf{t}$$

$$\Sigma = (\mathbf{A} + \beta \Phi^t \Phi)^{-1}$$

and

$$\Phi_{ni} = \phi_i(\mathbf{x}_n)$$

$$\mathbf{A} = \text{diag}(\alpha_i)$$

# RVMs for Regression

- The values for  $\alpha$  and  $\beta$  are determined using the evidence approximation, where we maximize

$$p(\mathbf{t} | \mathbf{X}, \alpha, \beta) = \int p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w} | \alpha) d\mathbf{w}$$

In general, this results in many of the precision parameters  $\alpha_i \rightarrow \infty$ , so that  $w_i \rightarrow 0$ .

Unfortunately, this is a non-convex problem.

# Example

