

LINEAR CLASSIFIERS

J. Elder

CSE 4404/5327 Introduction to Machine Learning and Pattern Recognition

Classification: Problem Statement

2

Probability & Bayesian Inference

- In regression, we are modeling the relationship between a continuous input variable x and a continuous target variable t .
- In classification, the input variable x may still be continuous, but the target variable is discrete.
- In the simplest case, t can have only 2 values.

e.g., Let $t = +1$ ~ ~ assigned to C_1

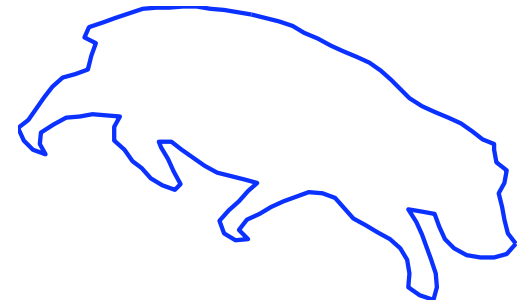
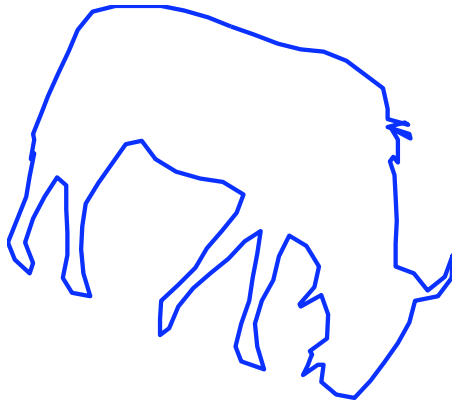
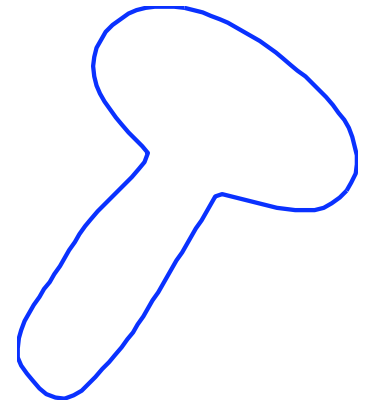
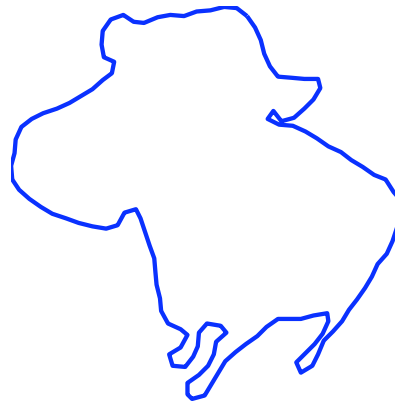
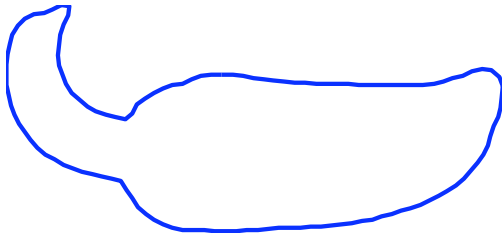
$t = -1$ ~ ~ assigned to C_2

Example Problem

3

Probability & Bayesian Inference

□ Animal or Vegetable?



Linear Models for Classification

4

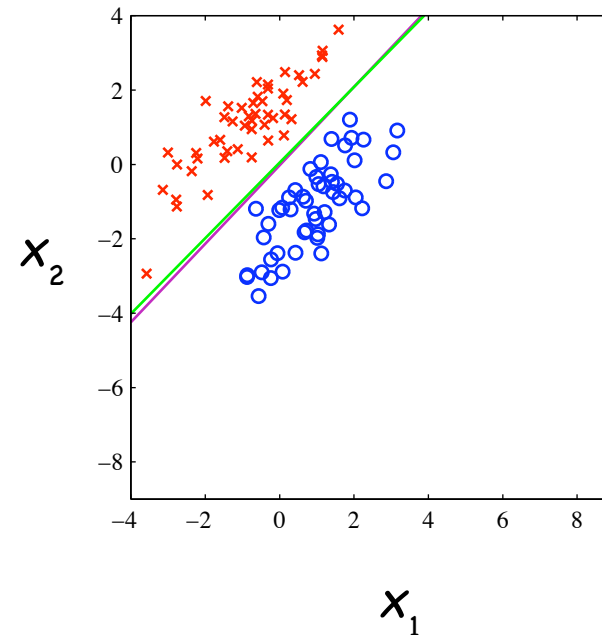
Probability & Bayesian Inference

- Linear models for classification separate input vectors into classes using linear (hyperplane) *decision boundaries*.

- Example:

2D Input vector \mathbf{x}

Two discrete classes C_1 and C_2



Two Class Discriminant Function

5

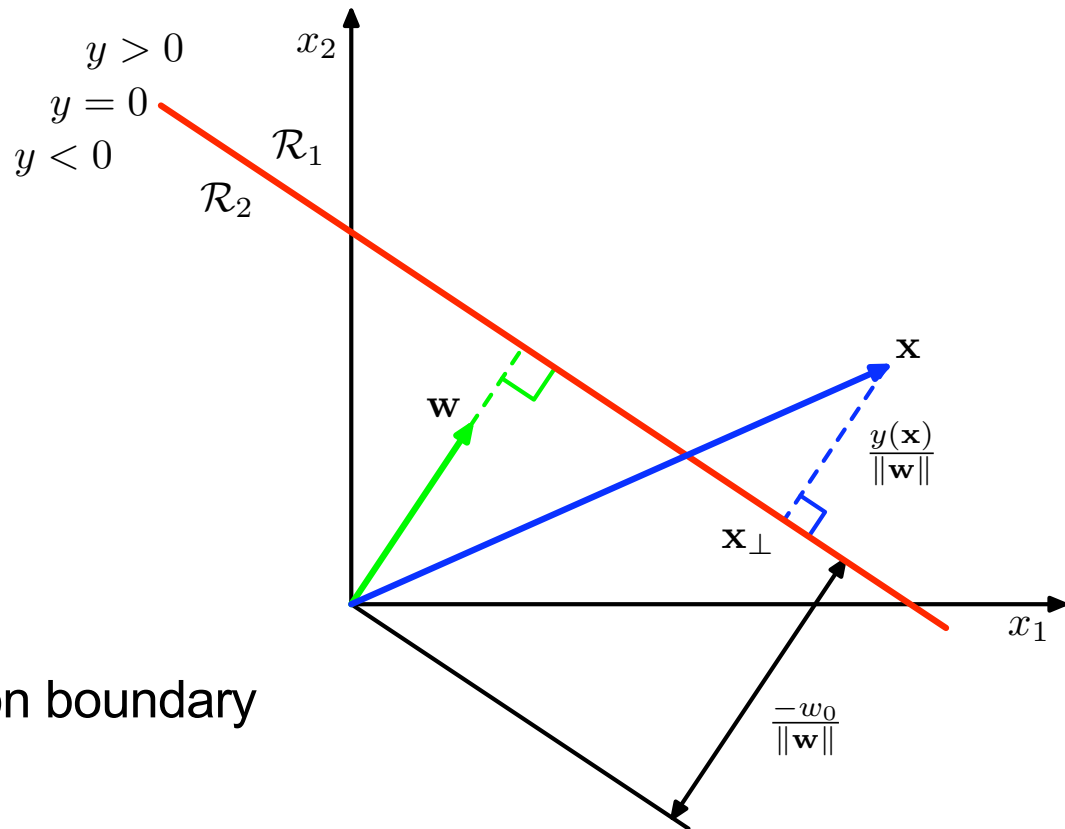
Probability & Bayesian Inference

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

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

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

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



Two-Class Discriminant Function

6

Probability & Bayesian Inference

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

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

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

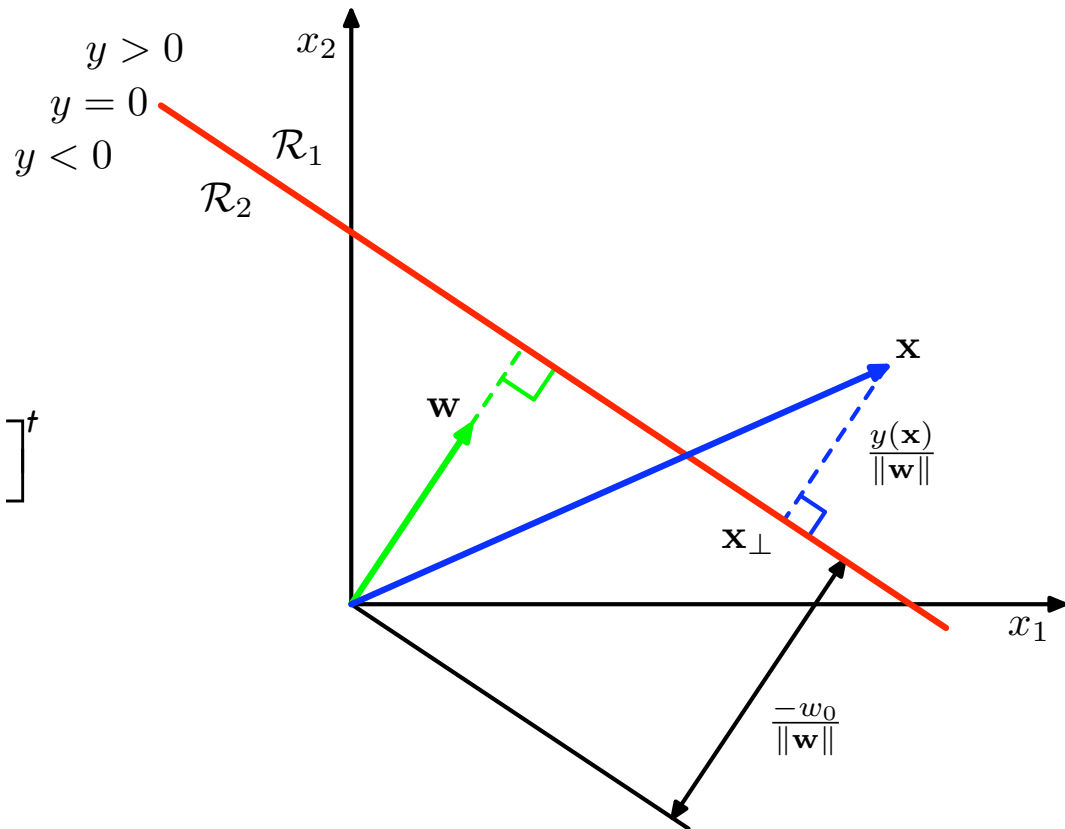
For convenience, let

$$\mathbf{w} = [w_1 \dots w_M]^t \Rightarrow [w_0 \ w_1 \dots w_M]^t$$

and

$$\mathbf{x} = [x_1 \dots x_M]^t \Rightarrow [1 \ x_1 \dots x_M]^t$$

So we can express $y(\mathbf{x}) = \mathbf{w}^t \mathbf{x}$



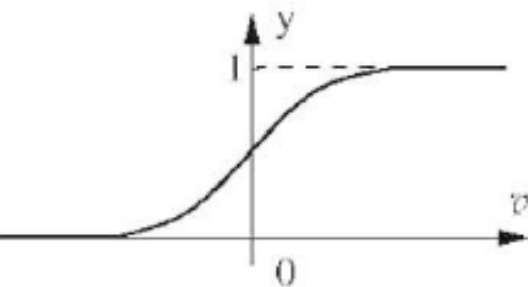
Generalized Linear Models

7

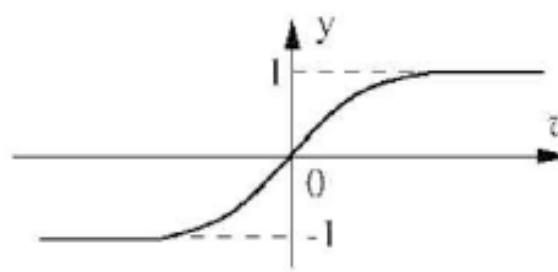
Probability & Bayesian Inference

- For classification problems, we want y to be a predictor of t . In other words, we wish to map the input vector into one of a number of discrete classes, or to posterior probabilities that lie between 0 and 1.
- For this purpose, it is useful to elaborate the linear model by introducing a nonlinear activation function f , which typically will constrain y to lie between -1 and 1 or between 0 and 1.

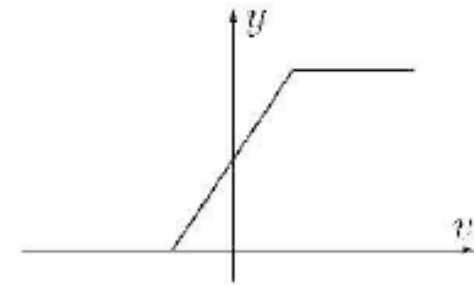
$$y(\mathbf{x}) = f(\mathbf{w}^t \mathbf{x} + w_0)$$



Log-sigmoid function



Tan-sigmoid function



Linear function

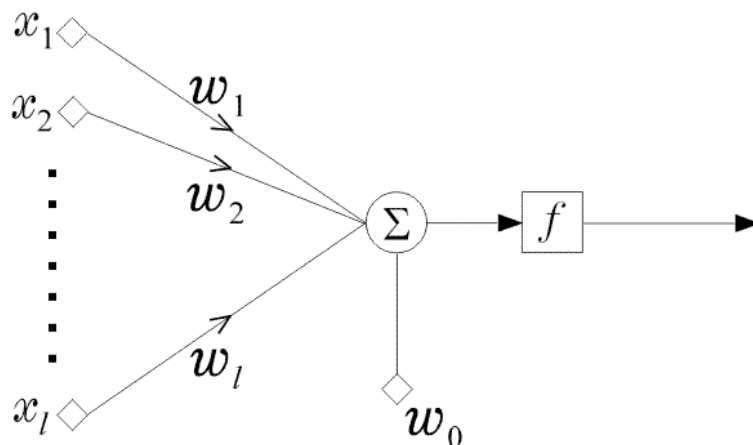
The Perceptron

8

Probability & Bayesian Inference

$$y(\mathbf{x}) = f(\mathbf{w}^t \mathbf{x} + w_0) \quad \begin{array}{l} y(\mathbf{x}) \geq 0 \rightarrow \mathbf{x} \text{ assigned to } C_1 \\ y(\mathbf{x}) < 0 \rightarrow \mathbf{x} \text{ assigned to } C_2 \end{array}$$

- A classifier based upon this simple generalized linear model is called a (single layer) **perceptron**.
- It can also be identified with an abstracted model of a neuron called the McCulloch Pitts model.



Parameter Learning

9

Probability & Bayesian Inference

- How do we learn the parameters of a perceptron?

Outline

10

Probability & Bayesian Inference

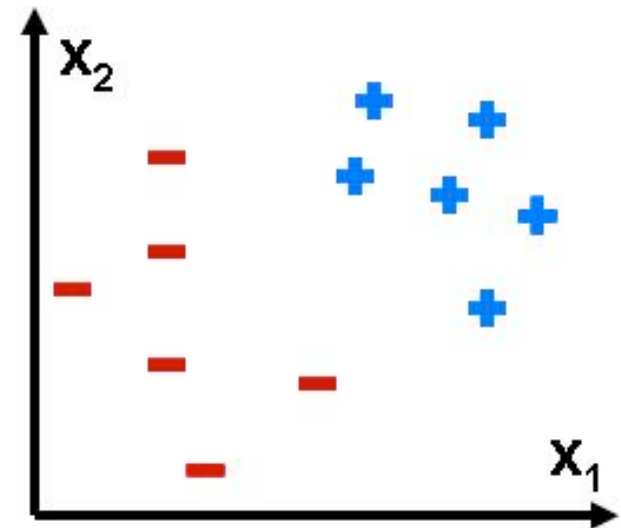
- **The Perceptron Algorithm**
- Least-Squares Classifiers
- Fisher's Linear Discriminant
- Logistic Classifiers
- Support Vector Machines

Case 1. Linearly Separable Inputs

11

Probability & Bayesian Inference

- For starters, let's assume that the training data is in fact perfectly linearly separable.
- In other words, there exists at least one hyperplane (one set of weights) that yields 0 classification error.
- We seek an algorithm that can automatically find such a hyperplane.



The Perceptron Algorithm

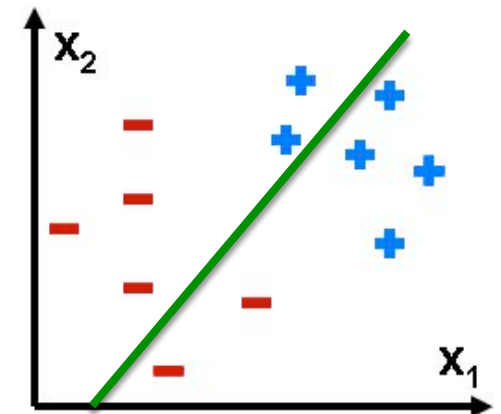
12

Probability & Bayesian Inference

- The perceptron algorithm was invented by Frank Rosenblatt (1962).
- The algorithm is iterative.
- The strategy is to start with a random **guess** at the weights \mathbf{w} , and to then iteratively change the weights to move the hyperplane in a direction that lowers the classification error.



Frank Rosenblatt (1928 – 1971)

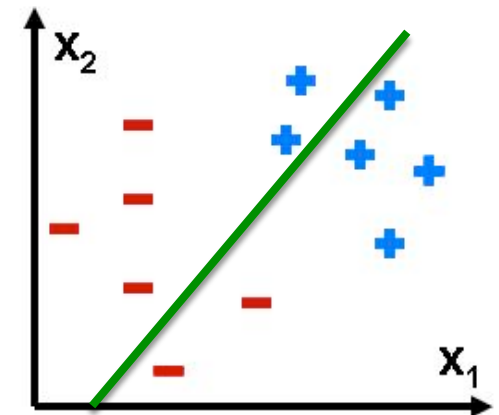


The Perceptron Algorithm

13

Probability & Bayesian Inference

- Note that as we change the weights continuously, the classification error changes in discontinuous, piecewise constant fashion.
- Thus we cannot use the classification error per se as our objective function to minimize.
- What would be a better objective function?



The Perceptron Criterion

14

Probability & Bayesian Inference

- Note that we seek \mathbf{w} such that

$$\mathbf{w}^t \mathbf{x} \geq 0 \text{ when } t = +1$$

$$\mathbf{w}^t \mathbf{x} < 0 \text{ when } t = -1$$

- In other words, we would like

$$\mathbf{w}^t \mathbf{x}_n t_n \geq 0 \quad \forall n$$

- Thus we seek to minimize

$$E_P(\mathbf{w}) = - \sum_{n \in \mathcal{M}} \mathbf{w}^t \mathbf{x}_n t_n$$

where \mathcal{M} is the set of misclassified inputs.

The Perceptron Criterion

15

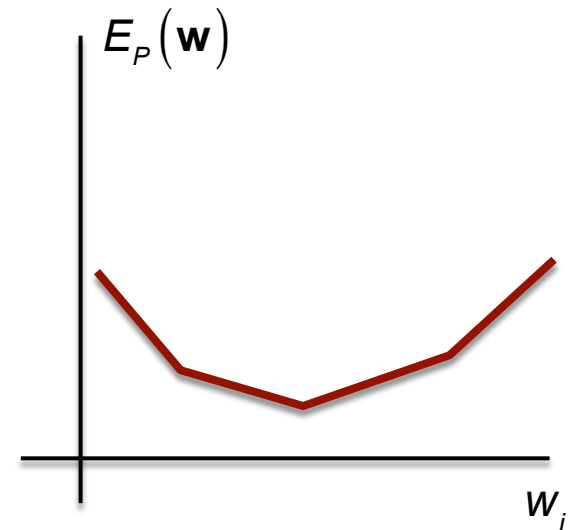
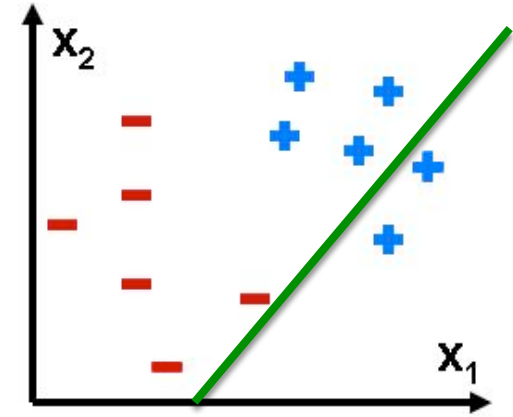
Probability & Bayesian Inference

$$E_P(\mathbf{w}) = - \sum_{n \in \mathcal{M}} \mathbf{w}^t \mathbf{x}_n t_n$$

where \mathcal{M} is the set of misclassified inputs.

□ Observations:

- $E_P(\mathbf{w})$ is always non-negative.
- $E_P(\mathbf{w})$ is continuous and piecewise linear, and thus easier to minimize.



The Perceptron Algorithm

16

Probability & Bayesian Inference

$$E_P(\mathbf{w}) = - \sum_{n \in \mathcal{M}} \mathbf{w}^t \mathbf{x}_n t_n$$

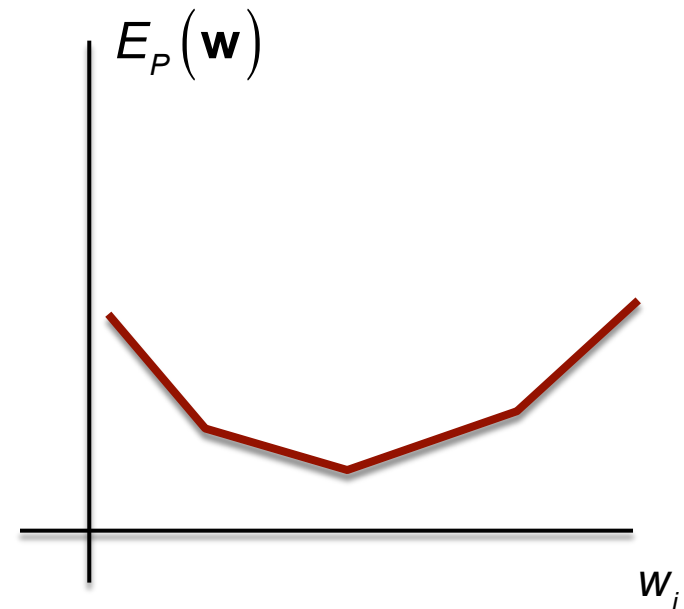
where \mathcal{M} is the set of misclassified inputs.

$$\frac{dE_P(\mathbf{w})}{d\mathbf{w}} = - \sum_{n \in \mathcal{M}} \mathbf{x}_n t_n$$

where the derivative exists.

□ Gradient descent:

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E_P(\mathbf{w}) = \mathbf{w}^{\tau} + \eta \sum_{n \in \mathcal{M}} \mathbf{x}_n t_n$$



The Perceptron Algorithm

17

Probability & Bayesian Inference

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E_P(\mathbf{w}) = \mathbf{w}^t + \eta \sum_{n \in \mathcal{M}} \mathbf{x}_n t_n$$

□ Why does this make sense?

- ▣ If an input from $C_1 (t = +1)$ is misclassified, we need to make its projection on \mathbf{w} more positive.
- ▣ If an input from $C_2 (t = -1)$ is misclassified, we need to make its projection on \mathbf{w} more negative.

The Perceptron Algorithm

18

Probability & Bayesian Inference

□ The algorithm can be implemented sequentially:

□ Repeat until convergence:

■ For each input (\mathbf{x}_n, t_n) :

■ If it is correctly classified, do nothing

■ If it is misclassified, update the weight vector to be

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} + \eta \mathbf{x}_n t_n$$

■ Note that this will lower the contribution of input n to the objective function:

$$-\left(\mathbf{w}^{(\tau)}\right)^t \mathbf{x}_n t_n \rightarrow -\left(\mathbf{w}^{(\tau+1)}\right)^t \mathbf{x}_n t_n = -\left(\mathbf{w}^{(\tau)}\right)^t \mathbf{x}_n t_n - \eta \left(\mathbf{x}_n t_n\right)^t \mathbf{x}_n t_n < -\left(\mathbf{w}^{(\tau)}\right)^t \mathbf{x}_n t_n.$$

Not Monotonic

- While updating with respect to a misclassified input n will lower the error for that input, the error for other misclassified inputs may increase.
- Also, new inputs that had been classified correctly may now be misclassified.
- The result is that the perceptron algorithm is not guaranteed to reduce the total error monotonically at each stage.

The Perceptron Convergence Theorem

20

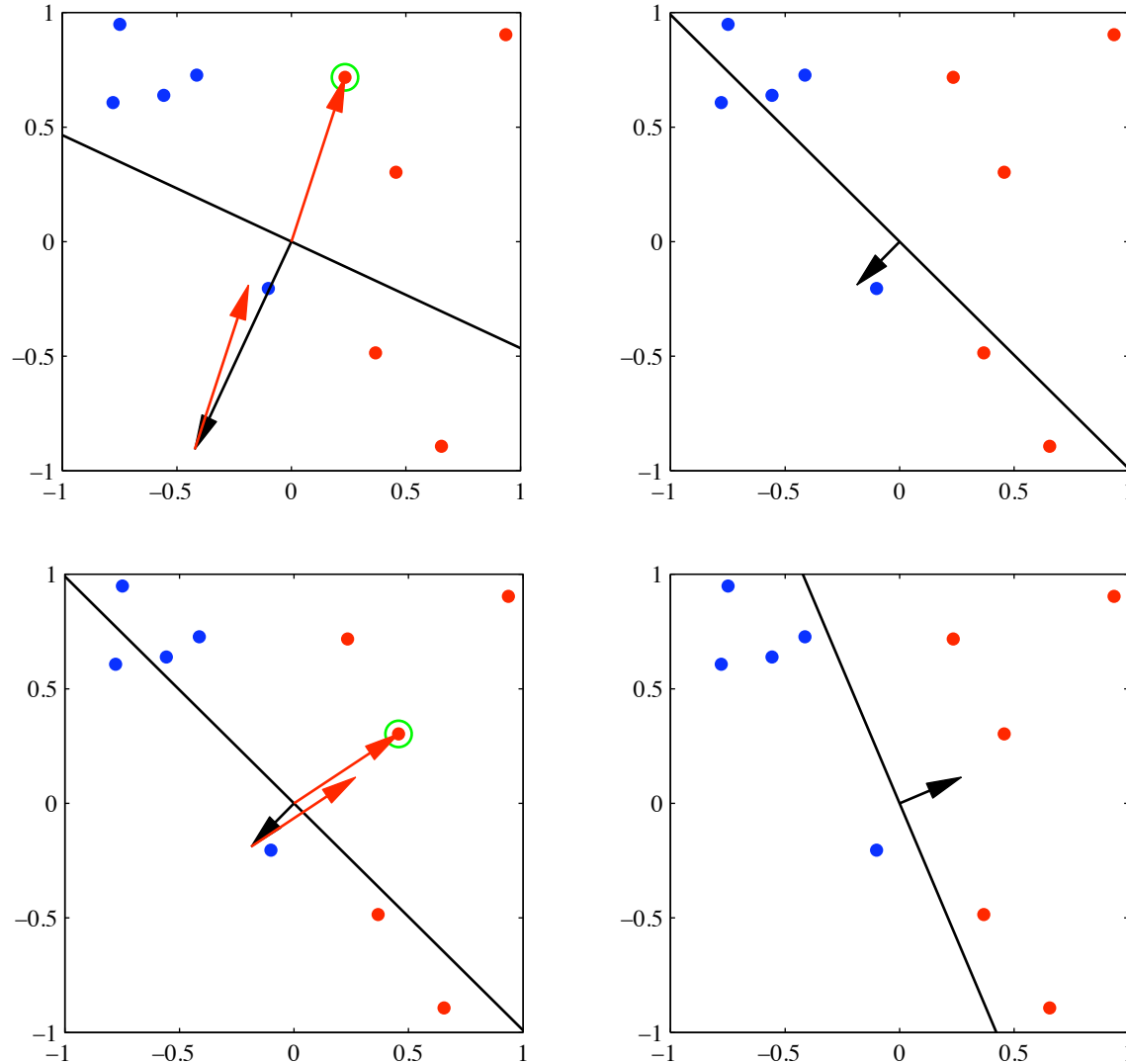
Probability & Bayesian Inference

- Despite this non-monotonicity, **if in fact the data are linearly separable, then the algorithm is guaranteed to find an exact solution in a finite number of steps** (Rosenblatt, 1962).

Example

21

Probability & Bayesian Inference

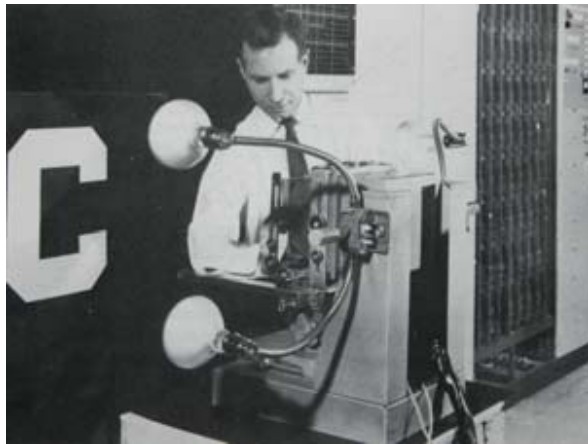


The First Learning Machine

22

Probability & Bayesian Inference

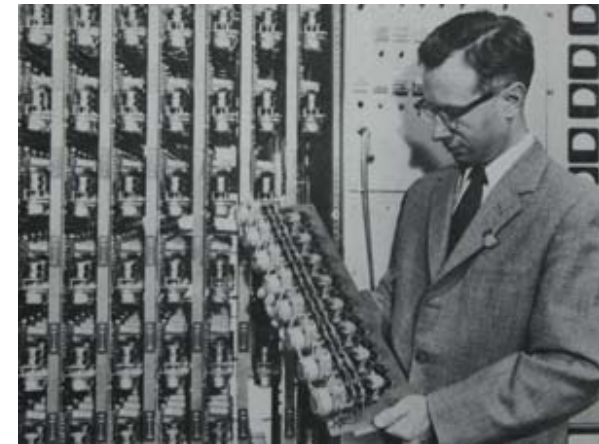
□ Mark 1 Perceptron Hardware (c. 1960)



Visual Inputs



Patch board allowing
configuration of inputs ϕ



Rack of adaptive weights \mathbf{w}
(motor-driven potentiometers)

Practical Limitations

- The Perceptron Convergence Theorem is an important result. However, there are practical limitations:
 - ▣ Convergence may be slow
 - ▣ If the data are not separable, the algorithm will not converge.
 - ▣ We will only know that the data are separable once the algorithm converges.
 - ▣ The solution is in general not unique, and will depend upon initialization, scheduling of input vectors, and the learning rate η .

Generalization to inputs that are not linearly separable.

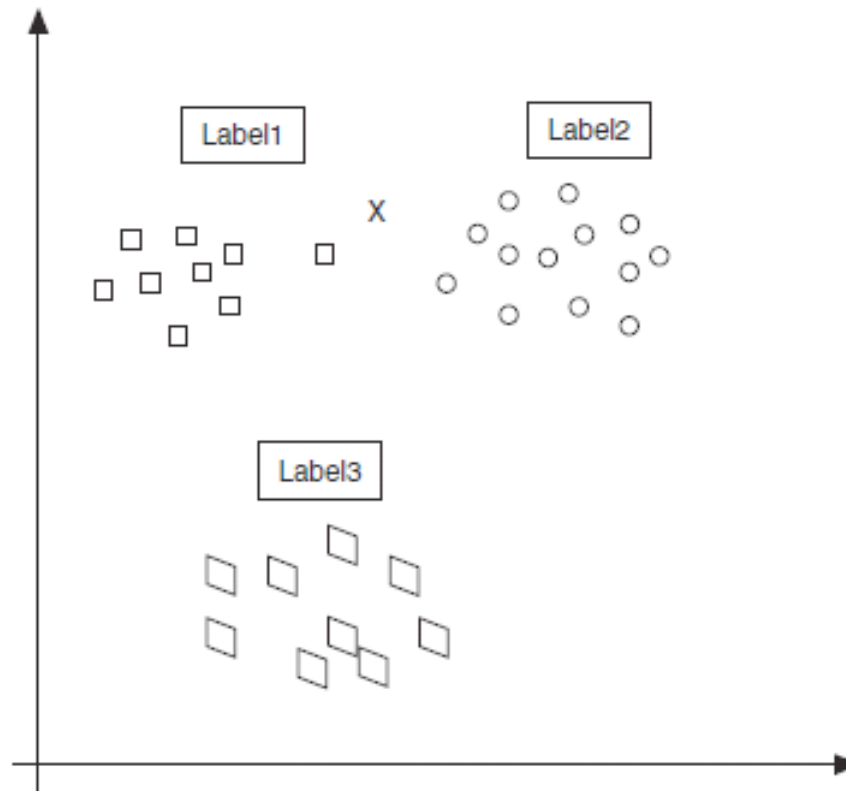
- The single-layer perceptron can be generalized to yield good linear solutions to problems that are not linearly separable.
- Example: The Pocket Algorithm (Gal 1990)
 - ▣ Idea:
 - Run the perceptron algorithm
 - Keep track of the weight vector \mathbf{w}^* that has produced the best classification error achieved so far.
 - It can be shown that \mathbf{w}^* will converge to an optimal solution with probability 1.

Generalization to Multiclass Problems

25

Probability & Bayesian Inference

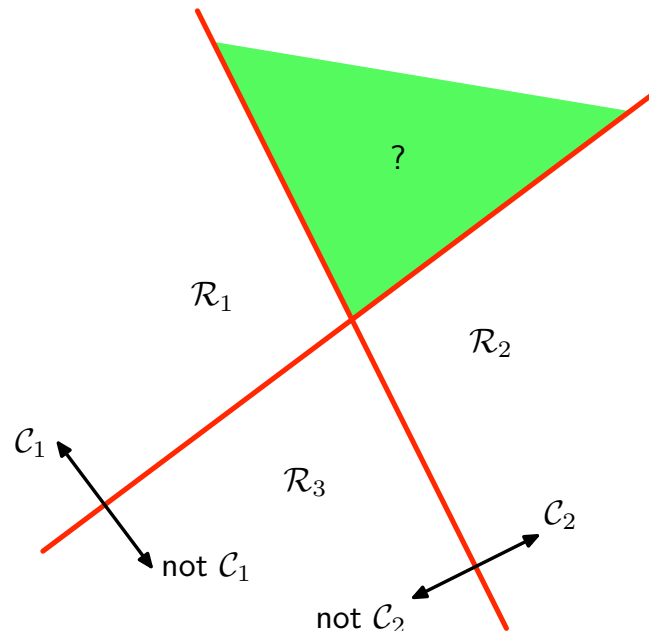
- How can we use perceptrons, or linear classifiers in general, to classify inputs when there are $K > 2$ classes?



$K > 2$ Classes

26

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

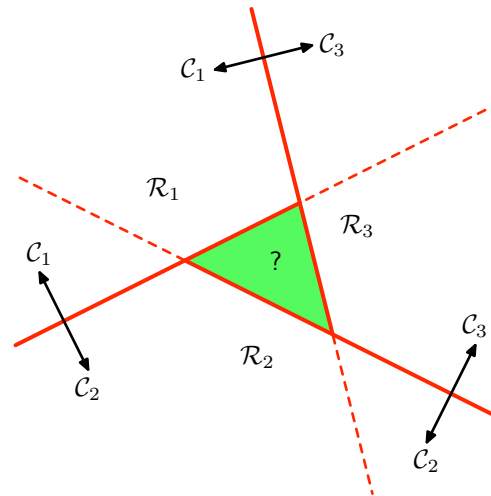


$K > 2$ Classes

27

Probability & Bayesian Inference

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



$K > 2$ Classes

28

Probability & Bayesian Inference

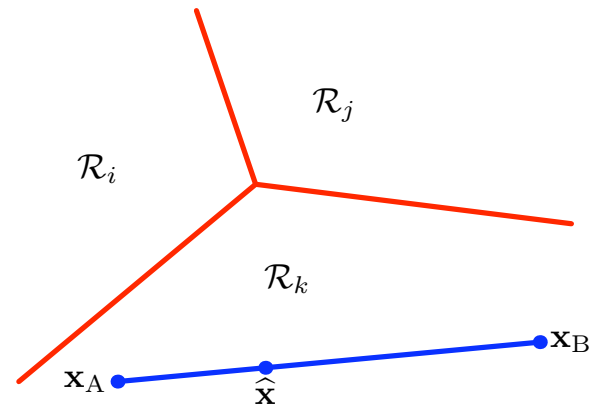
- Idea #3: Use K discriminant functions $y_k(\mathbf{x})$
- Use the **magnitude** of $y_k(\mathbf{x})$, not just the sign.

$$y_k(\mathbf{x}) = \mathbf{w}_k^t \mathbf{x}$$

\mathbf{x} assigned to C_k if $y_k(\mathbf{x}) > y_j(\mathbf{x}) \forall j \neq k$

Decision boundary $y_k(\mathbf{x}) = y_j(\mathbf{x}) \rightarrow (\mathbf{w}_k - \mathbf{w}_j)^t \mathbf{x} + (\mathbf{w}_{k0} - \mathbf{w}_{j0}) = 0$

Results in decision regions that are simply-connected and convex.



Example: Kesler's Construction

29

Probability & Bayesian Inference

- The perceptron algorithm can be generalized to K -class classification problems.
- Example:
 - ▣ Kesler's Construction:
 - Allows use of the perceptron algorithm to simultaneously learn K separate weight vectors \mathbf{w}_i .
 - Inputs are then classified in Class i if and only if
$$\mathbf{w}_i^t \mathbf{x} > \mathbf{w}_j^t \mathbf{x} \quad \forall j \neq i$$
 - The algorithm will converge to an optimal solution if a solution exists, i.e., if all training vectors can be correctly classified according to this rule.

1-of-K Coding Scheme

30

Probability & Bayesian Inference

- When there are $K > 2$ classes, target variables can be coded using the 1-of-K coding scheme:

Input from Class $C_i \Leftrightarrow t = [0 \ 0 \ \dots 1 \ \dots 0 \ 0]^t$



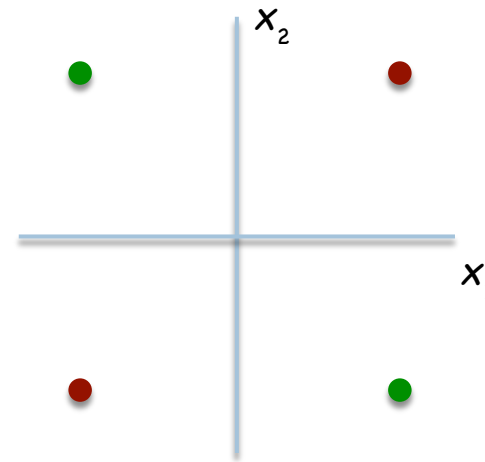
Element i

Computational Limitations of Perceptrons

31

Probability & Bayesian Inference

- Initially, the perceptron was thought to be a potentially powerful learning machine that could model human neural processing.
- However, Minsky & Papert (1969) showed that the single-layer perceptron could not learn a simple XOR function.
- This is just one example of a non-linearly separable pattern that cannot be learned by a single-layer perceptron.



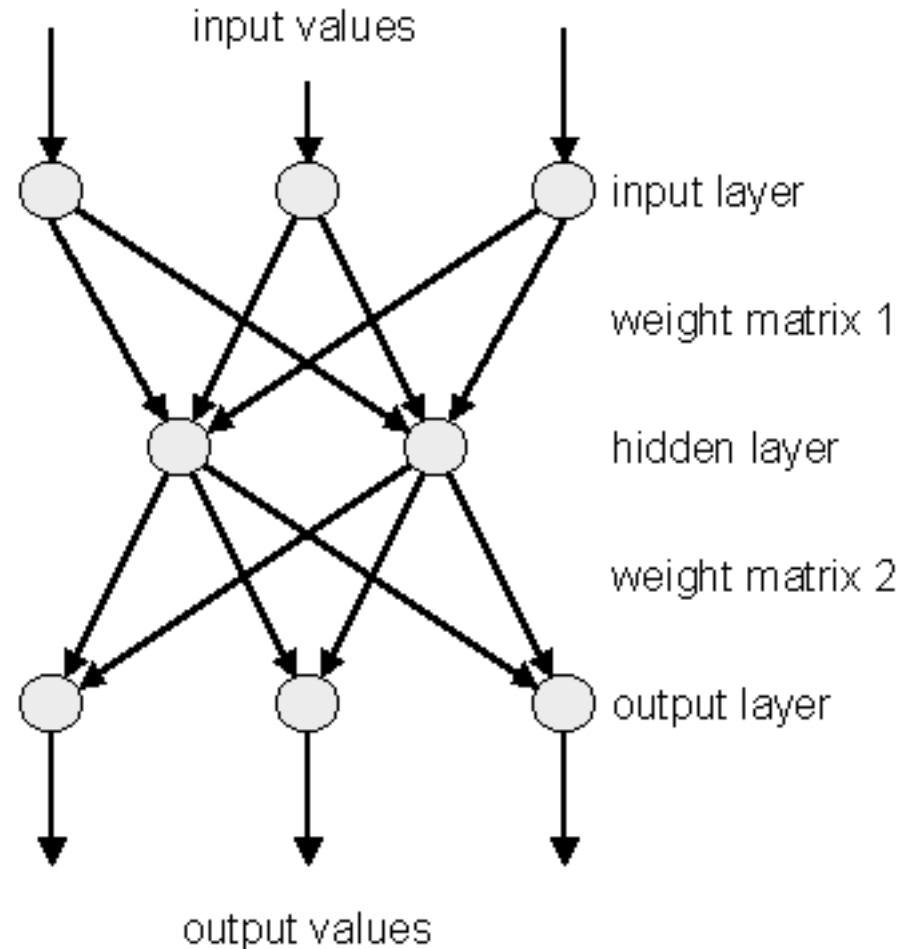
Marvin Minsky (1927 -)

Multi-Layer Perceptrons

32

Probability & Bayesian Inference

- Minsky & Papert's book was widely misinterpreted as showing that artificial neural networks were inherently limited.
- This contributed to a decline in the reputation of neural network research through the 70s and 80s.
- However, their findings apply only to single-layer perceptrons. Multi-layer perceptrons are capable of learning highly nonlinear functions, and are used in many practical applications.





End of Lecture 11

Outline

34

Probability & Bayesian Inference

- The Perceptron Algorithm
- **Least-Squares Classifiers**
- Fisher's Linear Discriminant
- Logistic Classifiers
- Support Vector Machines

Dealing with Non-Linearly Separable Inputs

35

Probability & Bayesian Inference

- The perceptron algorithm fails when the training data are not perfectly linearly separable.
- Let's now turn to methods for learning the parameter vector \mathbf{w} of a perceptron (linear classifier) even when the training data are not linearly separable.

The Least Squares Method

36

Probability & Bayesian Inference

- In the least squares method, we simply fit the (\mathbf{x}, t) observations with a hyperplane $y(\mathbf{x})$.
- Note that this is kind of a weird idea, since the t values are binary (when $K=2$), e.g., 0 or 1.
- However it can work pretty well.

Least Squares: Learning the Parameters

37

Probability & Bayesian Inference

Assume D – dimensional input vectors \mathbf{x} .

For each class $k \in 1 \dots K$:

$$y_k(\mathbf{x}) = \mathbf{w}_k^t \mathbf{x} + w_{k0}$$

$$\rightarrow \mathbf{y}(\mathbf{x}) = \tilde{\mathbf{W}}^t \tilde{\mathbf{x}}$$

where

$$\tilde{\mathbf{x}} = (1, \mathbf{x}^t)^t$$

$\tilde{\mathbf{W}}$ is a $(D+1) \times K$ matrix whose k th column is $\tilde{\mathbf{w}}_k = (w_0, \mathbf{w}_k^t)^t$

Learning the Parameters

38

Probability & Bayesian Inference

□ Method #2: Least Squares

$$\mathbf{y}(\mathbf{x}) = \tilde{\mathbf{W}}^t \tilde{\mathbf{x}}$$

Training dataset $(\mathbf{x}_n, \mathbf{t}_n)$, $n = 1, \dots, N$

where we use the 1-of- K coding scheme for \mathbf{t}_n

Let \mathbf{T} be the $N \times K$ matrix whose n^{th} row is \mathbf{t}_n^t

Let $\tilde{\mathbf{X}}$ be the $N \times (D + 1)$ matrix whose n^{th} row is $\tilde{\mathbf{x}}_n^t$

Let $R_D(\tilde{\mathbf{W}}) = \tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T}$

Then we define the error as $E_D(\tilde{\mathbf{W}}) = \frac{1}{2} \sum_{i,j} R_{ij}^2 = \frac{1}{2} \text{Tr} \left\{ R_D(\tilde{\mathbf{W}})^t R_D(\tilde{\mathbf{W}}) \right\}$

Setting derivative wrt $\tilde{\mathbf{W}}$ to 0 yields:

$$\tilde{\mathbf{W}} = (\tilde{\mathbf{X}}^t \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^t \mathbf{T} = \tilde{\mathbf{X}}^\dagger \mathbf{T}$$

Outline

39

Probability & Bayesian Inference

- The Perceptron Algorithm
- Least-Squares Classifiers
- **Fisher's Linear Discriminant**
- Logistic Classifiers
- Support Vector Machines

Fisher's Linear Discriminant

40

Probability & Bayesian Inference

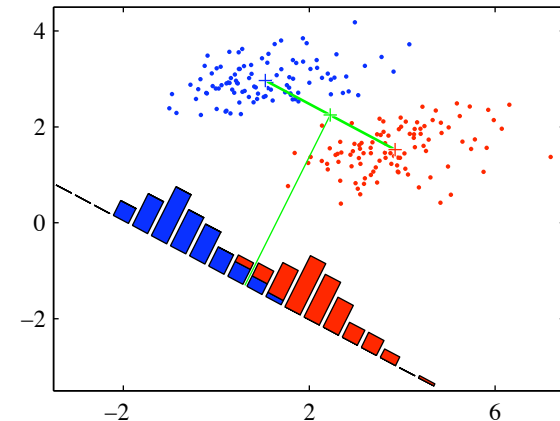
- Another way to view linear discriminants: find the 1D subspace that maximizes the separation between the two classes.

$$\text{Let } \mathbf{m}_1 = \frac{1}{N_1} \sum_{n \in C_1} \mathbf{x}_n, \quad \mathbf{m}_2 = \frac{1}{N_2} \sum_{n \in C_2} \mathbf{x}_n$$

For example, might choose \mathbf{w} to maximize $\mathbf{w}^t (\mathbf{m}_2 - \mathbf{m}_1)$, subject to $\|\mathbf{w}\| = 1$

This leads to $\mathbf{w} \propto \mathbf{m}_2 - \mathbf{m}_1$

However, if conditional distributions are not isotropic, this is typically not optimal.



Fisher's Linear Discriminant

41

Probability & Bayesian Inference

Let $m_1 = \mathbf{w}^t \mathbf{m}_1$, $m_2 = \mathbf{w}^t \mathbf{m}_2$ be the conditional means on the 1D subspace.

Let $s_k^2 = \sum_{n \in C_k} (y_n - m_k)^2$ be the within-class variance on the subspace for class C_k

The Fisher criterion is then $J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}$

This can be rewritten as

$$J(\mathbf{w}) = \frac{\mathbf{w}^t \mathbf{S}_B \mathbf{w}}{\mathbf{w}^t \mathbf{S}_W \mathbf{w}}$$

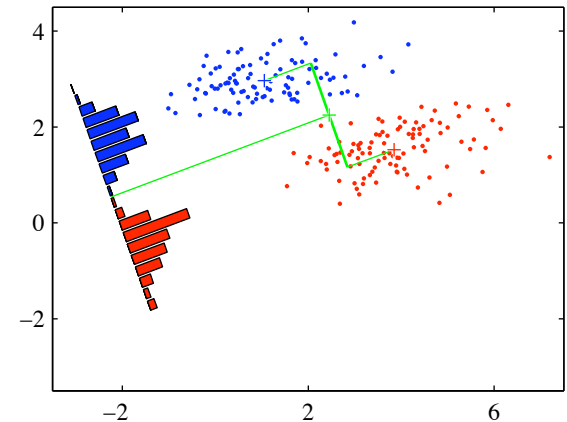
where

$\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^t$ is the between-class variance

and

$\mathbf{S}_W = \sum_{n \in C_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^t + \sum_{n \in C_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^t$ is the within-class variance

$J(\mathbf{w})$ is maximized for $\mathbf{w} \propto \mathbf{S}_W^{-1}(\mathbf{m}_2 - \mathbf{m}_1)$



Connection between Least-Squares and FLD

42

Probability & Bayesian Inference

Change coding scheme used in least-squares method to

$$t_n = \frac{N}{N_1} \text{ for } C_1$$

$$t_n = -\frac{N}{N_2} \text{ for } C_2$$

Then one can show that the ML \mathbf{w} satisfies

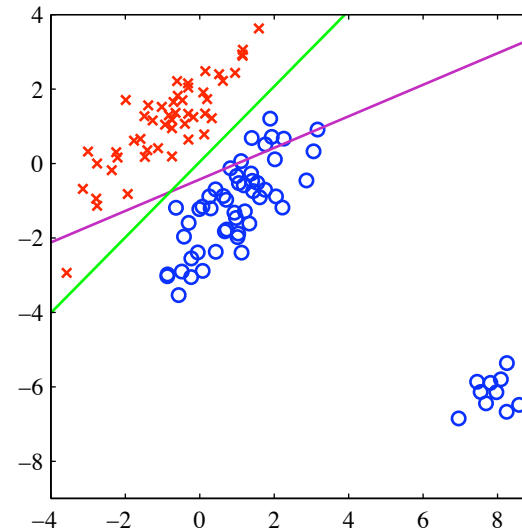
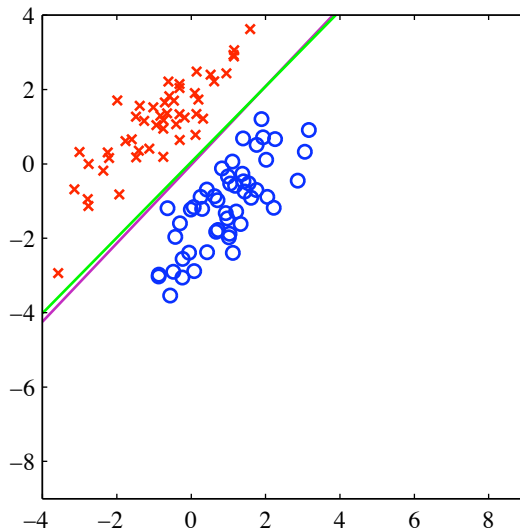
$$\mathbf{w} \propto \mathbf{S}_w^{-1} (\mathbf{m}_2 - \mathbf{m}_1)$$

Least Squares Classifier

43

Probability & Bayesian Inference

□ Problem #1: Sensitivity to outliers

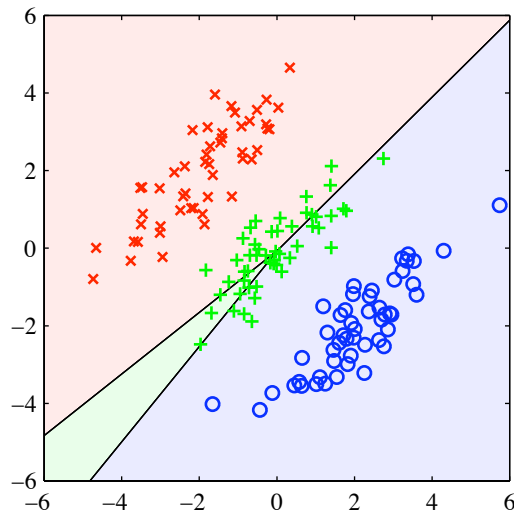


Least Squares Classifier

44

Probability & Bayesian Inference

- Problem #2: Linear activation function is not a good fit to binary data. This can lead to problems.



Outline

45

Probability & Bayesian Inference

- The Perceptron Algorithm
- Least-Squares Classifiers
- Fisher's Linear Discriminant
- **Logistic Classifiers**
- Support Vector Machines

Logistic Regression ($K = 2$)

46

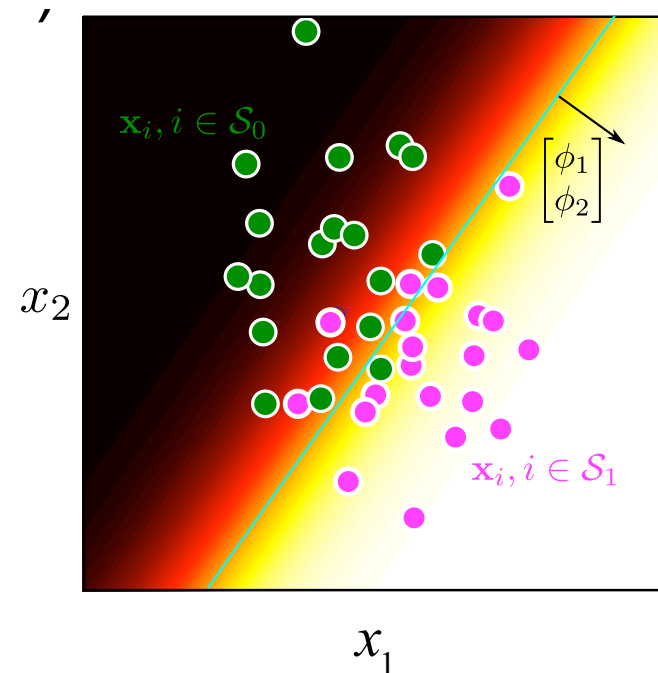
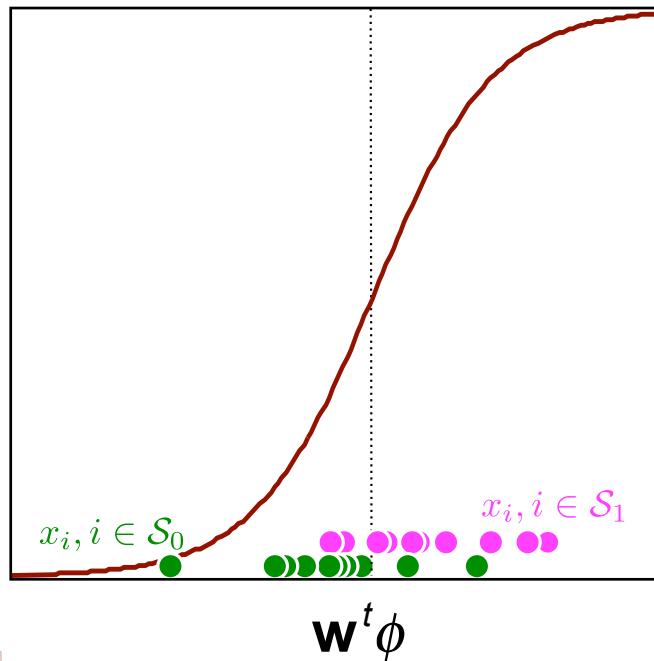
Probability & Bayesian Inference

$$p(C_1 | \phi) = y(\phi) = \sigma(\mathbf{w}^t \phi)$$

$$p(C_2 | \phi) = 1 - p(C_1 | \phi)$$

$$\text{where } \sigma(a) = \frac{1}{1 + \exp(-a)}$$

$$p(C_1 | \phi) = y(\phi) = \sigma(\mathbf{w}^t \phi)$$



Logistic Regression

47

Probability & Bayesian Inference

$$p(C_1 | \phi) = y(\phi) = \sigma(\mathbf{w}^t \phi)$$

$$p(C_2 | \phi) = 1 - p(C_1 | \phi)$$

where

$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$

□ Number of parameters

- ▣ Logistic regression: M

- ▣ Gaussian model: $2M + 2M(M+1)/2 + 1 = M^2 + 3M + 1$

ML for Logistic Regression

48

Probability & Bayesian Inference

$$p(\mathbf{t} | \mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n} \quad \text{where } \mathbf{t} = (t_1, \dots, t_N)^t \text{ and } y_n = p(C_1 | \phi_n)$$

We define the error function to be $E(\mathbf{w}) = -\log p(\mathbf{t} | \mathbf{w})$

Given $y_n = \sigma(a_n)$ and $a_n = \mathbf{w}^t \phi_n$, one can show that

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n$$

Unfortunately, there is no closed form solution for \mathbf{w} .



End of Lecture 12

□ Iterative Reweighted Least Squares

- Although there is no closed form solution for the ML estimate of \mathbf{w} , fortunately, the error function is convex.
- Thus an appropriate iterative method is guaranteed to find the exact solution.
- A good method is to use a local quadratic approximation to the log likelihood function (Newton-Raphson update):

$$\mathbf{w}^{(new)} = \mathbf{w}^{(old)} - \mathbf{H}^{-1} \nabla E(\mathbf{w})$$

where \mathbf{H} is the Hessian matrix of $E(\mathbf{w})$

ML for Logistic Regression

51

Probability & Bayesian Inference

$$\mathbf{w}^{(new)} = \mathbf{w}^{(old)} - \mathbf{H}^{-1} \nabla E(\mathbf{w})$$

where \mathbf{H} is the Hessian matrix of $E(\mathbf{w})$:

$$\mathbf{H} = \Phi^t \mathbf{R} \Phi$$

where \mathbf{R} is the $N \times N$ diagonal weight matrix with $R_{nn} = y_n (1 - y_n)$

(Note that, since $\mathbf{R}_{nn} \geq 0$, \mathbf{R} is positive semi-definite, and hence \mathbf{H} is positive semi-definite

Thus $E(\mathbf{w})$ is convex.)

Thus

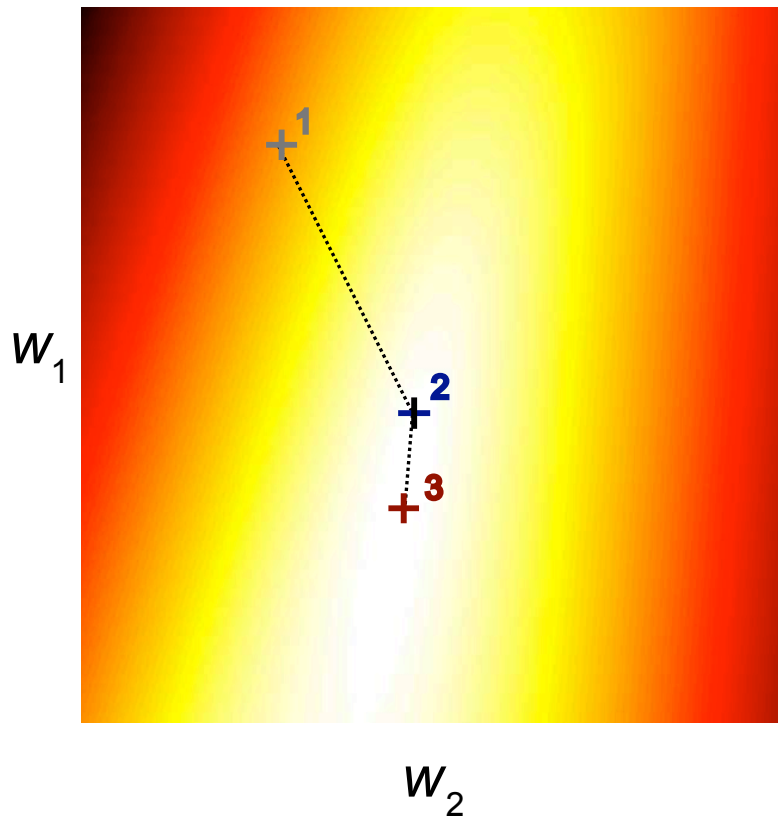
$$\mathbf{w}^{new} = \mathbf{w}^{(old)} - \left(\Phi^t \mathbf{R} \Phi \right)^{-1} \Phi^t (\mathbf{y} - \mathbf{t})$$

ML for Logistic Regression

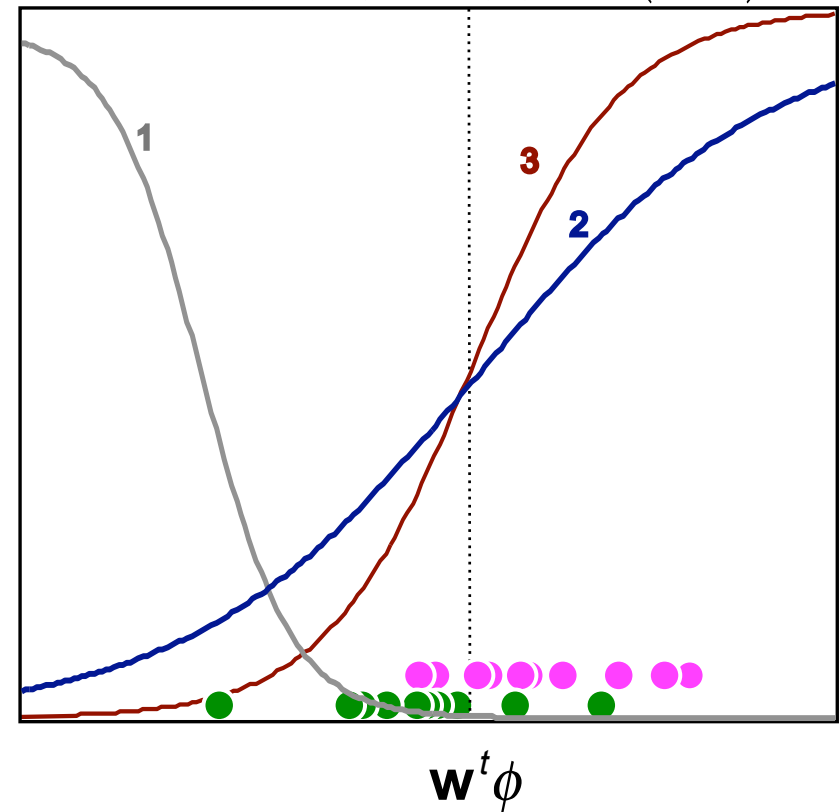
52

Probability & Bayesian Inference

□ Iterative Reweighted Least Squares



$$p(C_1 | \phi) = y(\phi) = \sigma(\mathbf{w}^t \phi)$$



Logistic Regression

53

Probability & Bayesian Inference

- For $K > 2$, we can generalize the activation function by modeling the posterior probabilities as

$$p(C_k | \phi) = y_k(\phi) = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

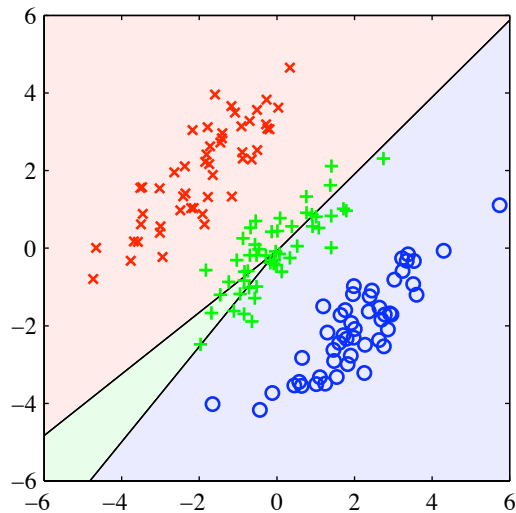
where the activations a_k are given by

$$a_k = \mathbf{w}_k^t \phi$$

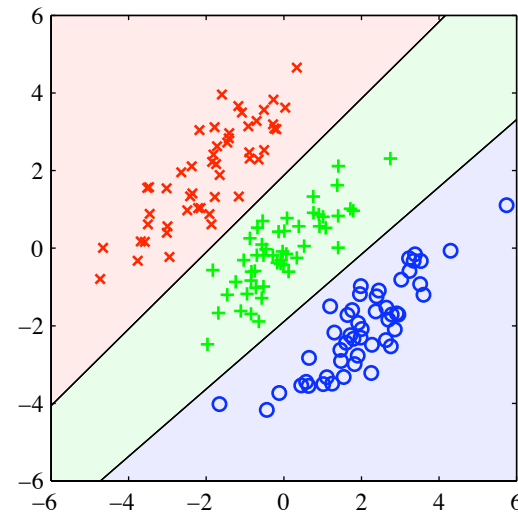
Example

54

Probability & Bayesian Inference



Least-Squares



Logistic

Outline

55

Probability & Bayesian Inference

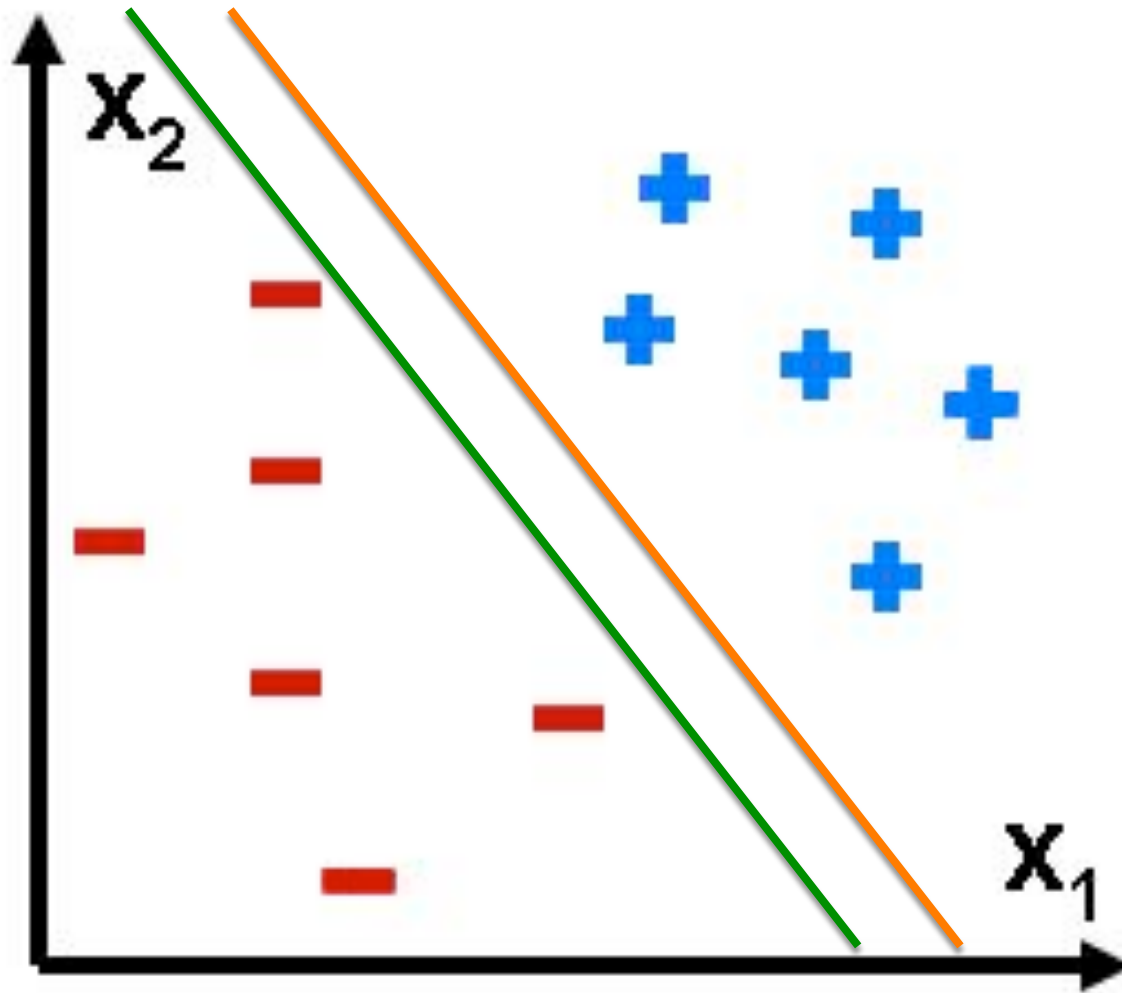
- The Perceptron Algorithm
- Least-Squares Classifiers
- Fisher's Linear Discriminant
- Logistic Classifiers
- **Support Vector Machines**

- The perceptron algorithm is guaranteed to provide a linear decision surface that separates the training data, if one exists.
- However, if the data are linearly separable, there are in general an infinite number of solutions, and the solution returned by the perceptron algorithm depends in a complex way on the initial conditions, the learning rate and the order in which training data are processed.
- While all solutions achieve a perfect score on the training data, they won't all necessarily generalize as well to new inputs.

Which solution would you choose?

57

Probability & Bayesian Inference



The Large Margin Classifier

- Unlike the Perceptron Algorithm, Support Vector Machines solve a problem that has a unique solution: they return the linear classifier with the maximum margin, that is, the hyperplane that separates the data and is farthest from any of the training vectors.
- Why is this good?

Support Vector Machines

59

Probability & Bayesian Inference

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.

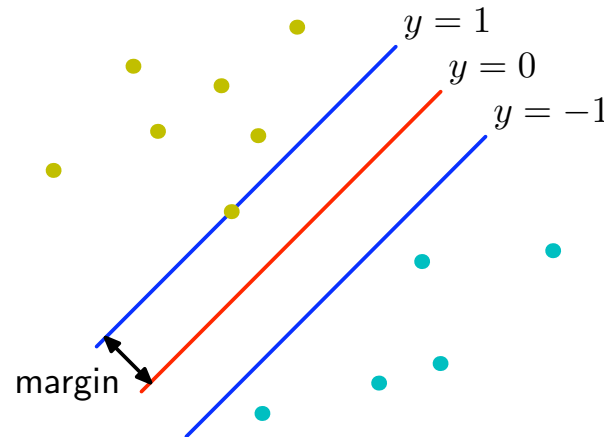
Then $\exists \mathbf{w}, b : t_n y(\mathbf{x}_n) > 0 \quad \forall n \in [1, \dots, N]$

Maximum Margin Classifiers

60

Probability & Bayesian Inference

- 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.



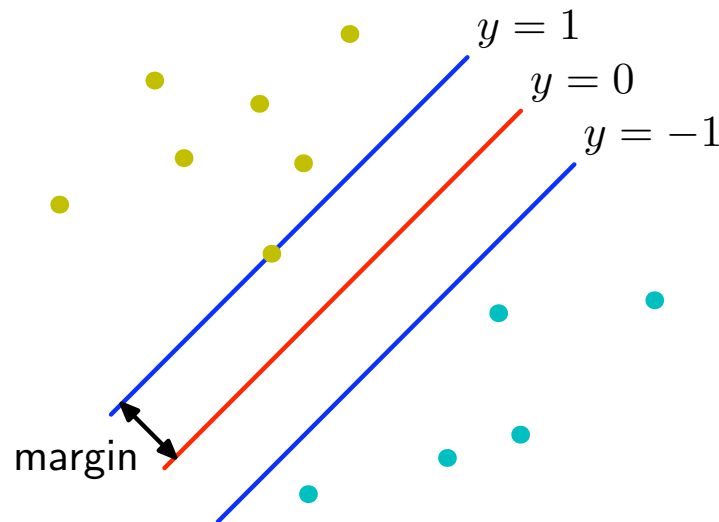
Probabilistic Motivation

61

Probability & Bayesian Inference

- The maximum margin classifier has a probabilistic motivation.

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



Two Class Discriminant Function

62

Probability & Bayesian Inference

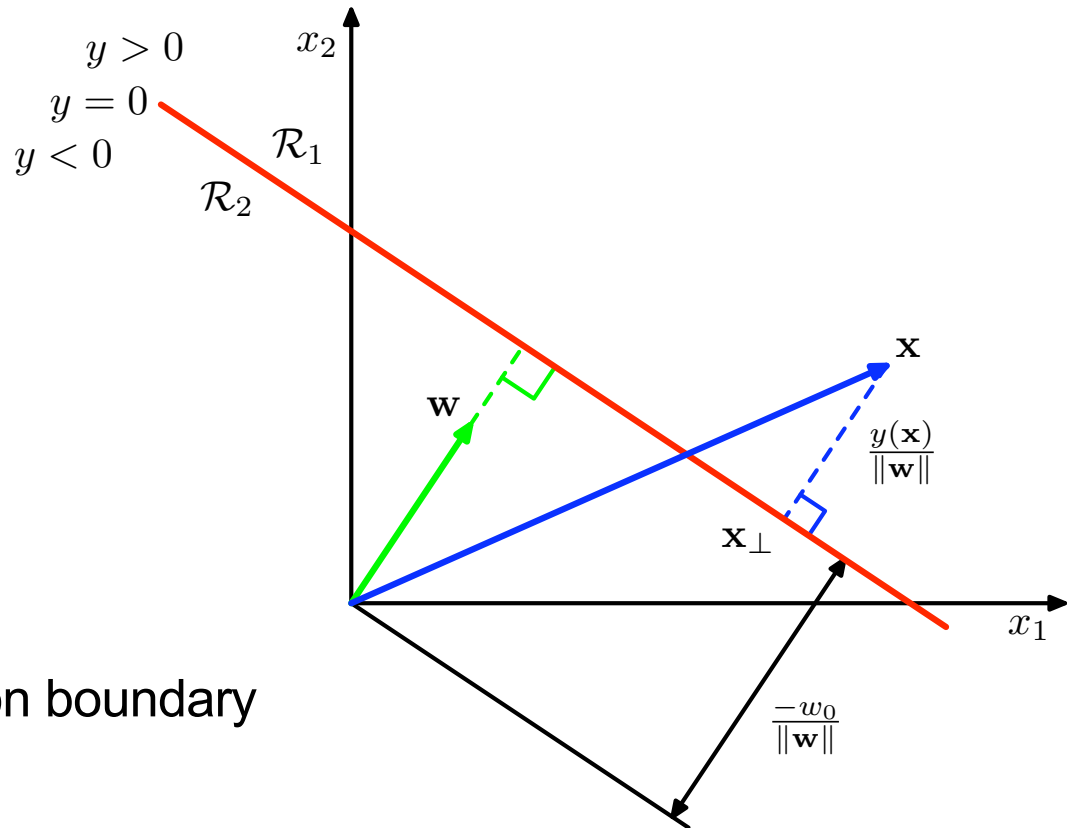
Recall:

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

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

63

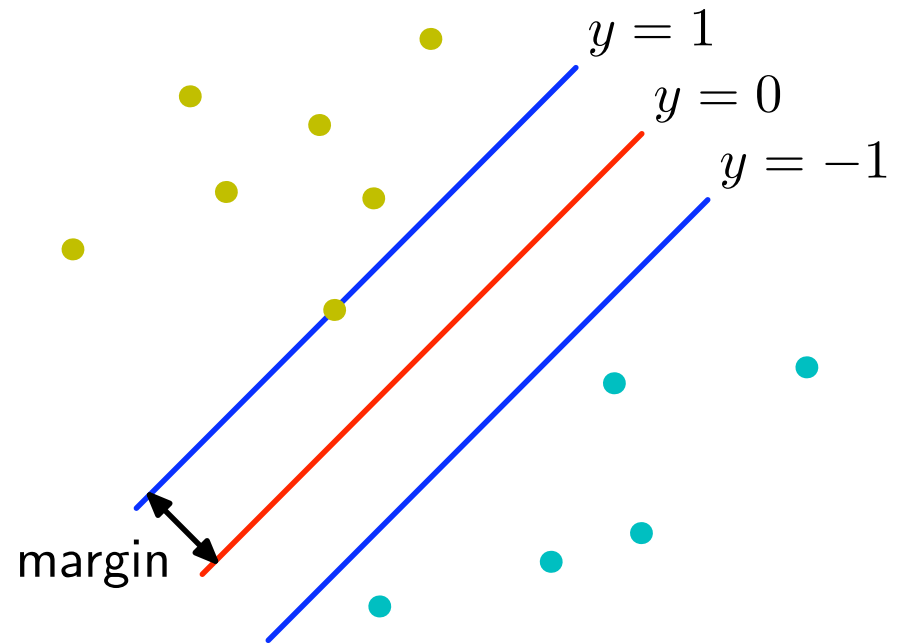
Probability & Bayesian Inference

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:

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



Maximum Margin Classifiers

64

Probability & Bayesian Inference

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

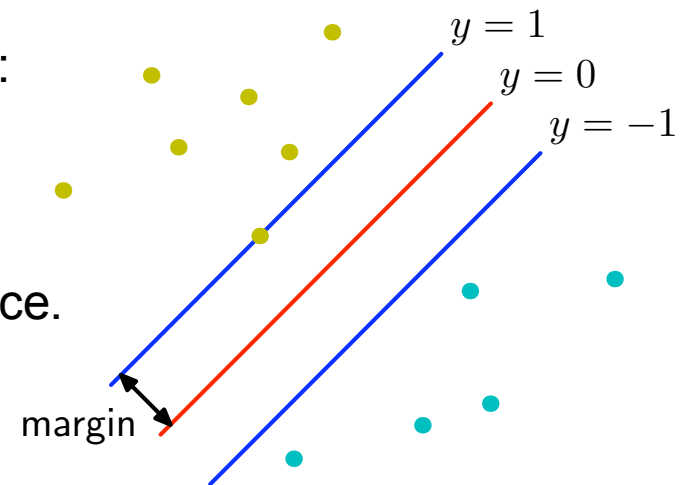
$$\frac{t_n y(\mathbf{x}_n)}{\|\mathbf{w}\|} = \frac{t_n (\mathbf{w}^t \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|}$$

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

65

Probability & Bayesian Inference

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

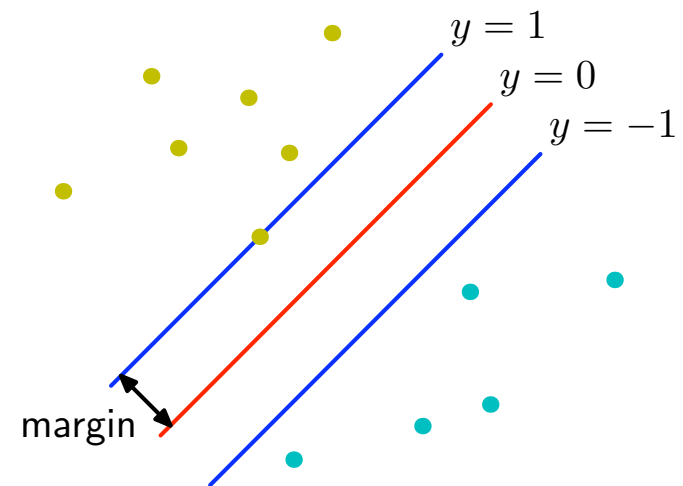
Points for which equality holds are said to be **active**.

All other points are **inactive**.

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

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

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.

Solving this problem will involve **Lagrange multipliers**.



Lagrange Multipliers

Lagrange Multipliers (Appendix C.4 in textbook)

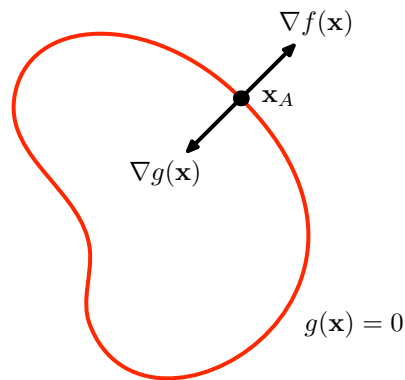
67

Probability & Bayesian Inference

- Used to find stationary points of a function subject to one or more constraints.

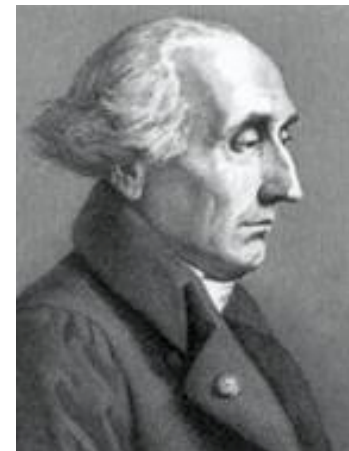
- Example (equality constraint):

Maximize $f(\mathbf{x})$ subject to $g(\mathbf{x}) = 0$.



- Observations:

1. At any point on the constraint surface, $\nabla g(\mathbf{x})$ must be orthogonal to the surface.
2. Let \mathbf{x}^* be a point on the constraint surface where $f(\mathbf{x})$ is maximized. Then $\nabla f(\mathbf{x})$ is also orthogonal to the constraint surface.
3. $\rightarrow \exists \lambda \neq 0$ such that $\nabla f(\mathbf{x}) + \lambda \nabla g(\mathbf{x}) = 0$ at \mathbf{x}^* .
 λ is called a **Lagrange multiplier**.



Joseph-Louis Lagrange
1736-1813

Lagrange Multipliers (Appendix C.4 in textbook)

68

Probability & Bayesian Inference

$\exists \lambda \neq 0$ such that $\nabla f(\mathbf{x}) + \lambda \nabla g(\mathbf{x}) = 0$ at \mathbf{x}^* .

□ Defining the **Lagrangian** function as:

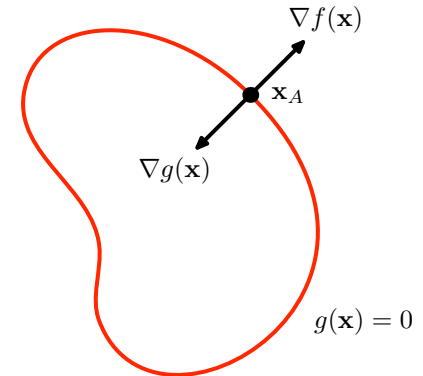
$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})$$

we then have

$$\nabla_{\mathbf{x}} L(\mathbf{x}, \lambda) = 0.$$

and

$$\frac{\partial L(\mathbf{x}, \lambda)}{\partial \lambda} = 0.$$



Example

69

Probability & Bayesian Inference

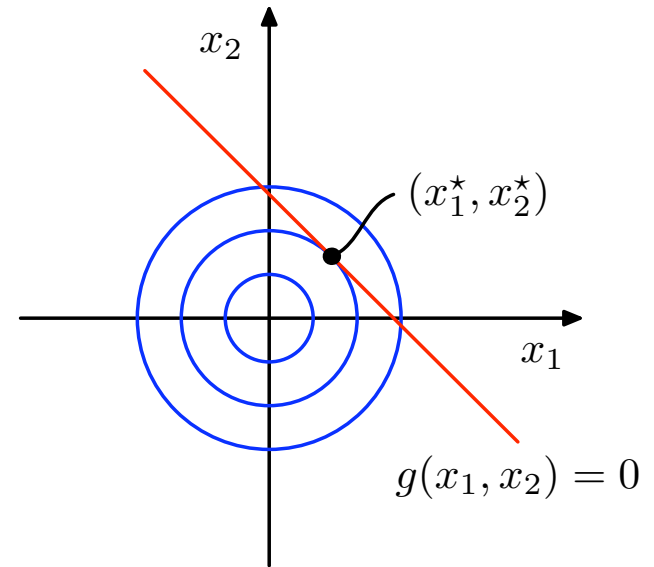
$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})$$

□ Find the stationary point of

$$f(x_1, x_2) = 1 - x_1^2 - x_2^2$$

subject to

$$g(x_1, x_2) = x_1 + x_2 - 1 = 0$$





End of Lecture 13

Inequality Constraints

71

Probability & Bayesian Inference

Maximize $f(\mathbf{x})$ subject to $g(\mathbf{x}) \geq 0$.

□ There are 2 cases:

1. \mathbf{x}^* on the interior (e.g., \mathbf{x}_B)

■ Here $g(\mathbf{x}) > 0$ and the stationary condition is simply

$$\nabla f(\mathbf{x}) = 0.$$

■ This corresponds to a stationary point of the Lagrangian where $\lambda = 0$.

2. \mathbf{x}^* on the boundary (e.g., \mathbf{x}_A)

■ Here $g(\mathbf{x}) = 0$ and the stationary condition is

$$\nabla f(\mathbf{x}) = -\lambda \nabla g(\mathbf{x}), \lambda > 0.$$

■ This corresponds to a stationary point of the Lagrangian where $\lambda > 0$.

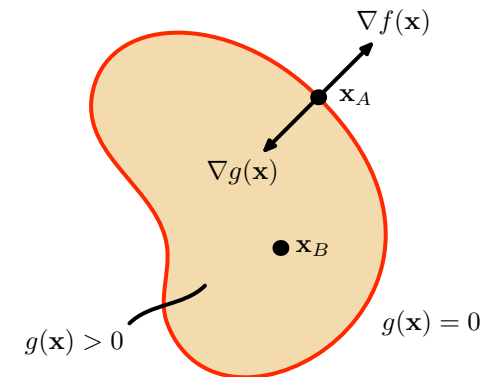
□ Thus the general problem can be expressed as

maximizing the Lagrangian subject to

1. $g(\mathbf{x}) \geq 0$
2. $\lambda \geq 0$
3. $\lambda g(\mathbf{x}) = 0$

}

Karush-Kuhn-Tucker (KKT) conditions



$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})$$

Minimizing vs Maximizing

72

Probability & Bayesian Inference

- If we want to **minimize** $f(\mathbf{x})$ subject to $g(\mathbf{x}) \geq 0$, then the Lagrangian becomes

$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) - \lambda g(\mathbf{x})$$

with $\lambda \geq 0$.

Extension to Multiple Constraints

73

Probability & Bayesian Inference

- Suppose we wish to maximize $f(\mathbf{x})$ subject to

$$g_j(\mathbf{x}) = 0 \text{ for } j = 1, \dots, J$$

$$h_k(\mathbf{x}) \geq 0 \text{ for } k = 1, \dots, K$$

- We then find the stationary points of

$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \sum_{j=1}^J \lambda_j g_j(\mathbf{x}) + \sum_{k=1}^K \mu_k h_k(\mathbf{x})$$

subject to

$$h_k(\mathbf{x}) \geq 0$$

$$\mu_k \geq 0$$

$$\mu_k h_k(\mathbf{x}) = 0$$



Back to our quadratic programming problem...

Quadratic Programming Problem

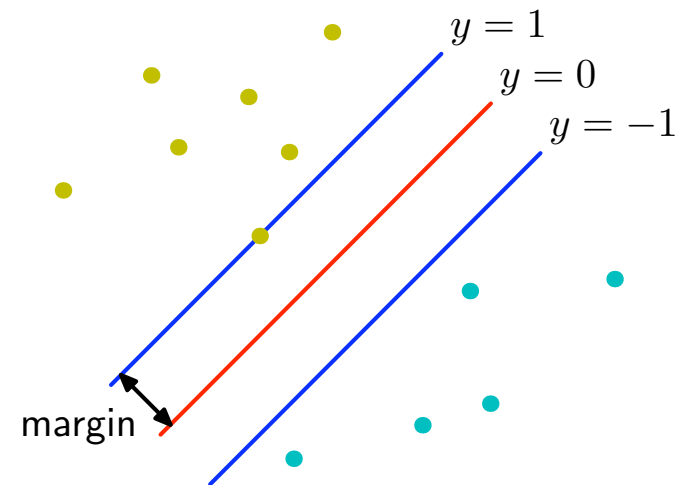
75

Probability & Bayesian Inference

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

Solve using Lagrange multipliers a_n :

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



Dual Representation

76

Probability & Bayesian Inference

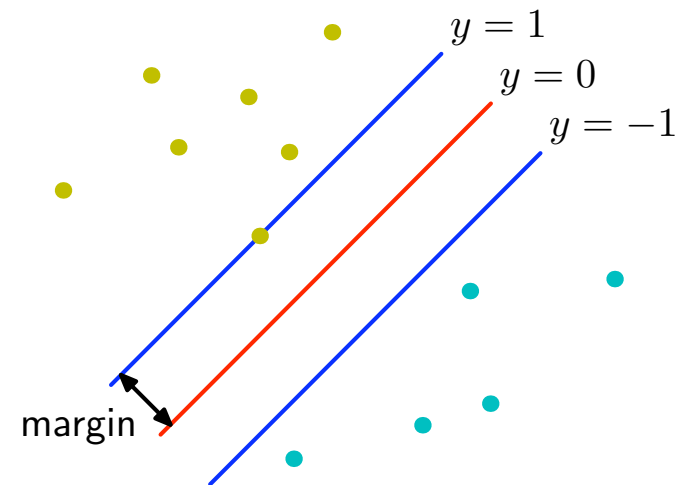
Solve using Lagrange multipliers a_n :

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

Setting derivatives with respect to \mathbf{w} and b to 0, 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

77

Probability & Bayesian Inference

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

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

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

Substituting leads to the dual representation of the maximum margin problem, in which we maximize:

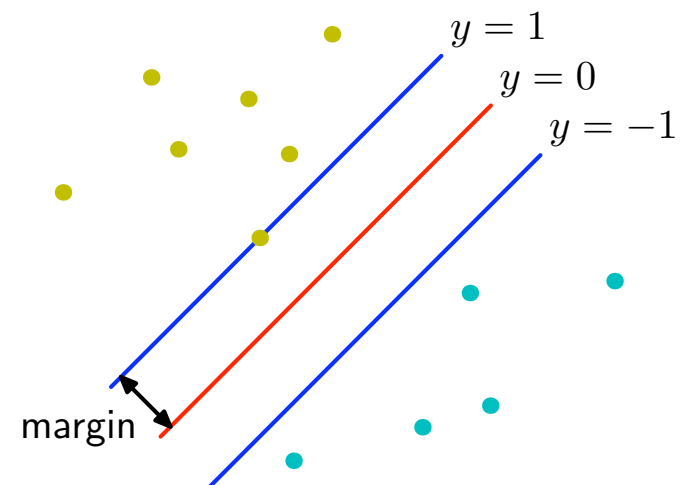
$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

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

78

Probability & Bayesian Inference

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:

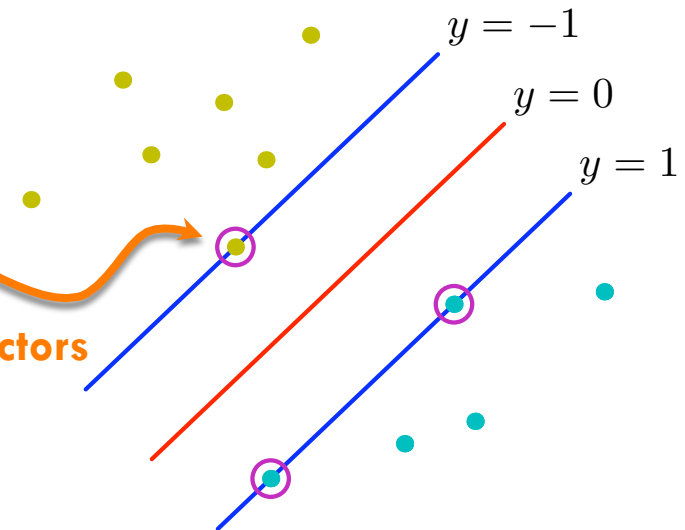
$$a_n \geq 0$$

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

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

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 by noting that any support vector \mathbf{x}_n satisfies $t_n y(\mathbf{x}_n) = 1$.

$$\text{Using } y(\mathbf{x}) = \sum_{n=1}^N a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b$$
$$\text{we have } t_n \left(\sum_{m=1}^N a_m t_m k(\mathbf{x}_n, \mathbf{x}_m) + b \right) = 1$$
$$\text{and so } b = t_n - \sum_{m=1}^N a_m t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

A more numerically accurate solution can be obtained by averaging over all support vectors:

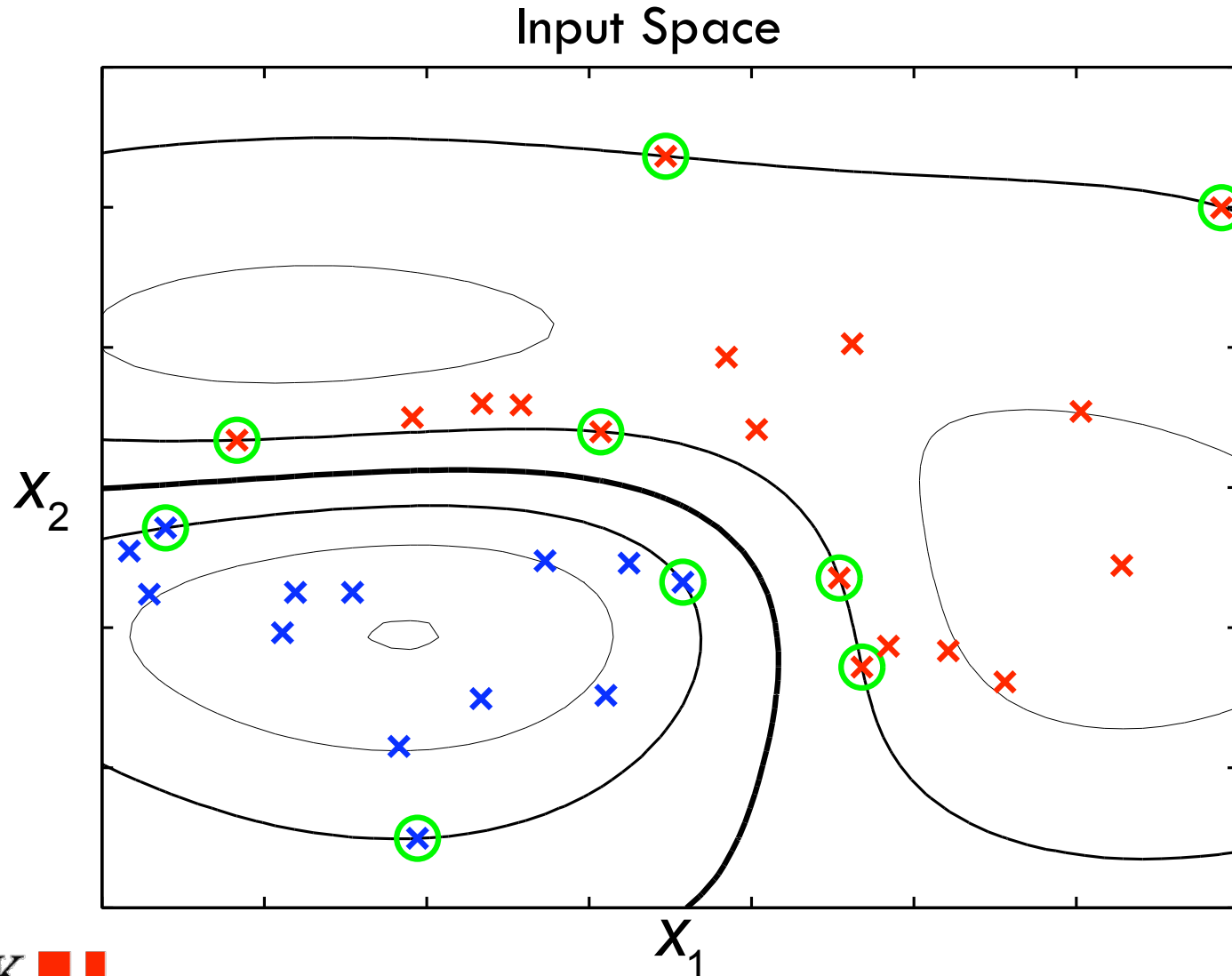
$$b = \frac{1}{N_S} \sum_{n \in S} \left(t_n - \sum_{m \in S} a_m t_m k(\mathbf{x}_n, \mathbf{x}_m) \right)$$

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

Example (Gaussian Kernel)

80

Probability & Bayesian Inference



Overlapping Class Distributions

81

Probability & Bayesian Inference

- The SVM for non-overlapping class distributions is determined by solving

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

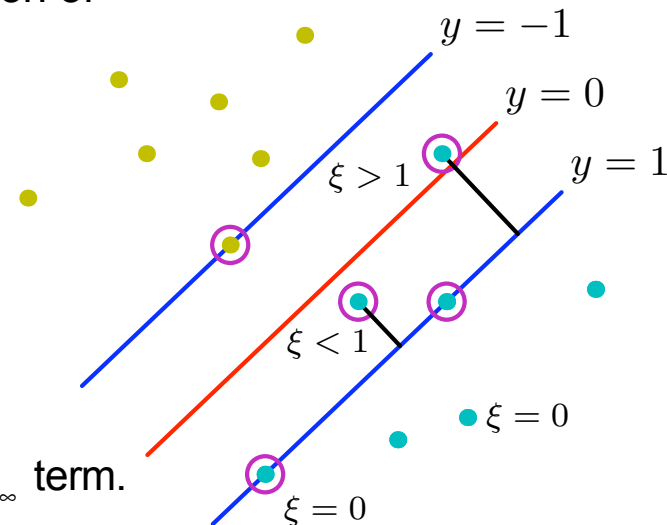
Alternatively, this can be expressed as the minimization of

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

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



Slack Variables

82

Probability & Bayesian Inference

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

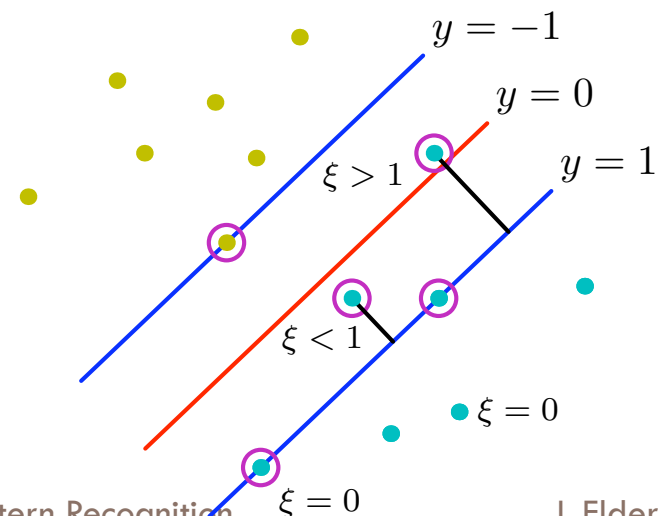
$\xi_n = |t_n - y(\mathbf{x}_n)|$ for all other points.

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

83

Probability & Bayesian Inference

This leads to a dual representation, where we maximize

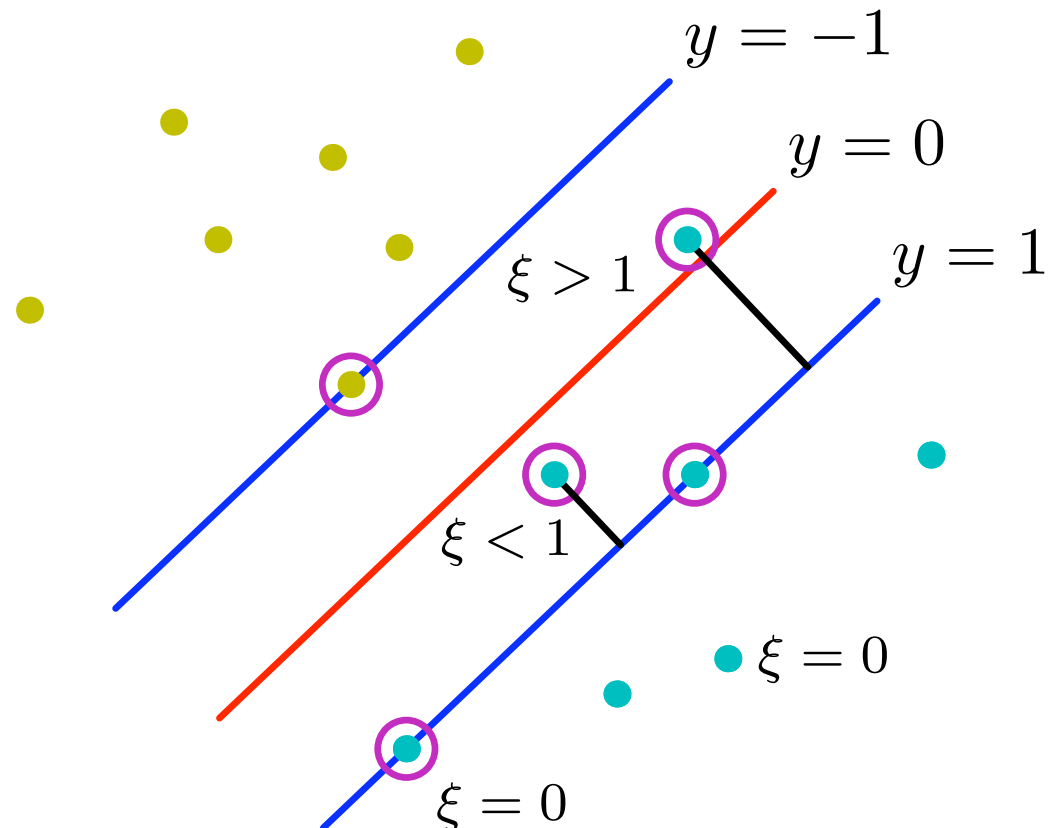
$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

with constraints

$$0 \leq a_n \leq C$$

and

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



Support Vectors

84

Probability & Bayesian Inference

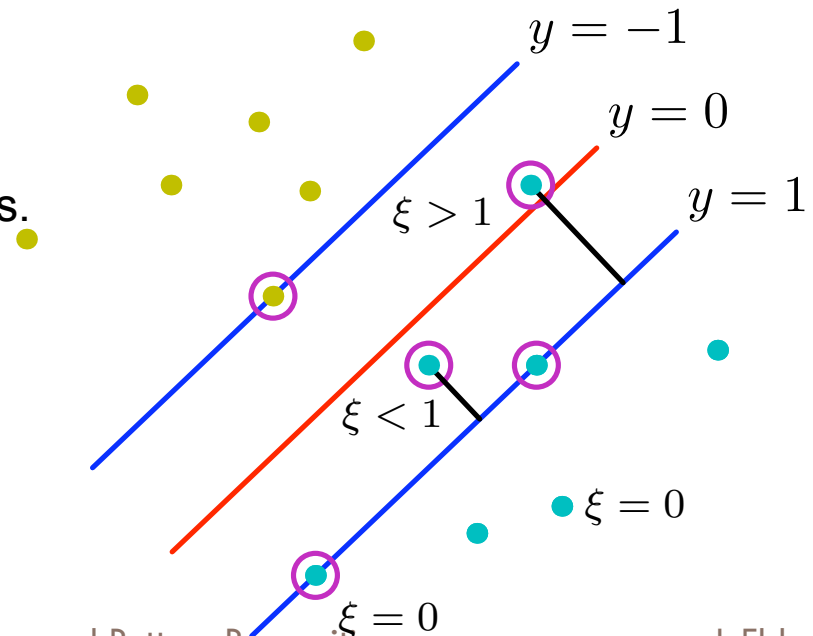
Again, a new point \mathbf{x} is classified by computing

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

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.

In other words, all points that are not on the right side of their margin are support vectors.



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

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

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

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

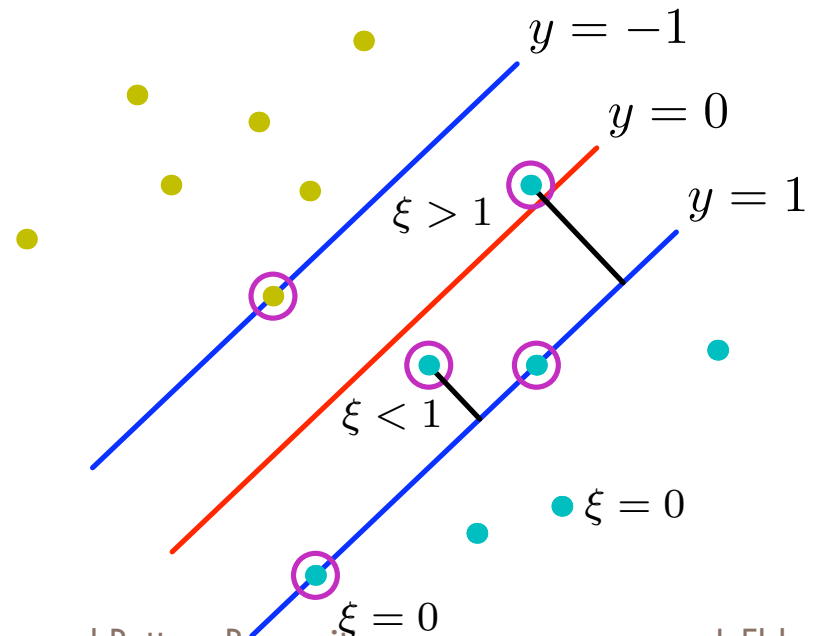
where

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

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

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

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

Chunking

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

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

- Conventional quadratic programming solution requires that matrices with N^2 elements be maintained in memory.

$$K \sim O(N^2), \text{ where } K_{nm} = k(\mathbf{x}_n, \mathbf{x}_m)$$

$$T \sim O(N^2), \text{ where } T_{nm} = t_n t_m$$

$$A \sim O(N^2), \text{ where } A_{nm} = a_n a_m$$

- This becomes infeasible when N exceeds $\sim 10,000$.

Chunking

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

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

- Chunking (Vapnik, 1982) exploits the fact that the value of the Lagrangian is unchanged if we remove the rows and columns of the kernel matrix where $a_n = 0$ or $a_m = 0$.

Chunking

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

$\xi_n = 0$ 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.

□ Chunking (Vapnik, 1982)

1. Select a small number (a 'chunk') of training vectors
2. Solve the QP problem for this subset
3. Retain only the support vectors
4. Consider another chunk of the training data
5. Ignore the subset of vectors in all chunks considered so far that lie on the correct side of the margin, since these do not contribute to the cost function
6. Add the remainder to the current set of support vectors and solve the new QP problem
7. Return to Step 4
8. Repeat until the set of support vectors does not change.

This method reduces memory requirements to $O(N_s^2)$, where N_s is the number of support vectors.

This may still be big!

Decomposition Methods

- It can be shown that the global QP problem is solved when, for all training vectors, satisfy the following optimality conditions:

$$a_i = 0 \Leftrightarrow t_i y(\mathbf{x}_i) \geq 1.$$

$$0 < a_i < C \Leftrightarrow t_i y(\mathbf{x}_i) = 1.$$

$$a_i = C \Leftrightarrow t_i y(\mathbf{x}_i) \leq 1.$$

- Decomposition methods decompose this large QP problem into a series of smaller subproblems.
- Decomposition (Osuna et al, 1997)
 - ▣ Partition the training data into a small working subset B and a fixed subset N.
 - ▣ Minimize the global objective function by adjusting the coefficients in B
 - ▣ Swap 1 or more vectors in B for an equal number in N that fail to satisfy the optimality conditions
 - ▣ Re-solve the global QP problem for B
- Each step is $O(B)^2$ in memory.
- Osuna et al (1997) proved that the objective function decreases on each step and will converge in a finite number of iterations.

Sequential Minimal Optimization

- Sequential Minimal Optimization (Platt 1998) takes decomposition to the limit.
- On each iteration, the working set consists of just two vectors.
- The advantage is that in this case, the QP problem can be solved analytically.
- Memory requirement are $O(N)$.
- Compute time is typically $O(N) - O(N^2)$.

- LIBSVM is a widely used library for SVMs developed by Chang & Lin (2001).
 - ▣ Can be downloaded from www.csie.ntu.edu.tw/~cjlin/libsvm
 - ▣ MATLAB interface
 - ▣ Uses SMO
 - ▣ Will use for Assignment 2.



End of Lecture 14

LIBSVM Example: Face Detection

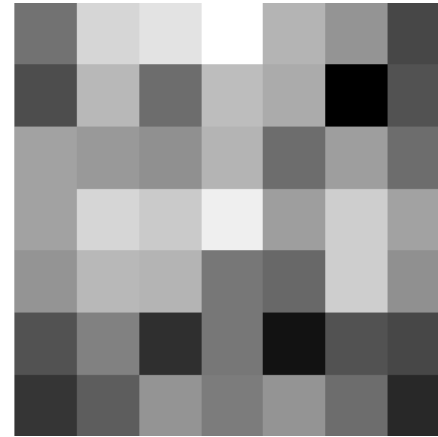
94

Probability & Bayesian Inference

Face



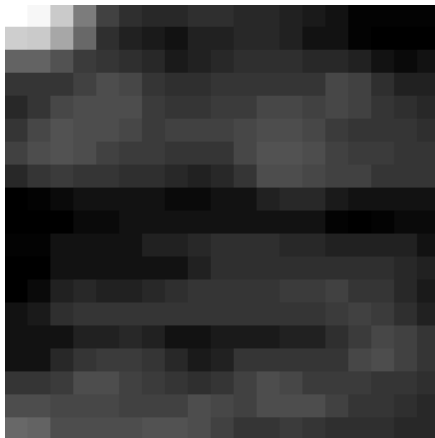
Preprocess:
Subsample &
Normalize



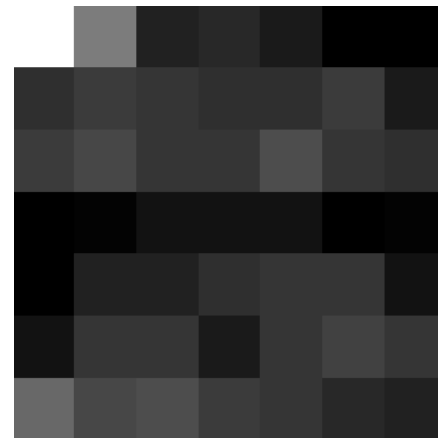
$$\mu = 0, \sigma^2 = 1$$

svmtrain

Non-Face



Preprocess:
Subsample &
Normalize



$$\mu = 0, \sigma^2 = 1$$

LIBSVM Example: MATLAB Interface

95

Probability & Bayesian Inference

Selects linear SVM



```
model=svmtrain(traint, trainx, '-t 0');
```

```
[predicted_label, accuracy, decision_values] = svmpredict(testt, testx, model);
```

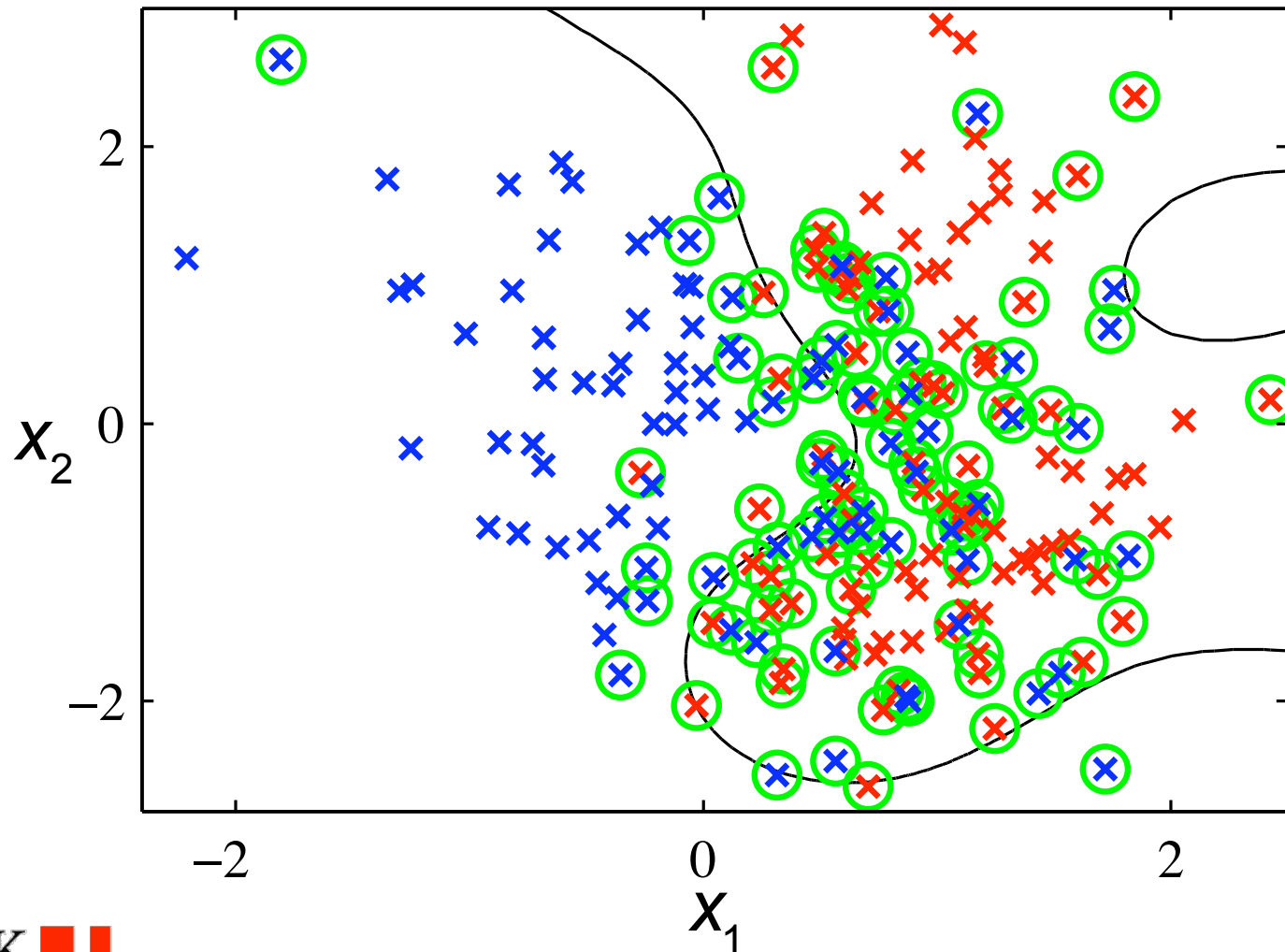
Accuracy = 70.0212% (661/944) (classification)

Example

96

Probability & Bayesian Inference

Input Space



Relation to Logistic Regression

97

Probability & Bayesian Inference

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

