

Linear Models for Classification

Henrik I Christensen

Robotics & Intelligent Machines @ GT
Georgia Institute of Technology,
Atlanta, GA 30332-0280
hic@gatech.gatech.edu

Outline

- 1 Introduction
- 2 Linear Discriminant Functions
- 3 LSQ for Classification
- 4 Fisher's Discriminant Method
- 5 Perceptrons
- 6 Summary

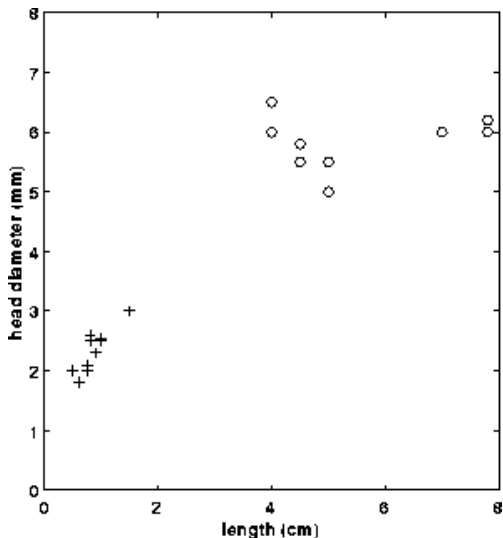
Introduction

- Today: linear classification of data
 - Basic pattern recognition
 - Separation of data: buy/sell
 - Segmentation of line data, ...

Getting Started with Data Analysis

- Important to get started with data analysis / classification
 - Get Wine Dataset
 - Plot Alcohol vs Hue to demonstrate that you can access the data
- Consideration of papers for class discussions (list on the web) - Which papers do you want to bid on.
- Start annotation of data / video data

Simple Example - Bolts or Needles



Classification

- Given
 - An input vector: X
 - A set of classes: $c_i \in \mathcal{C}, \quad i = 1, \dots, k$
- Mapping $m : X \rightarrow \mathcal{C}$
- Separation of space into decision regions
- Boundaries termed decision boundaries/surfaces

Basis Formulation

- It is a 1-of-K coding problem
- Target vector: $\mathbf{t} = (0, \dots, 1, \dots, 0)$
- Consideration of 3 different approaches
 - 1 Optimization of discriminant function
 - 2 Bayesian Formulation: $p(c_i|x)$
 - 3 Learning & Decision fusion

Outline

- 1 Introduction
- 2 Linear Discriminant Functions**
- 3 LSQ for Classification
- 4 Fisher's Discriminant Method
- 5 Perceptrons
- 6 Summary

Discriminant Functions

- Objective: input vector \mathbf{x} assigned to a class c_i
- Simple formulation:

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

- \mathbf{w} is termed a weight vector
- w_0 is termed a bias
- Two class example: c_1 if $y(\mathbf{x}) \geq 0$ otherwise c_2

Basic Design

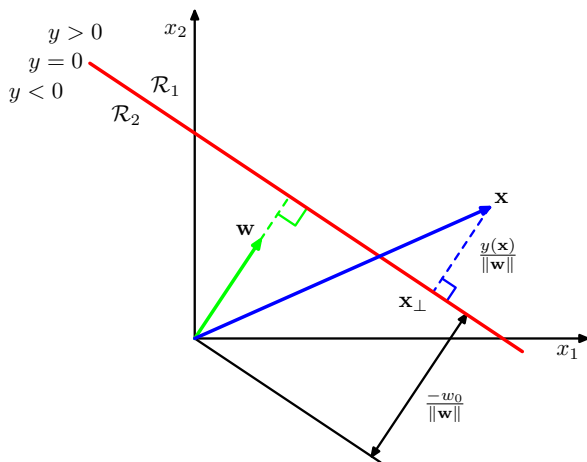
- Two points on decision surface \mathbf{x}_a and \mathbf{x}_b
- $y(\mathbf{x}_a) = y(\mathbf{x}_b) = 0 \Rightarrow \mathbf{w}^T(\mathbf{x}_a - \mathbf{x}_b) = 0$
- \mathbf{w} perpendicular to decision surface

$$\frac{\mathbf{w}^T \mathbf{x}}{\|\mathbf{w}\|} = -\frac{w_0}{\|\mathbf{w}\|}$$

- Define: $\tilde{\mathbf{w}} = (w_0, \mathbf{w})$ and $\tilde{\mathbf{x}} = (1, \mathbf{x})$ so that:

$$y(\mathbf{x}) = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}$$

Linear discriminant function



Multi Class Discrimination

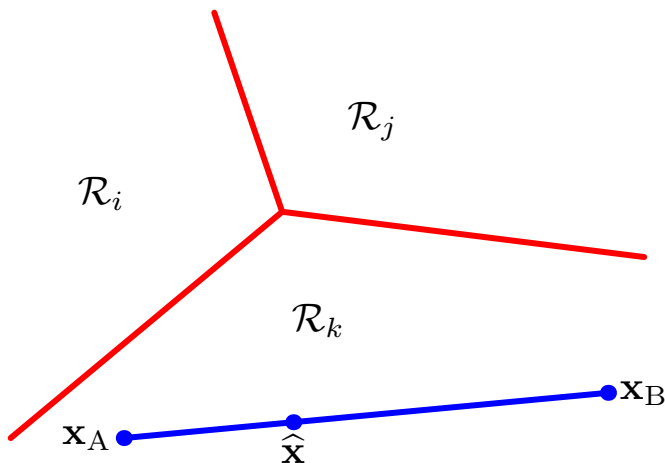
- Generation of multiple decision functions

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

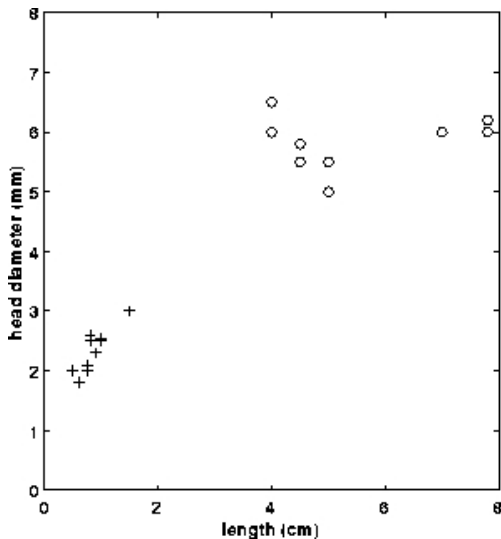
- Decision strategy

$$j = \arg \max_{i \in 1..k} y_i(\mathbf{x})$$

Multi-Class Decision Regions



Example - Bolts or Needles



Minimum distance classification

- Suppose we have computed the mean value for each of the classes
- $m_{needle} = [0.86, 2.34]^T$ and $m_{bolt} = [5.74, 5, 85]^T$
- We can then compute the minimum distance

$$d_j(x) = \|x - m_j\|$$

- $\operatorname{argmin}_i d_i(x)$ is the best fit
- Decision functions can be derived

Bolts / Needle Decision Functions

Needle $d_{needle}(x) = 0.86x_1 + 2.34x_2 - 3.10$

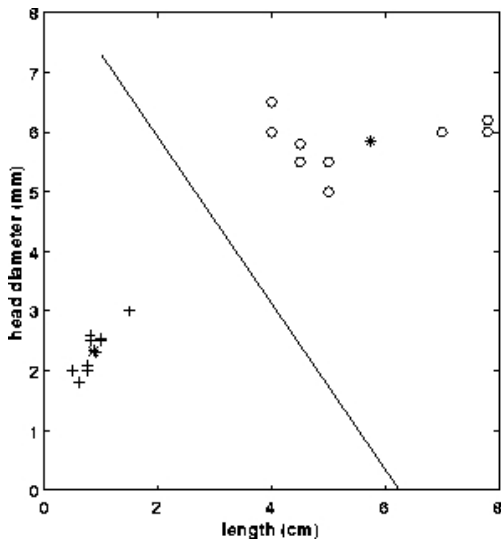
Bolt $d_{bolt}(x) = 5.74x_1 + 5.85x_2 - 33.59$

Decision boundary

$$d_i(x) - d_j(x) = 0$$

$$d_{needle/bolt}(x) = -4.88x_1 - 3.51x_2 + 30.49$$

Example decision surface



Outline

- 1 Introduction
- 2 Linear Discriminant Functions
- 3 LSQ for Classification**
- 4 Fisher's Discriminant Method
- 5 Perceptrons
- 6 Summary

Maximum Likelihood & Least Squares

- Assume observation from a deterministic function contaminated by Gaussian Noise

$$t = y(x, w) + \epsilon \quad p(\epsilon|\beta) = N(\epsilon|0, \beta^{-1})$$

the problem at hand is then

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

- From a series of observations we have the likelihood

$$p(\mathbf{t}|\mathbf{X}, w, \beta) = \prod_{i=1}^N N(t_i|w^T \phi(x_i), \beta^{-1})$$

Maximum Likelihood & Least Squares (2)

- This results in

$$\ln p(\mathbf{t}|\mathbf{w}, \beta) = \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(\mathbf{w})$$

- where

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N \{t_i - \mathbf{w}^T \phi(x_i)\}^2$$

is the sum of squared errors

Maximum Likelihood & Least Squares (3)

- Computing the extrema yields:

$$\mathbf{w}_{ML} = \left(\Phi^T \Phi \right)^{-1} \Phi^T \mathbf{t}$$

- where

$$\Phi = \begin{pmatrix} \phi_0(x_1) & \phi_1(x_1) & \cdots & \phi_{M-1}(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \cdots & \phi_{M-1}(x_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(x_N) & \phi_1(x_N) & \cdots & \phi_{M-1}(x_N) \end{pmatrix}$$

Line Estimation

- Least square minimization:

- Line equation: $y = ax + b$

- Error in fit: $\sum_i (y_i - ax_i - b)^2$

- Solution:

$$\begin{pmatrix} \bar{y}^2 \\ \bar{y} \end{pmatrix} = \begin{pmatrix} \bar{x}^2 & \bar{x} \\ \bar{x} & 1 \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix}$$

- Why might this be a non-robust solution?

LSQ on Lasers

- Line model: $r_i \cos(\phi_i - \theta) = \rho$
- Error model: $d_i = r_i \cos(\phi_i - \theta) - \rho$
- Optimize: $\operatorname{argmin}_{(\rho, \theta)} \sum_i (r_i \cos(\phi_i - \theta) - \rho)^2$
- Error model derived in Deriche *et al.* (1992)
- Well suited for “clean-up” of Hough lines

Total Least Squares

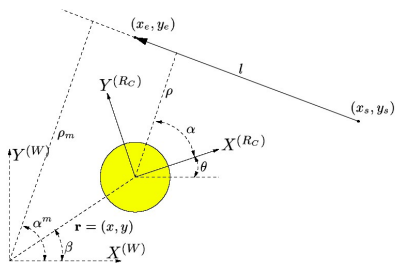
- Line equation: $ax + by + c = 0$
- Error in fit: $\sum_i (ax_i + by_i + c)^2$ where $a^2 + b^2 = 1$.
- Solution:

$$\begin{pmatrix} \bar{x}^2 - \bar{x}\bar{x} & \bar{x}\bar{y} - \bar{x}\bar{y} \\ \bar{x}\bar{y} - \bar{x}\bar{y} & \bar{y}^2 - \bar{y}\bar{y} \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \mu \begin{pmatrix} a \\ b \end{pmatrix}$$

where μ is a scale factor.

- $c = -a\bar{x} - b\bar{y}$

Line Representations



- The line representation is crucial
- Often a redundant model is adopted
- Line parameters vs end-points
- Important for fusion of segments.
- End-points are less stable

Sequential Adaptation

- In some cases one at a time estimation is more suitable
- Also known as gradient descent

$$\begin{aligned}\mathbf{w}^{(\tau+1)} &= \mathbf{w}^{(\tau)} - \eta \nabla E_n \\ &= \mathbf{w}^{(\tau)} - \eta (t_n - \mathbf{w}^{(\tau)T} \phi(x_n)) \phi(x_n)\end{aligned}$$

- Known as least-mean square (LMS). An issue is how to choose η ?

Regularized Least Squares

- As seen in lecture 2 sometime control of parameters might be useful.
- Consider the error function:

$$E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})$$

- which generates

$$\frac{1}{2} \sum_{i=1}^N \{t_i - \mathbf{w}^T \phi(x_i)\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

- which is minimized by

$$\mathbf{w} = \left(\lambda I + \Phi^T \Phi \right)^{-1} \Phi^T \mathbf{t}$$

Least Squares for Classification

- We could do LSQ for regression and we can perform an approximation to the classification vector \mathcal{C}
- Consider:

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

- Rewrite to

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

- Assuming we have a target vector \mathbf{T}

Least Squares for Classification

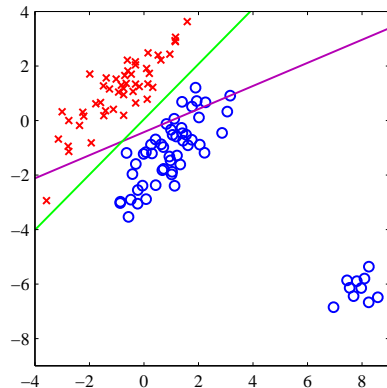
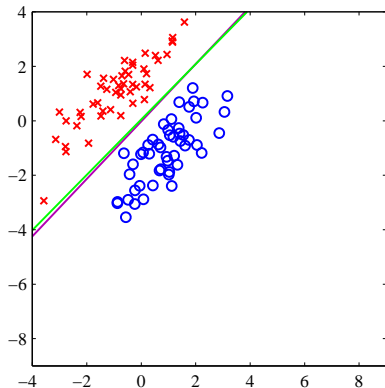
- The error is then:

$$E_D(\tilde{\mathbf{W}}) = \frac{1}{2} \text{Tr} \left\{ (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T})^T (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T}) \right\}$$

- The solution is then

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

LSQ and Outliers



Outline

- 1 Introduction
- 2 Linear Discriminant Functions
- 3 LSQ for Classification
- 4 Fisher's Discriminant Method**
- 5 Perceptrons
- 6 Summary

Fisher's linear discriminant

- Selection of a decision function that maximizes distance between classes
- Assume for a start

$$y = \mathbf{W}^T \mathbf{x}$$

- Compute m_1 and m_2

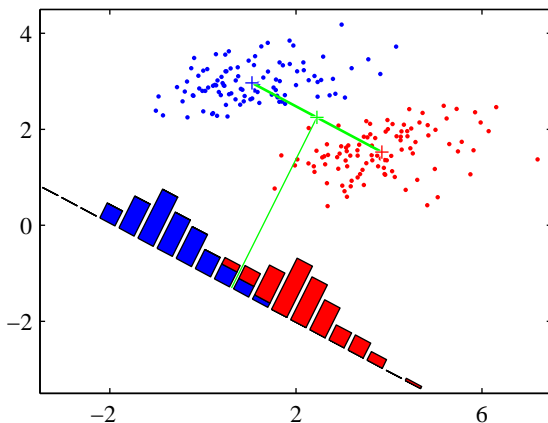
$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{i \in C_1} \mathbf{x}_i \quad \mathbf{m}_2 = \frac{1}{N_2} \sum_{j \in C_2} \mathbf{x}_j$$

- Distance:

$$m_2 - m_1 = \mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1)$$

- where $m_i = \mathbf{w} \mathbf{m}_i$

The suboptimal solution



The Fisher criterion

- Consider the expression

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

- where \mathbf{S}_B is the between class covariance and \mathbf{S}_W is the within class covariance, i.e.

$$\mathbf{S}_B = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T$$

and

$$\mathbf{S}_W = \sum_{i \in \mathcal{C}_1} (\mathbf{x}_i - \mathbf{m}_1)(\mathbf{x}_i - \mathbf{m}_1)^T + \sum_{i \in \mathcal{C}_2} (\mathbf{x}_i - \mathbf{m}_2)(\mathbf{x}_i - \mathbf{m}_2)^T$$

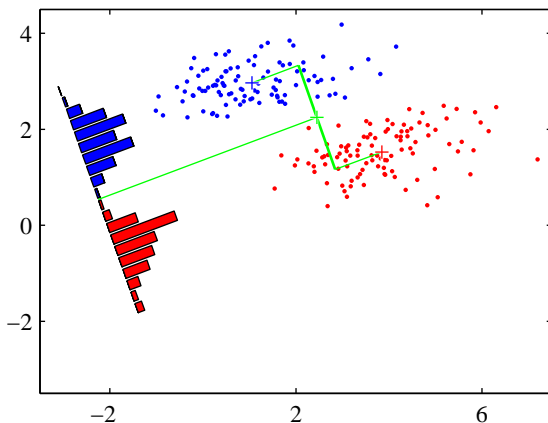
- Optimized when

$$(\mathbf{w}^T \mathbf{S}_B \mathbf{w}) \mathbf{S}_W \mathbf{w} = (\mathbf{w}^T \mathbf{S}_W \mathbf{w}) \mathbf{S}_B \mathbf{w}$$

or

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

The Fisher result



Generalization to $N > 2$

- Define a stacked weight factor

$$\mathbf{y} = \mathbf{W}^T \mathbf{x}$$

- The within class covariance generalizes to

$$\mathbf{S}_w = \sum_{k=1}^K \mathbf{S}_k$$

- The between class covariance is

$$\mathbf{S}_B = \sum_{k=1}^K N_k (\mathbf{m}_k - \mathbf{m})(\mathbf{m}_k - \mathbf{m})^T$$

- It can be shown that $J(\mathbf{w})$ is optimized by the eigenvectors to the equation

$$S = \mathbf{S}_W^{-1} \mathbf{S}_B$$

Outline

- 1 Introduction
- 2 Linear Discriminant Functions
- 3 LSQ for Classification
- 4 Fisher's Discriminant Method
- 5 Perceptrons**
- 6 Summary

Perceptron Algorithm

- Developed by Rosenblatt (1962)
- Formed an important basis for neural networks
- Use a non-linear transformation $\phi(\mathbf{x})$
- Construct a decision function

$$y(\mathbf{x}) = f\left(\mathbf{w}^T \phi(\mathbf{x})\right)$$

- where

$$f(a) = \begin{cases} +1, & a \geq 0 \\ -1, & a < 0 \end{cases}$$

The perceptron criterion

- Normally we want

$$\mathbf{w}^t \phi(\mathbf{x}_n) > 0$$

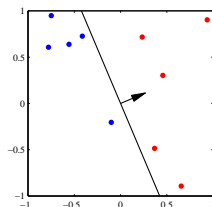
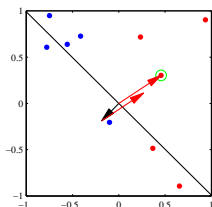
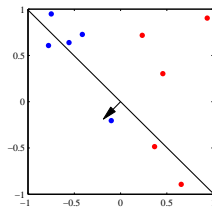
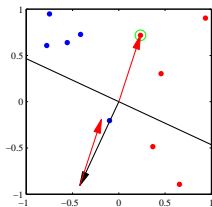
- Given the target vector definition

$$E_P(\mathbf{w}) = - \sum_{n \in \mathcal{M}} \mathbf{w}^T \phi_n t_n$$

- Where \mathcal{M} represents all the mis-classified samples
- We can use gradient descent

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_P(\mathbf{w}) = \mathbf{w}^{(\tau)} + \eta \phi_n t_n$$

Perceptron learning example



Outline

- 1 Introduction
- 2 Linear Discriminant Functions
- 3 LSQ for Classification
- 4 Fisher's Discriminant Method
- 5 Perceptrons
- 6 Summary**

Summary

- Basics for discrimination / classification
- Obviously not all problems are linear
- Optimization of the distance/overlap between classes
 - Minimizing the probability of error classification
- Basic formulation as an optimization problem
- How to optimize between cluster distance? Covariance Weighted
- Basic recursive formulation
- Could we make it more robust?

Deriche, R., Vaillant, R., & Faugeras, O. 1992. *From Noisy Edges Points to 3D Reconstruction of a Scene : A Robust Approach and Its Uncertainty Analysis*. Vol. 2. World Scientific. Series in Machine Perception and Artificial Intelligence. Pages 71–79.