

# Linear Models for Classification: Overview

Sargur N. Srihari

University at Buffalo, State University of New York  
USA

# Topics in Linear Models for Classification

- Overview

1. Discriminant Functions

2. Probabilistic Generative Models

3. Probabilistic Discriminative Models

4. The Laplace Approximation

# Topics in Overview

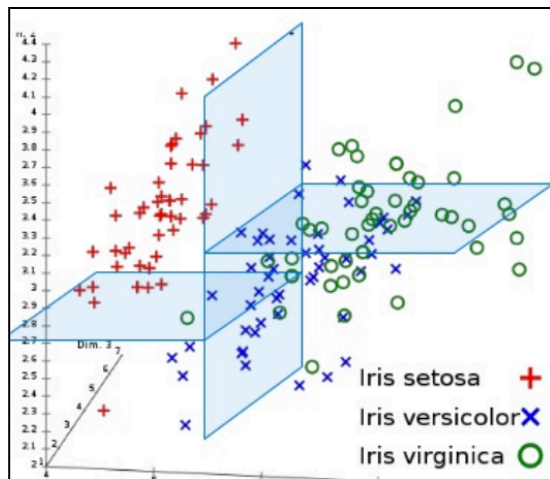
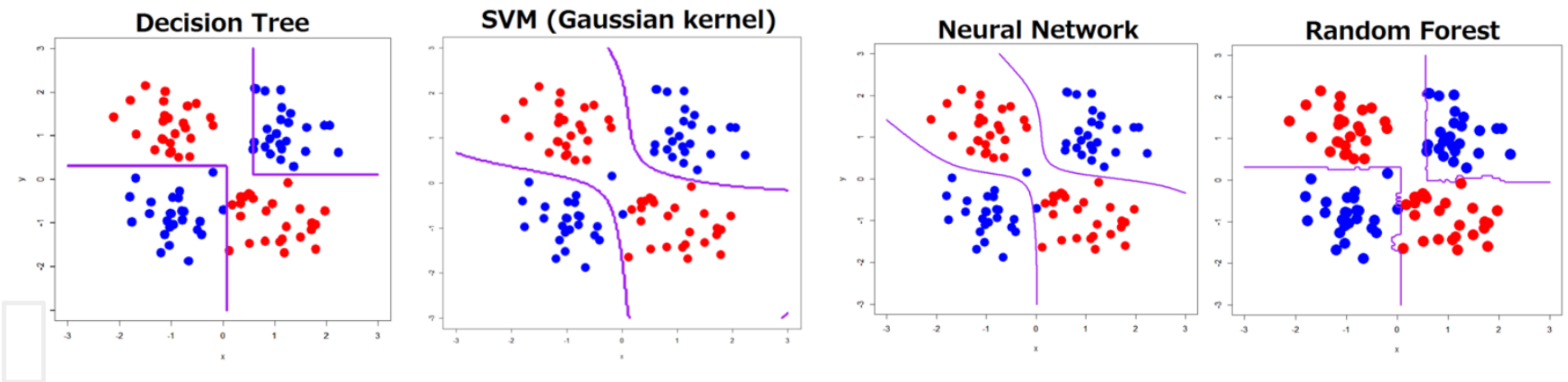
1. Regression vs Classification
2. Linear Classification Models
3. Converting probabilistic regression output to classification output
4. Three classes of classification models

# Regression vs Classification

- In *Regression* we assign input vector  $x$  to one or more continuous target variables  $t$ 
  - Linear regression has simple analytical and computational properties
- In *Classification* we assign input vector  $x$  to one of  $K$  discrete classes  $C_k, k = 1, \dots, K$ 
  - Common classification scenario: classes considered disjoint
    - Each input assigned to only one class
  - Input space is thereby divided into decision regions

# Boundaries of decision regions

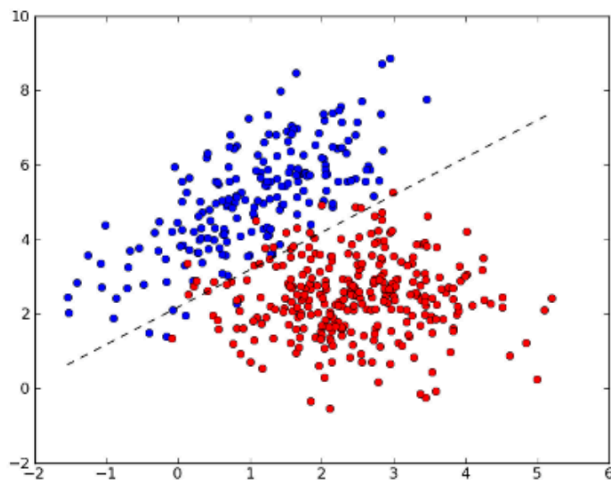
Boundaries are called *decision boundaries* or *decision surfaces*



# Linear Classification Models

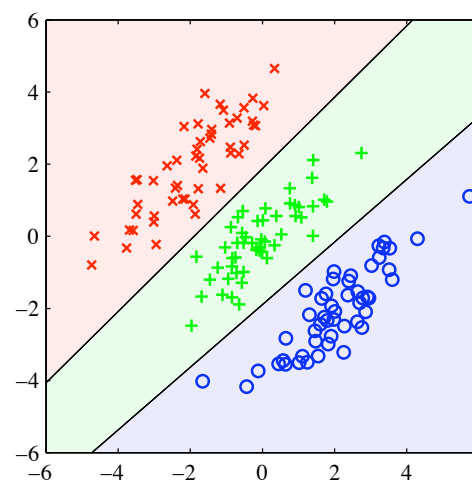
- Decision surfaces are linear functions of input  $x$ 
  - Defined by  $(D - 1)$  dimensional hyperplanes within  $D$  dimensional input space

$K=2$



Not linearly separable

$K=3$



Linearly separable

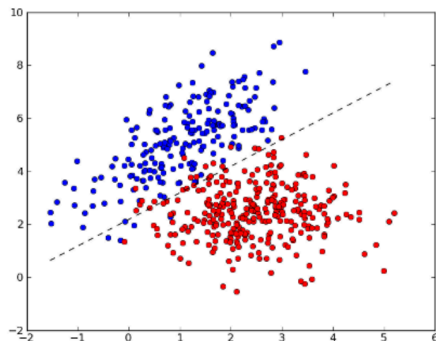
Straight line is 1-D  
decision boundary in 2-D  
space

A plane is 2-D surface in  
3-D space

Data sets whose classes can be separated exactly by linear decision surfaces are said to be Linearly separable

# Representing the target in Classification

- In *regression*:
  - target variable  $t$  is a real number (or vector of real numbers  $\mathbf{t}$ ) which we wish to predict
- In *classification*:
  - there are various ways of using target values to represent class labels, depending on whether
    - Model is probabilistic
    - Model is non-probabilistic



# Representing Class in Probabilistic Model

- Two class: Binary representation is convenient
  - Discrete  $t \in \{0, 1\}$ ,  $t = 1$  represents  $C_1$ ,  
 $t = 0$  means class  $C_2$ 
    - Can interpret value of  $t$  as probability that class is  $C_1$
    - Probabilities taking only extreme values of 0 and 1
- For  $K > 2$  : Use a 1-of- $K$  coding scheme.
  - $\mathbf{t}$  is a vector of length  $K$ 
    - Eg. if  $K = 5$ , a pattern of class 2 has  $\mathbf{t} = (0, 1, 0, 0, 0)^T$
    - Value of  $t_k$  interpreted as probability of class  $C_k$ 
      - If  $t_k$  assume real values then we allow different class probabilities



# Representing Class: Nonprobabilistic

- For non-probabilistic models, e.g, nearest neighbor
  - other choices of target variable representation used

# Three Approaches to Classification

## 1. Discriminant function

– Directly assign  $\mathbf{x}$  to a specific class

- E.g., Linear discriminant, Fisher Linear Disc, Perceptron

## 2. Probabilistic Models (2)

### 1. Discriminative approach

- Model  $p(C_k|\mathbf{x})$  in *inference* stage (direct)
- Use it to make *optimal* decisions
- E.g., Logistic Regression

### 2. Generative approach

- Model class-conditional density  $p(\mathbf{x}|C_k)$
- Together with  $p(C_k)$  use Bayes rule to compute posterior
- E.g., Naïve Bayes classifier

# Separating Inference from Decision

- Separating Inference from Decision is better:
  - Minimize risk (loss function can change in financial app)
  - Reject option (minimize expected loss)
  - Compensate for unbalanced data
    - use modified balanced data & scale by class fractions
  - Combine models

# Probabilistic Models: Generative/Discriminative

- Model  $p(C_k | \mathbf{x})$  in an *inference* stage and use it to make optimal decisions
- Approaches to computing the  $p(C_k | \mathbf{x})$

## 1. Generative

- Model class conditional densities by  $p(\mathbf{x} | C_k)$  together with prior probabilities  $p(C_k)$
- Then use Bayes rule to compute posterior

$$p(C_k | \mathbf{x}) = \frac{p(\mathbf{x} | C_k)p(C_k)}{p(\mathbf{x})}$$

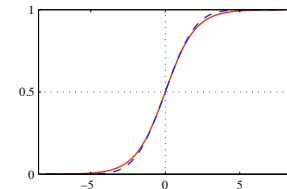
## 2. Discriminative

- Directly model conditional probabilities  $p(C_k | \mathbf{x})$

# From Regression to Classification

- Linear Regression: model prediction  $y(\mathbf{x}, \mathbf{w})$  is a linear function of parameters  $\mathbf{w}$ 
  - In simple case model is also a linear function of  $\mathbf{x}$ 
    - Thus has the form  $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$  where  $y$  is a real no.
- Classification: we need need to predict class labels or posterior probabilities in range  $(0,1)$ 
  - For this, we use a generalization where we transform the linear function of  $\mathbf{w}$  using a nonlinear function  $f(\cdot)$ , so that

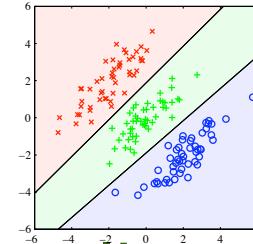
$$y(\mathbf{x}) = f(\mathbf{w}^T \mathbf{x} + w_0)$$



- $f(\cdot)$  is known as an *activation function*
- Whereas its inverse is called a *link function* in statistics
  - link function provides relationship between the linear predictor and the mean of the distribution function

# Decision surface of linear classifier

- Decision surfaces of  $y(\mathbf{x}) = f(\mathbf{w}^T \mathbf{x} + w_0)$  correspond to  $y(\mathbf{x}) = \text{constant}$  or  $\mathbf{w}^T \mathbf{x} + w_0 = \text{constant}$



- Surfaces are linear in  $\mathbf{x}$  even if  $f(\cdot)$  is nonlinear
  - For this reason they are called *generalized linear models*
- However no longer linear in parameters  $\mathbf{w}$  due to presence of  $f(\cdot)$ , therefore:
  - More complex models for classification than regression
- Linear classification algorithms we discuss are applicable even if we transform  $\mathbf{x}$  using a vector of basis functions  $\phi(\mathbf{x})$

# Overview of Linear Classifiers

## 1. Discriminant Functions

- Two class and Multi class
- Least squares for classification
- Fisher's linear discriminant
- Perceptron algorithm

## 2. Probabilistic Generative Models

- Continuous inputs and max likelihood
- Discrete inputs, Exponential Family

## 3. Probabilistic Discriminative Models

- Logistic regression for single and multi class
- Laplace approximation
- Bayesian logistic regression

# Well-known Probabilistic Models

