

# CS-453 Machine Learning

**Nazar Khan**

Department of Computer Science  
University of the Punjab

Linear Classification

## Classification

- ▶ In the previous topic, regression, the goal was to predict *continuous* target variable(s)  $t$  given input variables vector  $x$ .
- ▶ In *classification*, the goal is to predict *discrete* target variable(s)  $t$  given input variables vector  $x$ .
- ▶ Input space is divided into *decision regions*.
- ▶ Boundaries between regions are called *decision boundaries/surfaces*.
- ▶ Training corresponds to finding optimal decision boundaries given training data  $\{(x_1, t_1), \dots, (x_N, t_N)\}$ .

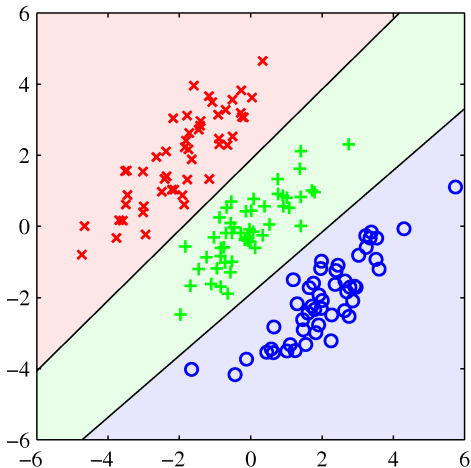
## Classification

- ▶ Assign  $x$  to 1-of- $K$  discrete classes  $C_k$ .
- ▶ Most commonly, the classes are distinct. That is,  $x$  is assigned to one and only one class.
- ▶ Convenient coding schemes for targets  $t$  are
  - ▶ 0/1 coding for binary classification.
  - ▶ 1-of- $K$  coding for multi-class classification. Example, for  $x$  belonging to class 3, the  $K \times 1$  target vector will be coded as  $t = (0, 0, 1, 0, \dots, 0)^T$ .

## Linear Classification

- ▶ Like regression, the simplest classification model is *linear classification*.
  - ▶ This means that the decision surfaces are linear functions of  $x$ , for example  $y(x, w) = w^T x + w_0 = 0$ .
  - ▶ That is, a linear decision surface is a  $D - 1$  dimensional hyperplane in  $D$ -dimensional space.
- ▶ Data in which classes can be *separated exactly* by *linear decision surfaces* is called *linearly separable*.

# Linear Classification



**Figure:** Linearly separable data and corresponding linear decision boundaries.

# Linear Classification

## *Generalized Linear Model*

- ▶ The simplest linear regression model computes continuous outputs  $y(x) = w^T x + w_0$ .
- ▶ By passing these continuous outputs through a non-linear function  $f(\cdot)$ , we can obtain discrete class labels.

$$y(x) = f(w^T x + w_0)$$

- ▶ This is known as a *generalised linear model* and  $f(\cdot)$  is known as the *activation function*.

# Linear Classification

## *Generalized Linear Model*

- ▶ Decision surfaces correspond to all inputs  $x$  where  $y(x) = \text{const}$ . This is equivalent to the condition  $w^T x + w_0 = \text{const}$ .
- ▶ Therefore, decision surfaces are linear functions of the input  $x$ , even if  $f(\cdot)$  is non-linear.
- ▶ As before, we can replace  $x$  by a non-linear transformation  $\phi(x)$  and learn non-linear boundaries in  $x$ -space by learning linear boundaries in  $\phi$ -space.

# Linear Discriminant Functions

## *Two class case*

- ▶ The simplest linear discriminant function is given by  $y(x) = w^T x + w_0$  where  $w$  is called the *weight vector* and  $w_0$  is called the *bias*.
- ▶ Classification is performed via the non-linear step

$$\text{class}(x) = \begin{cases} C_1 & \text{if } y(x) \geq 0 \\ C_2 & \text{if } y(x) < 0 \end{cases}$$

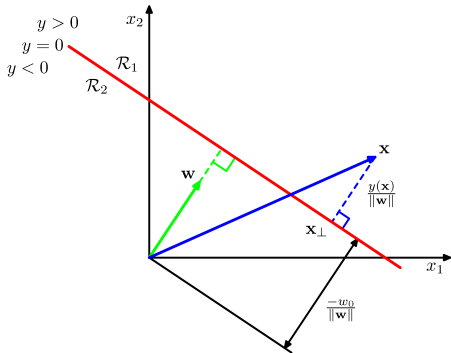
- ▶ We can view  $-w_0$  as a *threshold*.



# Linear Discriminant Functions

## Two class case

- ▶ Weight vector  $w$  is always orthogonal to the decision surface.



- ▶ Proof: For *any* two points  $x_A$  and  $x_B$  on the surface,  $y(x_A) = y(x_B) = 0 \Rightarrow w^T(x_A - x_B) = 0$ . Since vector  $x_A - x_B$  is along the surface,  $w$  must be orthogonal.

# Linear Discriminant Functions

## Two class case

- ▶ Normal distance of any point  $x$  from decision boundary can be computed as  $d = \frac{y(x)}{\|w\|}$ .
  - ▶ Proof:

$$\begin{aligned}
 x &= x_{\perp} + d \frac{w}{\|w\|} \\
 \Rightarrow \underbrace{w^T x + w_0}_{y(x)} &= \underbrace{w^T x_{\perp} + w_0}_{y(x_{\perp})=0} + d \underbrace{w^T \frac{w}{\|w\|}}_{\|w\|} \\
 \Rightarrow d &= \frac{y(x)}{\|w\|}
 \end{aligned}$$

- ▶ Normal distance to boundary from origin ( $x = 0$ ) is  $\frac{w_0}{\|w\|}$ .

## Linear Discriminant Functions

- ▶ For notational convenience, bias can be included as a component of the weight vector via

$$\tilde{w} = (w_0, w)$$

$$\tilde{x} = (1, x)$$

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

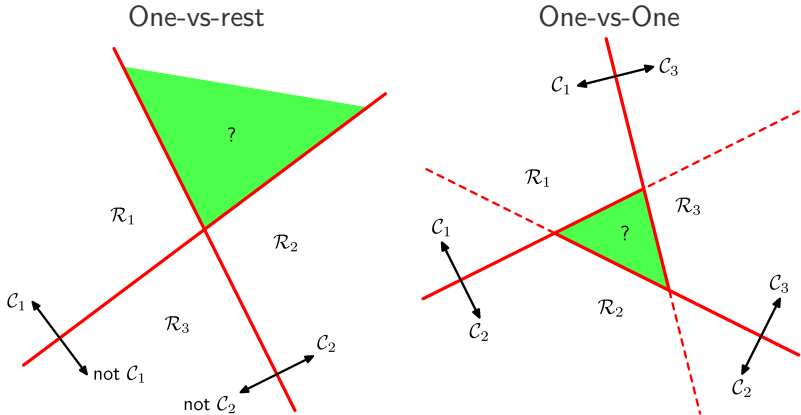
# Linear Discriminant Functions

## *Multiclass case*

- ▶ For  $K$  class classification with  $K > 2$ , we have 3 options
  1. Learn  $K - 1$  *one-vs-rest* binary classifiers.
  2. Learn  $K(K - 1)/2$  *one-vs-one* binary classifiers for every possible pair of classes. Each point can be classified based on majority vote among the discriminant functions.
  3. Learn  $K$  discriminant functions  $y_1, \dots, y_K$  and then  $\text{class}(x) = \arg \max_k y_k(x)$ .
  
- ▶ Options 1 and 2 lead to ambiguous classification regions.

# Linear Discriminant Functions

## Multiclass Ambiguity



**Figure:** Ambiguity of multiclass classification using two-class linear discriminant functions.

# Linear Discriminant Functions

## *Multiclass case*

- ▶ We can write the  $K$ -class discriminant function as

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

- ▶ For learning, we can write the error function as

$$\begin{aligned} E(\tilde{W}) &= \frac{1}{2} \sum_{n=1}^N \|y(x_n) - t_n\|^2 \\ &= \frac{1}{2} \sum_{n=1}^N (\tilde{W}^T \tilde{x}_n - t_n)^T (\tilde{W}^T \tilde{x}_n - t_n) \end{aligned}$$

# Linear Discriminant Functions

## *Least Squares Solution*

- ▶ Optimal discriminant function parameters  $\tilde{W}^*$  that minimize the SSE  $E(\tilde{W})$  are known as the *least-squares-solution*.
- ▶ Can be computed as

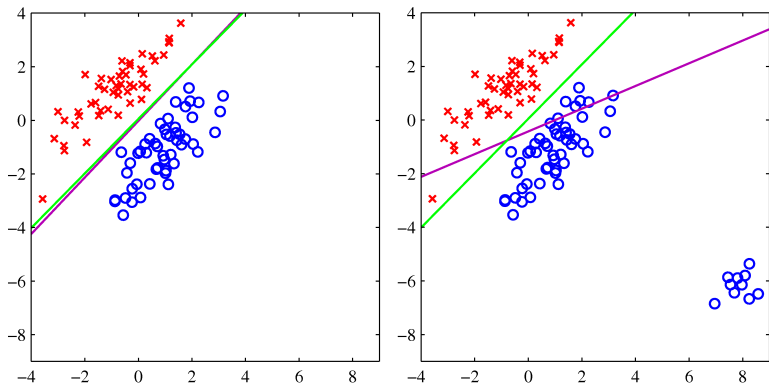
$$\tilde{W}^* = \tilde{X}^\dagger T$$

where  $\tilde{X}^\dagger$  is the pseudo-inverse of the design matrix  $\tilde{X}$  and  $T$  is the matrix of target vectors.

- ▶ As before, we can also work in  $\phi$ -space where we will use the corresponding matrix  $\tilde{\Phi}$  as the design matrix.

# Linear Discriminant Functions

## Least Squares Solution



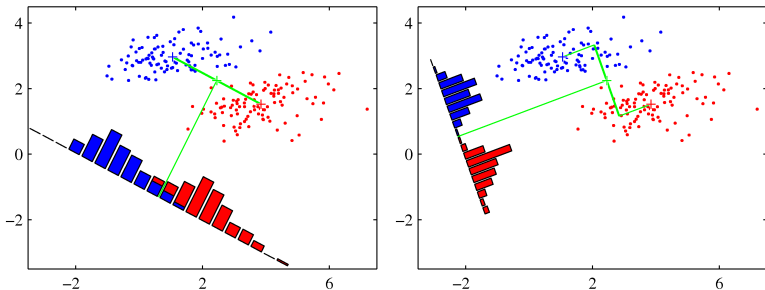
**Figure:** Least squares solution is sensitive to outliers.



# Fisher's Linear Discriminant

## Two class case

- ▶ Project all data onto a single vector  $w$ .
- ▶ Classify by thresholding projected coefficients.
- ▶ Optimal vector is one which
  - ▶ maximises between-class distance, and
  - ▶ minimises within-class distance.



**Figure:** Fisher's linear discriminant. Classify by thresholding projections onto a vector  $w$  that maximises inter-class distance and minimises intra-class distances.

# Fisher's Linear Discriminant

## Two class case

- ▶ Let  $m_k = \frac{\sum_{n \in \mathcal{C}_k} x_n}{N_k}$  be the mean vector of points belonging to class  $\mathcal{C}_k$ .
- ▶ Projection of this mean is then  $m_k = w^T m_k$ .
- ▶ Variance around projected mean can be written as  $s_k^2 = \sum_{n \in \mathcal{C}_k} (w^T x_n - w^T m_k)^2$ .
- ▶ Suitability of any projection direction  $w$  can then be written as

$$\begin{aligned}
 J(w) &= \frac{\text{Inter-class variance}}{\text{Intra-class variance}} \\
 &= \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} \\
 &= \frac{(w^T m_2 - w^T m_1)^2}{\sum_{n \in \mathcal{C}_1} (w^T x_n - w^T m_1)^2 + \sum_{n \in \mathcal{C}_2} (w^T x_n - w^T m_2)^2}
 \end{aligned}$$

# Fisher's Linear Discriminant

Two class case

$$\begin{aligned}
 J(\mathbf{w}) &= \frac{(\mathbf{w}^T(\mathbf{m}_2 - \mathbf{m}_1))(\mathbf{w}^T(\mathbf{m}_2 - \mathbf{m}_1))^T}{\sum_{k=1}^2 \sum_{n \in \mathcal{C}_k} (\mathbf{w}^T(\mathbf{x}_n - \mathbf{m}_k))^2} \\
 &= \frac{\mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1) (\mathbf{m}_2 - \mathbf{m}_1)^T \mathbf{w}}{\mathbf{w}^T \left( \sum_{k=1}^2 \sum_{n \in \mathcal{C}_k} (\mathbf{x}_n - \mathbf{m}_k) (\mathbf{x}_n - \mathbf{m}_k)^T \right) \mathbf{w}} \\
 &= \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \quad (\mathbf{S}_B \text{ and } \mathbf{S}_W \text{ are symmetric due to outer-products})
 \end{aligned}$$

$$\begin{aligned}
 \nabla_{\mathbf{w}} E(\mathbf{w}) &= \frac{\mathbf{w}^T \mathbf{S}_W \mathbf{w} \nabla_{\mathbf{w}} (\mathbf{w}^T \mathbf{S}_B \mathbf{w}) - \mathbf{w}^T \mathbf{S}_B \mathbf{w} \nabla_{\mathbf{w}} (\mathbf{w}^T \mathbf{S}_W \mathbf{w})}{(\mathbf{w}^T \mathbf{S}_W \mathbf{w})^2} \quad (\because \text{quotient rule}) \\
 &= \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w} (2\mathbf{S}_W \mathbf{w}) - \mathbf{w}^T \mathbf{S}_W \mathbf{w} (2\mathbf{S}_B \mathbf{w})}{(\mathbf{w}^T \mathbf{S}_W \mathbf{w})^2} \quad (\because \nabla_{\mathbf{v}} (\mathbf{v}^T \mathbf{M} \mathbf{v}) = (\mathbf{M} + \mathbf{M}^T) \mathbf{v})
 \end{aligned}$$

# Fisher's Linear Discriminant

## Two class case

- ▶ Objective  $J$  can be maximized by equating gradient to the 0 vector

$$w^T S_B w (S_W w) = w^T S_W w (S_B w)$$

- ▶ Since we only care about the direction of projection, we can drop the scalar factors to get

$$S_W w = S_B w$$

$$S_W w = (m_2 - m_1) \underbrace{(m_2 - m_1)^T w}_{\text{scalar}}$$

$$S_W w \propto (m_2 - m_1)$$

$$w \propto S_W^{-1} (m_2 - m_1)$$

# Perceptron Algorithm

## Two-class Classification

- ▶ Target  $t_n$  is taken to be either  $+1$  or  $-1$ .
- ▶ A perceptron classifies its input via the non-linear step function

$$y(\phi) = \begin{cases} 1 & \text{if } w^T \phi_n \geq 0 \\ -1 & \text{if } w^T \phi_n < 0 \end{cases}$$

- ▶ Extremely simplified model of biological neuron.
- ▶ *Perceptron criterion*:  $w^T \phi_n t_n > 0$  for correctly classified point.
- ▶ Error can be defined on the set  $\mathcal{M}(w)$  of misclassified points.

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

- ▶ Optimal  $w$  can be learned via gradient descent.
- ▶ For linearly separable data, perceptron learning is guaranteed to find the decision boundary in finite iterations.

## Gradient Descent

- ▶ Gradient is the direction (in input space) of maximum rate of increase of a function.
- ▶ To minimize, move in negative gradient direction.

$$\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} - \eta \nabla_{\mathbf{w}} E(\mathbf{w})$$

- ▶ Also known as *gradient descent*.
- ▶ Local versus global minima.
- ▶ Learning rate  $\eta$  should be decayed to *avoid oscillation* and to *converge* to a local minimum.
- ▶ Different types of gradient descent:
  - ▶ Batch ( $\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} - \eta \nabla_{\mathbf{w}} E$ )
  - ▶ Sequential ( $\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} - \eta \nabla_{\mathbf{w}} E_n$ )
  - ▶ Stochastic (same as sequential but  $n$  is chosen randomly).
  - ▶ Mini-batches ( $\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} - \eta \nabla_{\mathbf{w}} E_B$ )
- ▶ Most common variations are stochastic gradient descent (SGD) and SGD using mini-batches.