

CS 559: Machine Learning Fundamentals & Applications

Lecture 5: Linear Classification



Outline



5.0. Lecture 4 Review

5.1. Introduction

5.2. Discriminant Functions

5.2.1. Linear Discriminant Analysis (LDA)

5.2.2. Perceptron

5.3. Probabilistic Generative Models - MLE

5.4. Probabilistic Discriminative Models - Logistic Regression

5.0 Lecture 4 Review



Linear Regression

- One of simplest linear models in ML.
- Overcoming the conditions and assumptions may be challenge.

Model Selection

- Overfit vs. Underfit
- Regularization
- Bias-Variance Trade-off



5.0. Lecture 4 Review

5.1. Introduction

5.2. Discriminant Functions

5.2.1. Linear Discriminant Analysis (LDA)

5.2.2. Perceptron

5.3. Probabilistic Generative Models - MLE

5.4. Probabilistic Discriminative Models - Logistic Regression



5.1. Introduction

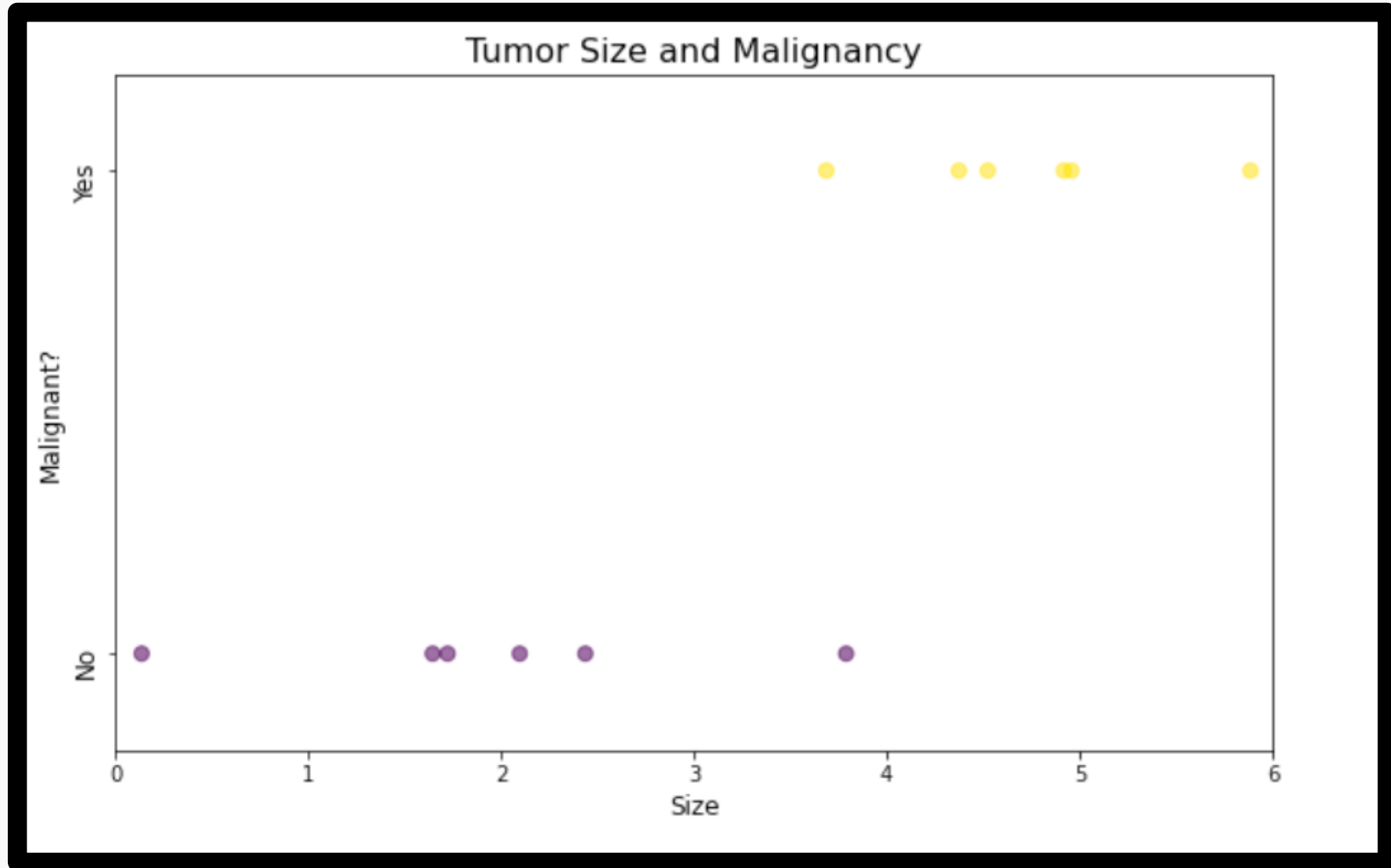
Goal: to take an input vector \mathbf{x} and to assign it to one of K discrete classes C_k where $k = 1, \dots, K$.

- Data Condition: The classes are disjoint, and the target is discrete.
- The input space is divided into ***decision regions*** – linear decision surfaces called ***hyperplanes*** in $D-1$ dimensions if the input space is in D -dimensions.
- For binary problem, target $t \in \{0,1\}$ or $\{-1,1\}$.

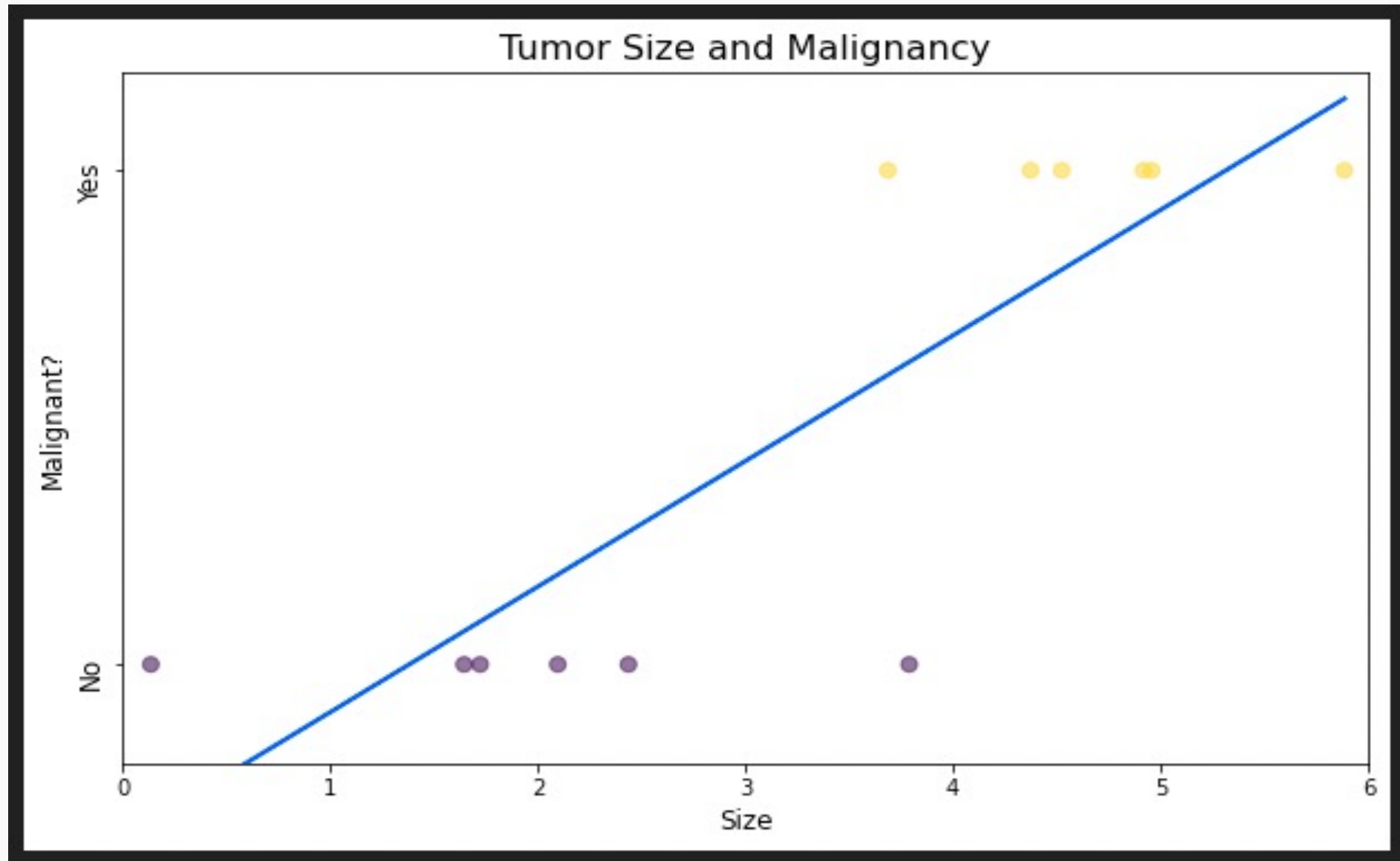
We approach in three different approaches:

1. A ***discriminant modeling***: directly assigns each vector \mathbf{x} to a specific class (LDA, Perceptron).
2. A ***probabilistic modeling***: determines the conditional probability distribution $p(C_k|\mathbf{x})$ and represent them as a parametric model.
 1. ***Probabilistic generative modeling***: a classification modeling using Bayes' theorem (MLE).
 2. ***Probabilistic discriminative modeling***: a classification modeling using MLE in a form of discriminative modeling (logistic regression).

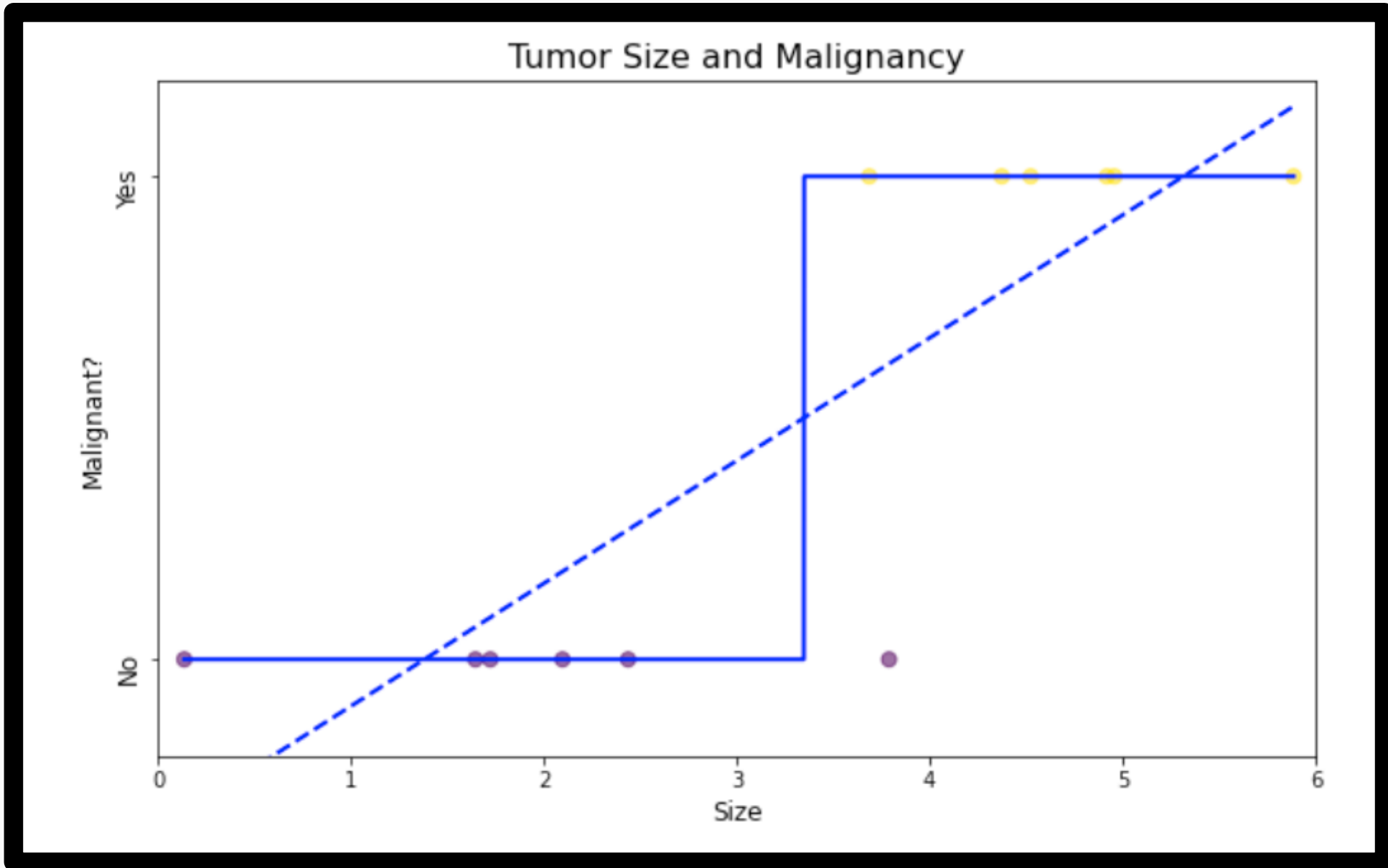
5.1. Introduction



5.1. Introduction



5.1. Introduction





5.0. Lecture 4 Review

5.1. Introduction

5.2. Discriminant Functions

5.2.1. Linear Discriminant Analysis (LDA)

5.2.2. Perceptron

5.3. Probabilistic Generative Models - MLE

5.4. Probabilistic Discriminative Models - Logistic Regression



5.2. Discriminant Functions

Consider a simple linear discriminant function of the input vector \mathbf{x} ,

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 \quad (5-1)$$

where \mathbf{w} is called a *weight vector*, and w_0 is a bias.

- If $w_0 < 0$, sometimes we call it a *threshold*.
- If $y(\mathbf{x}) \geq 0$, then \mathbf{x} is assigned to class C_1 . Otherwise, C_2 .
- The corresponding hyperplane is when $y(\mathbf{x}) = 0$.

5.2. Discriminant Functions

- If two points \mathbf{x}_A and \mathbf{x}_B both lying on the decision surface, $y(\mathbf{x}_A) = y(\mathbf{x}_B) = 0$.
- The normal distance from the origin to the decision surface is

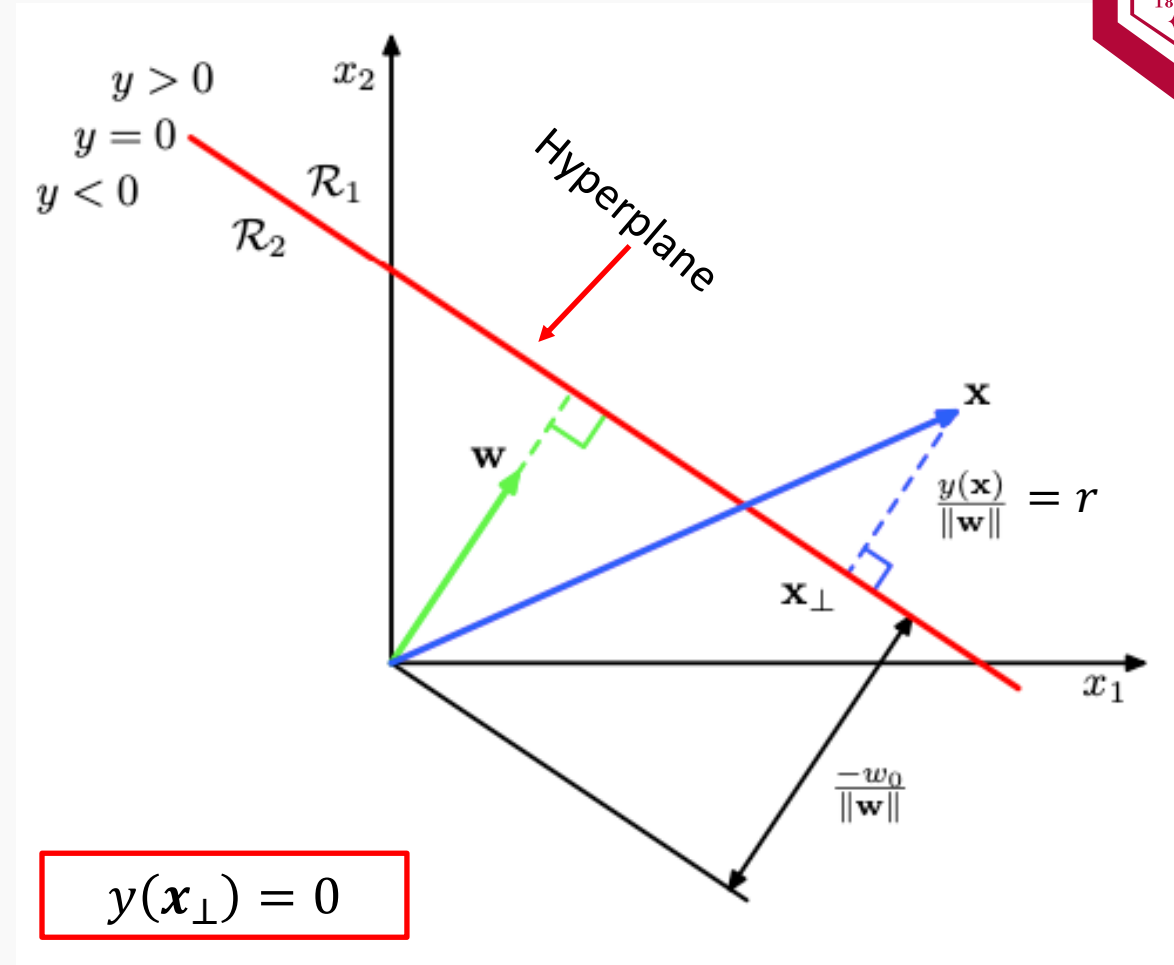
$$\frac{\mathbf{w}^T \mathbf{x}}{|\mathbf{w}|} = -\frac{w_0}{|\mathbf{w}|}.$$

- The perpendicular distance r of the point \mathbf{x} from the hyperplane can be expressed as

$$\mathbf{x} = \mathbf{x}_\perp + r \left(\frac{\mathbf{w}}{|\mathbf{w}|} \right)$$

where \mathbf{x}_\perp is the orthogonal projection onto the hyperplane. It also can be expressed in terms of y :

$$y(\mathbf{x}) = y(\mathbf{x}_\perp) + r \left(\frac{\mathbf{w}^T \mathbf{w}}{|\mathbf{w}|} \right) \rightarrow r = \frac{y(\mathbf{x})}{|\mathbf{w}|}.$$



$$y(\mathbf{x}_\perp) = 0$$

$$|\mathbf{w}| = \sqrt{\mathbf{w}^T \mathbf{w}} \rightarrow \frac{|\mathbf{w}|}{\mathbf{w}^T \mathbf{w}} = \frac{1}{|\mathbf{w}|}$$

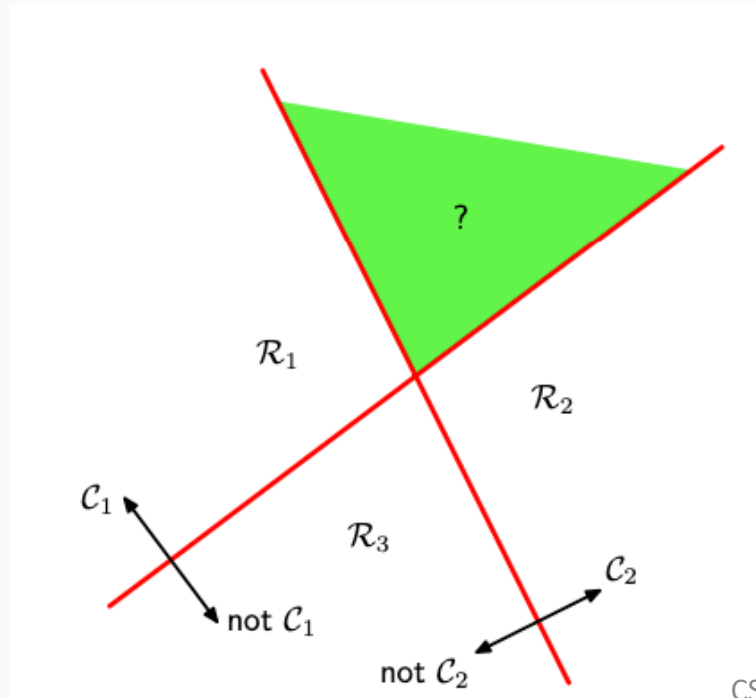
(5-2)

5.2. Discriminant Functions

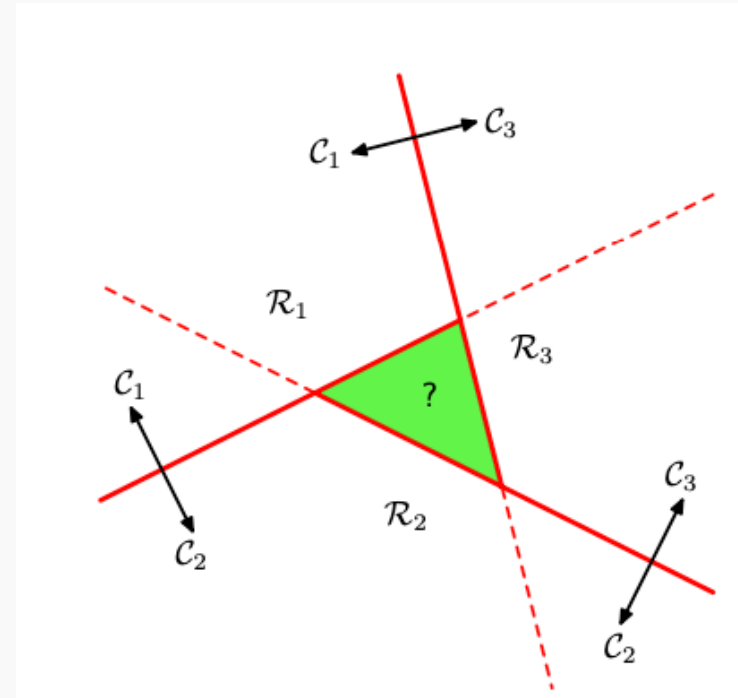
What happens when the target has more than two classes?

If we approach as binary classification problem, there are two possible approaches.

- *One-verse-rest* classifier: the $K-1$ classifier that solves as a two-class problem of separating points in a particular class C_k from points not in that class.
- *One-verse-one* classifier: each point is classified according to a majority vote amongst the discriminant functions.
- Either holds an *ambiguous region* where $y(\mathbf{x})$ is not possible for some \mathbf{x} .



CS559 - Lecture 5 - Linear Classification



5.2. Discriminant Functions

- Instead, we consider a single K -class linear discriminant function comprising K many functions in the form of

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

to assign a class C_k for a point \mathbf{x} if $y_k(\mathbf{x}) > y_j(\mathbf{x})$ for all $j \neq k$.

- The class boundary between C_k and C_j is when $y_k = y_j$ defined by

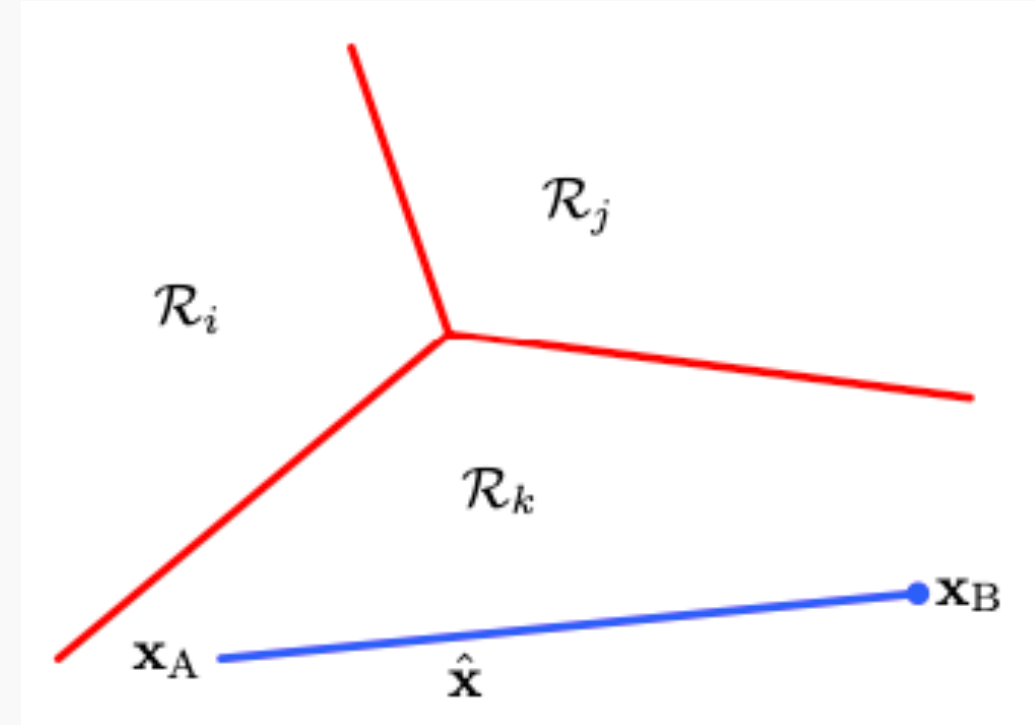
$$(\mathbf{w}_k - \mathbf{w}_j)^T \mathbf{x} + (w_{k0} - w_{j0}) = 0. \quad (5-3)$$

- In this way, we can connect the decision regions into a single point and be convex.
- Consider any point $\hat{\mathbf{x}}$ lies on the line between two points \mathbf{x}_A and \mathbf{x}_B in R_k ,

$$\hat{\mathbf{x}} = \lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B$$

where $0 \leq \lambda \leq 1$. From the linearity of the discriminant functions, it follows that

$$y_k(\hat{\mathbf{x}}) = \lambda y_k(\mathbf{x}_A) + (1 - \lambda) y_k(\mathbf{x}_B).$$





5.0. Lecture 4 Review

5.1. Introduction

5.2. Discriminant Functions

5.2.1. Linear Discriminant Analysis (LDA)

5.2.2. Perceptron

5.3. Probabilistic Generative Models - MLE

5.4. Probabilistic Discriminative Models - Logistic Regression



5.2.1. Fisher's Linear Discriminant Analysis (LDA)

- Consider a binary linear classifier in terms of *dimensionality reduction* for D -dimensional input vector \mathbf{x} having an output $\{C_1, C_2\}$.
- Let y be C_1 if $y \geq -w_0$ and otherwise C_2 .
- Suppose we project down to 1-dimension using
$$y = \mathbf{w}^T \mathbf{x}.$$
- We can control \mathbf{w} to project by maximizing the class separation to avoid the considerable loss information that may occur in projection D -dimensional space to 1-dimensional space.



5.2.1. Fisher's Linear Discriminant Analysis (LDA)

If there are N_1 points in C_1 and N_2 in C_2 , the mean vector of the two classes is

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

The separation mean of the projected means can measure the separation of the classes to maximize

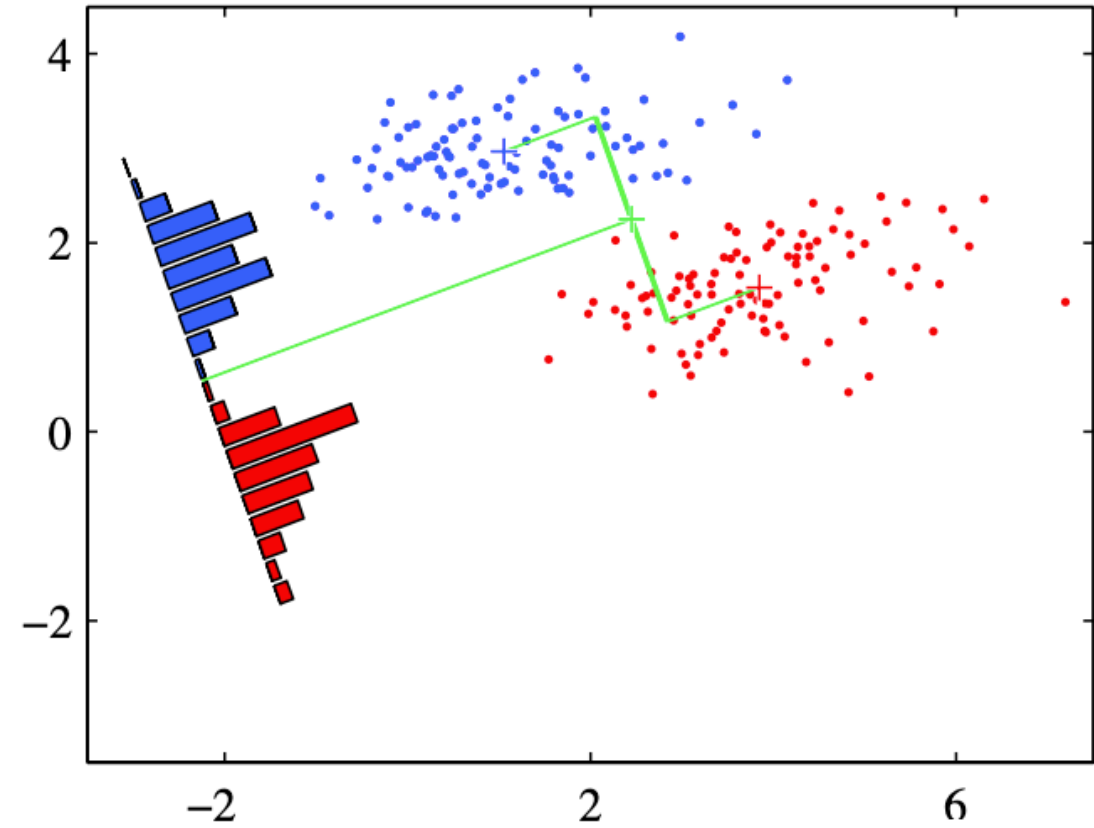
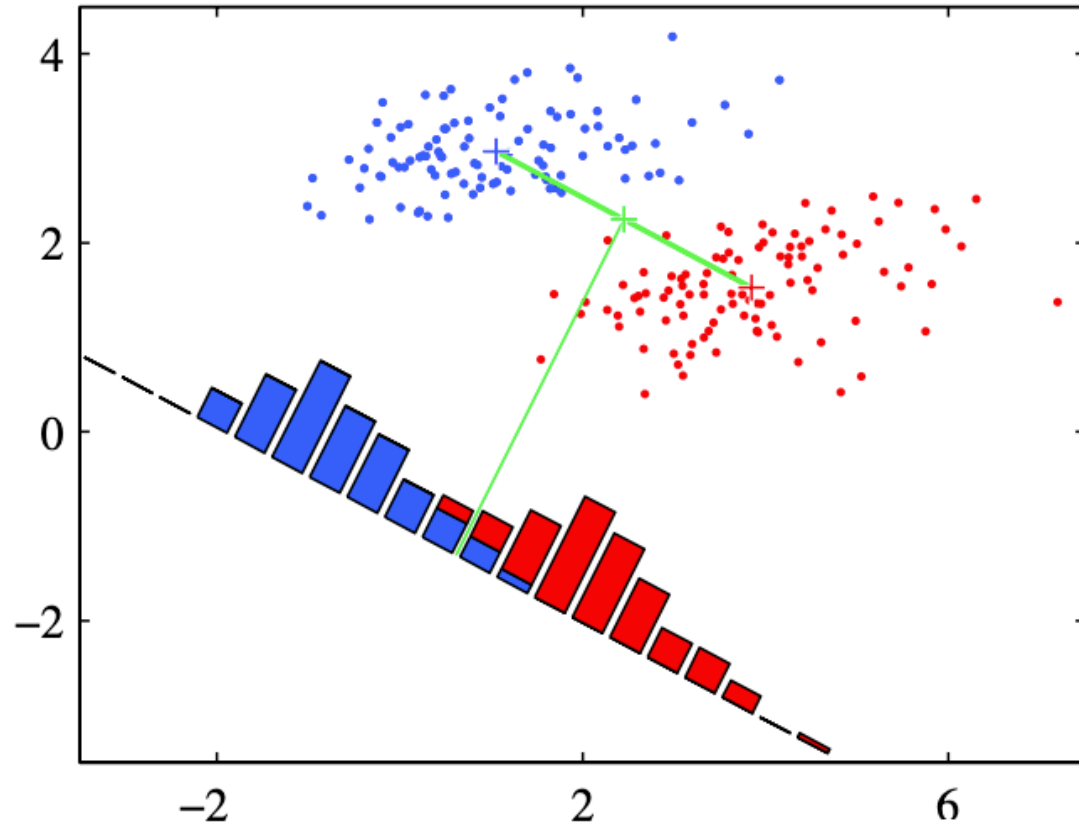
$$\mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1). \quad (5-13)$$

However, the separation becomes arbitrary large when the magnitude of \mathbf{w} is large.

To avoid this, \mathbf{w} will be forced to be a unit length $\sum_i w_i^2 = 1$.

- However, the strong nondiagonal **covariances** of the class distributions still may have **overlap** when projected onto the line joining their means.
- To minimize the overlap, the function must give the **largest separation** between the projected class means while the **smallest variance** within each class, also known as the within-class variance.

5.2.1. Fisher's Linear Discriminant Analysis (LDA)





5.2.1. Fisher's Linear Discriminant Analysis (LDA)

The **within-class variance** of the transformed data from class C_k is

$$s_k^2 = \sum_{n \in C_k} (y_n - m_k)^2 \quad \text{where } y_n = \mathbf{w}^T \mathbf{x}_n \quad (5-4)$$

and the **total within-class variance** is

$$s^2 = \sum_{k=1}^K s_k^2. \quad (5-5)$$

The **Fisher criterion**, $J(\mathbf{w})$, is the ratio of the between-class variance to the within-class variance (Eq. 5-5) and is defined as

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}. \quad (5-6)$$



5.2.1. Fisher's Linear Discriminant Analysis (LDA)

It can be rewritten as the dependence on \mathbf{w} 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 is the *between-class* covariance matrix

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

and \mathbf{S}_w is the **total within-class** covariance matrix is given by

$$\mathbf{S}_w = \sum_{n \in C_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T + (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^T.$$

$J(\mathbf{w})$ is maximized when $\nabla_{\mathbf{w}} J(\mathbf{w}) = 0$:

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

(5-7)



5.2.1. Fisher's Linear Discriminant Analysis (LDA)

To obtain \mathbf{w} , we can consider the following two properties.

- The *between-class* shows that $\mathbf{S}_B \mathbf{w}$ is always in the direction of $(\mathbf{m}_2 - \mathbf{m}_1)$.
- Considering the direction only (not the magnitude), the scalar factors $(\mathbf{w}^T \mathbf{S}_B \mathbf{w})$ and $(\mathbf{w}^T \mathbf{S}_w \mathbf{w})$ can be dropped in Eq. (5-7).

Then multiplying \mathbf{S}_w^{-1} on both sides, then \mathbf{w} can be obtained as

$$\mathbf{w} \propto \mathbf{S}_w^{-1}(\mathbf{m}_2 - \mathbf{m}_1). \quad (5-8)$$

Alternatively, we can find \mathbf{w} that maximizes the criterion using Linear Algebra via generalized eigenvalue problem:

- The eigenvalue, λ , can be computed as

$$\begin{aligned} \mathbf{S}_B \mathbf{w} &= \lambda \mathbf{S}_w \mathbf{w} \rightarrow (\mathbf{S}_w^{-1} \mathbf{S}_B - \lambda \mathbf{I}) = \mathbf{0} \\ \det(\mathbf{S}_w^{-1} \mathbf{S}_B - \lambda \mathbf{I}) &= 0 \end{aligned} \quad (5-9)$$

- The eigenvector for λ_{\max} is then solution.



5.2.1. Fisher's Linear Discriminant Analysis (LDA)

We can relate to the least squares.

Suppose we can model the class-conditional densities $p(y|C_k)$ using Gaussian distribution. The solution of \mathbf{w} can be approached via MLE.

Consider a binary classification and let the target $t_1 = \frac{N}{N_1}$ for class C_1 where N is the total number and N_1 is the number in C_1 . Similarly, let $t_2 = -\frac{N}{N_2}$.

Start from the sum-of-squares error function,

$$E = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n)^2.$$



5.2.1. Fisher's Linear Discriminant Analysis (LDA)

Set the derivative of the sum-of-squares error function w.r.t. \mathbf{w} equals to 0,

$$\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n) \mathbf{x}_n = 0.$$

The use of

$$\sum_{n=1}^N t_n = \frac{N_1 N}{N_1} - \frac{N_2 N}{N_2} = 0$$
$$\mathbf{m} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \frac{1}{N} (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2)$$

and

$$w_0 = -\mathbf{w}^T \mathbf{m},$$

the derivative of the sum-of-squares error function (using \mathbf{S}_w and \mathbf{S}_B from slide #19) becomes

$$\left(\mathbf{S}_w + \frac{N_1 N_2}{N} \mathbf{S}_B \right) \mathbf{w} = N(\mathbf{m}_1 - \mathbf{m}_2). \quad (5-9)$$



5.2.1. Fisher's Linear Discriminant Analysis (LDA)

For $K > 2$ classes, assume $D > K$.

Suppose there are $D' > 1$ linear new features $y_k = \mathbf{w}_k^T \mathbf{x}$ for $k = 1, \dots, D'$ where \mathbf{y} is the vector for y_k :

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

where the bias parameter w_0 is excluded.

The within-class covariance matrix is

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

where \mathbf{S}_k is

$$\mathbf{S}_k = \sum_{n \in C_k} (\mathbf{x}_n - \mathbf{m}_k)(\mathbf{x}_n - \mathbf{m}_k)^T = \sum_{n \in C_k} \left(\mathbf{x}_n - \frac{1}{N_k} \sum_{n \in C_k} \mathbf{x}_n \right) \left(\mathbf{x}_n - \frac{1}{N_k} \sum_{n \in C_k} \mathbf{x}_n \right)^T.$$



5.2.1. Fisher's Linear Discriminant Analysis (LDA)

The total covariance matrix

$$\mathbf{S}_T = \sum_{n=1}^N (\mathbf{x}_n - \mathbf{m})(\mathbf{x}_n - \mathbf{m})^T$$

where \mathbf{m} is the mean of total set

$$\mathbf{m} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \frac{1}{N} \sum_{k=1}^K N_k \mathbf{m}_k .$$

The total covariance is the sum of within-class variance and the between-class covariance

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



5.2.1. Fisher's Linear Discriminant Analysis (LDA)

In the projected D' -dimensional \mathbf{y} -space

$$\mathbf{s}_W = \sum_{k=1}^K \sum_{n \in C_k} (\mathbf{y}_n - \boldsymbol{\mu}_k)(\mathbf{y}_n - \boldsymbol{\mu}_k)^T$$

and

$$\mathbf{s}_B = \sum_{k=1}^K N_k (\boldsymbol{\mu}_k - \boldsymbol{\mu})(\boldsymbol{\mu}_k - \boldsymbol{\mu})^T$$

where

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n \in C_k} \mathbf{y}_n, \quad \boldsymbol{\mu} = \frac{1}{N} \sum_{k=1}^K N_k \boldsymbol{\mu}_k.$$

Then, the criterion can be expressed as

$$J(\mathbf{W}) = \text{Tr}\{\mathbf{s}_W^{-1} \mathbf{s}_B\} \rightarrow J(\mathbf{w}) = \text{Tr}\left\{(\mathbf{W} \mathbf{S}_W \mathbf{W}^T)^{-1} (\mathbf{W} \mathbf{S}_B \mathbf{W})\right\}.$$

The weight values to maximize the criterion can be determined by the eigenvectors of $\mathbf{s}_W^{-1} \mathbf{s}_B$ that correspond to the D' largest eigenvalues.



5.2.1. Fisher's Linear Discriminant Analysis (LDA)

| Class | (x_1, x_2) | m |
|-------|------------------------|----------|
| C_1 | (1,2), (2,3), (3,4.9) | (2, 3.3) |
| C_2 | (2,1), (3,2), (4, 3.9) | (3,2.3) |

The total within-class variance is

$$\mathbf{S}_w = \begin{bmatrix} 4 & 5.8 \\ 5.8 & 8.68 \end{bmatrix}$$

and its inverse is

$$\mathbf{S}_w^{-1} = \begin{bmatrix} 8.04 & -5.37 \\ -5.37 & 3.70 \end{bmatrix}.$$

Since \mathbf{w} is proportional to $\mathbf{S}_w^{-1}(\mathbf{m}_2 - \mathbf{m}_1)$,

$$\mathbf{w} \propto \mathbf{S}_w^{-1}(\mathbf{m}_2 - \mathbf{m}_1) = \begin{bmatrix} -13.41 \\ 9.07 \end{bmatrix}$$

and its unit vector is

$$\hat{\mathbf{w}} = \frac{\mathbf{w}}{|\mathbf{w}|} = \begin{bmatrix} -0.83 \\ 0.56 \end{bmatrix}$$

and therefore, the model is

$$y = -0.83x_1 + 0.56x_2.$$



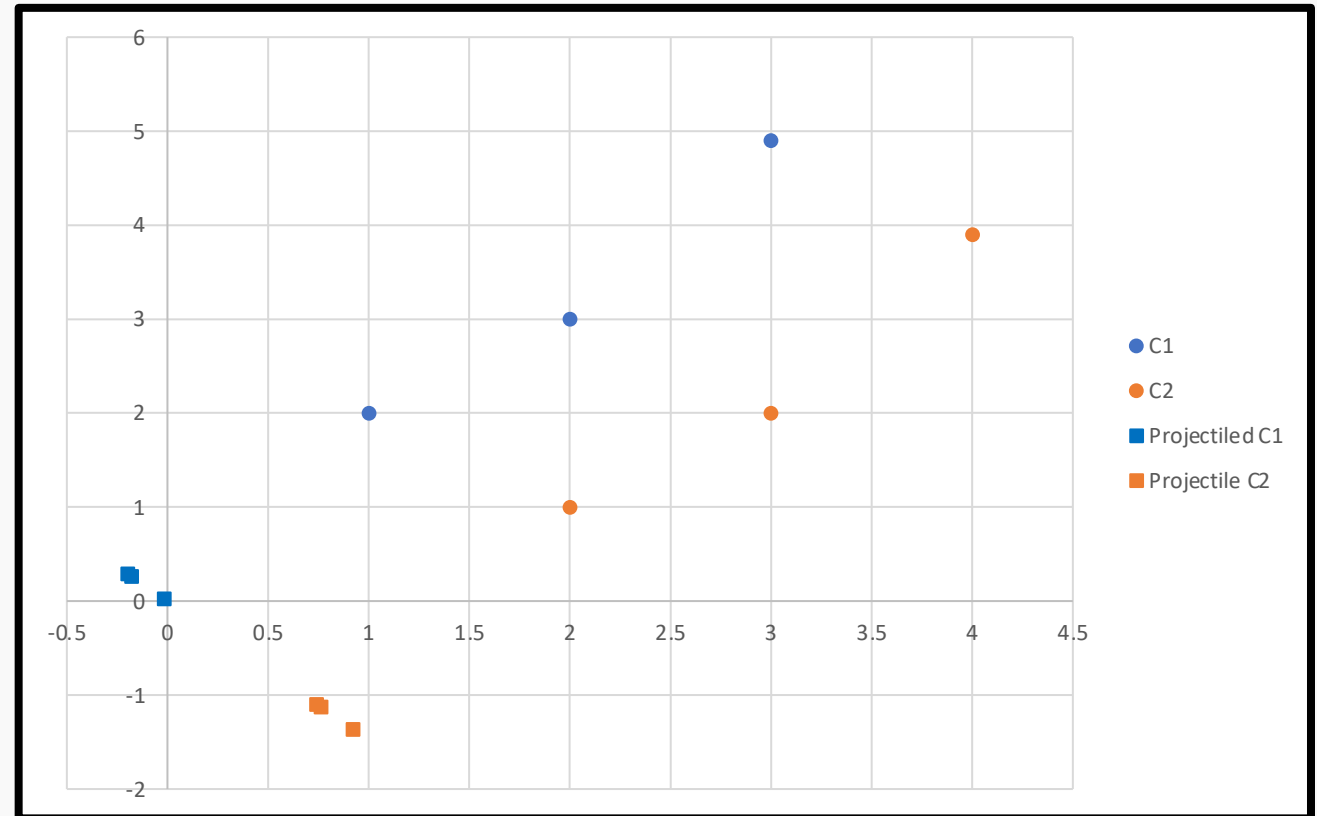
5.2.1. Fisher's Linear Discriminant Analysis (LDA)

Substituting data into the model reveals

$$y = [0.29 \quad 0.025 \quad 0.26 \quad -1.10 \quad -1.36 \quad -1.13].$$

For the complete orthogonal projection, set $y=0$ and find the model

$$x_2 = \frac{0.83}{0.56} x_1$$





5.0. Lecture 4 Review

5.1. Introduction

5.2. Discriminant Functions

5.2.1. Linear Discriminant Analysis (LDA)

5.2.2. Perceptron

5.3. Probabilistic Generative Models - MLE

5.4. Probabilistic Discriminative Models - Logistic Regression



5.2.2. Perceptron

Consider a binary class model of the form

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

where the nonlinear activation function $f(\cdot)$ is a step function of the form

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

The target values are

$$t = \begin{cases} 1, & \text{for } C_1 \\ -1, & \text{for } C_2 \end{cases}.$$



5.2.2. Perceptron

The learning is motivated from the error function minimization.

- But the change of \mathbf{w} can cause the change of decision boundary and lead to misclassification.
- Therefore, the approaching in gradient descent way is not possible. We need a different error function.
- Instead, we use the *perceptron criterion* to ensure the misclassification is minimized.

The perceptron criterion is given by

$$E_p(\mathbf{w}) = - \sum_{n \in \mathcal{M}} \mathbf{w}^T \phi_n t_n \quad (5-10)$$

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



5.2.2. Perceptron

- Suppose we seek for a linear pattern s.t. $\mathbf{w}^T \phi(\mathbf{x}_n) > 0$ for C_1 .
- The patterns with $\mathbf{w}^T \phi(\mathbf{x}_n)t_n > 0$ will also be followed. For \mathbf{x}_n in C_2 will also follow the same pattern since $\mathbf{w}^T \phi(\mathbf{x}_n) < 0$ and $t_n = -1$ for all correct classification.
- In other words, misclassified \mathbf{x}_n , that is $n \in \mathcal{M}$, will have pattern of $\mathbf{w}^T \phi(\mathbf{x}_n)t_n < 0$ for any classes.
- We need to therefore minimize the quantity $-\mathbf{w}^T \phi(\mathbf{x}_n)t_n$ by applying the stochastic gradient descent algorithm,

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_p(\mathbf{x}_n).$$

- For each iteration, $y(\mathbf{x}, \mathbf{w})$ is evaluated.
- If the pattern is correctly classified, the weight does not change.

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} \tag{5-11a}$$

- If the classification is incorrect,

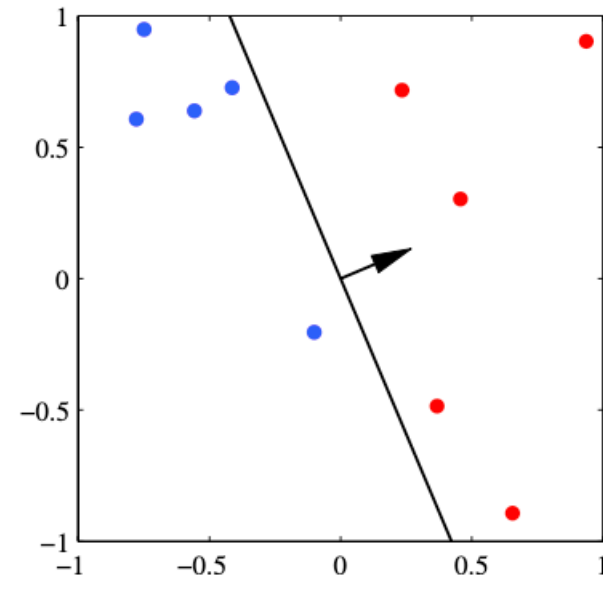
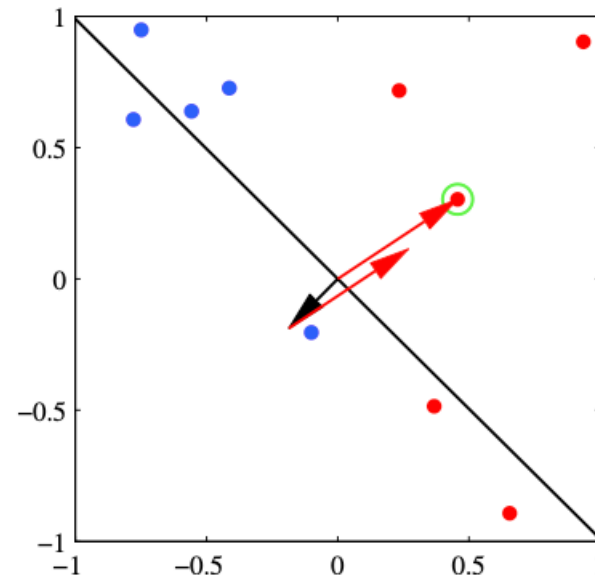
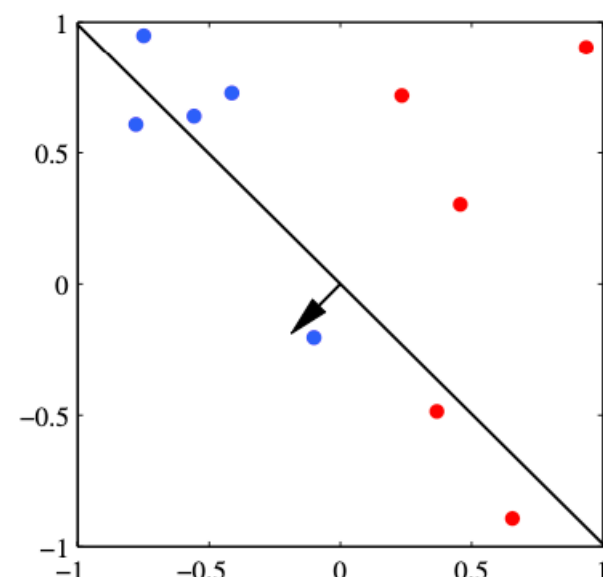
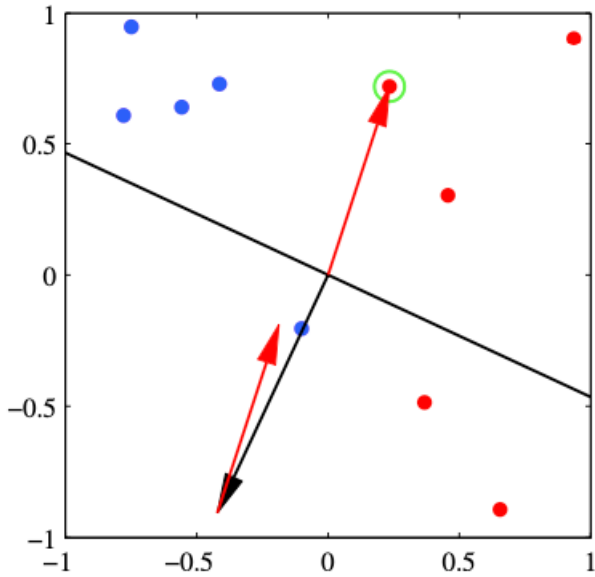
- For C_1 , we add the vector $\phi(\mathbf{x}_n)$ to \mathbf{w} .

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^\tau + \mathbf{x}_n \tag{5-11b}$$

- For C_2 , we subtract vector $\phi(\mathbf{x}_n)$ from \mathbf{w} .

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^\tau - \mathbf{x}_n \tag{5-11c}$$

5.2.2. Perceptron





5.2.2. Perceptron

Consider $\eta = 1$ and $||\phi_n t_n||^2 > 0$, the single update will reduce the error from the misclassification:

$$-\mathbf{w}^{(\tau+1)T} \phi_n t_n = -\mathbf{w}^{(\tau)T} \phi_n t_n - (\phi_n t_n)^T \phi_n t_n < -\mathbf{w}^{(\tau)T} \phi_n t_n.$$

- It will not impact the error from another misclassified pattern.
- It may change the result on previously corrected patterns.
- It does not guarantee the reduction of the total error function at each stage.
- However, if the training data is linearly separable, the perceptron guarantees the exact solution after the finite iteration – the *perceptron convergence theorem*.
- The perceptron may have many solutions depends on the parameter initialization.
- The learning algorithm will never converge if the training data is not linearly separable.



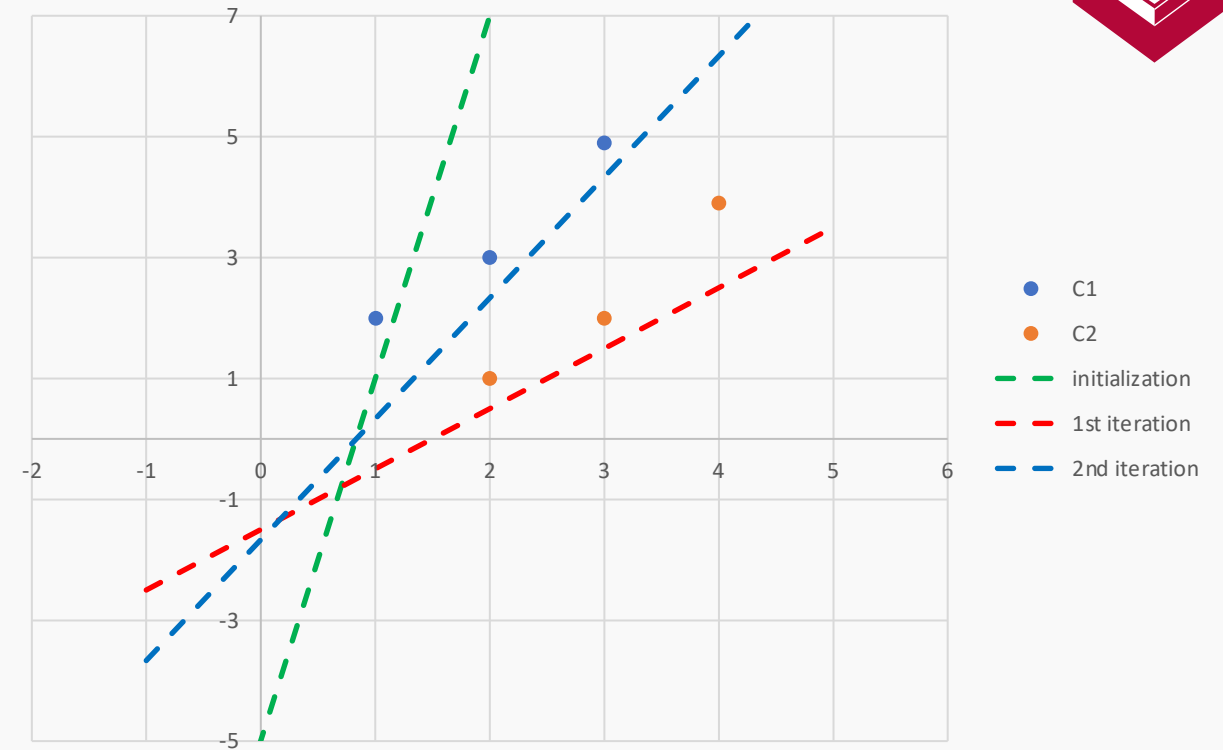
5.2.2. Perceptron

| Class | (x_1, x_2) |
|-------|------------------------|
| 1 | (1,2), (2,3), (3,4.9) |
| -1 | (2,1), (3,2), (4, 3.9) |

Initial $\mathbf{w} = [5, -6, 1]$, $y = w_0 + w_1x_1 + w_2x_2$.

| | w_0 | w_1 | w_2 |
|----------------|-------|-------|-------|
| Initialization | 5 | -6 | 1 |
| 1st Iteration | 6 | -4 | 4 |
| 2nd Iteration | 5 | -6 | 3 |

$\hat{\mathbf{w}} = [0.598, -0.717, 0.585]$



| | 1 | 2 | 3 | 4 | 5 | 6 |
|----------------|---|----|----|----|----|----|
| Initialization | 1 | -1 | -1 | -1 | -1 | -1 |
| 1st Iteration | 1 | 1 | 1 | 1 | 1 | 1 |
| 2nd Iteration | 1 | 1 | 1 | -1 | -1 | -1 |



5.0. Lecture 4 Review

5.1. Introduction

5.2. Discriminant Functions

5.2.1. Linear Discriminant Analysis (LDA)

5.2.2. Perceptron

5.3. Probabilistic Generative Models - MLE

5.4. Probabilistic Discriminative Models - Logistic Regression

5.3. Probabilistic Generative Models

Consider a binary class classification, C_1 and C_2 .

The posterior probability for C_1 ,

$$p(C_1|\mathbf{x}) = \frac{p(\mathbf{x}|C_1)p(C_1)}{p(\mathbf{x}|C_1)p(C_1) + p(\mathbf{x}|C_2)p(C_2)} = \frac{1}{1 + \exp(-a)} = \sigma(a) \quad (5-12)$$

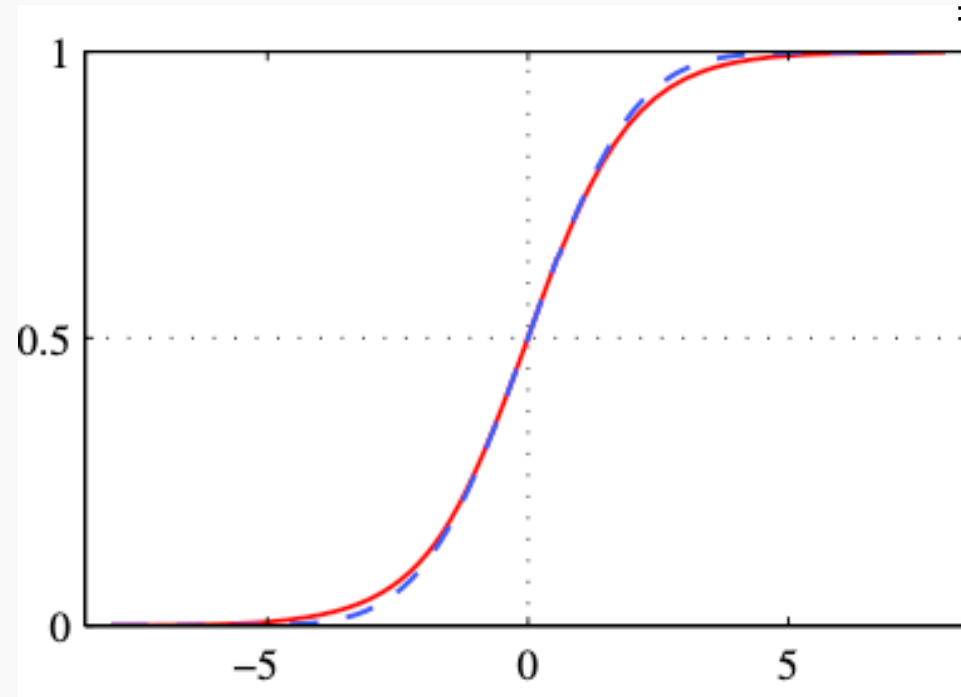
where

$$a = \ln \frac{p(\mathbf{x}|C_1)p(C_1)}{p(\mathbf{x}|C_2)p(C_2)}.$$

Note: $\exp(-a) = \exp\left(-\ln \frac{p(\mathbf{x}|C_1)p(C_1)}{p(\mathbf{x}|C_2)p(C_2)}\right)$

$$= \left(\frac{p(\mathbf{x}|C_1)p(C_1)}{p(\mathbf{x}|C_2)p(C_2)}\right)^{-1}$$

$\sigma(a)$ is called the *logistic sigmoid*.





5.3. Probabilistic Generative Models

It follows the symmetry property

$$\sigma(-a) = 1 - \sigma(a).$$

The inverse of the logistic sigmoid is known as the *logit* function in the form of

$$a = \ln\left(\frac{\sigma}{1 - \sigma}\right)$$

and represents to log of the ratio of $\ln[p(C_1|\mathbf{x})/p(C_2|\mathbf{x})]$ for the two classes, also known as the *log odds*.



5.3. Probabilistic Generative Models

For the case of $K > 2$ classes, the posterior probability of C_k is in the form

$$p(C_k|\mathbf{x}) = \frac{p(\mathbf{x}|C_k)p(C_k)}{\sum_j p(\mathbf{x}|C_j)p(C_j)} = \frac{\exp a_k}{\sum_j \exp(a_j)}$$

which is known as the *normalized exponential* and can be regarded as a multiclass generalization of the logistic sigmoid and the quantities a_k are defined by

$$a_k = \ln p(\mathbf{x}|C_k)p(C_k).$$

This is also known as the ***softmax function*** that smooths the max function:

$$\begin{aligned} p(C_k|\mathbf{x}) &\approx 1 \\ p(C_j|\mathbf{x}) &\approx 0 \end{aligned}$$

for all $j \neq k$.



5.3. Probabilistic Generative Models

Assume the class-conditional densities are Gaussian and all classes share the same covariance matrix. The density for class C_k is

$$p(\mathbf{x}|C_k) = \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\mathbf{\Sigma}|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \mathbf{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\}.$$

Consider for binary classes,

$$p(C_1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

Derivation:

Let the prior probability be $p(C_1) = \pi_1$ and $p(C_2) = \pi_2$.



5.3. Probabilistic Generative Models

$$\begin{aligned} p(C_1|\mathbf{x}) &= \frac{p(\mathbf{x}|C_1)p(C_1)}{p(\mathbf{x}|C_1)p(C_1) + p(\mathbf{x}|C_2)p(C_2)} \\ &= \frac{\pi_1 \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_1)\right\}}{\pi_1 \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_1)\right\} + \pi_2 \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_2)\right\}} \end{aligned}$$

The exponential term can be simplified as

$$-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_1) = -\frac{1}{2}(\mathbf{x} \boldsymbol{\Sigma}^{-1} \mathbf{x} - 2\boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1).$$

The first term in the right side is not associate with C_1 and it can be vanished. The equation above can be expressed as

$$\boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1.$$

5.3. Probabilistic Generative Models

The posterior probability then becomes

$$\begin{aligned}
 &= \frac{\pi_1 \exp \left\{ \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 \right\}}{\pi_1 \exp \left\{ \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 \right\} + \pi_2 \exp \left\{ \boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2 \right\}} \\
 &= \frac{1}{1 + \frac{\pi_2}{\pi_1} \exp \left\{ \boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2 - \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \frac{1}{2} \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 \right\}} \\
 &= \frac{1}{1 + \exp \left\{ \boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2 - \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \frac{1}{2} \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 + \ln \frac{\pi_2}{\pi_1} \right\}} \\
 &= \frac{1}{1 + \exp \left\{ - \left[(\boldsymbol{\mu}_1^T - \boldsymbol{\mu}_2^T) \boldsymbol{\Sigma}^{-1} \mathbf{x} + \frac{1}{2} (\boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2) + \ln \frac{\pi_1}{\pi_2} \right] \right\}} \\
 &= \frac{1}{1 + \exp(-a)}
 \end{aligned}$$

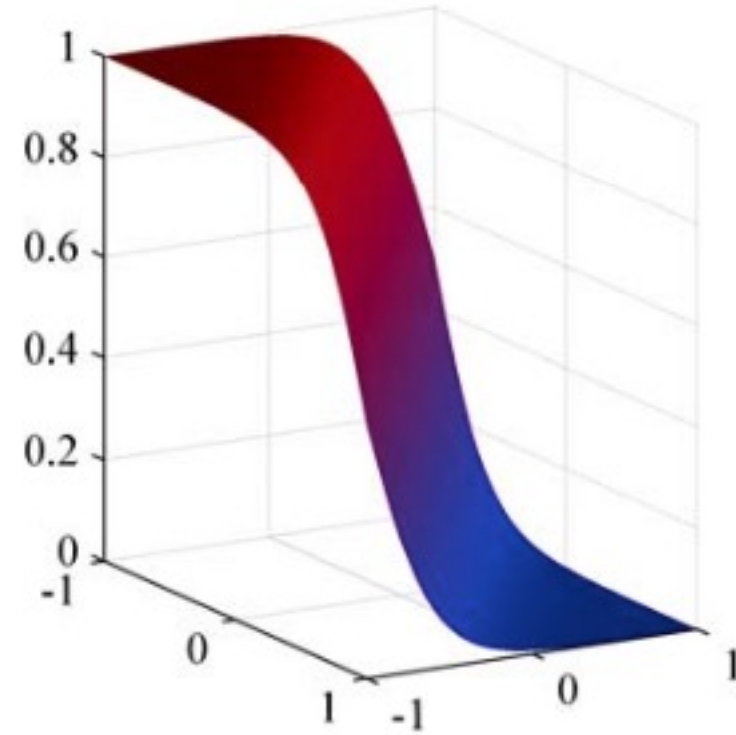
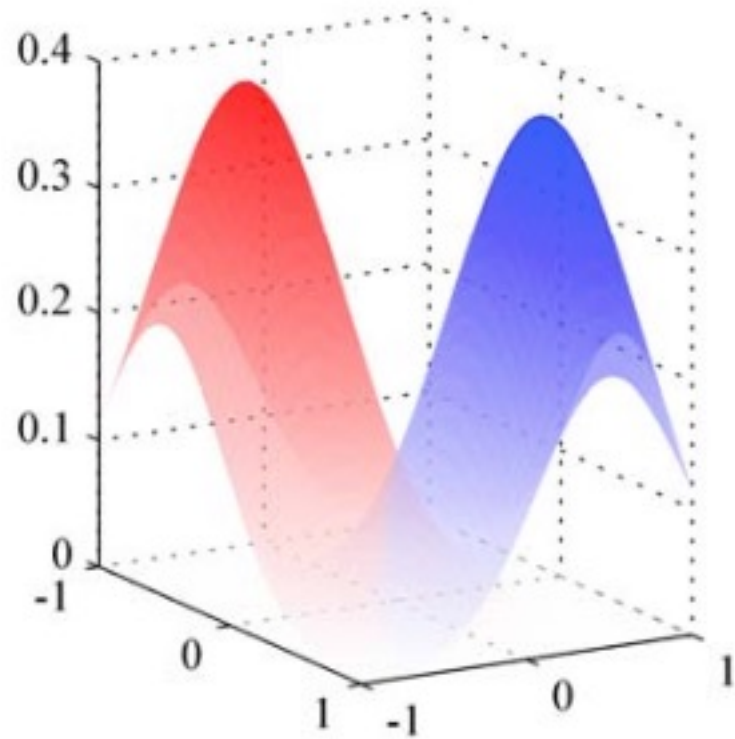
where

$$\begin{aligned}
 \mathbf{w}^T &= (\boldsymbol{\mu}_1^T - \boldsymbol{\mu}_2^T) \boldsymbol{\Sigma}^{-1}, \\
 w_0 &= \frac{1}{2} (\boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2) + \ln \frac{\pi_1}{\pi_2},
 \end{aligned} \tag{5-13}$$

and

$$a = \mathbf{w}^T \mathbf{x} + w_0.$$

5.3. Probabilistic Generative Models



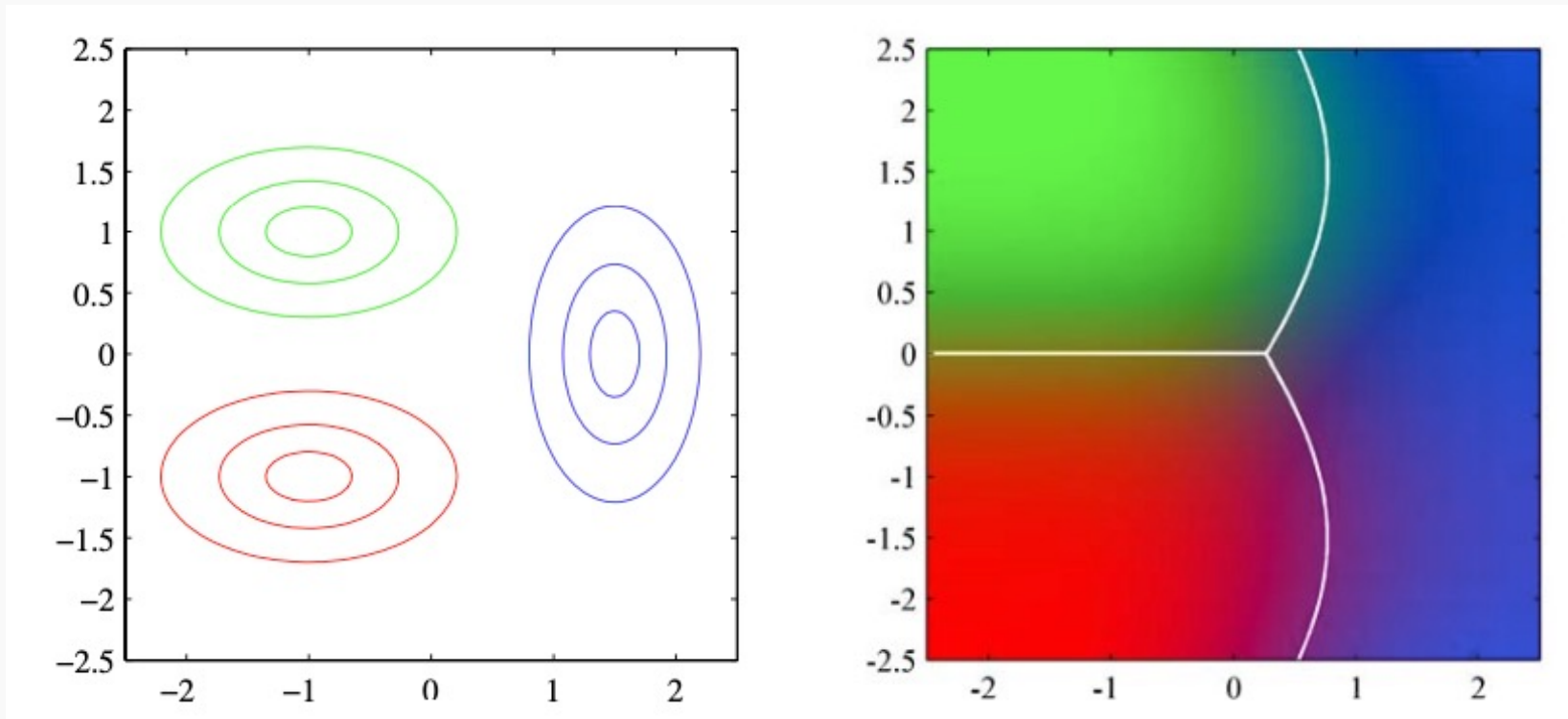
5.3. Probabilistic Generative Models

For the case of $K > 2$ classes,

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

where $\mathbf{w}_k = \Sigma^{-1} \boldsymbol{\mu}_k$ and $w_{k0} = -\frac{1}{2} \boldsymbol{\mu}_k^T \Sigma^{-1} \boldsymbol{\mu}_k + \ln p(C_k)$.

If each class has its own covariance matrix Σ_k , then $\mathbf{x} \Sigma_k^{-1} \mathbf{x}$ term will not be vanished.





5.3. Probabilistic Generative Models

The parameters can be estimated using MLE if the data set comprises observes of \mathbf{x} along with their corresponding class labels.

Consider a binary class dataset, $\{\mathbf{x}_n, t_n\}$ where $n = 1, \dots, N$.

- Each class having a Gaussian class-conditional density with a shared Σ .
- $t_n = \{1, 0\}$ for C_1 and C_2 , respectively.
- Let the prior class probability be $p(C_1) = \pi$ and $p(C_2) = 1 - \pi$.



5.3. Probabilistic Generative Models

For a data point \mathbf{x}_n from class C_1 , $t_n = 1$ and hence

$$p(\mathbf{x}_n, C_1) = p(C_1)p(\mathbf{x}_n|C_1) = \pi\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}).$$

Similarly,

$$p(\mathbf{x}_n, C_2) = p(C_2)p(\mathbf{x}_n|C_2) = (1 - \pi)\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}).$$

Then the likelihood function is given by

$$\begin{aligned} p(\mathbf{t}|\pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) &= \prod_{n=1}^N [\pi\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_1, \boldsymbol{\Sigma})]^{t_n} [(1 - \pi)\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_2, \boldsymbol{\Sigma})]^{1-t_n} \\ &\propto \left(\pi \exp \left[-\frac{1}{2}(\mathbf{x}_n - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu}_1) \right] \right)^{t_n} \left((1 - \pi) \exp \left[-\frac{1}{2}(\mathbf{x}_n - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu}_2) \right] \right)^{1-t_n}, \end{aligned}$$

where log-likelihood terms that depend on π are

$$\sum_{n=1}^N \{t_n \ln \pi + (1 - t_n) \ln(1 - \pi)\}.$$

(5-14)



5.3. Probabilistic Generative Models

Setting the derivative w.r.t. $\pi = 0$,

$$\pi = \frac{1}{N} \sum_{n=1}^N t_n = \frac{N_1}{N}$$

where N_1 is the number of data points in class C_1 .

Consider the log-likelihood function terms that depend on $\boldsymbol{\mu}_1$,

$$\sum_{n=1}^N t_n \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}) = -\frac{1}{2} \sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_1) + \text{Const.}$$

The derivative w.r.t. $\boldsymbol{\mu}_1$ will have

$$\boldsymbol{\mu}_1 = \frac{1}{N_1} \sum_{n=1}^N t_n \mathbf{x}_n ,$$

the mean of all the input vectors \mathbf{x}_n assigned to C_1 .

Similarly,

$$\boldsymbol{\mu}_2 = \frac{1}{N_2} \sum_{n=1}^N (1 - t_n) \mathbf{x}_n .$$



5.3. Probabilistic Generative Models

The covariance then can be estimated from the derivative of log-likelihood w.r.t. Σ ,

$$\Sigma = \frac{N_1}{N} \mathbf{S}_1 + \frac{N_2}{N} \mathbf{S}_2$$

where

$$\mathbf{S}_1 = \frac{1}{N_1} \sum_{n \in C_1} (\mathbf{x}_n - \boldsymbol{\mu}_1)(\mathbf{x}_n - \boldsymbol{\mu}_1)^T$$

and

$$\mathbf{S}_2 = \frac{1}{N_2} \sum_{n \in C_2} (\mathbf{x}_n - \boldsymbol{\mu}_2)(\mathbf{x}_n - \boldsymbol{\mu}_2)^T.$$

5.3. Probabilistic Generative Models

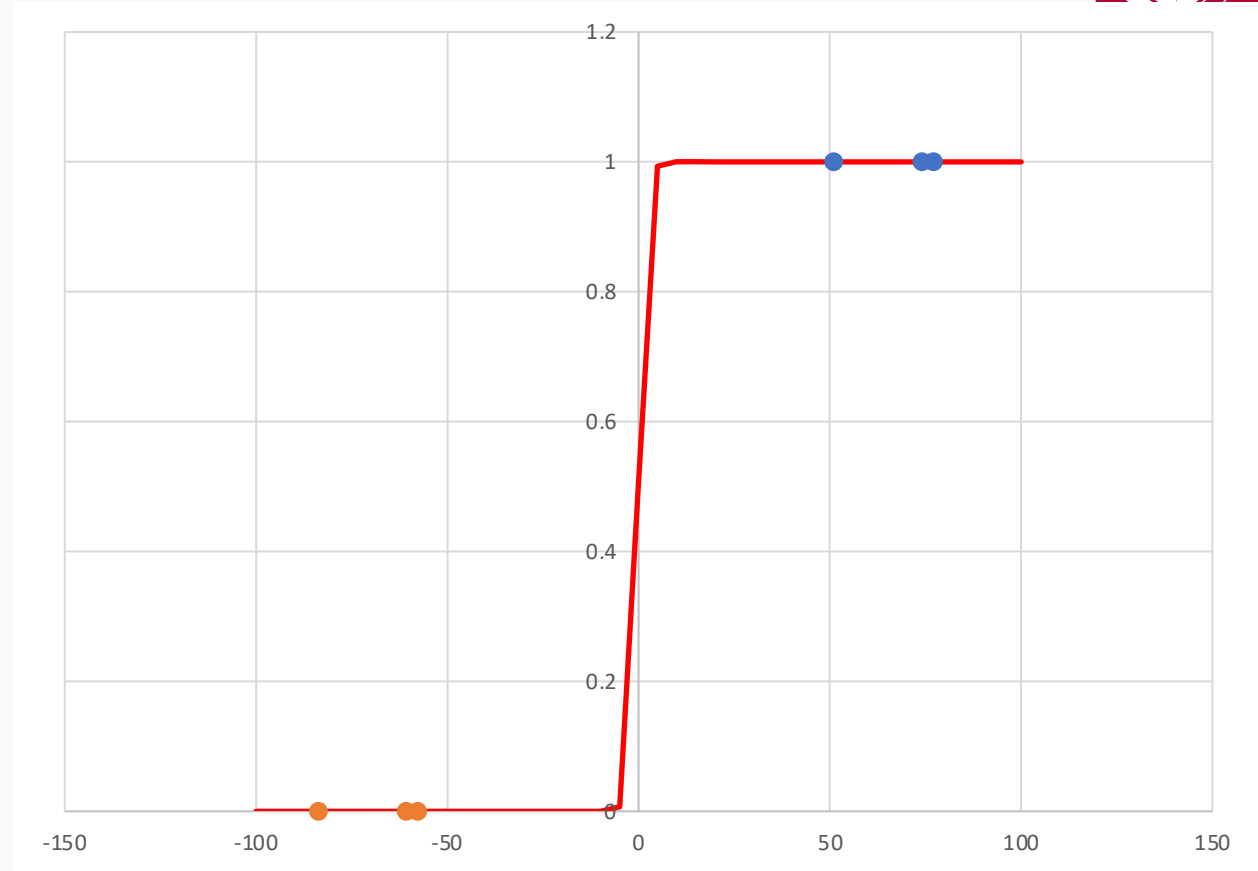
| Class | (x_1, x_2) |
|-------|------------------------|
| 1 | (1,2), (2,3), (3,4.9) |
| -1 | (2,1), (3,2), (4, 3.9) |

$$\Sigma = \begin{bmatrix} 0.667 & 0.967 \\ 0.967 & 1.467 \end{bmatrix}$$

$$\Sigma^{-1} = \begin{bmatrix} 48.22 & -32.22 \\ -32.22 & 22.22 \end{bmatrix}$$

$$\mathbf{w} = (\mu_1 - \mu_2)\Sigma^{-1} = [48.667, -80.444, 54.444]$$

$$\hat{\mathbf{w}} = \frac{\mathbf{w}}{|\mathbf{w}|} = [0.448, -0.828, 0.560]$$





5.3. Probabilistic Generative Models

- Consider a case of discrete feature values x_i .
- For the simplicity, let $x_i \in \{0,1\}$ and there are D many inputs.
- Assume the *naïve Bayes*, **the feature values are independent**.
- The class-conditional distribution is

$$p(\mathbf{x}|C_k) = \prod_{i=1}^D \mu_{ki}^{x_i} (1 - \mu_{ki})^{1-x_i}.$$

- Substituting into the quantities a_k ,

$$a_k(\mathbf{x}) = \sum_{n=1}^D \{x_i \ln \mu_{ki} + (1 - x_i) \ln(1 - \mu_{ki})\} + \ln p(C_k)$$

- which is the linear function of the input values x_i .



5.3. Probabilistic Generative Models

Each class C_k is described by its own linear model

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

where $k = 1, \dots, K$, $\tilde{\mathbf{w}}_k = (w_{k0}, \mathbf{w}_k^T)^T$, and $\tilde{\mathbf{x}}$ is the corresponding augmented input vector $(1, \mathbf{x}^T)^T$. In the matrix formation, the expression becomes much simpler as below

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

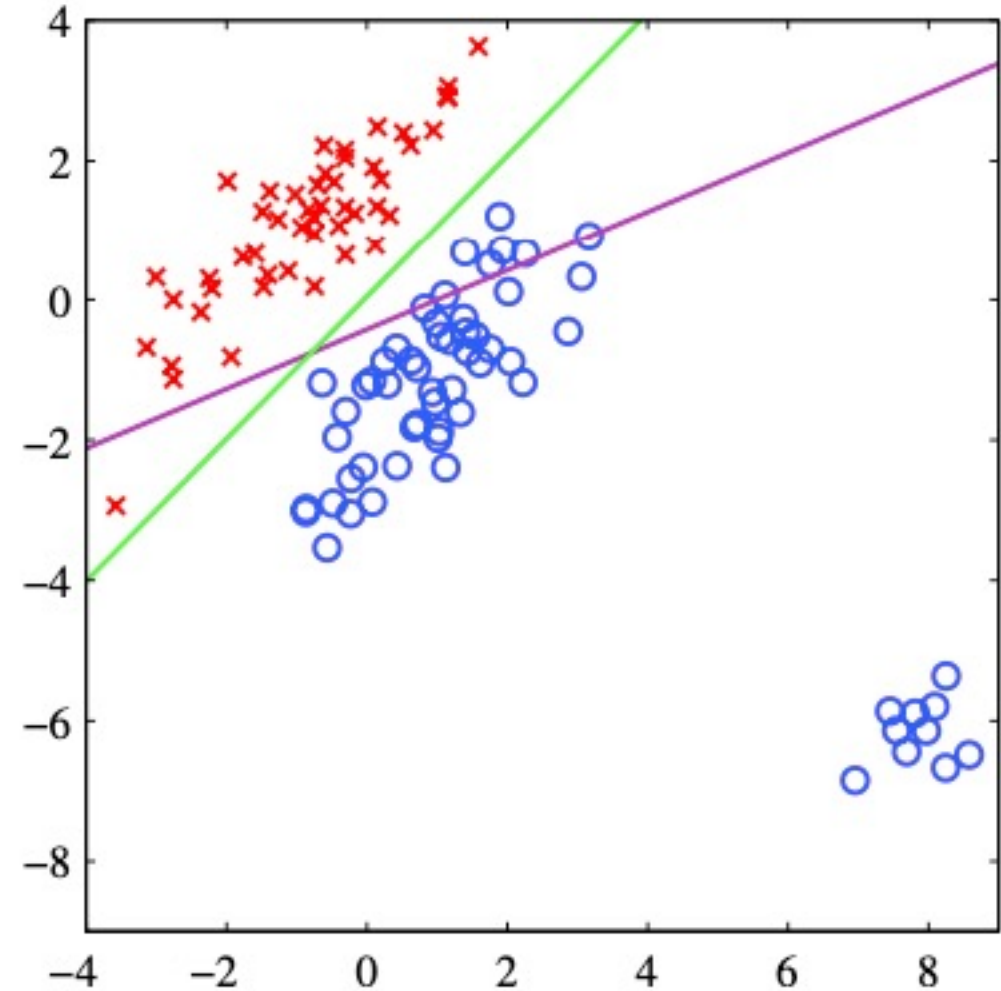
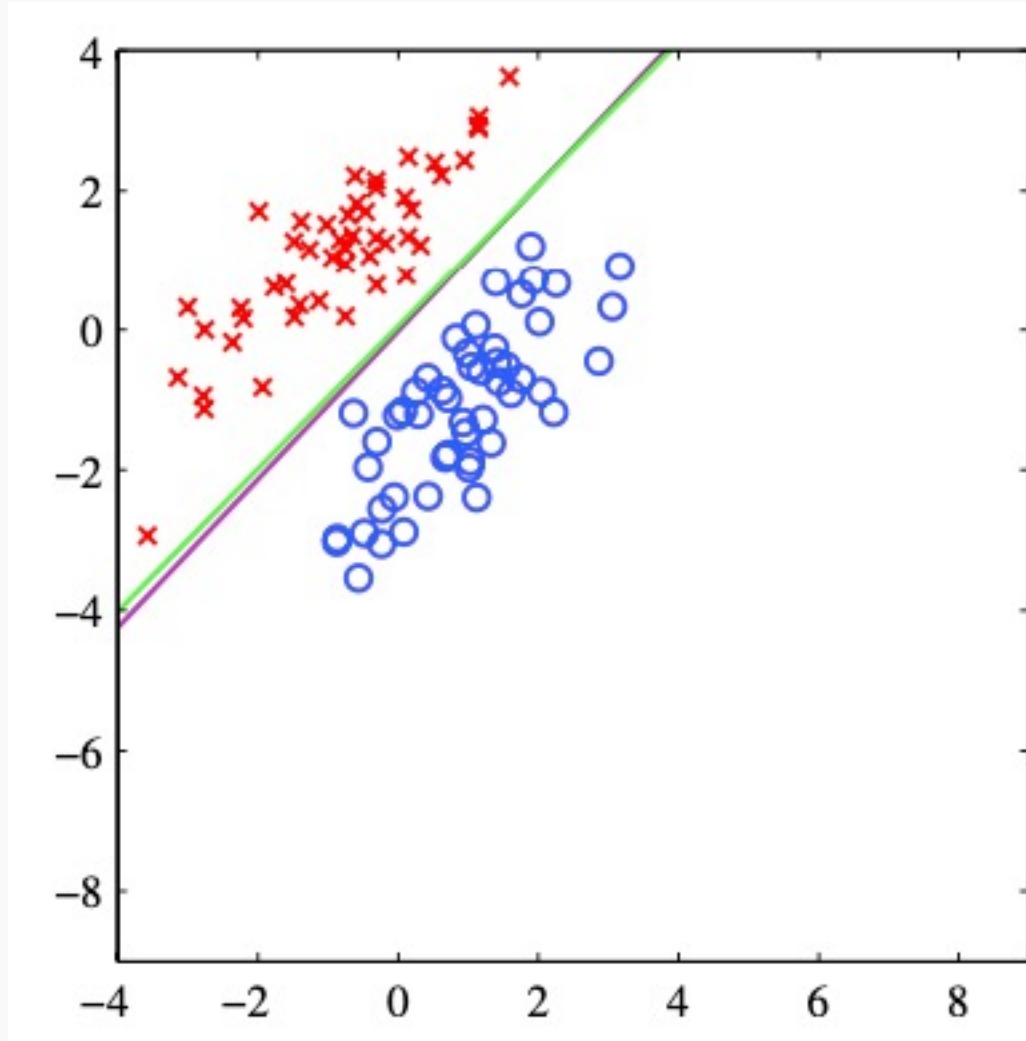
Consider a training data set $\{\mathbf{x}_n, \mathbf{t}_n\}$ where $n = 1, \dots, N$, the *sum-of-squares* error function

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

can set the derivative w.r.t. $\tilde{\mathbf{W}}$ to zero to have

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

5.3. Probabilistic Generative Models





5.0. Lecture 4 Review

5.1. Introduction

5.2. Discriminant Functions

5.2.1. Linear Discriminant Analysis (LDA)

5.2.2. Perceptron

5.3. Probabilistic Generative Models - MLE

5.4. Probabilistic Discriminative Models - Logistic Regression



5.4 Probabilistic Discriminative Models – Logistic Regression

Consider a binary class problem.

Assume the posterior probability of class C_1 can be written as a logistic sigmoid acting on a linear function of the feature vector ϕ ,

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

with $p(C_2|\phi) = 1 - p(C_1|\phi)$. This model is known as *logistic regression*.

The parameters of the logistic regression model can be determined using maximum likelihood,

$$\frac{d\sigma}{da} = \sigma(1 - \sigma).$$



5.4 Probabilistic Discriminative Models – Logistic Regression

For a data set $\{\phi_n, t_n\}$, where

- $t_n \in \{0,1\}$
- $\phi_n = \phi(\mathbf{x}_n)$ with $n = 1, \dots, N$.

The likelihood function can be written

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n}$$

where $\mathbf{t} = (t_1, \dots, t_N)^T$ and $y_n = p(C_1|\phi_n)$.

The *cross-entropy* error function is the negative log-likelihood in the form

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w}) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$$

where $y_n = \sigma(\mathbf{w}^T \phi_n)$ and the gradient of the error function w.r.t. \mathbf{w} is

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

The maximum likelihood solution occurs when the hyperplane corresponding to $\sigma = 0.5$, equivalent to $\mathbf{w}^T \phi = 0$.



5.4 Probabilistic Discriminative Models – Logistic Regression

For the precise measurement, the error function can be minimized by an efficient iterative technique based on the *Newton-Raphson* iterative optimization, which uses a local quadratic approximation to the log-likelihood function in the form of

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

where \mathbf{H} is the Hessian matrix whose elements are the second derivatives of $E(\mathbf{w})$ w.r.t. \mathbf{w} :

$$\begin{aligned} \mathbf{H} &= \nabla \nabla E(\mathbf{w}) = \nabla \left(\sum_{n=1}^N (\mathbf{w}^T \phi_n - t_n) \phi_n \right) \\ &= \nabla (\phi^T \phi \mathbf{w} - \phi^T \mathbf{t}) = \phi^T \phi. \end{aligned}$$



5.4 Probabilistic Discriminative Models – Logistic Regression

Alternatively, we can update the cross-entropy error function as

$$\mathbf{H} = \nabla \left(\boldsymbol{\phi}^T (\mathbf{y} - \mathbf{t}) \right) = \sum_{n=1}^N y_n (1 - y_n) \boldsymbol{\phi}_n \boldsymbol{\phi}_n^T = \boldsymbol{\phi}^T \mathbf{R} \boldsymbol{\phi}$$

where \mathbf{R} is the $N \times N$ weight matrix with elements $R_{nn} = y_n(1 - y_n)$. The error function is not quadratic longer but is not constant.

The Newton-Raphson update formula for the logistic regression model then becomes

$$\begin{aligned} \mathbf{w}^{new} &= \mathbf{w}^{old} - (\boldsymbol{\phi}^T \mathbf{R} \boldsymbol{\phi})^{-1} \boldsymbol{\phi}^T (\mathbf{y} - \mathbf{t}) \\ &= (\boldsymbol{\phi}^T \mathbf{R} \boldsymbol{\phi})^{-1} \{ \boldsymbol{\phi}^T \mathbf{R} \boldsymbol{\phi} \mathbf{w}^{(old)} - \boldsymbol{\phi}^T (\mathbf{y} - \mathbf{t}) \} \\ &= (\boldsymbol{\phi}^T \mathbf{R} \boldsymbol{\phi})^{-1} \boldsymbol{\phi}^T \mathbf{R} \mathbf{z} \end{aligned}$$

where \mathbf{z} is an N -dimensional vector with elements

$$\mathbf{z} = \boldsymbol{\phi} \mathbf{w}^{(old)} - \mathbf{R}^{-1} (\mathbf{y} - \mathbf{t}).$$



5.4 Probabilistic Discriminative Models – Logistic Regression

| | | | | | | | |
|------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| w0 | 1 | -1 | 1 | 0.577350269 | -0.57735027 | 0.577350269 | |
| likelihood | -2.93486703 | | | | | | |
| z | 2 | 2 | 2.9 | 0 | 0 | 0.9 | |
| sigmoid | 0.880797078 | 0.880797078 | 0.947846437 | 0.5 | 0.5 | 0.710949503 | |
| y | 1 | 1 | 1 | 0 | 0 | 1 | |
| grad | -1.4203901 | -4.82972856 | -3.42113599 | | | | |
| H | 0.964921047 | 2.535282692 | 2.318643613 | H_inv | 9.296241883 | -2.81430131 | -0.36016317 |
| | 2.535282692 | 7.507874965 | 6.772426941 | | -2.81430131 | 1.979349879 | -0.9930948 |
| | 2.318643613 | 6.772426941 | 6.927476275 | | -0.36016317 | -0.9930948 | 1.235767788 |
| Delta | 1.620188714 | -2.16480455 | 1.080220871 | | | | |
| w_update | 2.620188714 | -3.16480455 | 2.080220871 | 0.568951952 | -0.6872107 | 0.451702475 | |

After 6 iterations:

| | | | | | | | |
|------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| w_update | 5.713665154 | -8.53726067 | 5.768785273 | 0.484958884 | -0.72461726 | 0.489637316 | |
| likelihood | -0.00960079 | | | | | | |
| z | 8.713975027 | 5.945499627 | 8.368930973 | -5.59207092 | -8.36054632 | -5.93711497 | |
| sigmoid | 0.999835753 | 0.99738924 | 0.99976809 | 0.00371346 | 0.000233862 | 0.002632685 | |
| y | 1 | 1 | 1 | 0 | 0 | 0 | |
| grad | -0.00357309 | -0.01257775 | -0.00515152 | | | | |
| H | 0.009559251 | 0.024671452 | 0.02368409 | H_inv | 949.3811091 | -344.421922 | 17.81857682 |
| | 0.024671452 | 0.071581703 | 0.069124327 | | -344.421922 | 263.5807941 | -135.555362 |
| | 0.02368409 | 0.069124327 | 0.074231849 | | 17.81857682 | -135.555362 | 134.0146649 |
| Delta | 0.848035592 | -1.38628593 | 0.950934258 | | | | |

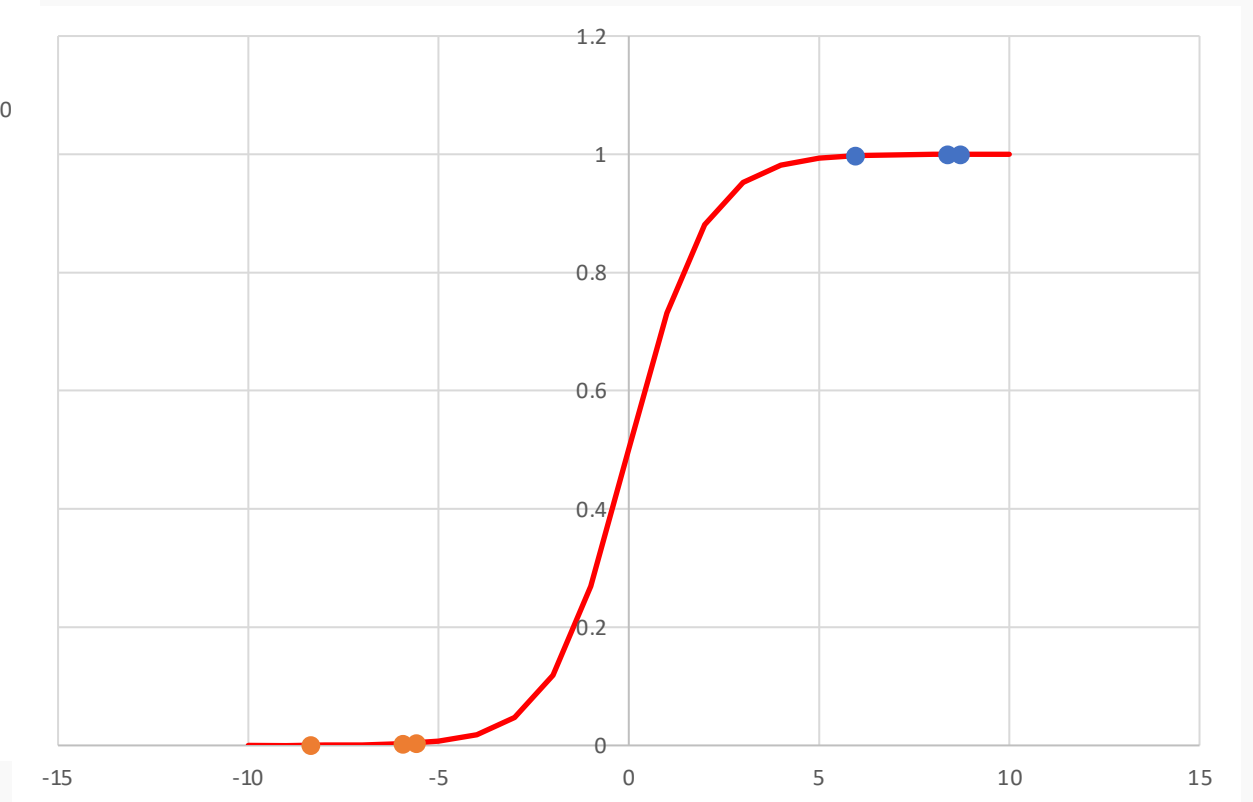
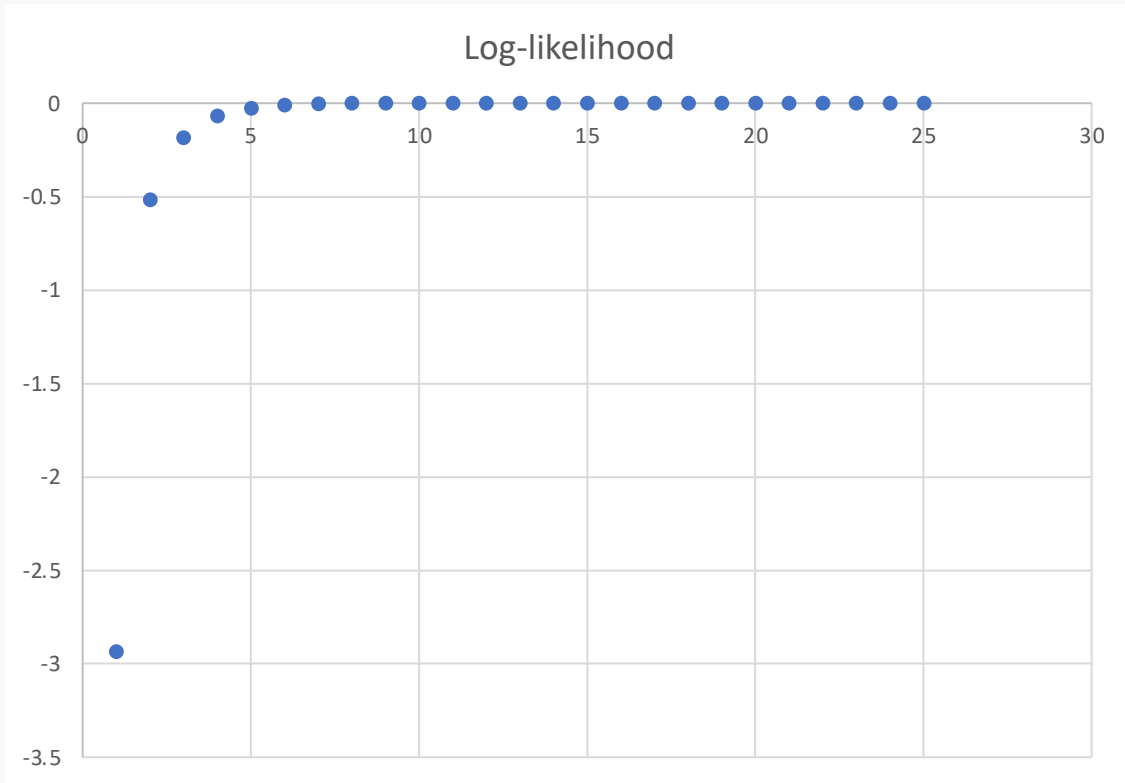
$$\nabla E = (t - \sigma(\mathbf{w}^T \mathbf{x}))\mathbf{x}$$

$$\mathbf{H} = \begin{bmatrix} \frac{\partial^2 \nabla E}{\partial x_0^2} & \frac{\partial^2 \nabla E}{\partial x_0 \partial x_1} & \frac{\partial^2 \nabla E}{\partial x_0 \partial x_2} \\ \frac{\partial^2 \nabla E}{\partial x_1 \partial x_0} & \frac{\partial^2 \nabla E}{\partial x_1^2} & \frac{\partial^2 \nabla E}{\partial x_1 \partial x_2} \\ \frac{\partial^2 \nabla E}{\partial x_2 \partial x_0} & \frac{\partial^2 \nabla E}{\partial x_2 \partial x_1} & \frac{\partial^2 \nabla E}{\partial x_2^2} \end{bmatrix}$$

$$\Delta = \mathbf{H}^{-1} \nabla E$$

$$\mathbf{w}^{(new)} = \mathbf{w}^{(old)} + \Delta$$

5.4 Probabilistic Discriminative Models – Logistic Regression





5.4 Probabilistic Discriminative Models – Logistic Regression

The posterior probabilities are given by a softmax transformation of linear functions of the feature variables

$$p(C_k|\phi) = y_k(\phi) = \frac{\exp(\mathbf{w}_k^T \phi)}{\sum_j \exp(\mathbf{w}_k^T \phi)}.$$

The derivative of y_k w.r.t. the activation function $a_k = \mathbf{w}_k^T \phi$ is

$$\frac{\partial y_k}{\partial a_j} = y_k(I_{kj} - y_j)$$

where I_{kj} are the elements of the identity matrix.



5.4 Probabilistic Discriminative Models – Logistic Regression

The target variables t_{nk} form a $N \times K$ matrix \mathbf{T} , can express the likelihood function as

$$p(\mathbf{T}|\mathbf{w}_1, \dots, \mathbf{w}_K) = \prod_{n=1}^N \prod_{k=1}^K p(C_k|\phi_n)^{t_{nk}} = \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}}$$

where $y_{nk} = y_k(\phi_n)$.

The negative logarithm gives the error function as

$$E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\ln p(\mathbf{T}|\mathbf{w}_1, \dots, \mathbf{w}_K) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk}$$

which is known as the *cross-entropy* error function for the multiclass classification problem.



5.4 Probabilistic Discriminative Models – Logistic Regression

The gradient of the error function w.r.t. one of the parameter vectors \mathbf{w}_j result the derivative of the softmax function,

$$\nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = \sum_{n=1}^N (y_{nj} - t_{nj}) \phi_n$$

and the parameter vectors can be found using via the iterative least square following Newton-Raphson method where the Hessian matrix is

$$\nabla_{\mathbf{w}_k} \nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = - \sum_{n=1}^N y_{nk} (I_{kj} - y_{nj}) \phi_n \phi_n^T .$$