

# Linear Classification

CSE 4309 – Machine Learning

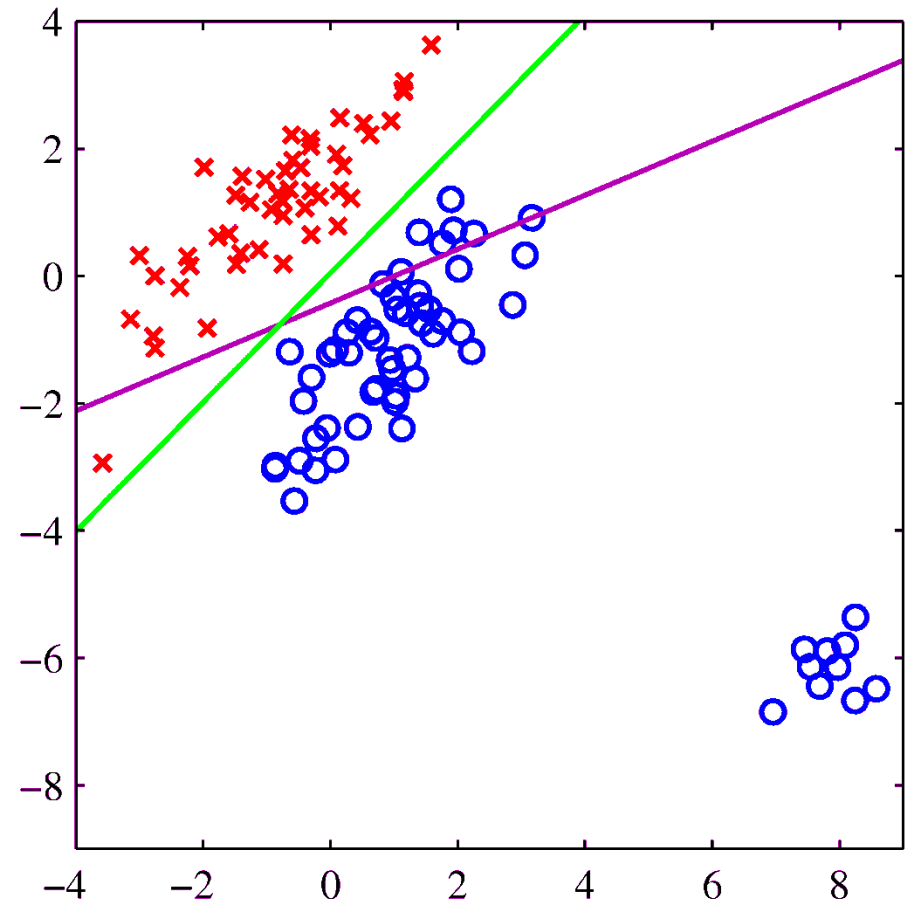
Vassilis Athitsos

Computer Science and Engineering Department

University of Texas at Arlington

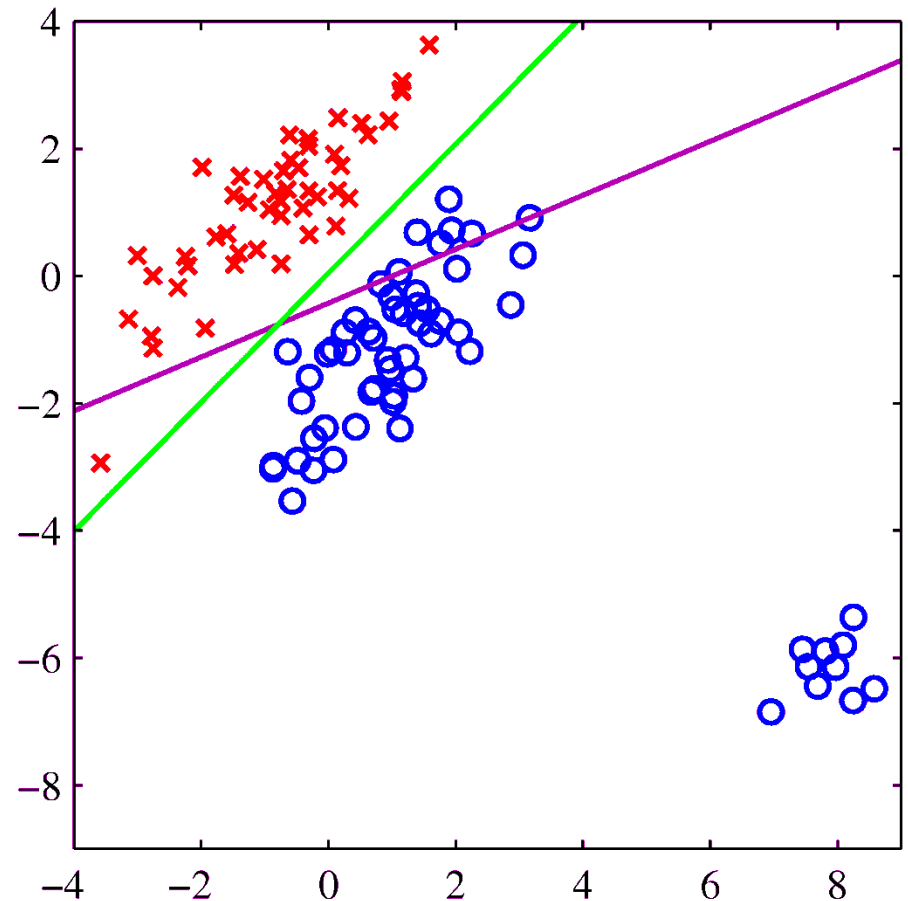
# Example of Linear Classification

- Red points: patterns belonging to class  $C_1$ .
- Blue points: patterns belonging to class  $C_2$ .
- Goal: find a **linear decision boundary** separating  $C_1$  from  $C_2$ .
- Points on one side of the line will be classified as belonging to  $C_1$ , points on the other side will be classified as  $C_2$ .
- The red line is one example of such a decision boundary.
  - It misclassified a few patterns.
- The green line is another example.



# Linear Classification

- Mathematically, assuming input patterns are  $D$ -dimensional vectors:
  - We are looking for a **decision boundary** in the form of a  $(D-1)$ -dimensional hyperplane separating the two classes.
  - Points on one side of the hyperplane will be classified as belonging to  $C_1$ , points on the other side will be classified as  $C_2$ .
- If inputs are 2-dimensional vectors, the decision boundary is a line.
- If inputs are 3-dimensional vectors, the decision boundary is a 2-dimensional surface.



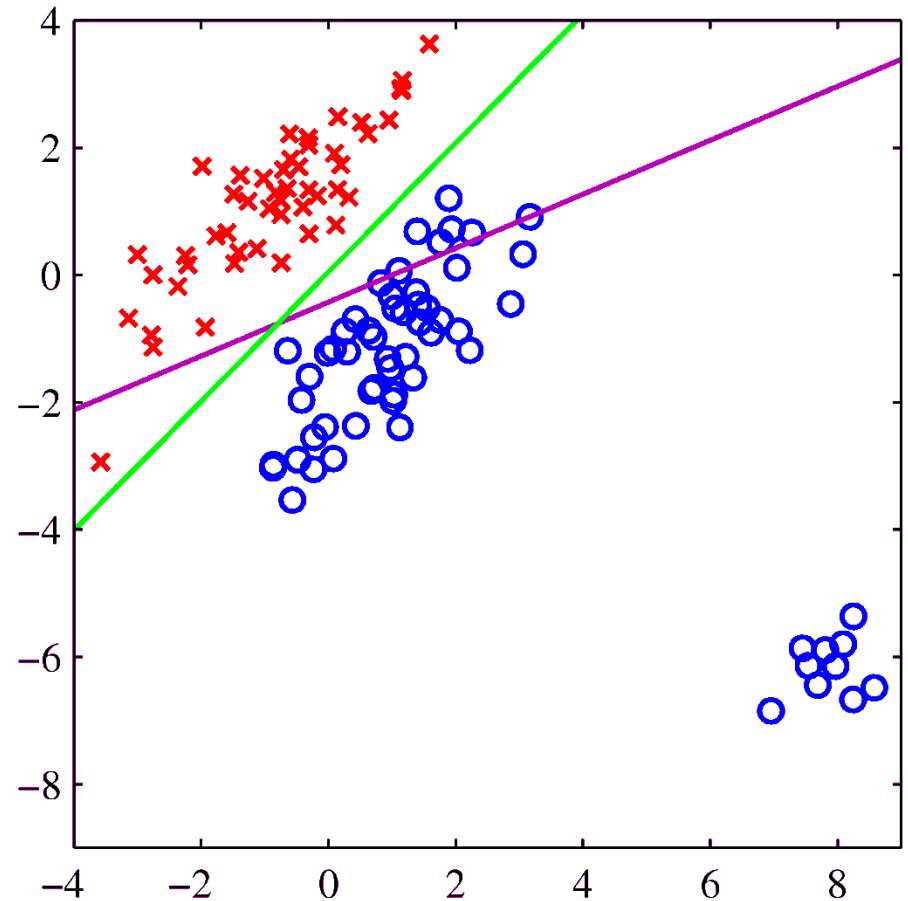
# Linear Classification

- Input space:  $\mathbb{R}^D$
- The decision boundary is a (D-1)-dimensional hyperplane defined using this equation:

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

- In the equation above:  $w_0$  and  $c$  are real numbers,  $\mathbf{w}$  and  $\mathbf{x}$  are column vectors:

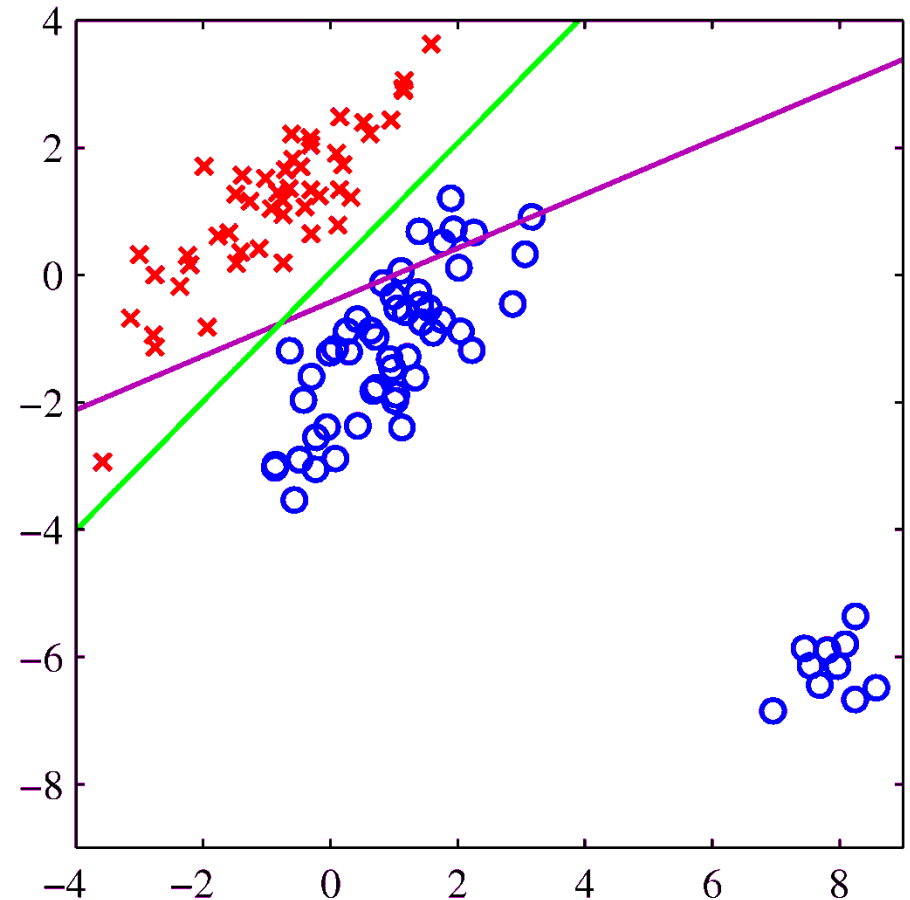
$$\mathbf{w} = \begin{bmatrix} w_1 \\ \dots \\ w_D \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} x_1 \\ \dots \\ x_D \end{bmatrix}$$



# Linear Classification

- In order to be able to use vector and matrix notation, we extend vector  $\mathbf{w}$  and  $\mathbf{x}$ , as follows:

$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \dots \\ w_D \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ \dots \\ x_D \end{bmatrix}$$



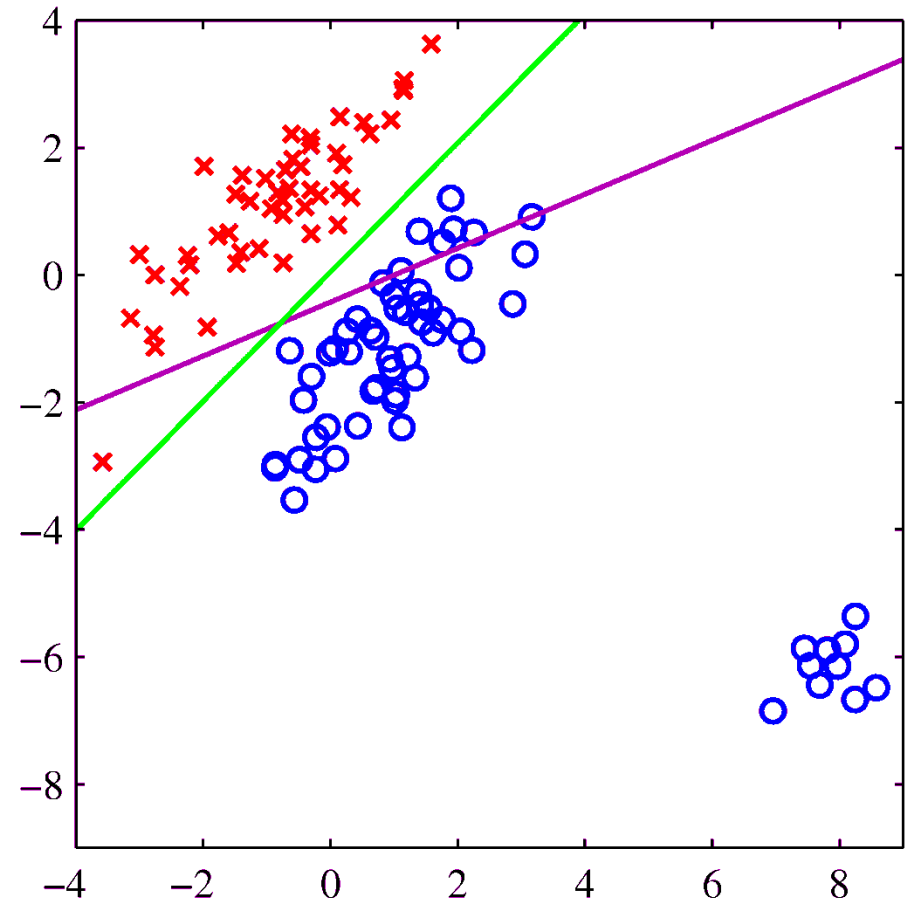
- Then, the decision boundary is defined by equation:

$$\mathbf{w}^T \mathbf{x} = c$$

# Linear Classification

- Furthermore, we can use basis functions  $\varphi_i(x)$ , which we pick manually.
- As before,  $\varphi_0(x) = 1$ .

$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \dots \\ w_D \end{bmatrix} \quad \varphi(x) = \begin{bmatrix} 1 \\ \varphi_1(x) \\ \dots \\ \varphi_D(x) \end{bmatrix}$$



- Then, the decision boundary is defined by equation:

$$\mathbf{w}^T \varphi(x) = c$$

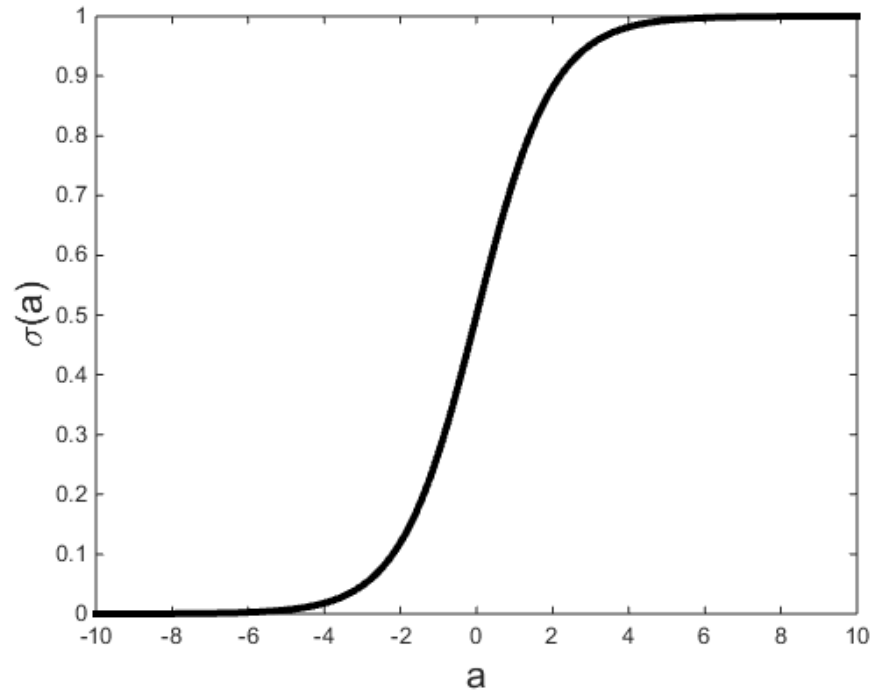
# Logistic Regression

- Despite the word "regression" in its name, **logistic regression** is actually a linear classification method.
- In logistic regression, the goal is to compute a decision boundary  $\mathbf{w}^T \varphi(x) = 0$ .
- The classification function itself is defined as:

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

where  $\sigma$  is the sigmoidal function we have seen before:

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

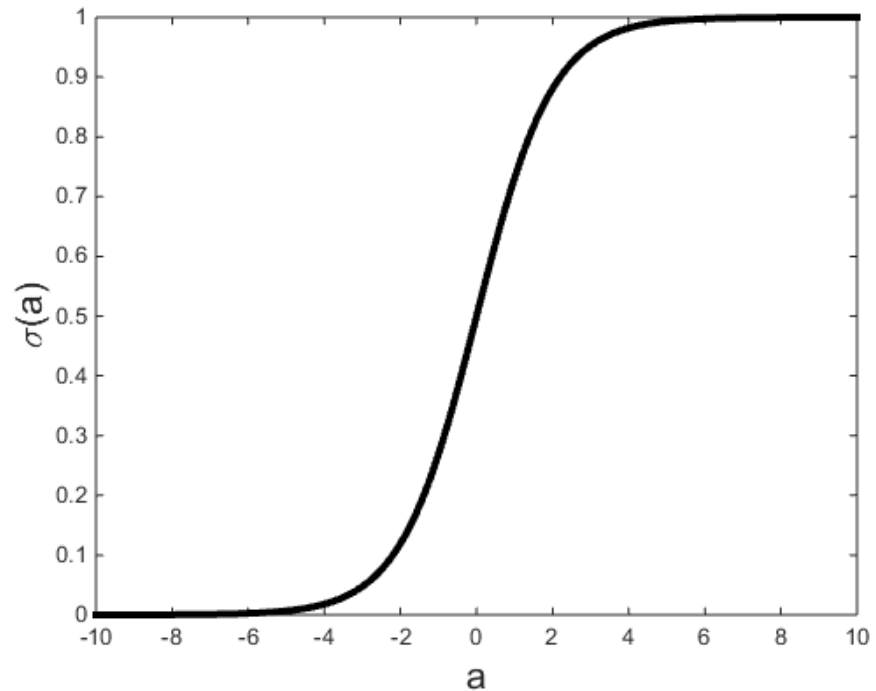


# Logistic Regression

- In logistic regression, the goal is to compute a decision boundary  $\mathbf{w}^T \varphi(x) = 0$ .
- The classification function itself is defined as:

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

- We want to separate two classes, denoted as  $C_0$  and  $C_1$ .
- We have a set of training inputs:  $X = \{x_1, \dots, x_N\}$
- We also have a set of corresponding outputs:  $\mathbf{t} = \{t_1, \dots, t_N\}$
- Each  $t_n$  is either 0 or 1:
  - 0 corresponds to class  $C_0$ .
  - 1 corresponds to class  $C_1$ .





# Logistic Regression

- Classification function:

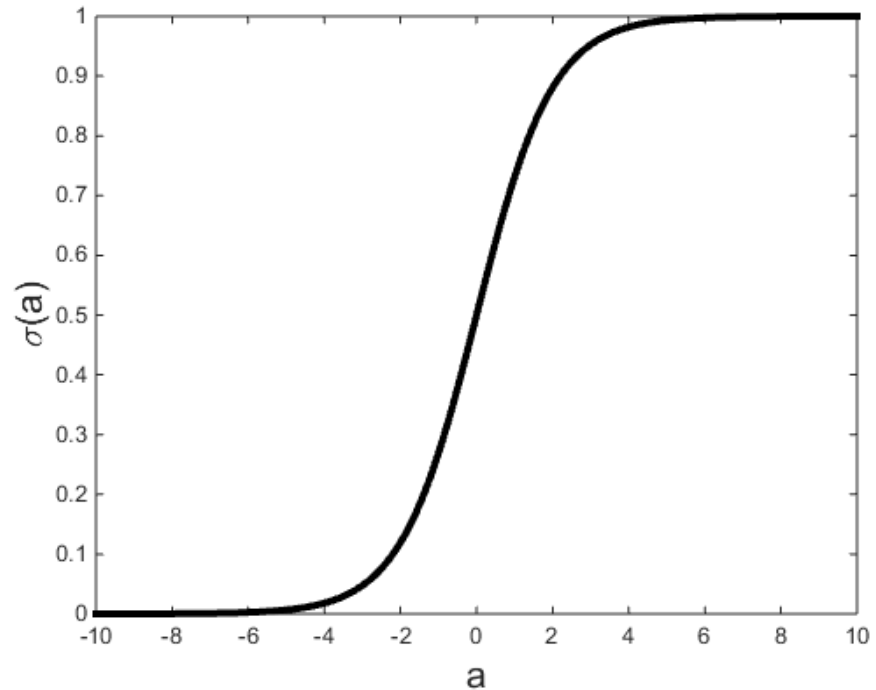
$$y(x) = \sigma(\mathbf{w}^T \varphi(x))$$

- We treat the output of  $\sigma(\mathbf{w}^T \varphi(x))$  as an estimate of probability  $p(C_1|x)$ .

- Then,  $p(C_0|x) = 1 - p(C_1|x)$ .

- When  $\mathbf{w}^T \varphi(x) > 0$ ,  $p(C_1|x) > 0.5$ , so  $C_1$  is more likely than  $C_0$ .

- When  $\mathbf{w}^T \varphi(x) < 0$ ,  $p(C_0|x) < 0.5$ , so  $C_0$  is more likely than  $C_1$ .



# Finding the Most Likely Solution

- Our rationale is almost identical to the rationale we used for finding the most likely solution for linear regression.
- Suppose we have a set of training inputs:  $X = \{x_1, \dots, x_N\}$
- We also have a set of corresponding outputs:  $\mathbf{t} = \{t_1, \dots, t_N\}$
- We assume that training inputs are independent of each other.
- We assume that outputs are conditionally independent of each other, given their inputs and our estimate of  $\mathbf{w}$ .
- Then:

$$p(\mathbf{w}|X, \mathbf{t}) = \frac{p(\mathbf{t}|X, \mathbf{w}) * p(\mathbf{w}|X)}{p(\mathbf{t}|X)}$$

# Finding the Most Likely Solution

$$p(\mathbf{w}|X, \mathbf{t}) = \frac{p(\mathbf{t}|X, \mathbf{w}) * p(\mathbf{w}|X)}{p(\mathbf{t}|X)}$$

- We assume that, given  $X$ , all values of  $\mathbf{w}$  are equally likely.
- Then,  $\frac{p(\mathbf{w}|X)}{p(\mathbf{t}|X)}$  is a constant that does not depend on  $\mathbf{w}$ .
- Therefore, finding the  $\mathbf{w}$  that maximizes  $p(\mathbf{w}|X, \mathbf{t})$  is the same as finding the  $\mathbf{w}$  that maximizes  $p(\mathbf{t}|X, \mathbf{w})$ .
- $p(\mathbf{t}|X, \mathbf{w})$  is the **likelihood** of the training data.
- So, to find the most likely answer  $\mathbf{w}$ , we must find the value of  $\mathbf{w}$  that maximizes the likelihood of the training data.

# Probability of $t_n$

- We want to find the  $\mathbf{w}$  that maximizes  $p(\mathbf{t} | X, \mathbf{w})$ .

$$p(\mathbf{t} | X, \mathbf{w}) = \prod_{n=1}^n p(t_n | x_n, \mathbf{w})$$

- How can we compute  $p(t_n | x_n, \mathbf{w})$ ?
- Remember: we assume that  $y(\mathbf{x}, \mathbf{w})$  is an estimate of  $p(C_1 | x)$ .
- Under this assumption, for a given  $x_n$ :
  - The probability that  $t_n$  (the correct class for  $x_n$ ) is 1 is  $y(\mathbf{x}, \mathbf{w})$ .
  - The probability that  $t_n$  (the correct class for  $x_n$ ) is 0 is  $\{1 - y(\mathbf{x}, \mathbf{w})\}$ .

# Probability of $t_n$

- To simplify notation, define  $y_n = y(x_n, \mathbf{w})$ .
- Under this assumption, for a given  $x_n$ :
  - $p(t_n = 1|x_n, w) = y_n$ .
  - $p(t_n = 0|x_n, w) = (1 - y_n)$ .
- So:
  - If  $t_n = 1$ , then  $p(t_n|x_n, w) = y_n$ .
  - If  $t_n = 0$ , then  $p(t_n|x_n, w) = (1 - y_n)$ .
- Regardless of whether  $t_n$  is 1 or 0, the following holds:

$$p(t_n|x_n, w) = (y_n)^{t_n}(1 - y_n)^{1-t_n}$$

# Probability of $t_n$

- Regardless of whether  $t_n$  is 1 or 0, the following holds:

$$p(t_n|x_n, w) = (y_n)^{t_n}(1 - y_n)^{1-t_n}$$

- To make sense of this formula, we just need to verify that it is correct when  $t_n = 1$  and when  $t_n = 0$ .
- If  $t_n = 1$ , then  $p(t_n|x_n, w) = y_n$ . The above formula becomes:

$$\begin{aligned} p(t_n|x_n, w) &= (y_n)^{t_n}(1 - y_n)^{1-t_n} \\ &= (y_n)^1(1 - y_n)^{1-1} \\ &= y_n(1 - y_n)^0 = y_n \end{aligned}$$

- So, the formula is verified for the case where  $t_n = 1$ .

# Probability of $t_n$

- Regardless of whether  $t_n$  is 1 or 0, the following holds:

$$p(t_n|x_n, w) = (y_n)^{t_n}(1 - y_n)^{1-t_n}$$

- If  $t_n = 0$ , then  $p(t_n|x_n, w) = (1 - y_n)$ . The above formula becomes:

$$\begin{aligned} p(t_n|x_n, w) &= (y_n)^{t_n}(1 - y_n)^{1-t_n} \\ &= (y_n)^0(1 - y_n)^{1-0} \\ &= 1(1 - y_n)^1 = (1 - y_n) \end{aligned}$$

- So, the formula is verified for the case where  $t_n = 0$ .

# Probability of $\mathbf{t}$

- We have shown that:

$$p(t_n | x_n, w) = (y_n)^{t_n} (1 - y_n)^{1-t_n}$$

- Therefore:

$$\begin{aligned} p(\mathbf{t} | X, \mathbf{w}) &= \prod_{n=1}^n p(t_n | x_n, \mathbf{w}) \\ &= \prod_{n=1}^n [(y_n)^{t_n} (1 - y_n)^{1-t_n}] \end{aligned}$$



# Log Likelihood of $\mathbf{t}$

$$p(\mathbf{t} | X, \mathbf{w}) = \prod_{n=1}^n [(y_n)^{t_n} (1 - y_n)^{1-t_n}]$$

- Therefore, the log likelihood  $\ln(p(\mathbf{t} | X, \mathbf{w}))$  is:

$$\ln(p(\mathbf{t} | X, \mathbf{w})) = \sum_{n=1}^n [\ln((y_n)^{t_n} (1 - y_n)^{1-t_n})] \Rightarrow$$

$$\ln(p(\mathbf{t} | X, \mathbf{w})) = \sum_{n=1}^n [t_n \ln(y_n) + (1 - t_n) \ln(1 - y_n)]$$

# Gradient of the Log Likelihood of $\mathbf{t}$

- We want to find the value of  $\mathbf{w}$  that maximizes

$$\ln(p(\mathbf{t} | X, \mathbf{w})) = \sum_{n=1}^n [t_n \ln(y_n) + (1 - t_n) \ln(1 - y_n)]$$

- The textbook defines  $E(\mathbf{w}) = -\ln(p(\mathbf{t} | X, \mathbf{w}))$ .
- That negative sign is not important for our purposes. It just means that to maximize  $\ln(p(\mathbf{t} | X, \mathbf{w}))$  we need to minimize  $E(\mathbf{w})$ .
- Thus, we need to find the  $\mathbf{w}$  such that the gradient  $\nabla E(\mathbf{w})$  becomes 0.

# Gradient of the Log Likelihood of $\mathbf{t}$

- We will not derive this, but it turns out that:

$$\nabla E(\mathbf{w}) = \sum_{n=1}^n [(y_n - t_n)\varphi_n]$$

- There is no closed-form solution for computing a value of  $\mathbf{w}$  so that  $\nabla E(\mathbf{w}) = 0$ .
- We have two alternatives for minimizing  $E(\mathbf{w})$ .
  - An online (sequential learning) method, where we update  $\mathbf{w}$  every time we get a new training example  $(x_n, t_n)$ .
  - A batch processing method called **iterative reweighted least squares**.

# Sequential Learning

- We follow the same approach as we did in sequential learning for linear regression.
- We first, somehow, get an initial estimate  $\mathbf{w}^{(0)}$ .
  - For example, we can initialize  $\mathbf{w}^{(0)}$  to random values between -1 and 1.
  - You can also initialize  $\mathbf{w}^{(0)}$  to be the zero vector (all entries equal to zero). That also works, I have verified it in my code.
- Then, we observe training examples, one by one.
- Every time we observe a new training example, we update the estimate.

# Sequential Learning

- The  $n^{\text{th}}$  training example contributes to the overall gradient  $\nabla E(\mathbf{w})$  a term  $\nabla E_n(\mathbf{w})$  defined as:

$$\nabla E_n(\mathbf{w}) = (y_n - t_n)\varphi(x_n)$$

- When we observe the  $n^{\text{th}}$  training example, we update the estimate from  $\mathbf{w}^{(\tau)}$  to  $\mathbf{w}^{(\tau+1)}$  as follows:

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n(\mathbf{w}) = \mathbf{w}^{(\tau)} - \eta (y_n - t_n) \varphi(x_n)$$

- $\eta$  is called the **learning rate**. It is picked manually.
- This whole process is called **stochastic gradient descent**.

# Sequential Learning - Intuition

- When we observe the  $n^{\text{th}}$  training example, we update the estimate from  $\mathbf{w}^{(\tau)}$  to  $\mathbf{w}^{(\tau+1)}$  as follows:

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n(\mathbf{w}) = \mathbf{w}^{(\tau)} - (y_n - t_n) \varphi(x_n)$$

- What is the intuition behind this update?
- The gradient  $\nabla E_n(\mathbf{w})$  is a vector that points in the direction where  $E_n(\mathbf{w})$  increases.
- Therefore, subtracting a very small amount of  $\nabla E_n(\mathbf{w})$  from  $\mathbf{w}$  should make  $E_n(\mathbf{w})$  a little bit smaller.

# Sequential Learning - Intuition

- When we observe the  $n^{\text{th}}$  training example, we update the estimate from  $\mathbf{w}^{(\tau)}$  to  $\mathbf{w}^{(\tau+1)}$  as follows:

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n(\mathbf{w}) = \mathbf{w}^{(\tau)} - (y_n - t_n) \varphi(x_n)$$

- Choosing a good value for  $\eta$  is important.
  - If  $\eta$  is too small, the minimization may happen too slowly and require too many training examples.
  - If  $\eta$  is large, the update may overfit the most recent training example, and overall  $\mathbf{w}$  may fluctuate too much from one update to the next.
- Unfortunately, picking a good  $\eta$  is more of an art than a science, and involves trial-and-error.

# Iterative Reweighted Least Squares

- Remember, we want the value of  $\mathbf{w}$  that minimizes  $E(\mathbf{w})$ , which is defined as:

$$E(\mathbf{w}) = - \sum_{n=1}^n [t_n \ln(y_n) + (1 - t_n) \ln(1 - y_n)]$$

- It can be proven that  $E(\mathbf{w})$  has a unique minimum.
- There is no closed form solution for that minimum.
- However, there is an iterative method minimizing  $E(\mathbf{w})$ .
- This method is called **iterative reweighted least squares**.



# Iterative Reweighted Least Squares

- **Iterative reweighted least squares** is an application of a more general optimization method, called **Newton-Raphson**.
- In the Newton-Raphson method, to minimize  $E(\mathbf{w})$  we do a sequence of updates on the value of  $\mathbf{w}$ , using this formula:

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

- In the above formula,  $\mathbf{H}$  is the **Hessian** matrix, whose elements are the second derivatives of  $E(\mathbf{w})$  with respect to  $\mathbf{w}$ .

# Iterative Reweighted Least Squares

- To minimize  $E(\mathbf{w})$  we iteratively update  $\mathbf{w}$ , using this formula:

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

- To use the above formula, we must compute matrix  $\mathbf{H}$ , which is the **Hessian** matrix. Its elements are the second derivatives of  $E(\mathbf{w})$  with respect to  $\mathbf{w}$ .
- We have already seen that  $\nabla E(\mathbf{w}) = \sum_{n=1}^n [(y_n - t_n)\varphi_n]$ .
- This can be rewritten as:  $\nabla E(\mathbf{w}) = \mathbf{\Phi}^T(\mathbf{y} - \mathbf{t})$ , where:

$$\mathbf{\Phi} = \begin{bmatrix} \varphi_0(x_1), \dots, \varphi_D(x_1) \\ \varphi_0(x_2), \dots, \varphi_D(x_2) \\ \dots \\ \varphi_0(x_N), \dots, \varphi_D(x_N) \end{bmatrix} \quad \mathbf{t} = \begin{bmatrix} t_1 \\ t_2 \\ \dots \\ t_N \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_N \end{bmatrix}$$

# Iterative Reweighted Least Squares

- So, we have:

$$\nabla E(\mathbf{w}) = \mathbf{\Phi}^T(\mathbf{y} - \mathbf{t})$$

$$H = \nabla \nabla E(\mathbf{w}) = \mathbf{\Phi}^T \mathbf{R} \mathbf{\Phi}$$

- In the above formula, we have introduced a new symbol  $\mathbf{R}$ , defined as a diagonal matrix as follows:
  - If  $i \neq j$ , then  $R_{ij} = 0$ , where  $R_{ij}$  is the value of  $\mathbf{R}$  at row  $i$  and column  $j$ .
  - Every diagonal value  $R_{nn}$  is defined as  $R_{nn} = y_n(1 - y_n)$

# Iterative Reweighted Least Squares

- Putting all those formulas together, we get:

$$\nabla E(\mathbf{w}) = \mathbf{\Phi}^T(\mathbf{y} - \mathbf{t})$$

$$H = \nabla \nabla E(\mathbf{w}) = \mathbf{\Phi}^T \mathbf{R} \mathbf{\Phi}$$

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

- This is now a closed form solution for how to iteratively update  $\mathbf{w}$ .
  - You stop the iterative updating when the value of  $\mathbf{w}$  converges, and changes little or not at all from one iteration to the next.

# Logistic Regression: Recap

- Classification function:  $y(x) = \sigma(\mathbf{w}^T \varphi(x))$ 
  - We treat  $y(x)$  as an estimate of probability  $p(C_1|x)$ .
- Error measure:  $E(\mathbf{w}) = -\ln(p(\mathbf{t} | X, \mathbf{w}))$
- There is no closed form formula for finding the best  $\mathbf{w}$ .
- However, a unique best  $\mathbf{w}$  exists, and can be found using iterative reweighted least squares:

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

- We can also minimize  $E(\mathbf{w})$  using sequential learning:

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n(\mathbf{w}) = \mathbf{w}^{(\tau)} - \eta (y_n - t_n) \varphi(x_n)$$

# Fisher's Linear Discriminant

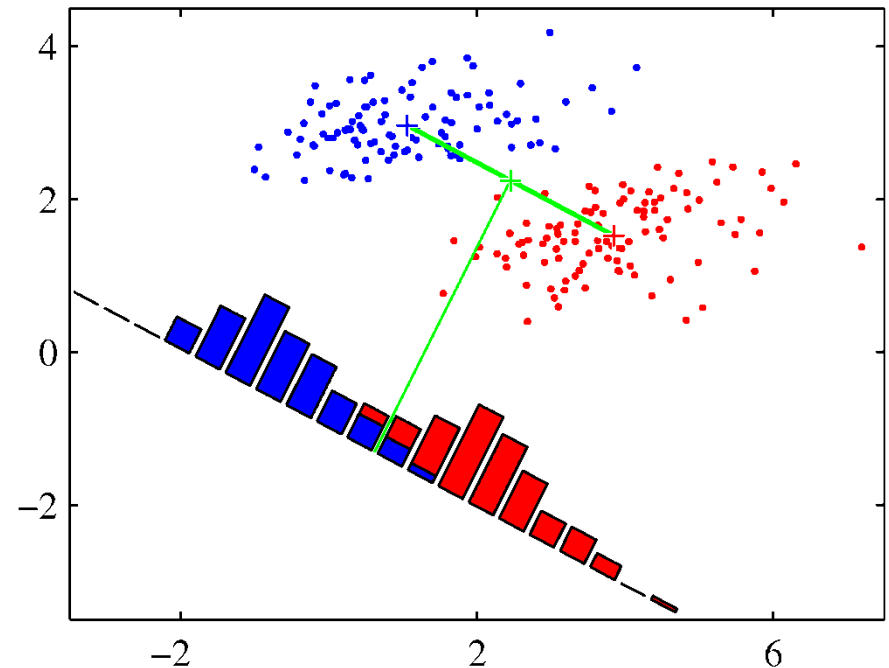
- Remember, in linear classification, the decision boundary is defined by equation:

$$\mathbf{w}^T \mathbf{x} = 0$$

- The formula above is the simple version, without using basis functions.
- Let's define function  $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ .
- Geometrically, what does function  $y(\mathbf{x})$  do?
  - $y(\mathbf{x}): \mathbb{R}^D \rightarrow \mathbb{R}$
  - $y(\mathbf{x})$  maps (projects)  $D$ -dimensional vectors  $\mathbf{x}$  to points on a line.
- If  $y(\mathbf{x}) > 0$  then  $\mathbf{x}$  is classified as belonging to class  $C_1$ .
- If  $y(\mathbf{x}) < 0$  then  $\mathbf{x}$  is classified as belonging to class  $C_0$ .

# Fisher's Linear Discriminant

- Let's define function  $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ .
- $y(\mathbf{x})$  maps (projects)  $D$ -dimensional vectors  $\mathbf{x}$  to points on a line.
  - If  $y(\mathbf{x}) > 0$  then  $\mathbf{x}$  is classified as belonging to class  $C_1$ .
  - If  $y(\mathbf{x}) < 0$  then  $\mathbf{x}$  is classified as belonging to class  $C_0$ .
- The figure shows an example:
  - Input points  $\mathbf{x}$  are projected on a line.
  - If they project onto one side of the line, they are classified as blue.
  - If they project on the other side, they are classified as red.
- In this particular example, some points are misclassified.



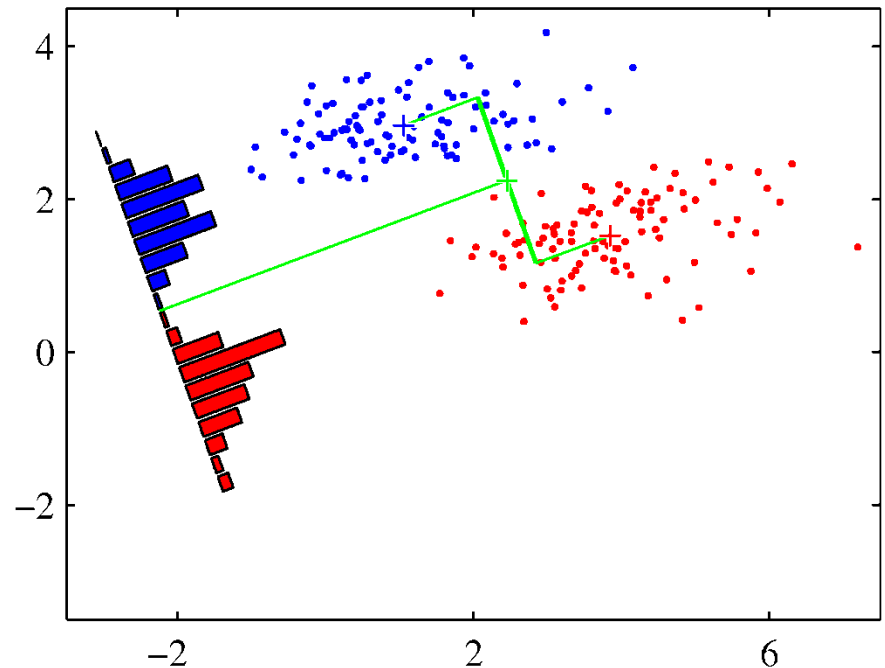
# Fisher's Linear Discriminant

- Let's define function  $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ .
- $y(\mathbf{x})$  maps (projects)  $D$ -dimensional vectors  $\mathbf{x}$  to points on a line.
  - If  $y(\mathbf{x}) > 0$  then  $\mathbf{x}$  is classified as belonging to class  $C_1$ .
  - If  $y(\mathbf{x}) < 0$  then  $\mathbf{x}$  is classified as belonging to class  $C_0$ .

- What  $\mathbf{w}$  is the best solution?
  - Logistic regression provided one answer, by minimizing a specific error measure:

$$E(\mathbf{w}) = -\ln(p(\mathbf{t} | X, \mathbf{w}))$$

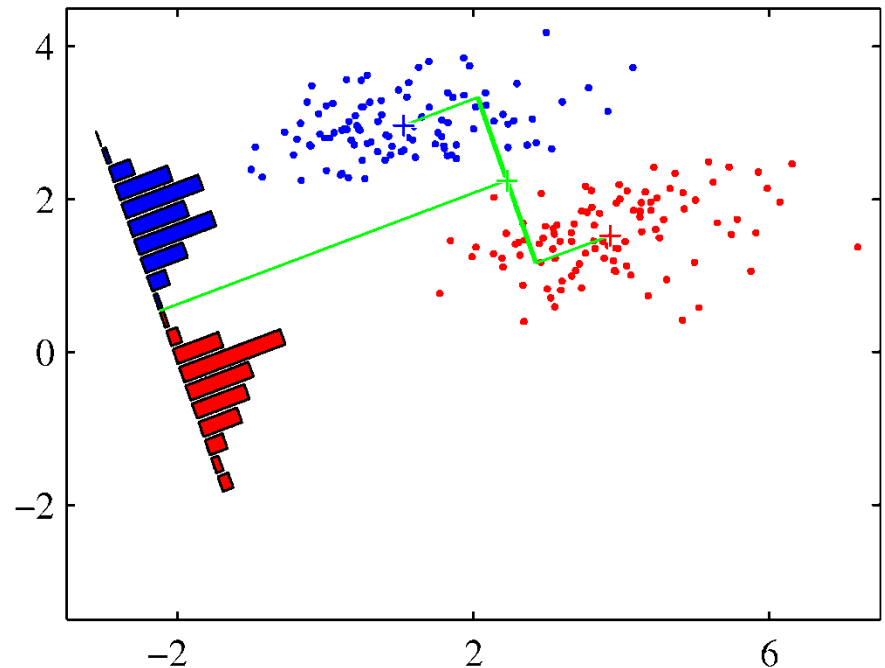
- Fisher's linear discriminant is an alternative method, that computes a value for  $\mathbf{w}$  using a different criterion.





# Fisher's Linear Discriminant

- Goal in Fisher's Linear Discriminant: find the line projection that maximizes the separation of the classes.
- Key question: how do we measure separation of the classes?
- One simple (but not very useful) answer: maximize the separation of the means of the classes.



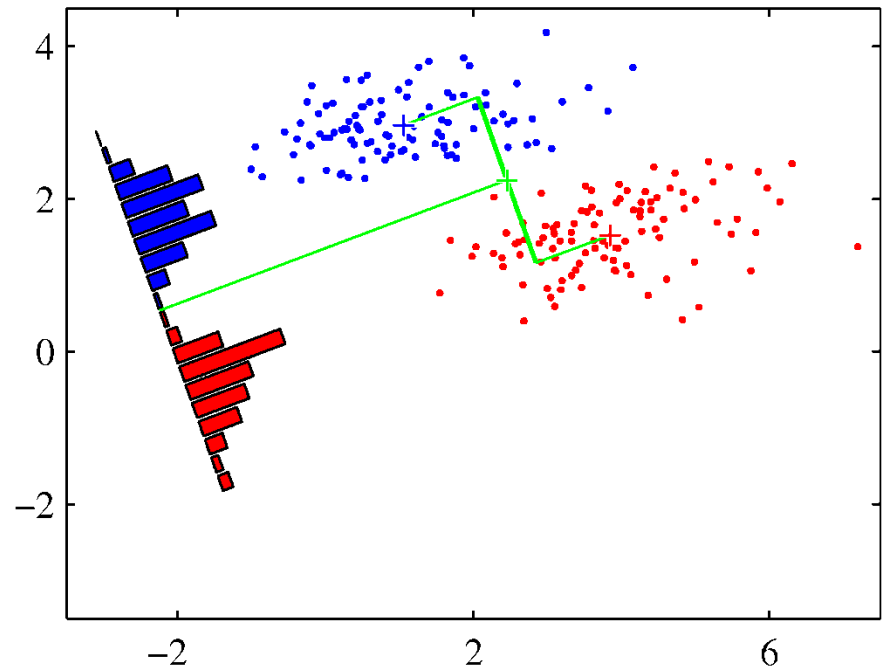
# Maximizing Separation of Means

- One simple (but not very useful) answer: maximize the separation of the means of the classes.

$$\mathbf{m}_0 = \frac{1}{N_0} \sum_{n \in C_0} x_n$$

$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{n \in C_1} x_n$$

- Goal: maximize  $\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_0$
- Same as maximizing  $\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_0)$
- Do you see any problem with this goal?



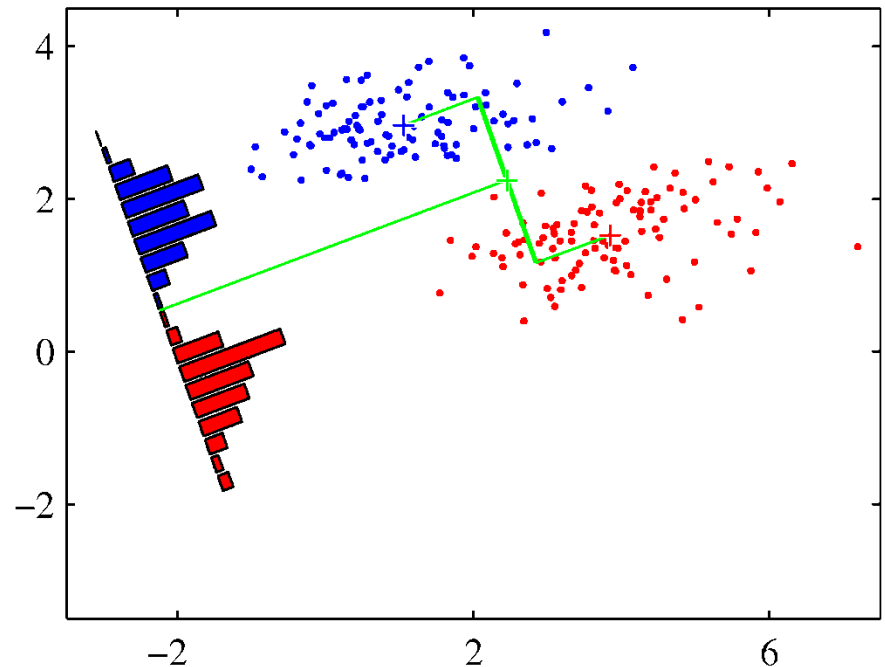
# Maximizing Separation of Means

- One simple (but not very useful) answer: maximize the separation of the means of the classes.

$$\mathbf{m}_0 = \frac{1}{N_0} \sum_{n \in C_0} x_n$$

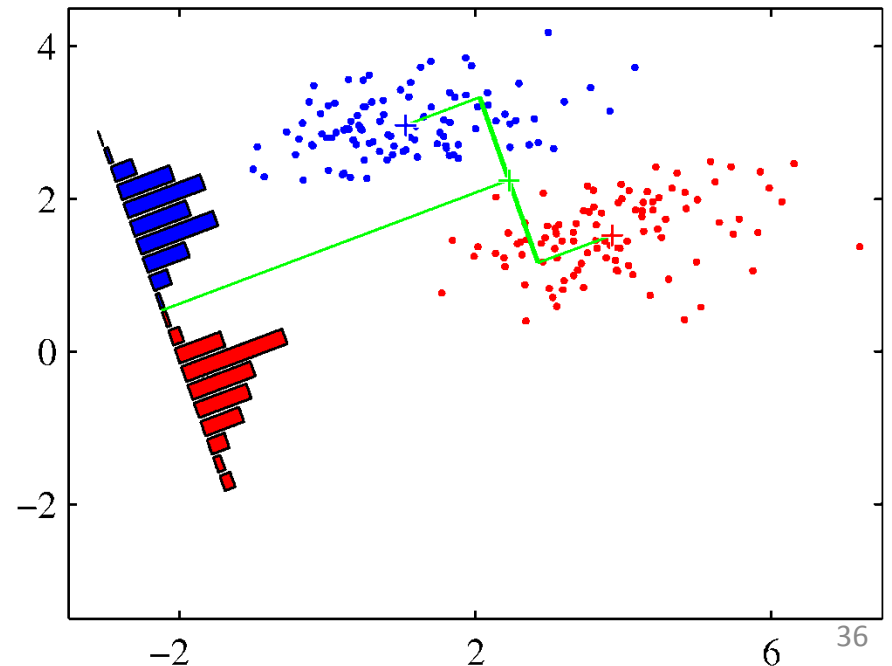
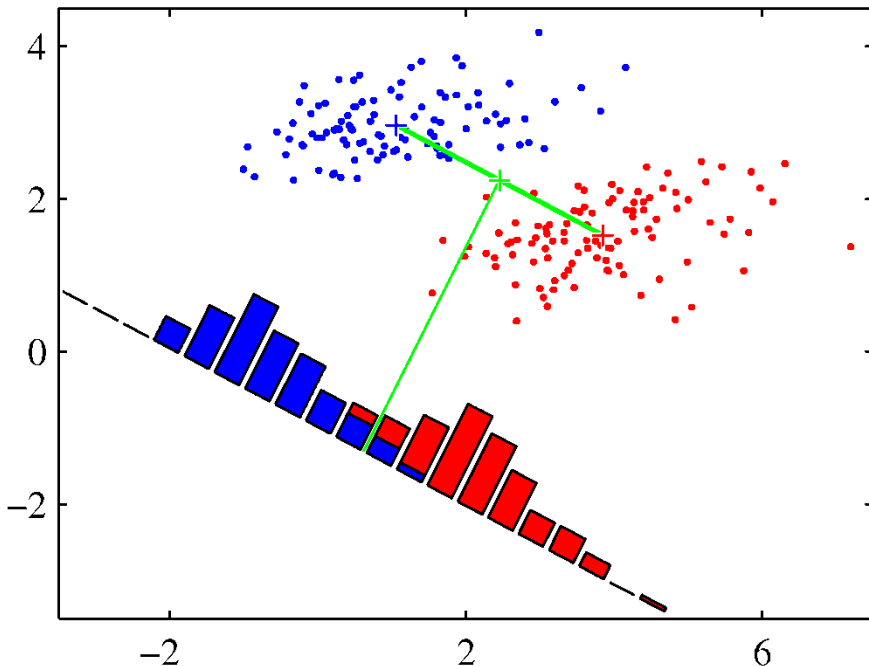
$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{n \in C_1} x_n$$

- Goal: maximize  $\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_0$
- Same as maximizing  $\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_0)$
- For any  $\mathbf{w}$ , replacing  $\mathbf{w}$  by  $1000 * \mathbf{w}$ , we increase the separation of the means by a factor of 1000, but classification accuracy stays the same.
- Maximizing  $\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_0)$  does not help improve accuracy.



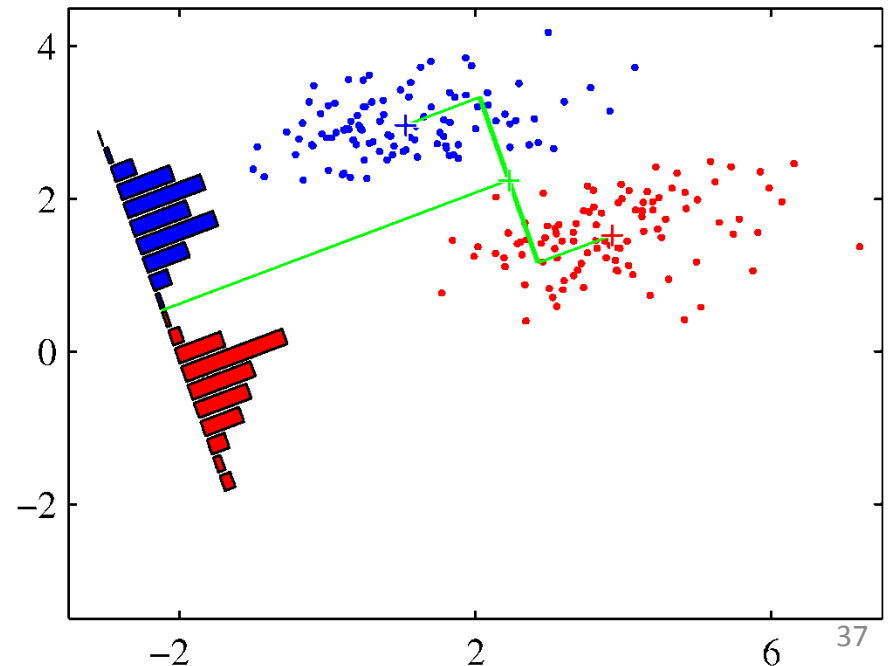
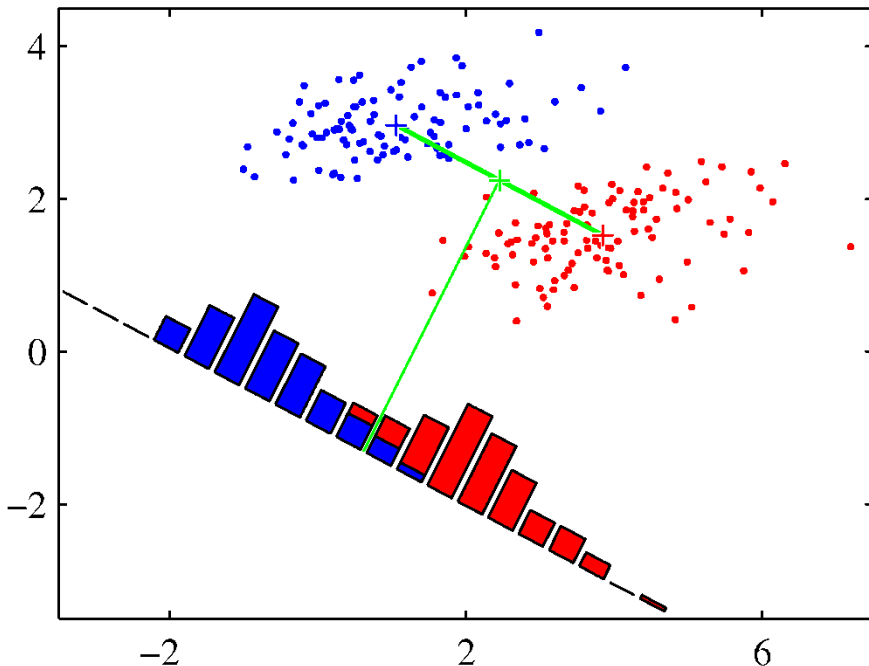
# Maximizing Separation of Means

- A second not very useful answer: maximize the separation of the means of the classes, subject to  $\|\mathbf{w}\| = 1$ .
- Now we cannot replace  $\mathbf{w}$  by  $1000 * \mathbf{w}$ , because that would violate the constraint that  $\|\mathbf{w}\| = 1$ .



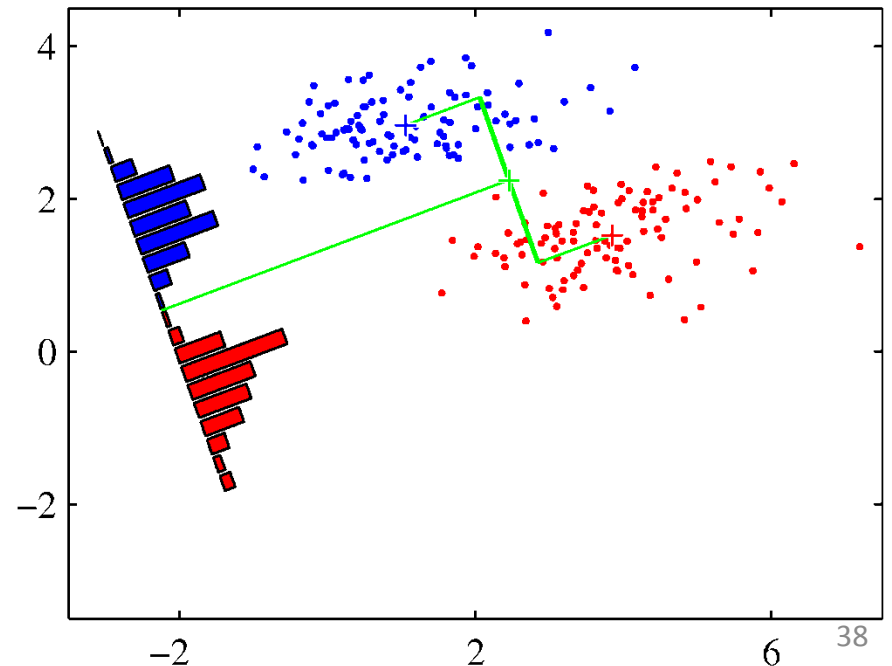
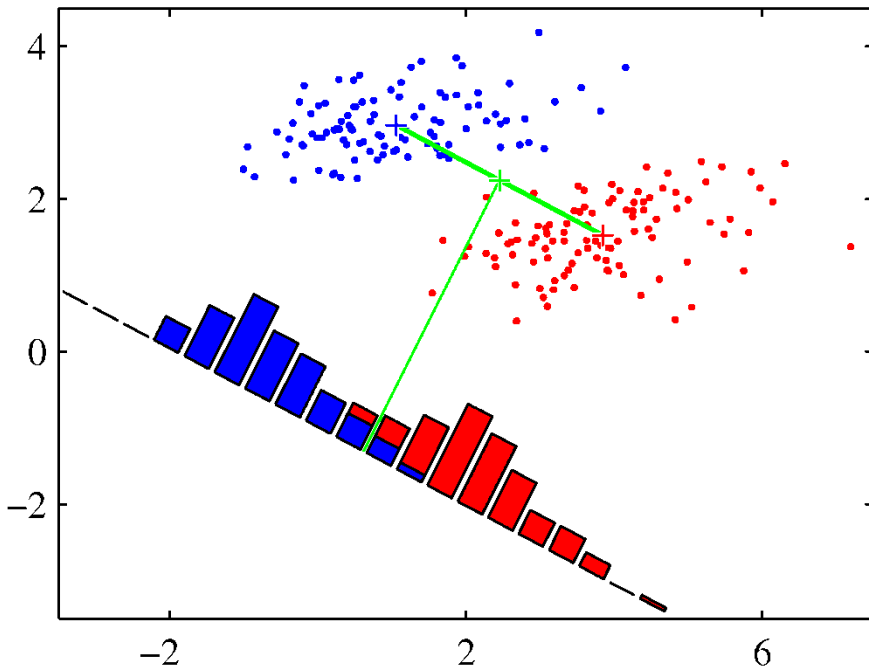
# Maximizing Separation of Means

- A second not very useful answer: maximize the separation of the means of the classes, subject to  $\|\mathbf{w}\| = 1$ .
- The solution according to this criterion is shown on the left figure.
  - Do you see any problems?



# Maximizing Separation of Means

- A second not very useful answer: maximize the separation of the means of the classes, subject to  $\|\mathbf{w}\| = 1$ .
- The solution according to this criterion is shown on the left figure.
  - The projection on the right figure gives more accurate classification.
  - Our criterion fails to identify that far better solution.



# Between Class Variance, Within-Class Variance,

- $\mathbf{w}^T(\mathbf{m}_1 - \mathbf{m}_0)$  is called the **between-class variance**.
  - It measures how far (on average) a point from class  $C_0$  is from a point in class  $C_1$ .
- We can also define the **within-class variance**  $(s_k)^2$ :
  - How far (on average) a point from class  $C_k$  is from another point in the same class  $C_k$ .

$$(s_k)^2 = \sum_{n \in C_k} (\mathbf{w}^T \mathbf{x}_n - \mathbf{w}^T \mathbf{m}_k)^2$$

# Fisher Criterion

- Fisher criterion: maximize the ratio  $J(\mathbf{w})$  of between-class variance over within-class variance.

$$J(\mathbf{w}) = \frac{(\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_0)^2}{(s_0)^2 + (s_1)^2}$$

- Intuition: we want, at the same time, a projection where:
  - Between-class variance is large (inputs belonging to different classes should map to points that are far from each other).
  - Within-class variance is small (inputs belonging to the same class should map to points that are close to each other).



# Maximizing $J(\mathbf{w})$

- We will use the following definitions:
  - $\mathbf{S}_B$  is the **between-class covariance matrix**.

$$\mathbf{S}_B = (\mathbf{m}_1 - \mathbf{m}_0)(\mathbf{m}_1 - \mathbf{m}_0)^T$$

- $\mathbf{S}_W$  is the total **within-class covariance matrix**.

$$\mathbf{S}_W = \left( \sum_{n \in C_0} (\mathbf{x}_n - \mathbf{m}_0)(\mathbf{x}_n - \mathbf{m}_0)^T \right) + \left( \sum_{n \in C_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T \right)$$

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

# Maximizing $J(\mathbf{w})$

- We want to find the  $\mathbf{w}$  that maximizes:

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

- Differentiating  $J(\mathbf{w})$  with respect to  $\mathbf{w}$ , and solving equation  $\nabla J(\mathbf{w}) = \mathbf{0}$ , we get this solution:

$$\mathbf{w} \propto (\mathbf{S}_W)^{-1}(\mathbf{m}_1 - \mathbf{m}_0)$$

- The symbol  $\propto$  means that  $\mathbf{w}$  is proportional to  $(\mathbf{S}_W)^{-1}(\mathbf{m}_1 - \mathbf{m}_0)$ .
- Therefore,  $\mathbf{w} = (\mathbf{S}_W)^{-1}(\mathbf{m}_1 - \mathbf{m}_0)$  is a solution.
- For any real number  $c$ ,  $\mathbf{w} = c(\mathbf{S}_W)^{-1}(\mathbf{m}_1 - \mathbf{m}_0)$  is also a solution.