

Linear classifiers

CE-717: Machine Learning
Sharif University of Technology

M. Soleymani

Fall 2016

Topics

- ▶ Discriminant functions
 - ▶ Linear classifiers
 - ▶ Perceptron
 - ▶ Fisher
- } → SVM will be covered in the later lectures
- ▶ Multi-class classification

Classification problem

- ▶ **Given: Training set**

- ▶ labeled set of N input-output pairs $D = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$
- ▶ $y \in \{1, \dots, K\}$

- ▶ **Goal: Given an input \mathbf{x} , assign it to one of K classes**

- ▶ **Examples:**

- ▶ Spam filter
- ▶ Handwritten digit recognition
- ▶ ...

Discriminant functions

- ▶ **Discriminant function** can directly assign each vector x to a specific class k
- ▶ A popular way of representing a classifier
 - ▶ Many classification methods are based on discriminant functions
- ▶ Assumption: the classes are taken to be disjoint
 - ▶ The input space is thereby divided into **decision regions**
 - ▶ boundaries are called **decision boundaries** or decision surfaces.

Discriminant Functions

- ▶ **Discriminant functions:** A discriminant function $f_i(\mathbf{x})$ for each class \mathcal{C}_i ($i = 1, \dots, K$):

- ▶ \mathbf{x} is assigned to class \mathcal{C}_i if:

$$f_i(\mathbf{x}) > f_j(\mathbf{x}) \quad \forall j \neq i$$

- ▶ Thus, we can easily divide the feature space into K decision regions

$$\forall \mathbf{x}, f_i(\mathbf{x}) > f_j(\mathbf{x}) \quad \forall j \neq i \Rightarrow \mathbf{x} \in \mathcal{R}_i$$

\mathcal{R}_i : Region of the i -th class

- ▶ **Decision surfaces (or boundaries)** can also be found using discriminant functions

- ▶ Boundary of the \mathcal{R}_i and \mathcal{R}_j separating samples of these two categories:

$$\forall \mathbf{x}, f_i(\mathbf{x}) = f_j(\mathbf{x})$$

Discriminant Functions: Two-Category

- ▶ Decision surface: $f(\mathbf{x}) = 0$
- ▶ For two-category problem, we can only find a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$
 - ▶ $f_1(\mathbf{x}) = f(\mathbf{x})$
 - ▶ $f_2(\mathbf{x}) = -f(\mathbf{x})$
- ▶ First, we explain two-category classification problem and then discuss the multi-category problems.
 - ▶ Binary classification: a target variable $y \in \{0,1\}$ or $y \in \{-1,1\}$

Linear classifiers

- ▶ Decision boundaries are linear in x , or linear in some given set of functions of x
- ▶ Linearly separable data: data points that can be exactly classified by a linear decision surface.
- ▶ Why linear classifier?
 - ▶ Even when they are not optimal, we can use their simplicity
 - ▶ are relatively easy to compute
 - ▶ In the absence of information suggesting otherwise, linear classifiers are an attractive candidates for initial, trial classifiers.

Two Category

- ▶ $f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x} + w_0 = w_0 + w_1 x_1 + \dots + w_d x_d$
 - ▶ $\mathbf{x} = [x_1 \ x_2 \ \dots \ x_d]$
 - ▶ $\mathbf{w} = [w_1 \ w_2 \ \dots \ w_d]$
 - ▶ w_0 : bias

if $\mathbf{w}^T \mathbf{x} + w_0 \geq 0$ then \mathcal{C}_1

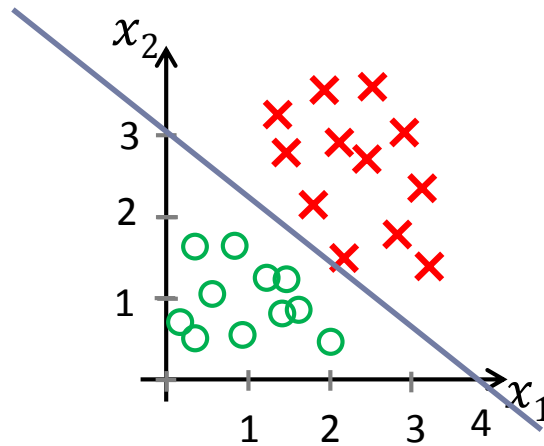
else \mathcal{C}_2

Decision surface (boundary): $\mathbf{w}^T \mathbf{x} + w_0 = 0$

\mathbf{w} is orthogonal to every vector lying within the decision surface

Example

$$3 - \frac{3}{4}x_1 - x_2 = 0$$



if $w^T x + w_0 \geq 0$ then C_1
else C_2

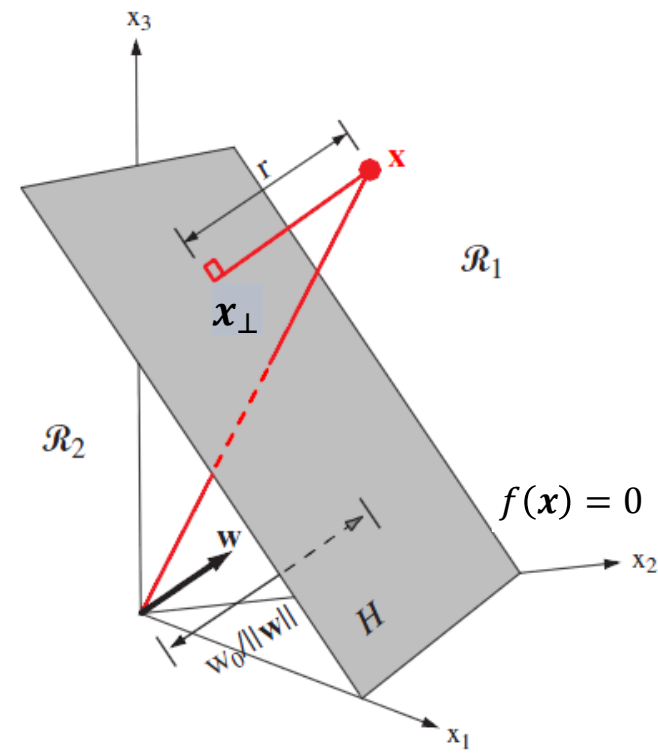
Linear classifier: Two Category

- ▶ Decision boundary is a $(d - 1)$ -dimensional hyperplane H in the d -dimensional feature space
 - ▶ The orientation of H is determined by the normal vector $[w_1, \dots, w_d]$
 - ▶ w_0 determine the location of the surface.
 - ▶ The normal distance from the origin to the decision surface is $\frac{w_0}{\|w\|}$

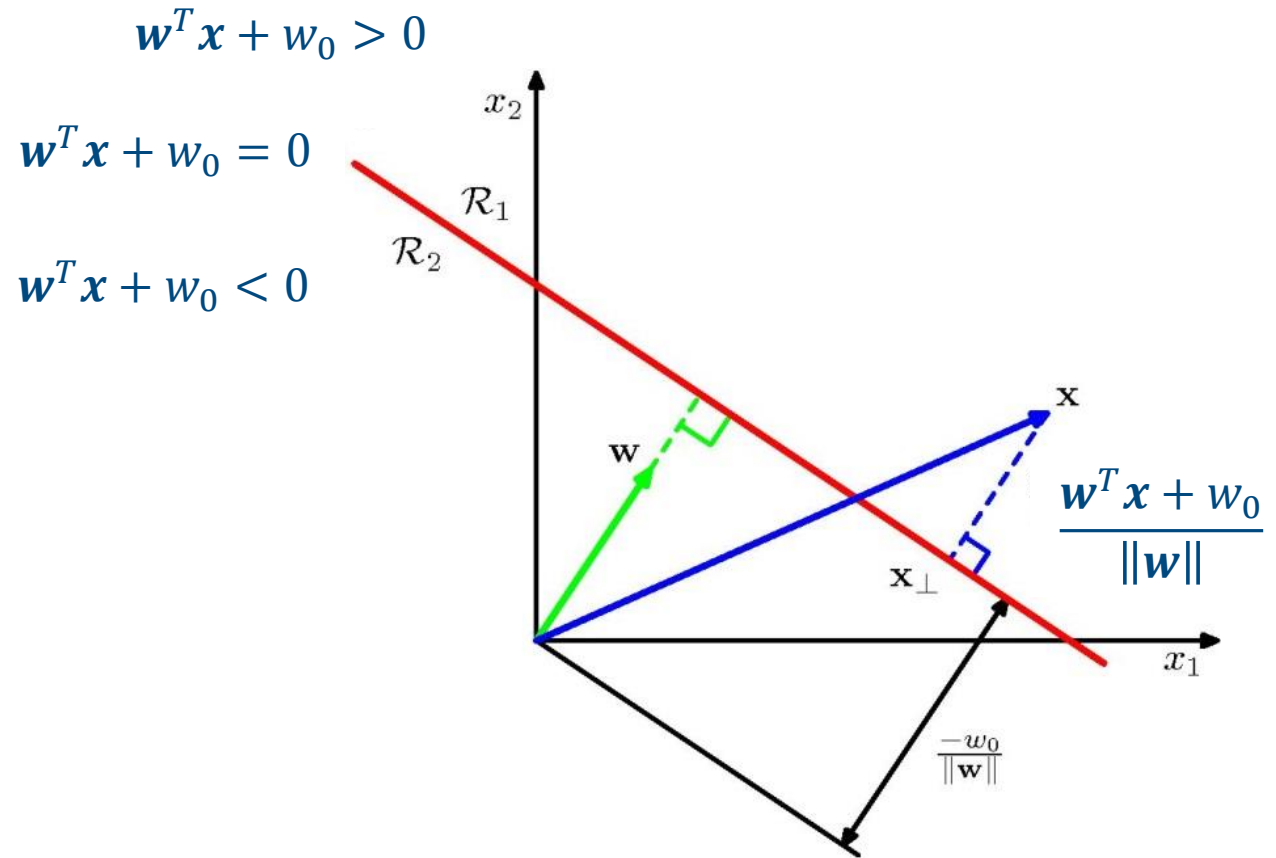
$$\mathbf{x} = \mathbf{x}_\perp + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

$$\mathbf{w}^T \mathbf{x} + w_0 = r \|\mathbf{w}\| \Rightarrow r = \frac{\mathbf{w}^T \mathbf{x} + w_0}{\|\mathbf{w}\|}$$

gives a signed measure of the perpendicular distance r of the point \mathbf{x} from the decision surface

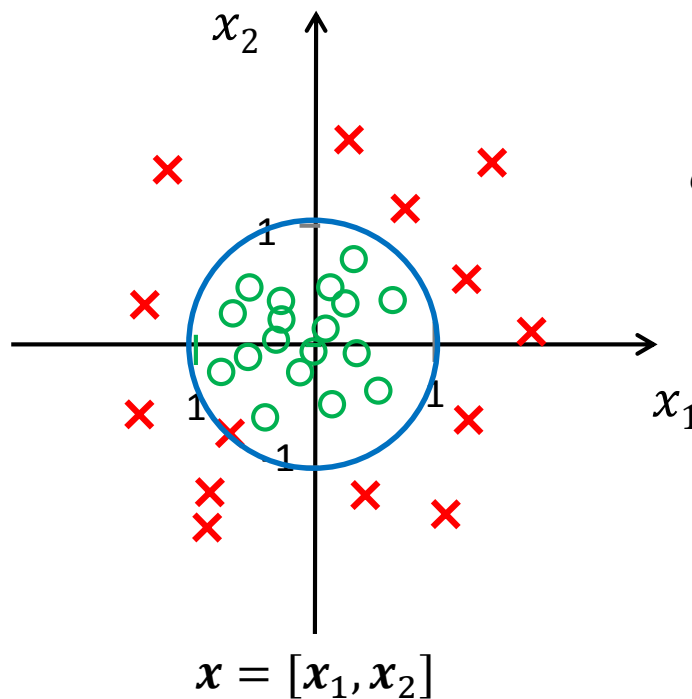


Linear boundary: geometry



Non-linear decision boundary

- ▶ Choose non-linear features
- ▶ Classifier still linear in parameters \mathbf{w}



$$-1 + x_1^2 + x_2^2 = 0$$

$$\phi(\mathbf{x}) = [1, x_1, x_2, x_1^2, x_2^2, x_1 x_2]$$

$$\mathbf{w} = [w_0, w_1, \dots, w_m] = [-1, 0, 0, 1, 1, 0]$$

if $\mathbf{w}^T \phi(\mathbf{x}) \geq 0$ then $y = 1$
else $y = -1$

Cost Function for linear classification

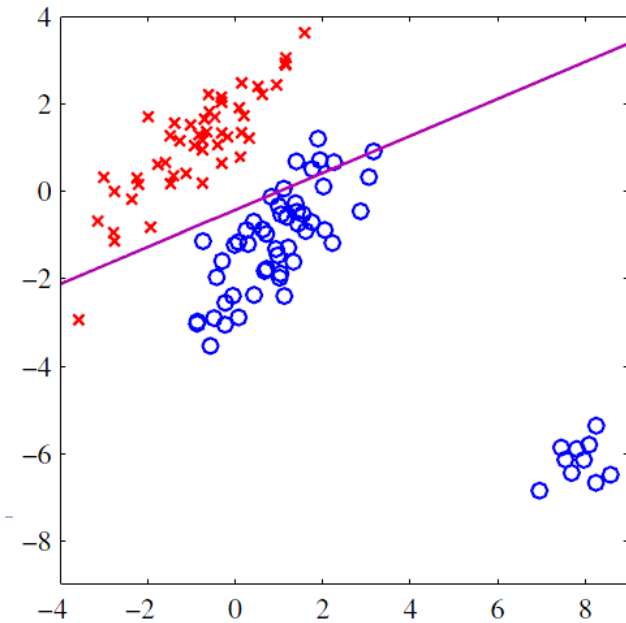
- ▶ Finding linear classifiers can be formulated as an optimization problem:
 - ▶ Select how to measure the prediction loss
 - ▶ Based on the training set $D = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$, a cost function $J(\mathbf{w})$ is defined
 - ▶ Solve the resulting optimization problem to find parameters:
 - ▶ Find optimal $\hat{f}(x) = f(x; \hat{\mathbf{w}})$ where $\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} J(\mathbf{w})$
- ▶ Criterion or cost functions for classification:
 - ▶ We will investigate several cost functions for the classification problem

SSE cost function for classification $K = 2$

SSE cost function is not suitable for classification:

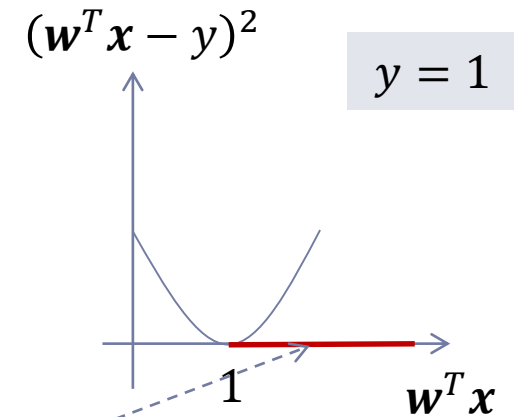
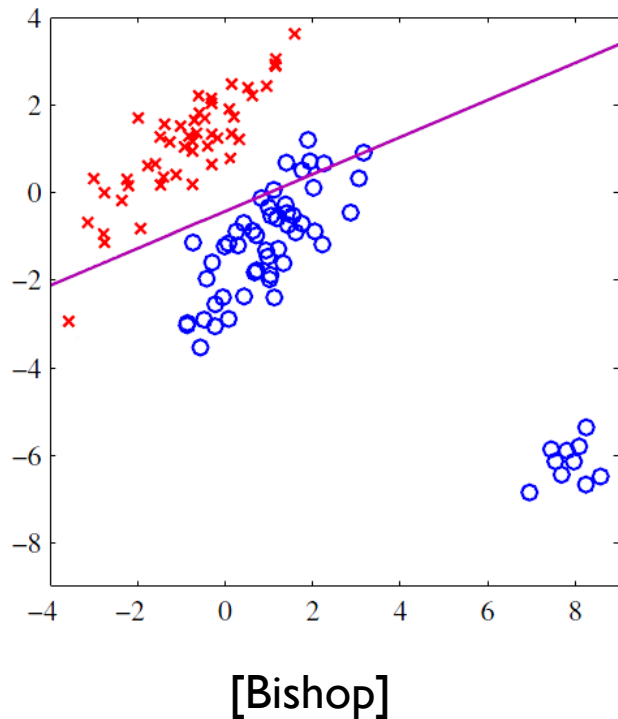
- ▶ Least square loss penalizes 'too correct' predictions (that they lie a long way on the correct side of the decision)
- ▶ Least square loss also lack robustness to noise

$$J(\mathbf{w}) = \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}^{(i)} - y^{(i)})^2$$

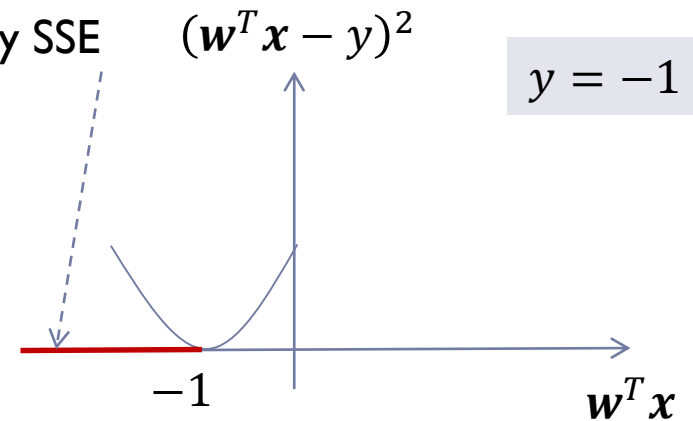


SSE cost function for classification

$K = 2$



Correct predictions that
are penalized by SSE

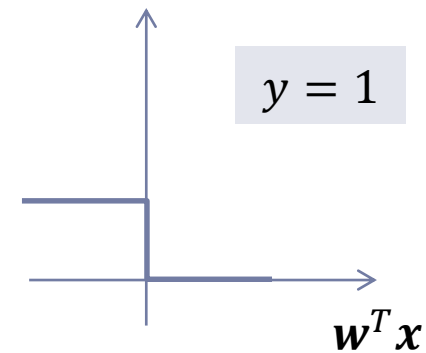


SSE cost function for classification $K = 2$

- ▶ Is it more suitable if we set $f(\mathbf{x}; \mathbf{w}) = g(\mathbf{w}^T \mathbf{x})$?

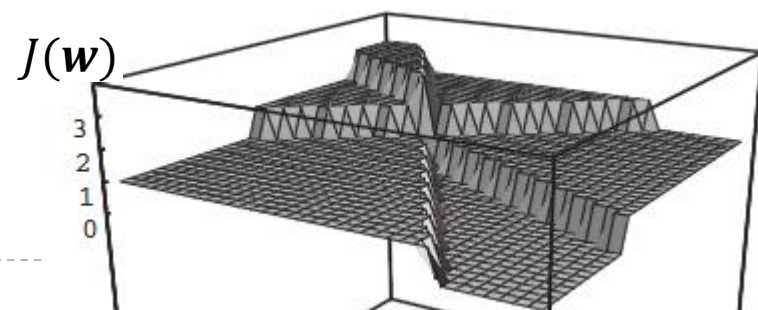
$$J(\mathbf{w}) = \sum_{i=1}^N (\text{sign}(\mathbf{w}^T \mathbf{x}^{(i)}) - y^{(i)})^2 \quad (\text{sign}(\mathbf{w}^T \mathbf{x}) - y)^2$$

$$\text{sign}(z) = \begin{cases} -1, & z < 0 \\ 1, & z \geq 0 \end{cases}$$



- ▶ $J(\mathbf{w})$ is a piecewise constant function shows the number of misclassifications

Training error incurred in classifying training samples



Perceptron algorithm

- ▶ Linear classifier
- ▶ Two-class: $y \in \{-1, 1\}$
 - ▶ $y = -1$ for C_2 , $y = 1$ for C_1
- ▶ Goal: $\forall i, \mathbf{x}^{(i)} \in C_1 \Rightarrow \mathbf{w}^T \mathbf{x}^{(i)} > 0$
 $\forall i, \mathbf{x}^{(i)} \in C_2 \Rightarrow \mathbf{w}^T \mathbf{x}^{(i)} < 0$
- ▶ $f(\mathbf{x}; \mathbf{w}) = \text{sign}(\mathbf{w}^T \mathbf{x})$

Perceptron criterion

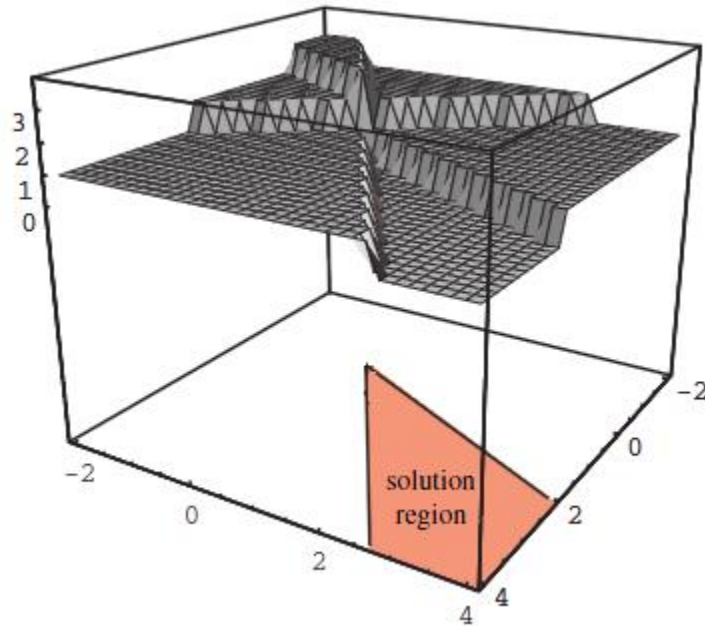
$$J_P(\mathbf{w}) = - \sum_{i \in \mathcal{M}} \mathbf{w}^T \mathbf{x}^{(i)} y^{(i)}$$

\mathcal{M} : subset of training data that are misclassified

Many solutions? Which solution among them?

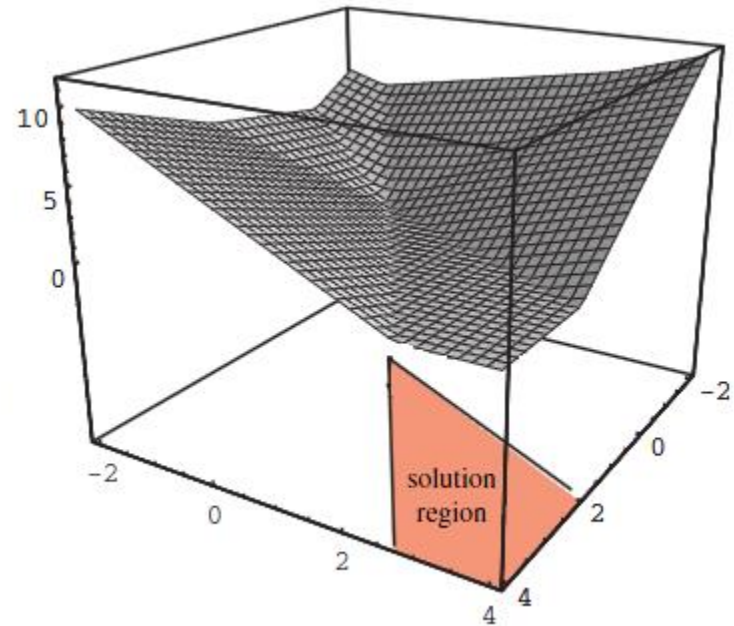
Cost function

$J(\mathbf{w})$



of misclassifications
as a cost function

$J_P(\mathbf{w})$



Perceptron's
cost function

There may be many solutions in these cost functions

Batch Perceptron

“Gradient Descent” to solve the optimization problem:

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta \nabla_{\mathbf{w}} J_P(\mathbf{w}^t)$$

$$\nabla_{\mathbf{w}} J_P(\mathbf{w}) = - \sum_{i \in \mathcal{M}} \mathbf{x}^{(i)} y^{(i)}$$

Batch Perceptron converges in finite number of steps for linearly separable data:

Initialize \mathbf{w}

Repeat

$$\mathbf{w} = \mathbf{w} + \eta \sum_{i \in \mathcal{M}} \mathbf{x}^{(i)} y^{(i)}$$

Until $\eta \sum_{i \in \mathcal{M}} \mathbf{x}^{(i)} y^{(i)} < \theta$

Stochastic gradient descent for Perceptron

- ▶ Single-sample perceptron:

- ▶ If $\mathbf{x}^{(i)}$ is misclassified:

$$\mathbf{w}^{t+1} = \mathbf{w}^t + \eta \mathbf{x}^{(i)} y^{(i)}$$

- ▶ Perceptron convergence theorem: for linearly separable data

- ▶ If training data are linearly separable, the single-sample perceptron is also guaranteed to find a solution in a finite number of steps

Fixed-Increment single sample Perceptron

Initialize $\mathbf{w}, t \leftarrow 0$

repeat

$t \leftarrow t + 1$

$i \leftarrow t \bmod N$

if $\mathbf{x}^{(i)}$ is misclassified then

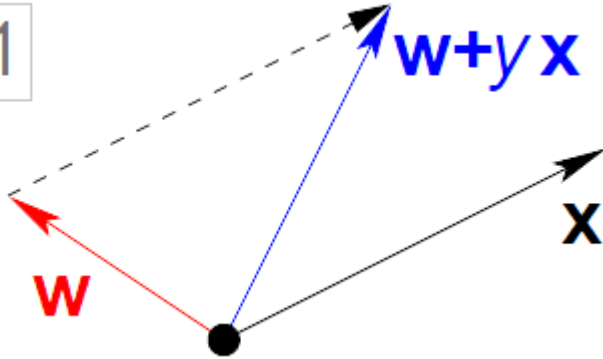
$\mathbf{w} = \mathbf{w} + \mathbf{x}^{(i)} y^{(i)}$

Until all patterns properly classified

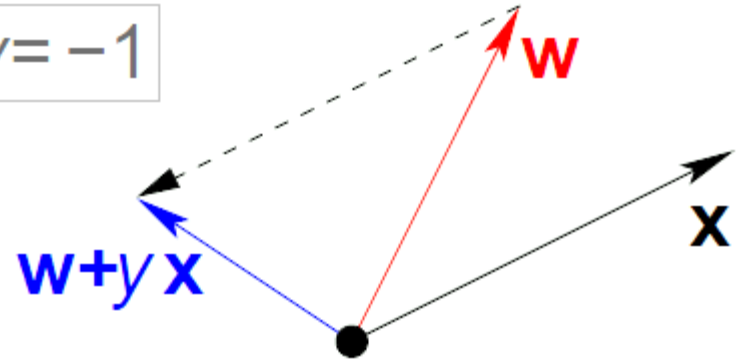
η can be set to 1 and
proof still works →

Example

$y = +1$

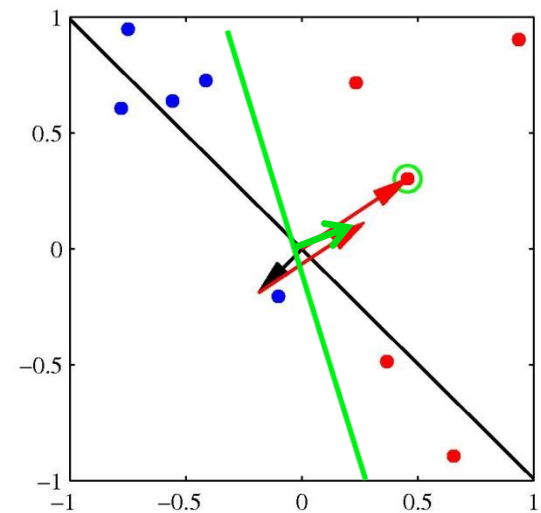
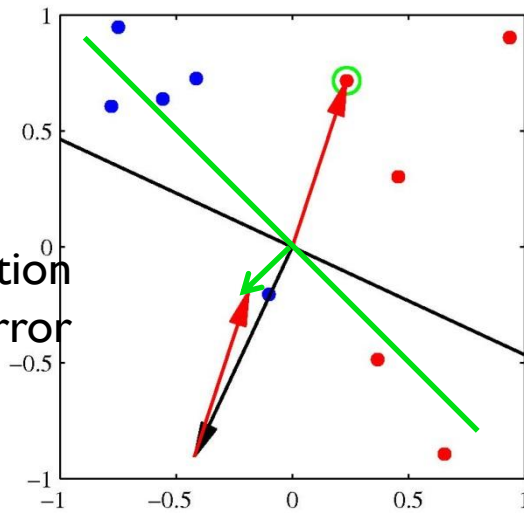


$y = -1$

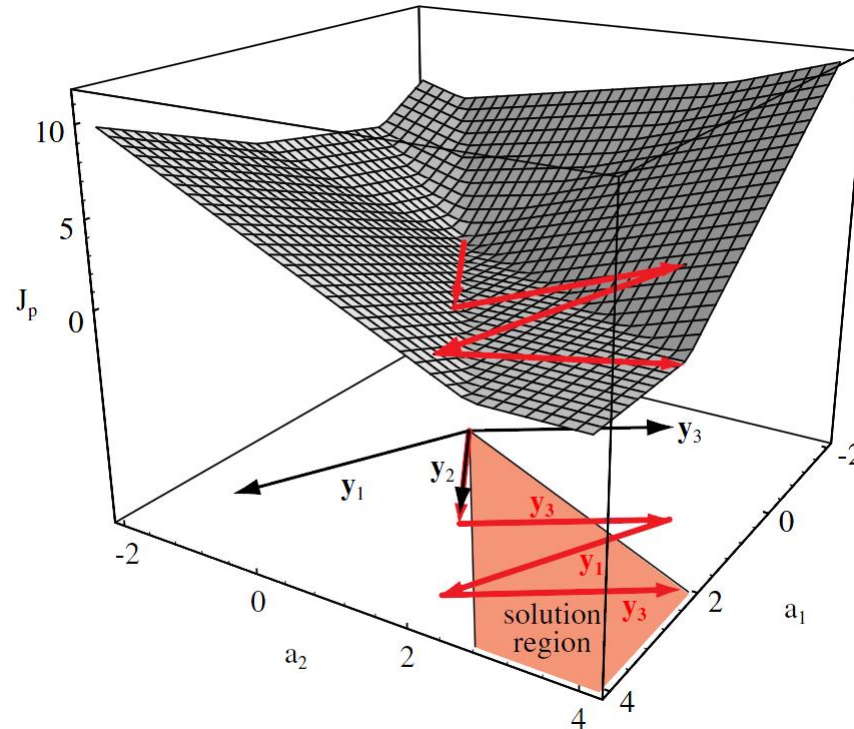


Perceptron: Example

Change w in a direction that corrects the error



Convergence of Perceptron



[Duda, Hart & Stork, 2002]

- ▶ For data sets that are not linearly separable, the single-sample perceptron learning algorithm will never converge

Pocket algorithm

- ▶ For the data that are not linearly separable due to noise:
 - ▶ Keeps in its pocket the best \mathbf{w} encountered up to now.

```
Initialize  $\mathbf{w}$ 
for  $t = 1, \dots, T$ 
   $i \leftarrow t \bmod N$ 
  if  $\mathbf{x}^{(i)}$  is misclassified then
     $\mathbf{w}^{new} = \mathbf{w} + \mathbf{x}^{(i)} y^{(i)}$ 
    if  $E_{train}(\mathbf{w}^{new}) < E_{train}(\mathbf{w})$  then
       $\mathbf{w} = \mathbf{w}^{new}$ 
end
```

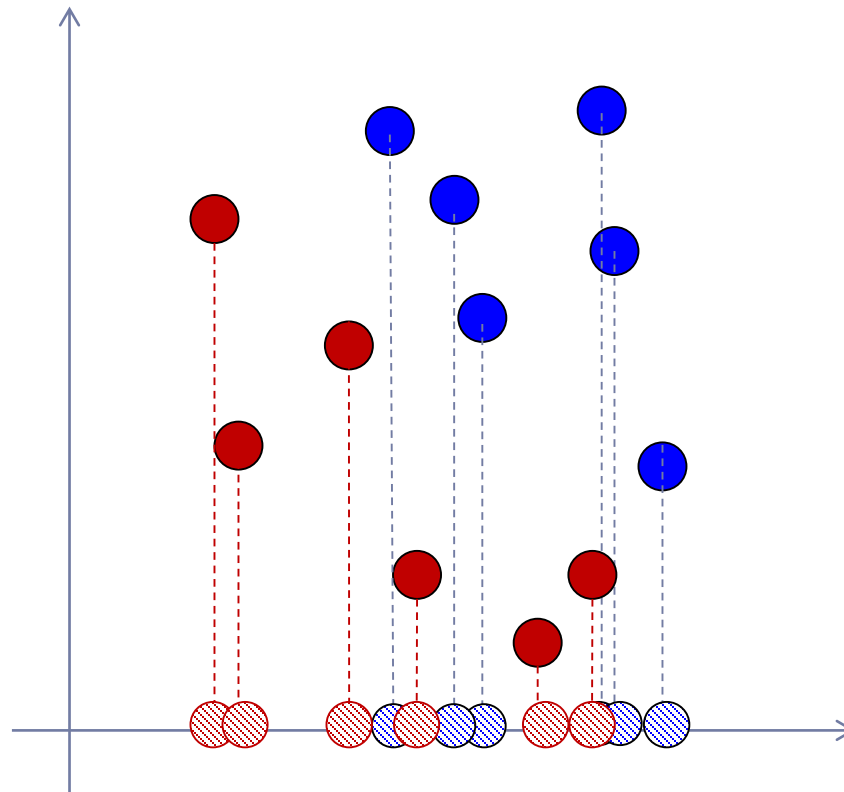
$$E_{train}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N [\text{sign}(\mathbf{w}^T \mathbf{x}^{(n)}) \neq y^{(n)}]$$

Linear Discriminant Analysis (LDA)

- ▶ Fisher's Linear Discriminant Analysis :
 - ▶ Dimensionality reduction
 - ▶ Finds linear combinations of features with large ratios of between-groups scatters to within-groups scatters (as discriminant new variables)
 - ▶ Classification
 - ▶ Predicts the class of an observation x by first projecting it to the space of discriminant variables and then classifying it in this space

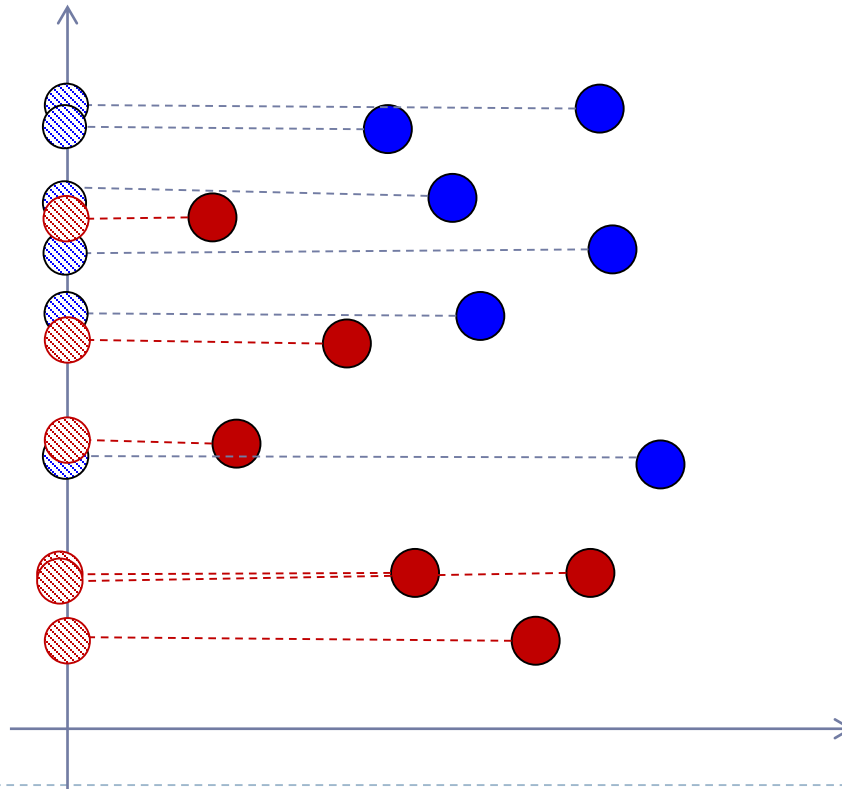
Good Projection for Classification

- ▶ What is a good criterion?
 - ▶ Separating different classes in the projected space



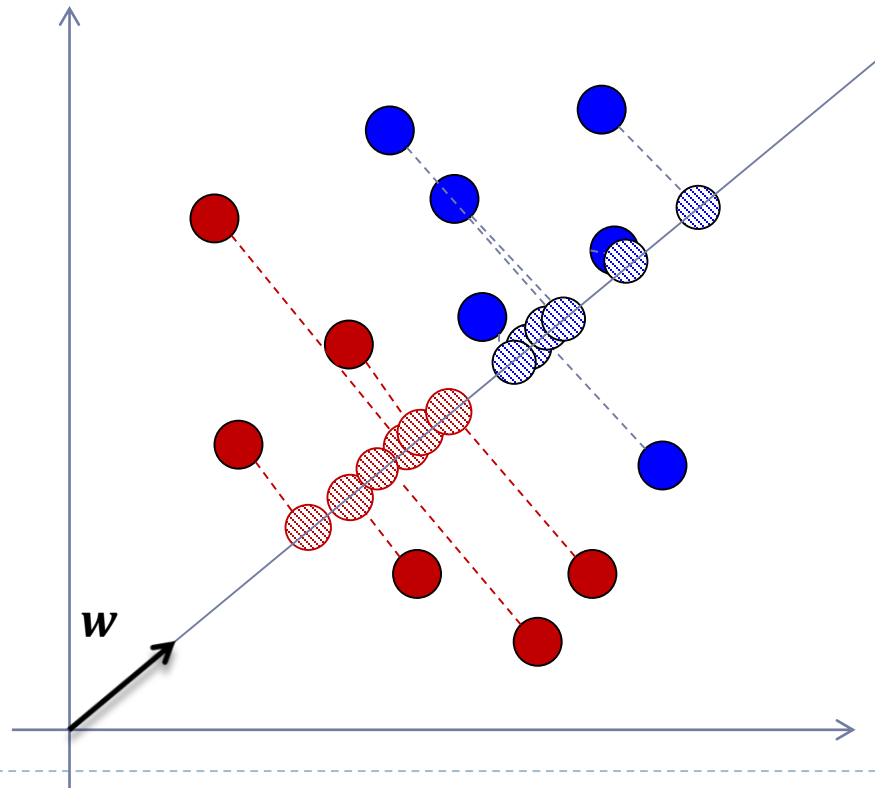
Good Projection for Classification

- ▶ What is a good criterion?
 - ▶ Separating different classes in the projected space



Good Projection for Classification

- ▶ What is a good criterion?
 - ▶ Separating different classes in the projected space

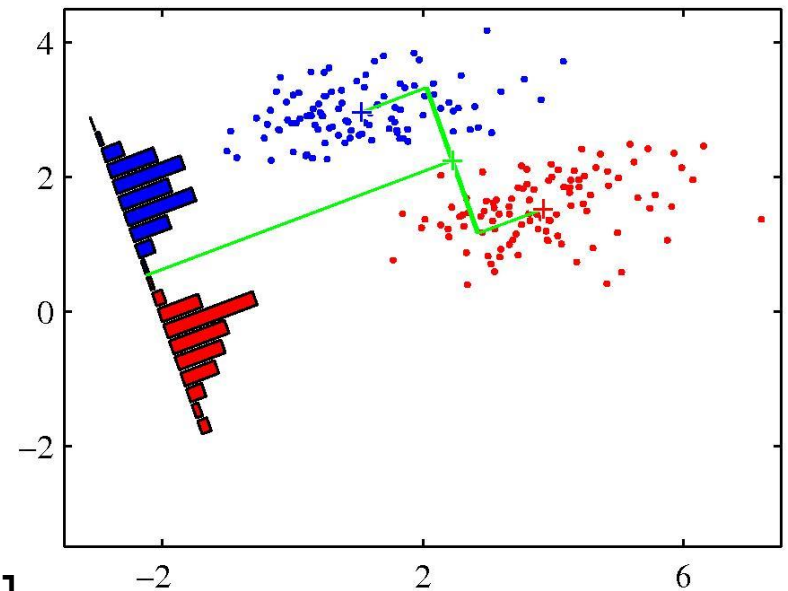
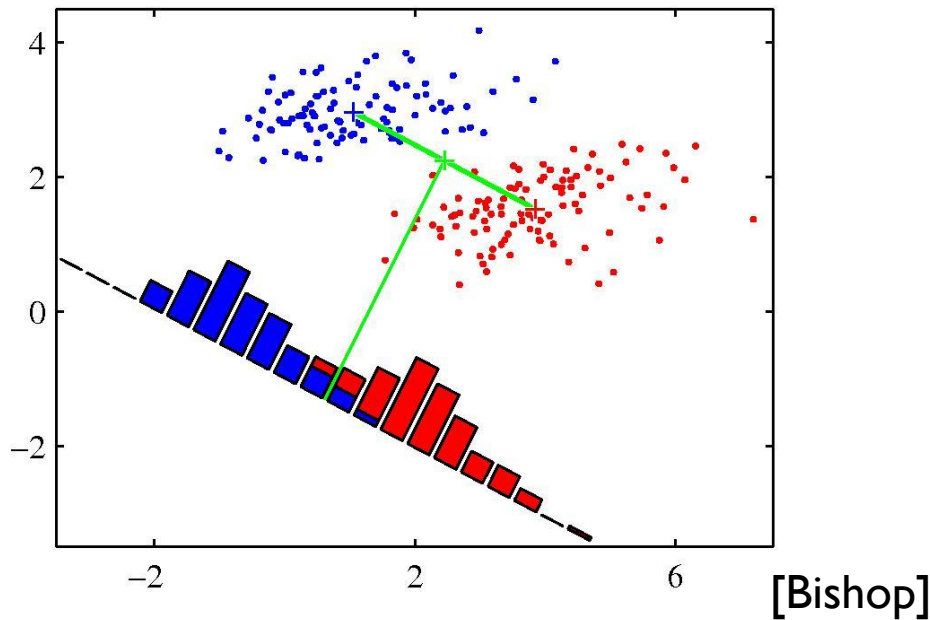


LDA Problem

- ▶ Problem definition:
 - ▶ $C = 2$ classes
 - ▶ $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$ training samples with N_1 samples from the first class (C_1) and N_2 samples from the second class (C_2)
 - ▶ Goal: finding the best direction \mathbf{w} that we hope to enable accurate classification
- ▶ The projection of sample \mathbf{x} onto a line in direction \mathbf{w} is $\mathbf{w}^T \mathbf{x}$
- ▶ What is the measure of the separation between the projected points of different classes?

Measure of Separation in the Projected Direction

- ▶ Is the direction of the line joining the class means a good candidate for w ?



Measure of Separation in the Projected Direction

- ▶ The direction of the line jointing the class means is the solution of the following problem:
 - ▶ Maximizes the separation of the projected class means

$$\begin{aligned} \max_{\mathbf{w}} J(\mathbf{w}) &= (\mu'_1 - \mu'_2)^2 \\ \text{s. t. } \|\mathbf{w}\| &= 1 \end{aligned}$$

$$\begin{aligned} \mu'_1 &= \mathbf{w}^T \boldsymbol{\mu}_1 & \boldsymbol{\mu}_1 &= \frac{\sum_{x^{(i)} \in \mathcal{C}_1} x^{(i)}}{N_1} \\ \mu'_2 &= \mathbf{w}^T \boldsymbol{\mu}_2 & \boldsymbol{\mu}_2 &= \frac{\sum_{x^{(i)} \in \mathcal{C}_2} x^{(i)}}{N_2} \end{aligned}$$

- ▶ What is the problem with the criteria considering only $|\mu'_1 - \mu'_2|$?
 - ▶ It does not consider the variances of the classes in the projected direction

LDA Criteria

- ▶ Fisher idea: maximize a function that will give
 - ▶ large separation between the projected class means
 - ▶ while also achieving a small variance within each class, thereby minimizing the class overlap.

$$J(\mathbf{w}) = \frac{|\mu'_1 - \mu'_2|^2}{s_1'^2 + s_2'^2}$$

LDA Criteria

- ▶ The scatters of the original data are:

$$s_1^2 = \sum_{\mathbf{x}^{(i)} \in \mathcal{C}_1} \|\mathbf{x}^{(i)} - \boldsymbol{\mu}_1\|^2$$

$$s_2^2 = \sum_{\mathbf{x}^{(i)} \in \mathcal{C}_2} \|\mathbf{x}^{(i)} - \boldsymbol{\mu}_2\|^2$$

- ▶ The scatters of projected data are:

$$s_1'^2 = \sum_{\mathbf{x}^{(i)} \in \mathcal{C}_1} (\mathbf{w}^T \mathbf{x}^{(i)} - \mathbf{w}^T \boldsymbol{\mu}_1)^2$$

$$s_2'^2 = \sum_{\mathbf{x}^{(i)} \in \mathcal{C}_2} (\mathbf{w}^T \mathbf{x}^{(i)} - \mathbf{w}^T \boldsymbol{\mu}_1)^2$$

LDA Criteria

$$J(\mathbf{w}) = \frac{|\mu'_1 - \mu'_2|^2}{s_1'^2 + s_2'^2}$$

$$\begin{aligned} |\mu'_1 - \mu'_2|^2 &= |\mathbf{w}^T \boldsymbol{\mu}_1 - \mathbf{w}^T \boldsymbol{\mu}_2|^2 \\ &= \mathbf{w}^T (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \mathbf{w} \end{aligned}$$

$$\begin{aligned} s_1'^2 &= \sum_{\mathbf{x}^{(i)} \in \mathcal{C}_1} (\mathbf{w}^T \mathbf{x}^{(i)} - \mathbf{w}^T \boldsymbol{\mu}_1)^2 \\ &= \mathbf{w}^T \left(\sum_{\mathbf{x}^{(i)} \in \mathcal{C}_1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_1)(\mathbf{x}^{(i)} - \boldsymbol{\mu}_1)^T \right) \mathbf{w} \end{aligned}$$

LDA Criteria

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

Between-class scatter matrix $\leftarrow \mathbf{S}_B = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T$

Within-class scatter matrix $\leftarrow \mathbf{S}_W = \mathbf{S}_1 + \mathbf{S}_2$

$$\mathbf{S}_1 = \sum_{\mathbf{x}^{(i)} \in \mathcal{C}_1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_1)(\mathbf{x}^{(i)} - \boldsymbol{\mu}_1)^T$$

$$\mathbf{S}_2 = \sum_{\mathbf{x}^{(i)} \in \mathcal{C}_2} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_2)(\mathbf{x}^{(i)} - \boldsymbol{\mu}_2)^T$$

LDA Derivation

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

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \frac{\frac{\partial \mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\partial \mathbf{w}} \times \mathbf{w}^T \mathbf{S}_W \mathbf{w} - \frac{\partial \mathbf{w}^T \mathbf{S}_W \mathbf{w}}{\partial \mathbf{w}} \times \mathbf{w}^T \mathbf{S}_B \mathbf{w}}{(\mathbf{w}^T \mathbf{S}_W \mathbf{w})^2} = \frac{(2\mathbf{S}_B \mathbf{w}) \mathbf{w}^T \mathbf{S}_W \mathbf{w} - (2\mathbf{S}_W \mathbf{w}) \mathbf{w}^T \mathbf{S}_B \mathbf{w}}{(\mathbf{w}^T \mathbf{S}_W \mathbf{w})^2}$$

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{S}_B \mathbf{w} = \lambda \mathbf{S}_W \mathbf{w}$$

LDA Derivation

$$\mathbf{S}_B \mathbf{w} = \lambda \mathbf{S}_W \mathbf{w} \quad \text{if } \mathbf{S}_W \text{ is full-rank} \quad \longrightarrow \quad \mathbf{S}_W^{-1} \mathbf{S}_B \mathbf{w} = \lambda \mathbf{w}$$

- ▶ $\mathbf{S}_B \mathbf{w}$ (for any vector \mathbf{w}) points in the same direction as $\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2$:

$$\mathbf{S}_B \mathbf{w} = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \mathbf{w} \propto (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

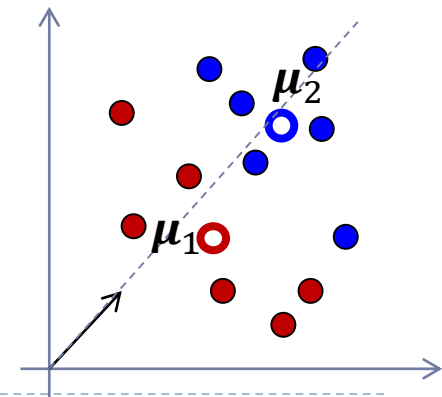
$$\mathbf{w} \propto \mathbf{S}_W^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

- ▶ Thus, we can solve the eigenvalue problem immediately

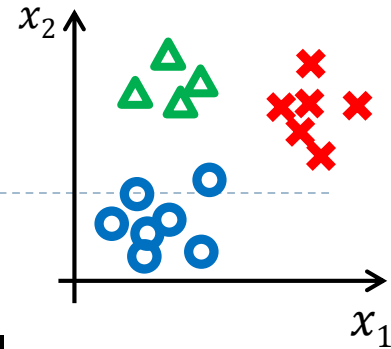
LDA Algorithm

- ▶ Find μ_1 and μ_2 as the mean of class 1 and 2 respectively
- ▶ Find S_1 and S_2 as scatter matrix of class 1 and 2 respectively
- ▶ $S_W = S_1 + S_2$
- ▶ $S_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T$

- ▶ Feature Extraction
 - ▶ $w = S_W^{-1}(\mu_1 - \mu_2)$ as the eigenvector corresponding to the largest eigenvalue of $S_W^{-1}S_B$
- ▶ Classification
 - ▶ $w = S_W^{-1}(\mu_1 - \mu_2)$
 - ▶ Using a threshold on $w^T x$, we can classify x



Multi-class classification



- ▶ Solutions to multi-category problems:

- ▶ Extend the learning algorithm to support multi-class:

- ▶ A function $f_i(\mathbf{x})$ for each class i is found

- $\hat{y} = \operatorname{argmax}_{i=1,\dots,c} f_i(\mathbf{x})$

\mathbf{x} is assigned to class C_i if $f_i(\mathbf{x}) > f_j(\mathbf{x}) \quad \forall j \neq i$

- ▶ Converting the problem to a set of two-class problems:

Converting multi-class problem to a set of two-class problems

- ▶ **“one versus rest”** or **“one against all”**

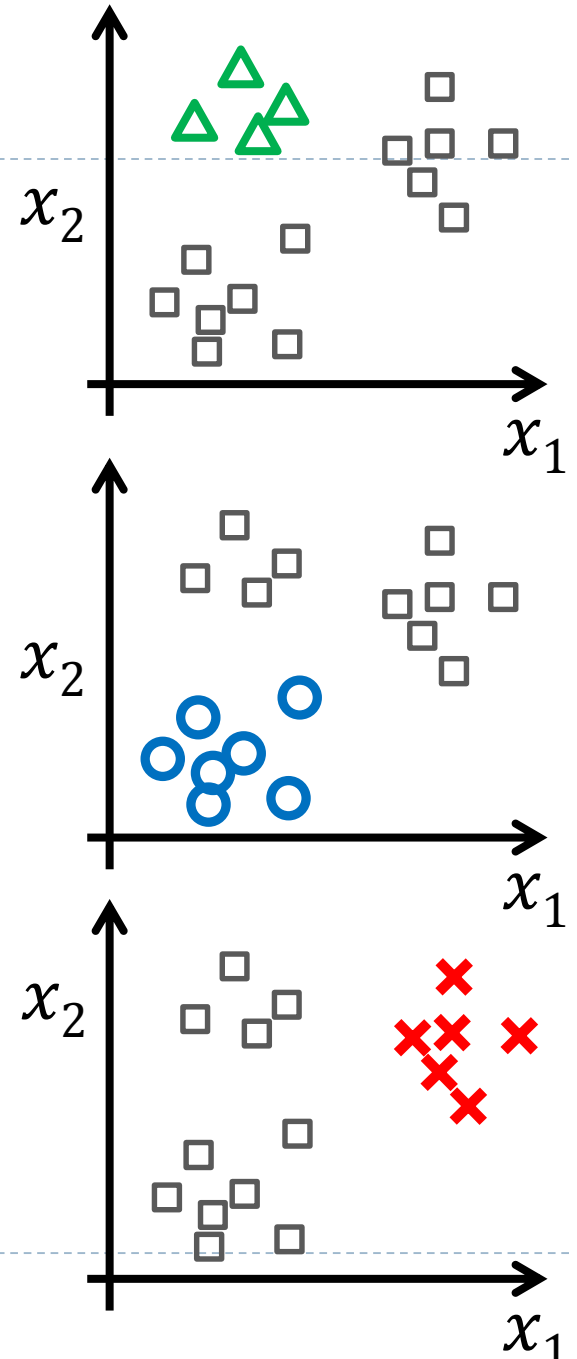
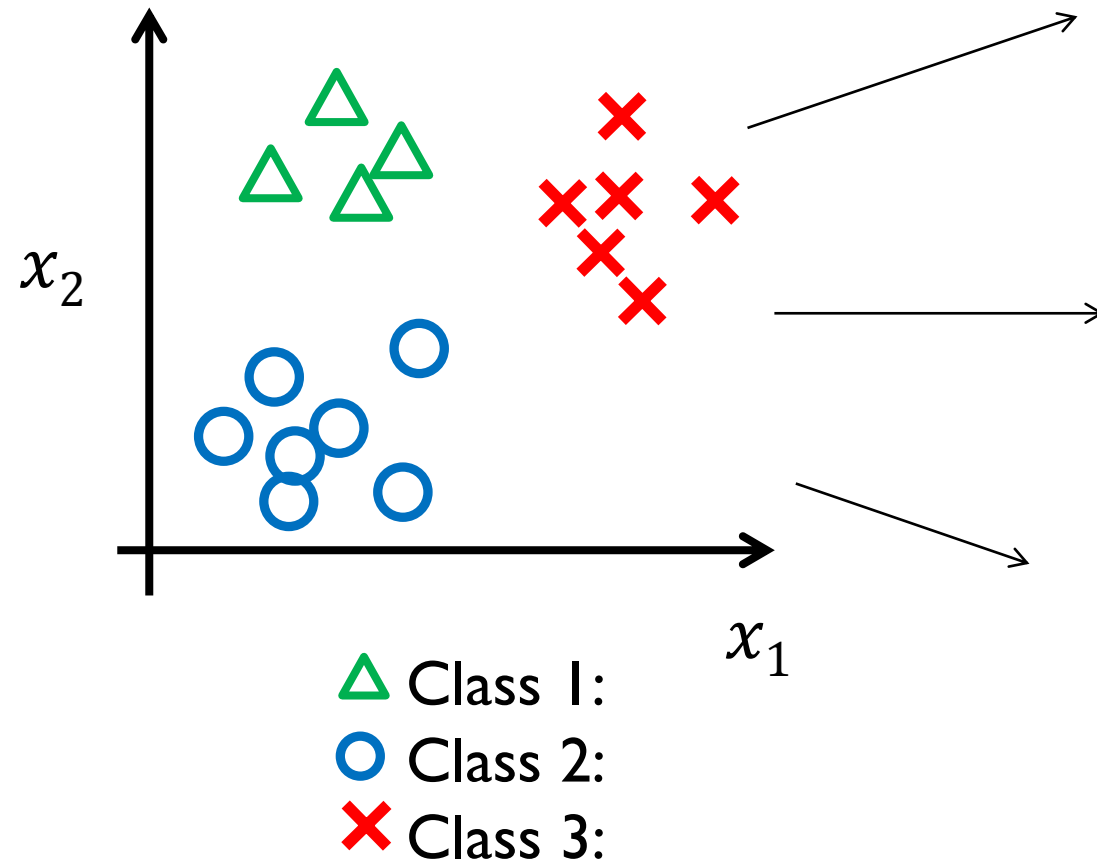
- ▶ For each class C_i , a linear discriminant function that separates samples of C_i from all the other samples is found.
 - ▶ Totally linearly separable

- ▶ **“one versus one”**

- ▶ $c(c - 1)/2$ linear discriminant functions are used, one to separate samples of a pair of classes.
 - ▶ Pairwise linearly separable

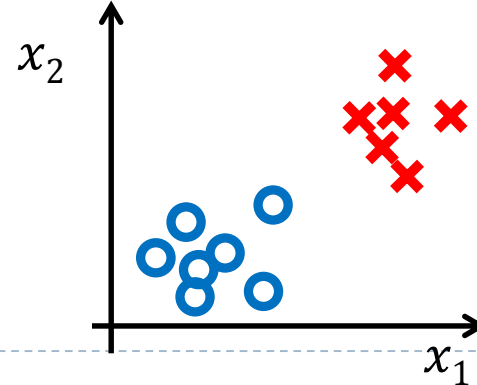
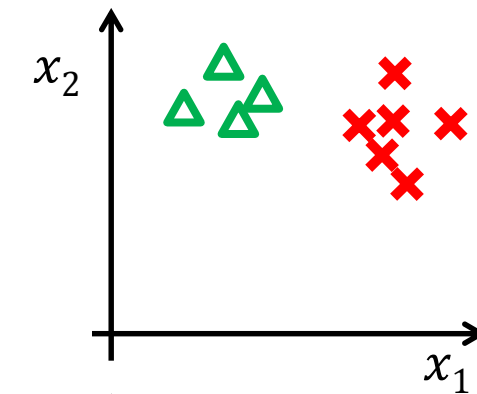
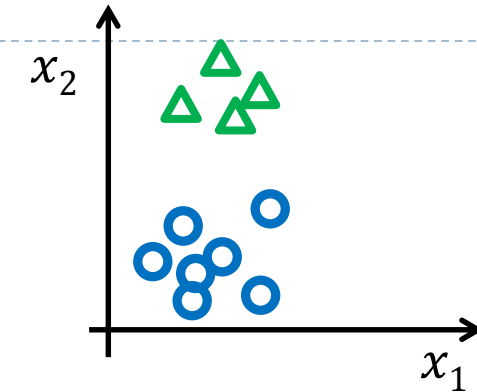
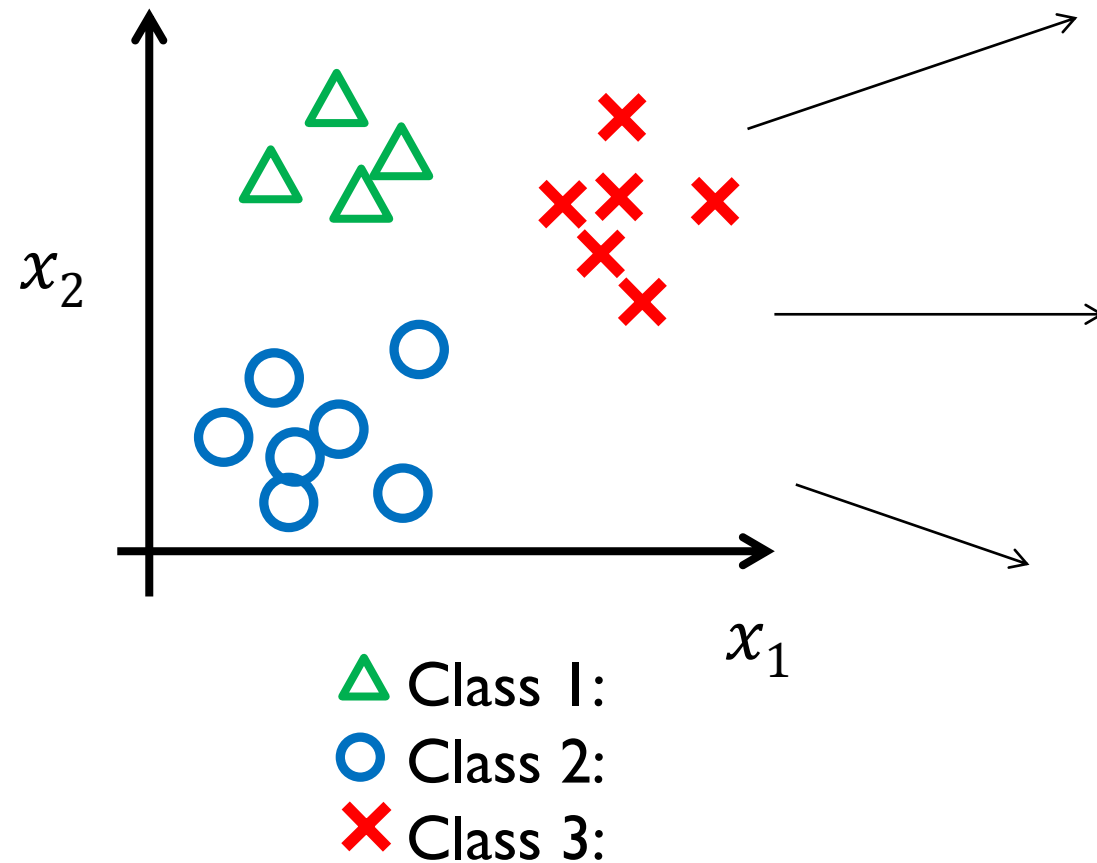
Multi-class classification

► One-vs-all (one-vs-rest)



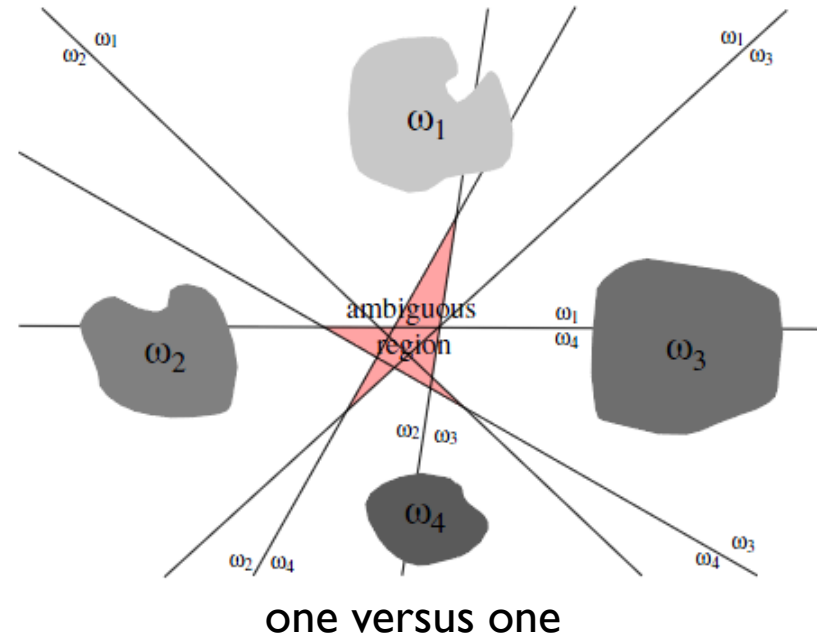
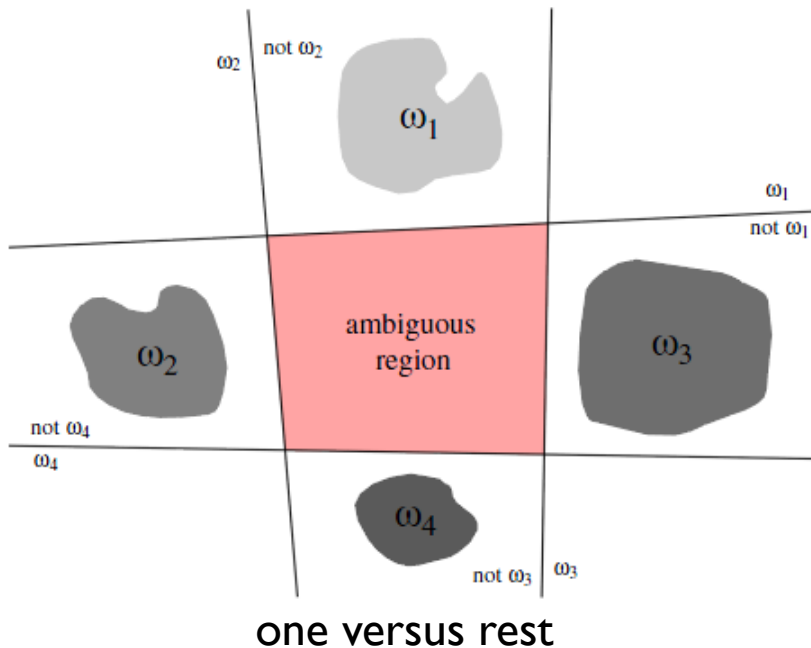
Multi-class classification

► One-vs-one



Multi-class classification: ambiguity

- ▶ Converting the multi-class problem to a set of two-class problems can lead to **regions in which the classification is undefined**



[Duda, Hart & Stork, 2002]

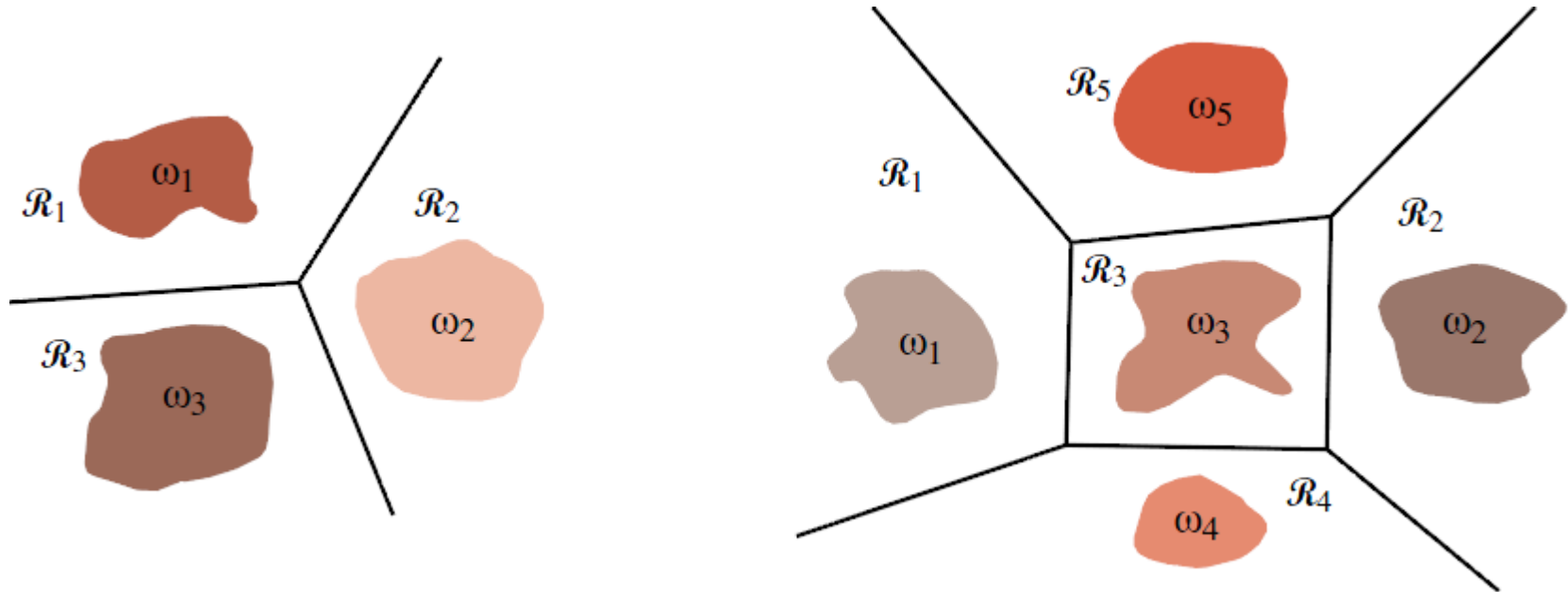
Multi-class classification: linear machine

- ▶ A discriminant function $f_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}$ for each class \mathcal{C}_i ($i = 1, \dots, K$):
 - ▶ \mathbf{x} is assigned to class \mathcal{C}_i if:

$$f_i(\mathbf{x}) > f_j(\mathbf{x}) \quad \forall j \neq i$$

- ▶ **Decision surfaces** (boundaries) can also be found using discriminant functions
 - ▶ Boundary of the contiguous \mathcal{R}_i and \mathcal{R}_j : $\forall \mathbf{x}, f_i(\mathbf{x}) = f_j(\mathbf{x})$
 - ▶ $(\mathbf{w}_i - \mathbf{w}_j)^T \mathbf{x} + (w_{i0} - w_{j0}) = 0$

Multi-class classification: linear machine



[Duda, Hart & Stork, 2002]

Perceptron: multi-class

$$\hat{y} = \operatorname{argmax}_{i=1,\dots,c} \mathbf{w}_i^T \mathbf{x}$$

$$J_P(\mathbf{W}) = - \sum_{i \in \mathcal{M}} \left(\mathbf{w}_{y^{(i)}} - \mathbf{w}_{\hat{y}^{(i)}} \right)^T \mathbf{x}^{(i)}$$

\mathcal{M} : subset of training data that are misclassified

$$\mathcal{M} = \{i \mid \hat{y}^{(i)} \neq y^{(i)}\}$$

Initialize $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_c]$, $k \leftarrow 0$

repeat

$k \leftarrow (k + 1) \bmod N$

if $\mathbf{x}^{(i)}$ is misclassified then

$$\mathbf{w}_{\hat{y}^{(i)}} = \mathbf{w}_{\hat{y}^{(i)}} - \mathbf{x}^{(i)}$$

$$\mathbf{w}_{y^{(i)}} = \mathbf{w}_{y^{(i)}} + \mathbf{x}^{(i)}$$

Until all patterns properly classified

Resources

- ▶ C. Bishop, “Pattern Recognition and Machine Learning”, Chapter 4.1.