



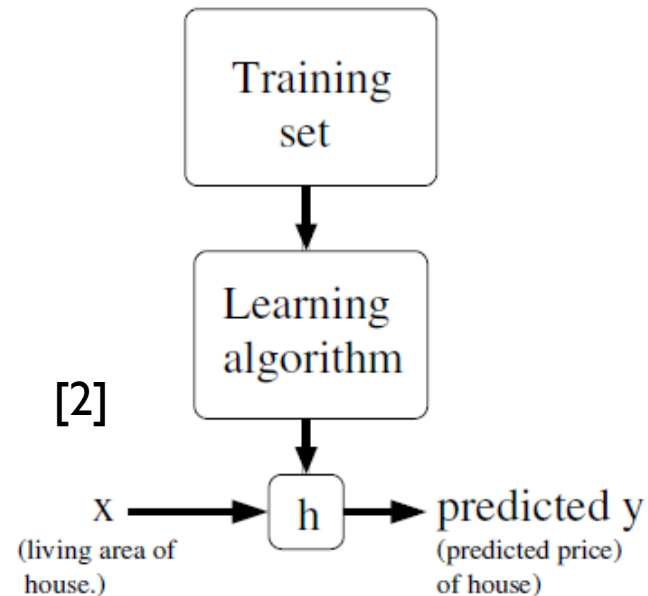
# Linear Classifiers

CE-477: Machine Learning - CS-828: Theory of Machine Learning  
Sharif University of Technology  
Fall 2024

Fatemeh Seyyedsalehi

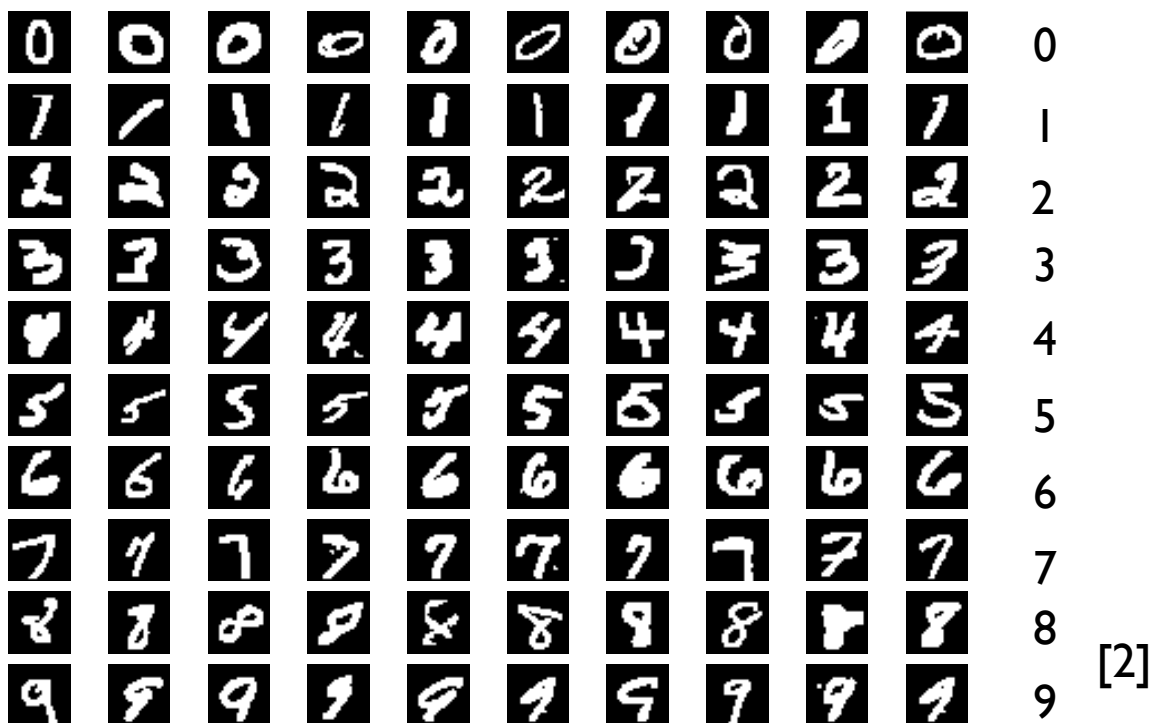
# A supervised problem

- ▶ Our goal is to learn a function  $h : \mathcal{X} \rightarrow \mathcal{Y}$
- ▶  $h(x)$  should be a good predictor for the corresponding  $y$
- ▶  $h$  is called a **hypothesis**



# Classification problem

- ▶ The values  $y$  takes on only a small number of discrete values.
  - ▶ A spam classifier for emails (0, 1)
  - ▶ Handwritten digit recognition



# Classification problem

- ▶ Given: Training set

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

- ▶  $y \in \{1, \dots, K\}$

- ▶ Two classes (binary):  $y \in \{0, 1\}$

- ▶ Multi classes: representing outputs by one-hot vectors,

$$y = [0, 1, 0, 0, 0]$$

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

- ▶ Discriminant function: takes an input vector  $\mathbf{x}$  and directly assigns it to one of  $K$  classes, denoted  $C_k$ .

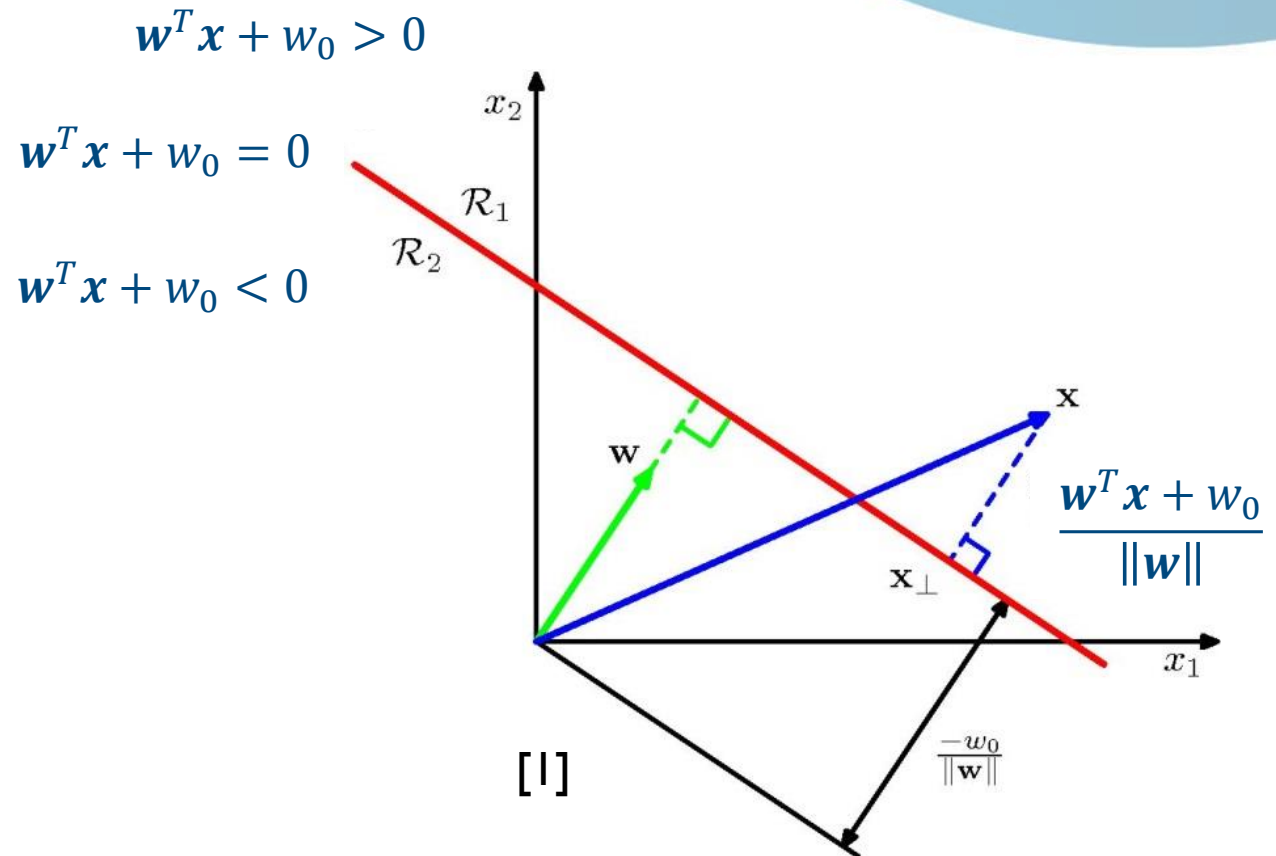
# Linear classifiers

- ▶ The hypothesis space:
  - ▶ The input space is divided into *decision regions* whose boundaries are called *decision boundaries* or *decision surfaces*.
  - ▶ Decision surfaces are linear functions of the input vector  $x$ 
    - ▶ Defined by  $(d-1)$ -dimensional hyperplanes within the  $d$ -dimensional input space.
- ▶ Linearly separable data: data points that can be exactly classified by a linear decision surface.
- ▶ Even when they are not optimal, we can use the simplicity of linear classifiers
  - ▶ Easy to compute
  - ▶ In the absence of information suggesting otherwise, linear classifiers are an attractive candidates for initial, trial classifiers.

# Binary classification

- ▶  $h(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x} + w_0 = w_0 + w_1 x_1 + \dots + w_d x_d$ 
  - ▶  $\mathbf{x} = [x_1 \ x_2 \ \dots \ x_d]$
  - ▶  $\mathbf{w} = [w_1 \ w_2 \ \dots \ w_d]$
  - ▶  $w_0$ : bias
- ▶ The linear discriminant function:
  - ▶ if  $\mathbf{w}^T \mathbf{x} + w_0 \geq 0$  then  $\mathcal{C}_1$  else  $\mathcal{C}_2$
- ▶ Decision surface (boundary):  $\mathbf{w}^T \mathbf{x} + w_0 = 0$

# Linear boundary: geometry



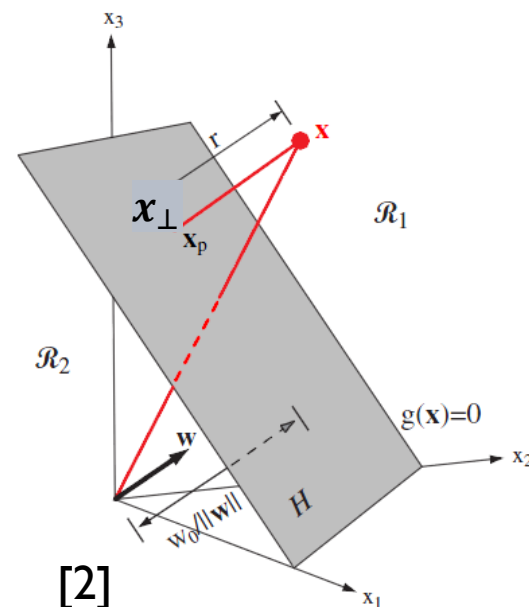
# Linear classifier: Two classes

- ▶ 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  $[w_1, \dots, w_d]$
- ▶  $w_0$  determines the location of the surface.
  - ▶ The normal distance from the origin to the decision surface is  $\frac{w_0}{\|w\|}$

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

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

gives a signed measure of the perpendicular distance  $r$  of the point  $x$  from the decision surface





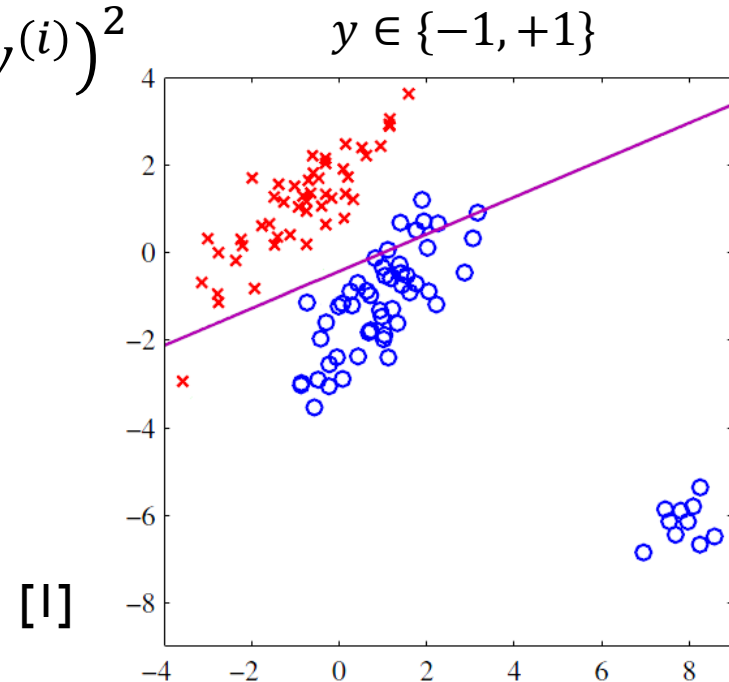
# Cost Function for linear classifiers

- ▶ 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 best parameters:
    - ▶ Find optimal  $\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} J(\mathbf{w})$
- ▶ Criterion or cost functions for classification:
  - ▶ We will investigate several cost functions for the classification problem

# Cost Function for linear classifiers

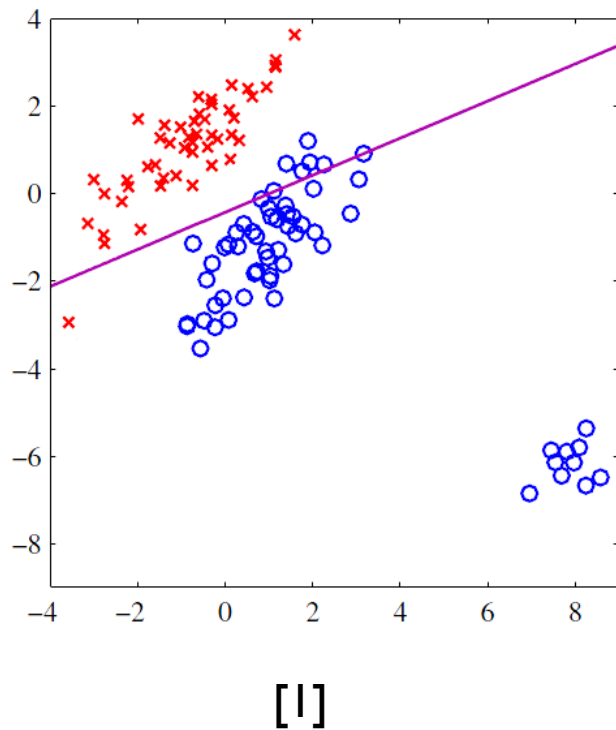
- ▶ SSE as the linear regression
  - ▶ Is not suitable for classification
  - ▶ Penalizes 'too correct' predictions (that they lie a long way on the correct side of the decision)

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



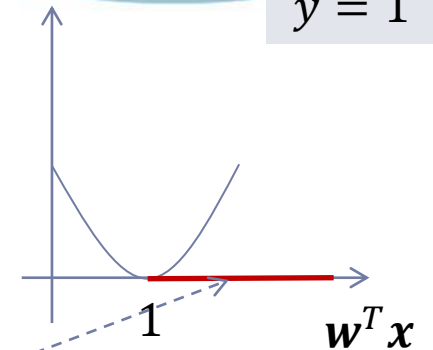
# SSE cost function for classification

$K = 2$

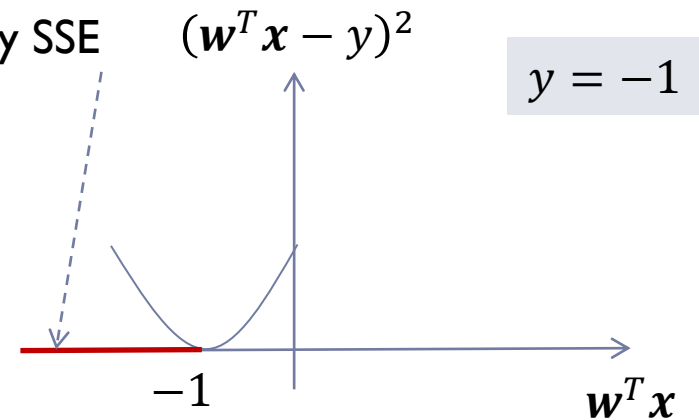


$$(w^T x - y)^2$$

$y = 1$



Correct predictions that  
are penalized by SSE

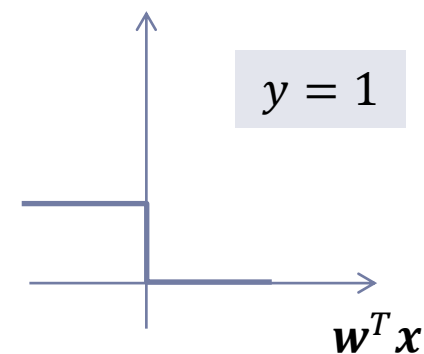


# Cost Function for linear classifiers

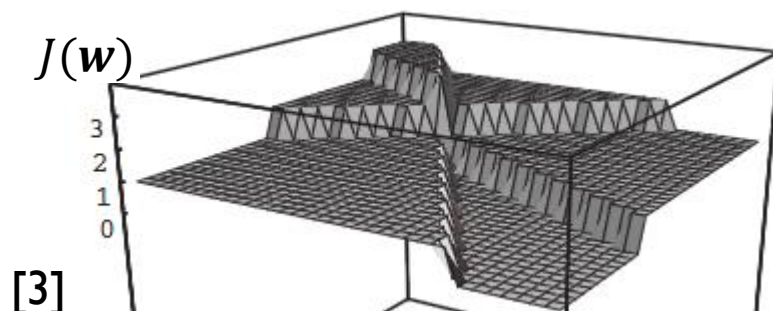
- Is it more suitable if we set  $h(\mathbf{x}) = g(\mathbf{w}^T \mathbf{x})$ ?

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

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



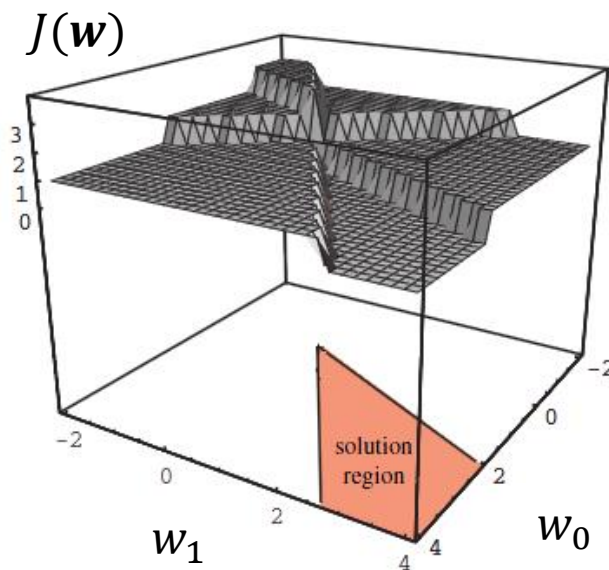
- $J(\mathbf{w})$  is a piecewise constant function shows the number of misclassifications



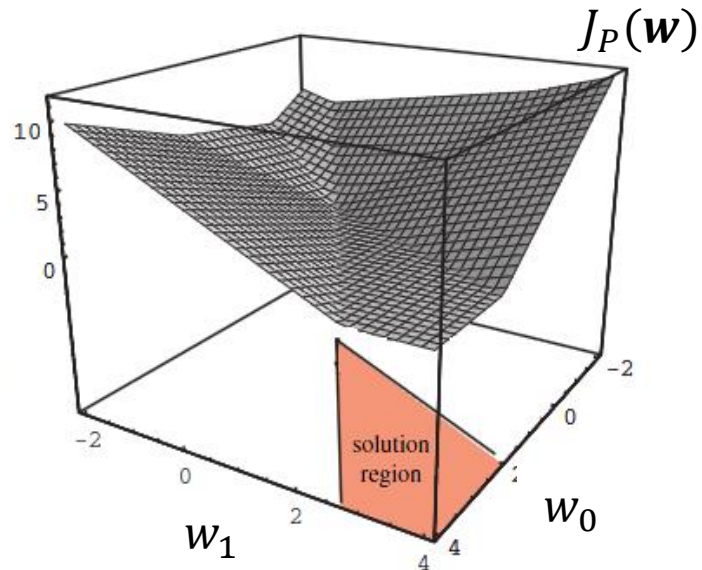
# Perceptron criterion

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

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



# of misclassifications  
as a cost function



[3]

Perceptron's  
cost function

# 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:

# Stochastic gradient descent for Perceptron

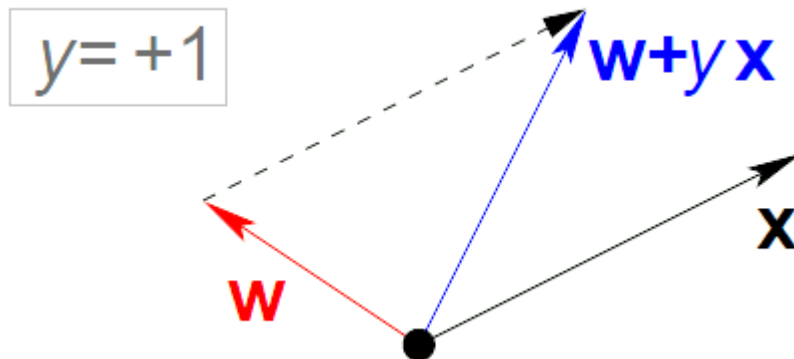
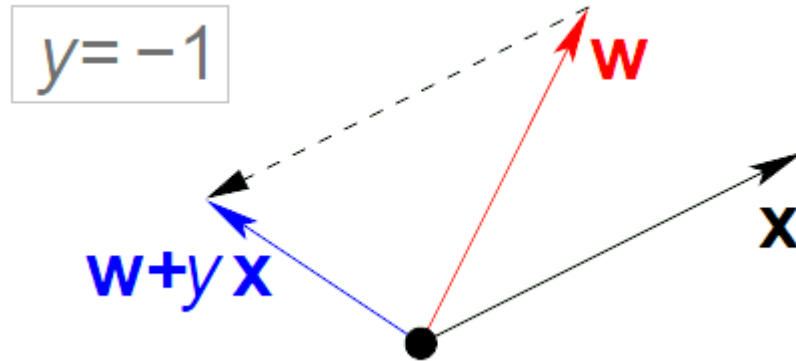
- ▶ Single-sample perceptron:

- ▶ 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

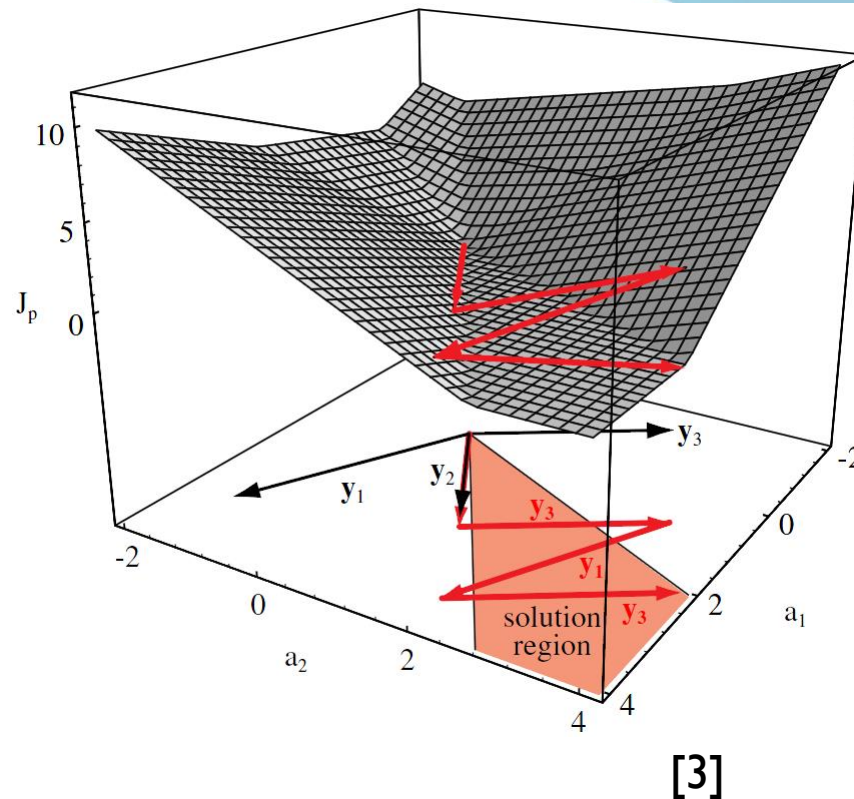
# Stochastic gradient descent for Perceptron



[2]



# Convergence of Perceptron



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

# Pocket algorithm

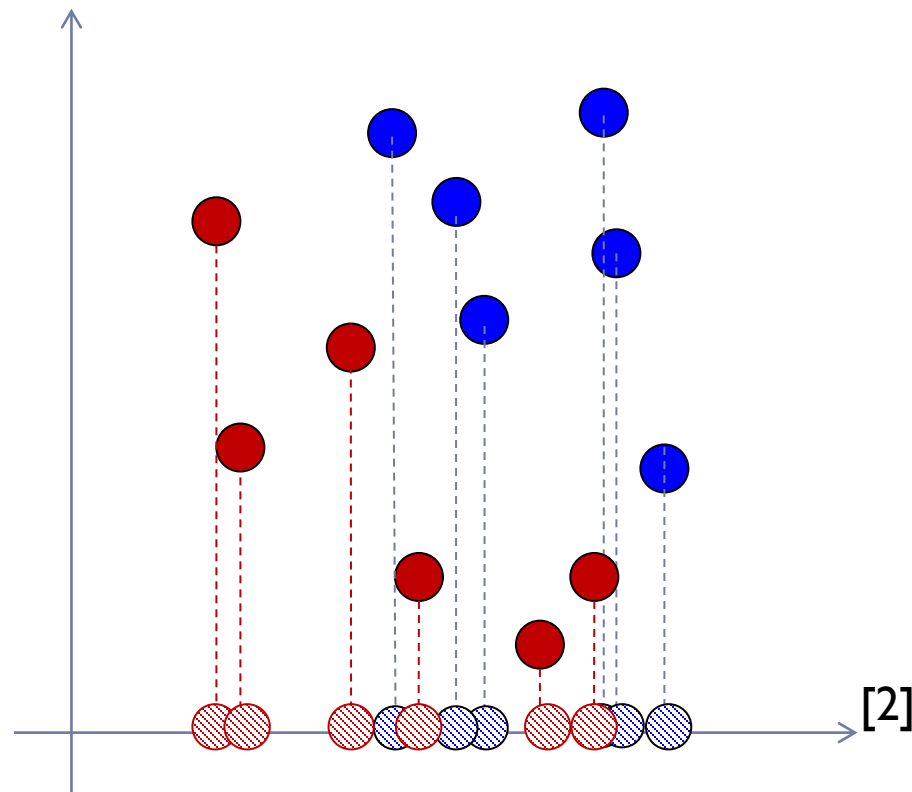
- ▶ For the data that are not linearly separable due to noise:
  - ▶ Keeps in its pocket the best  $\mathbf{w}$  encountered up to now.

```
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
```

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

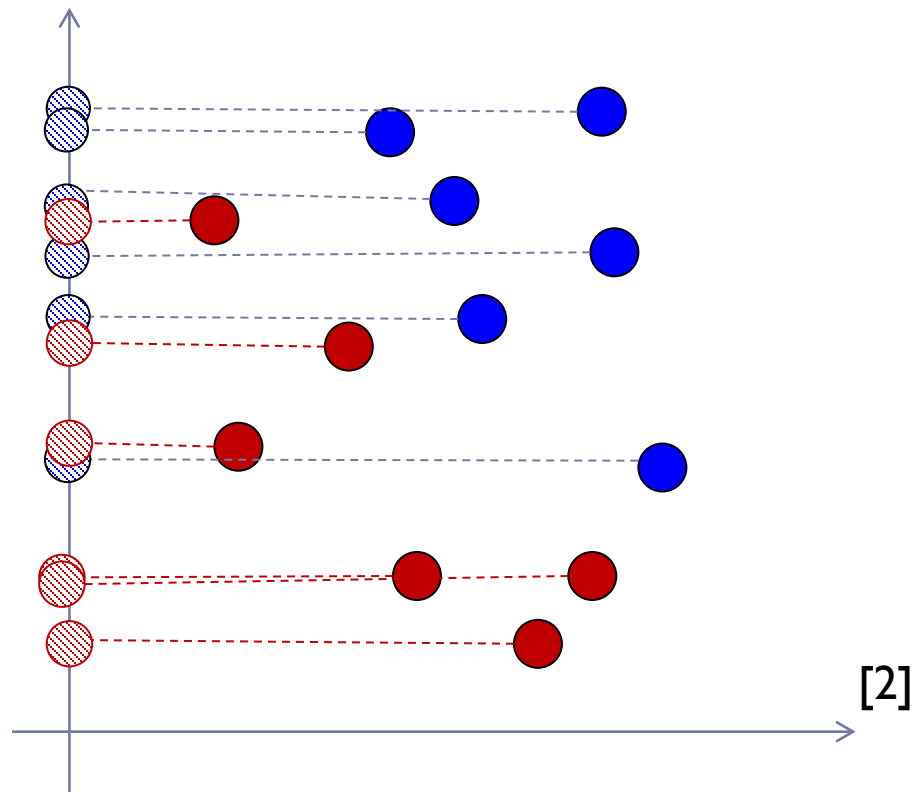
# Good Projection for Classification

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



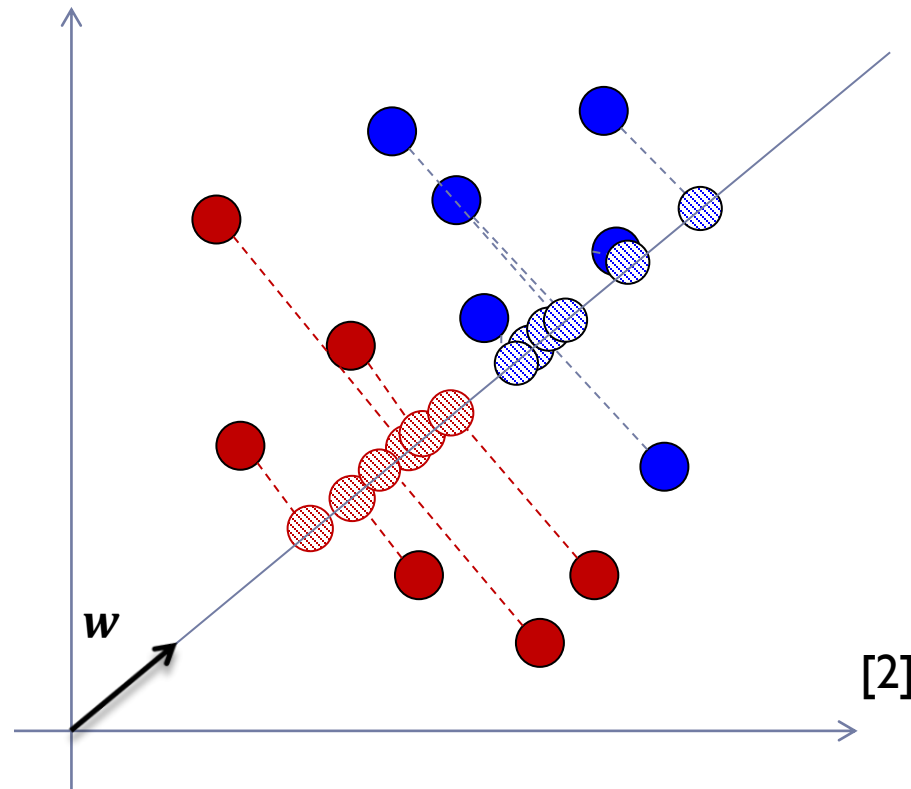
# Good Projection for Classification

- ▶ 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

- ▶ Fisher's Linear Discriminant Analysis
- ▶ Problem definition:
  - ▶  $K = 2$  classes
  - ▶  $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$  training samples with  $N_1$  samples from the first class ( $\mathcal{C}_1$ ) and  $N_2$  samples from the second class ( $\mathcal{C}_2$ )
  - ▶ 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

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

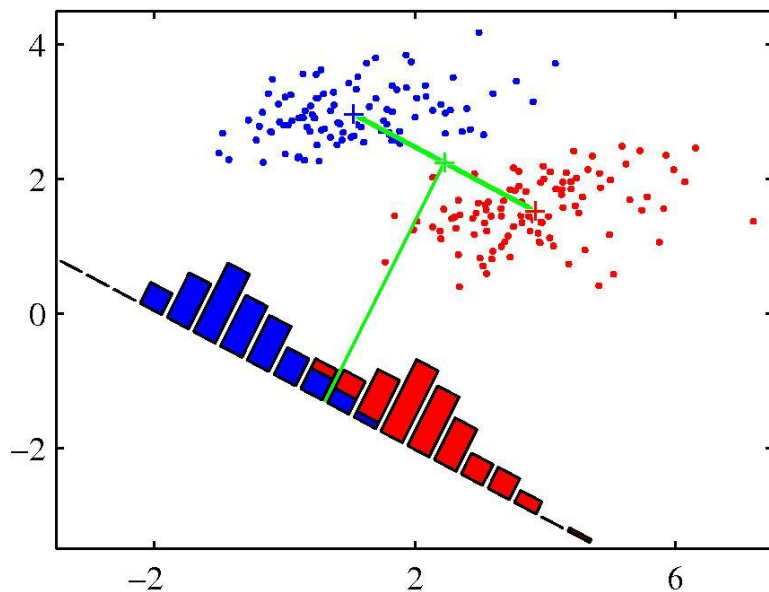
$$\begin{aligned} \max_{\mathbf{w}} J(\mathbf{w}) &= (\mu'_1 - \mu'_2)^2 \\ \text{s. t. } \|\mathbf{w}\| &= 1 \end{aligned}$$

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

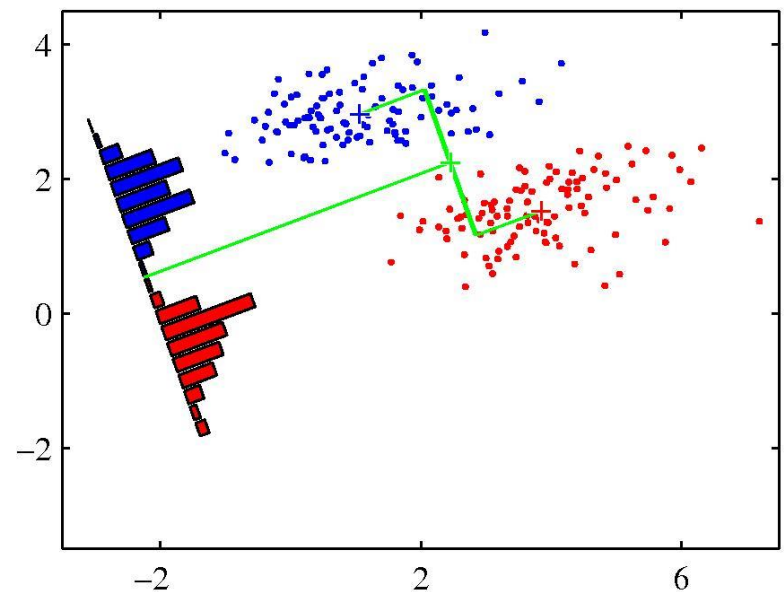
- ▶ What is the problem with the criteria considering only  $|\mu'_1 - \mu'_2|$ ?
  - ▶ It does not consider the variances of the classes in the projected direction

# Measure of Separation in the Projected Direction

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



[1]





# 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 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

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

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

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

# 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 Criteria

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

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

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

# LDA Derivation

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

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

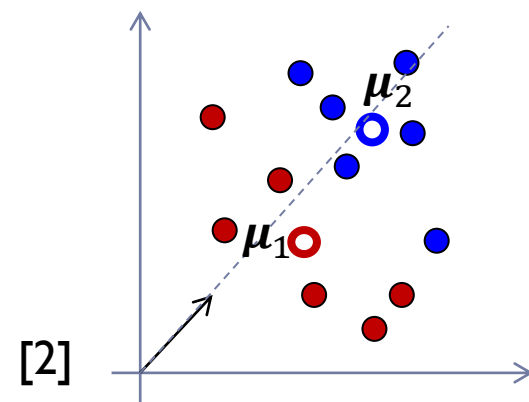
$$\mathbf{S}_B \mathbf{w} = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \mathbf{w} \propto (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

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

- ▶ Thus, we can solve the eigenvalue problem immediately

# LDA Algorithm

- ▶ Find  $\mu_1$  and  $\mu_2$  as the mean of class 1 and 2 respectively
- ▶ Find  $S_1$  and  $S_2$  as scatter matrix of class 1 and 2 respectively
- ▶  $S_W = S_1 + S_2$
- ▶ Classification
  - ▶  $w = S_W^{-1}(\mu_1 - \mu_2)$
  - ▶ Using a threshold on  $w^T x$ , we can classify  $x$



# References

- ▶ [1]: C. Bishop, “Pattern Recognition and Machine Learning”, Chapter 4.1.
- ▶ [2]: Mahdiah Soleymani, Machine learning, Sharif university of technology
- ▶ [3]: Pattern classification, Duda, Hart & Stork, 2002