

# ECE595 / STAT598: Machine Learning I

## Lecture 17 Perceptron 2: Algorithm and Property

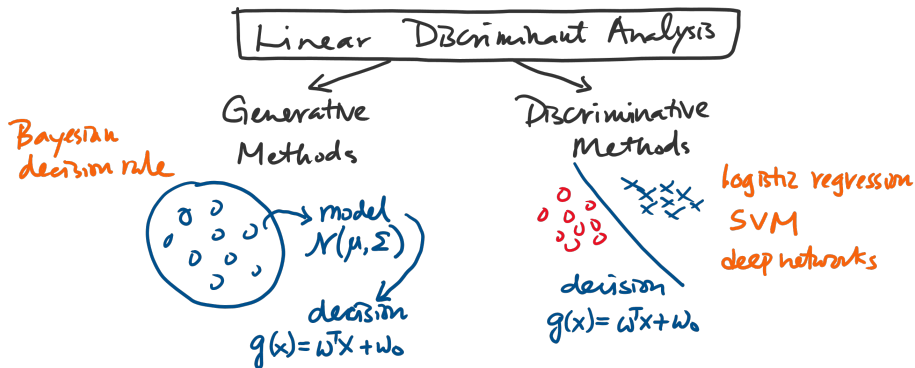
Spring 2020

Stanley Chan

School of Electrical and Computer Engineering  
Purdue University



# Overview



- In linear discriminant analysis (LDA), there are generally two types of approaches
- **Generative approach:** Estimate model, then define the classifier
- **Discriminative approach:** Directly define the classifier

# Outline

## Discriminative Approaches

- Lecture 16 Perceptron 1: Definition and Basic Concepts
- **Lecture 17 Perceptron 2: Algorithm and Property**
- Lecture 18 Multi-Layer Perceptron: Back Propagation

## This lecture: Perceptron 2

- Perceptron Algorithm
  - Loss Function
  - Algorithm
- Optimality
  - Uniqueness
  - Batch and Online Mode
- Convergence
  - Main Results
  - Implication

## Perceptron with Hard Loss

- Historically, we have perceptron algorithm way earlier than CVX.
- Before the age of CVX, people solve perceptron using gradient descent.
- Let us be explicit about which loss:

$$J_{\text{hard}}(\boldsymbol{\theta}) = \sum_{j=1}^N \max \left\{ -y_j h_{\boldsymbol{\theta}}(\mathbf{x}_j), 0 \right\}$$

$$J_{\text{soft}}(\boldsymbol{\theta}) = \sum_{j=1}^N \max \left\{ -y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j), 0 \right\}$$

- Goal:** To get a solution for  $J_{\text{hard}}(\boldsymbol{\theta})$
- Approach:** Gradient descent on  $J_{\text{soft}}(\boldsymbol{\theta})$

## Re-defining the Loss

- **Main idea:** Use the fact that

$$J_{\text{soft}}(\boldsymbol{\theta}) = \sum_{j=1}^N \max \left\{ -y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j), 0 \right\}$$

is the same as this loss function

$$J(\boldsymbol{\theta}) = - \sum_{j \in \mathcal{M}(\boldsymbol{\theta})} y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j).$$

- $\mathcal{M}(\boldsymbol{\theta}) \subseteq \{1, \dots, N\}$  is the set of misclassified samples.
- Run gradient descent on  $J(\boldsymbol{\theta})$ , but fixing  $\mathcal{M}(\boldsymbol{\theta}) \leftarrow \mathcal{M}(\boldsymbol{\theta}^k)$  for iteration  $k$ .

## Equivalent Perceptron Loss

- We want to show that the perceptron loss function is equivalent to

$$\underbrace{\sum_{j=1}^N \max \left\{ -y_j g_{\theta}(\mathbf{x}_j), 0 \right\}}_{J_{\text{soft}}(\theta)} = - \underbrace{\sum_{j \in \mathcal{M}(\theta)} y_j g_{\theta}(\mathbf{x}_j)}_{J(\theta)}$$

- If  $\mathbf{x}_j$  is misclassified ( $j \in \mathcal{M}(\theta)$ )
  - then by definition of  $\mathcal{M}(\theta)$  we have  $\text{sign} \{g_{\theta}(\mathbf{x}_j)\} \neq y_j$
  - So  $-y_j g_{\theta}(\mathbf{x}_j) > 0$
  - Therefore,  $\max \{-y_j g_{\theta}(\mathbf{x}_j), 0\} = -y_j g_{\theta}(\mathbf{x}_j)$ .
- If  $\mathbf{x}_j$  is correctly classified ( $j \notin \mathcal{M}(\theta)$ )
  - then by definition of  $\mathcal{M}(\theta)$  we have  $\text{sign} \{g_{\theta}(\mathbf{x}_j)\} = y_j$
  - So  $-y_j g_{\theta}(\mathbf{x}_j) < 0$
  - Therefore,  $\max \{-y_j g_{\theta}(\mathbf{x}_j), 0\} = 0$ .

## Equivalent Perceptron Loss

- Therefore, we conclude that

$$\mathcal{M}(\boldsymbol{\theta}) = \{j \mid y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j) < 0\}$$

## Equivalent Perceptron Loss

- Therefore, we conclude that

$$\mathcal{M}(\boldsymbol{\theta}) = \{j \mid y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j) < 0\}$$

- and

$$\begin{aligned} J_{\text{soft}}(\boldsymbol{\theta}) &= \sum_{j \in \mathcal{M}(\boldsymbol{\theta})} \max\{-y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j), 0\} + \sum_{j \notin \mathcal{M}(\boldsymbol{\theta})} \max\{-y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j), 0\} \\ &= \sum_{j \in \mathcal{M}(\boldsymbol{\theta})} -y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j) + \sum_{j \notin \mathcal{M}(\boldsymbol{\theta})} 0 \\ &= \sum_{j \in \mathcal{M}(\boldsymbol{\theta})} -y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j) = J(\boldsymbol{\theta}). \end{aligned}$$



## Equivalent Perceptron Loss

- Therefore, we conclude that

$$\mathcal{M}(\boldsymbol{\theta}) = \{j \mid y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j) < 0\}$$

- and

$$\begin{aligned} J_{\text{soft}}(\boldsymbol{\theta}) &= \sum_{j \in \mathcal{M}(\boldsymbol{\theta})} \max\{-y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j), 0\} + \sum_{j \notin \mathcal{M}(\boldsymbol{\theta})} \max\{-y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j), 0\} \\ &= \sum_{j \in \mathcal{M}(\boldsymbol{\theta})} -y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j) + \sum_{j \notin \mathcal{M}(\boldsymbol{\theta})} 0 \\ &= \sum_{j \in \mathcal{M}(\boldsymbol{\theta})} -y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j) = J(\boldsymbol{\theta}). \end{aligned}$$

## Equivalent Perceptron Loss

- Therefore, we conclude that

$$\mathcal{M}(\boldsymbol{\theta}) = \{j \mid y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j) < 0\}$$

- and

$$\begin{aligned} J_{\text{soft}}(\boldsymbol{\theta}) &= \sum_{j \in \mathcal{M}(\boldsymbol{\theta})} \max\{-y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j), 0\} + \sum_{j \notin \mathcal{M}(\boldsymbol{\theta})} \max\{-y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j), 0\} \\ &= \sum_{j \in \mathcal{M}(\boldsymbol{\theta})} -y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j) + \sum_{j \notin \mathcal{M}(\boldsymbol{\theta})} 0 \\ &= \sum_{j \in \mathcal{M}(\boldsymbol{\theta})} -y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j) = J(\boldsymbol{\theta}). \end{aligned}$$

## Equivalent Perceptron Loss

- Therefore, we conclude that

$$\mathcal{M}(\boldsymbol{\theta}) = \{j \mid y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j) < 0\}$$

- and

$$\begin{aligned} J_{\text{soft}}(\boldsymbol{\theta}) &= \sum_{j \in \mathcal{M}(\boldsymbol{\theta})} \max\{-y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j), 0\} + \sum_{j \notin \mathcal{M}(\boldsymbol{\theta})} \max\{-y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j), 0\} \\ &= \sum_{j \in \mathcal{M}(\boldsymbol{\theta})} -y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j) + \sum_{j \notin \mathcal{M}(\boldsymbol{\theta})} 0 \\ &= \sum_{j \in \mathcal{M}(\boldsymbol{\theta})} -y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j) = J(\boldsymbol{\theta}). \end{aligned}$$

## Equivalent Perceptron Loss

- Therefore, we conclude that

$$\mathcal{M}(\boldsymbol{\theta}) = \{j \mid y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j) < 0\}$$

- and

$$\begin{aligned} J_{\text{soft}}(\boldsymbol{\theta}) &= \sum_{j \in \mathcal{M}(\boldsymbol{\theta})} \max\{-y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j), 0\} + \sum_{j \notin \mathcal{M}(\boldsymbol{\theta})} \max\{-y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j), 0\} \\ &= \sum_{j \in \mathcal{M}(\boldsymbol{\theta})} -y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j) + \sum_{j \notin \mathcal{M}(\boldsymbol{\theta})} 0 \\ &= \sum_{j \in \mathcal{M}(\boldsymbol{\theta})} -y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j) = J(\boldsymbol{\theta}). \end{aligned}$$

- Minimizing  $J(\boldsymbol{\theta})$  is less obvious because  $\mathcal{M}(\boldsymbol{\theta})$  depends on  $\boldsymbol{\theta}$ .

## Equivalent Perceptron Loss

- Therefore, we conclude that

$$\mathcal{M}(\boldsymbol{\theta}) = \{j \mid y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j) < 0\}$$

- and

$$\begin{aligned} J_{\text{soft}}(\boldsymbol{\theta}) &= \sum_{j \in \mathcal{M}(\boldsymbol{\theta})} \max\{-y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j), 0\} + \sum_{j \notin \mathcal{M}(\boldsymbol{\theta})} \max\{-y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j), 0\} \\ &= \sum_{j \in \mathcal{M}(\boldsymbol{\theta})} -y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j) + \sum_{j \notin \mathcal{M}(\boldsymbol{\theta})} 0 \\ &= \sum_{j \in \mathcal{M}(\boldsymbol{\theta})} -y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j) = J(\boldsymbol{\theta}). \end{aligned}$$

- Minimizing  $J(\boldsymbol{\theta})$  is less obvious because  $\mathcal{M}(\boldsymbol{\theta})$  depends on  $\boldsymbol{\theta}$ .
- But it gives a very easy algorithm.

# Perceptron Algorithm

- The loss is

$$J(\boldsymbol{\theta}) = - \sum_{j \in \mathcal{M}(\boldsymbol{\theta})} y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j),$$

- At iteration  $k$ , fix  $\mathcal{M}_k = \mathcal{M}(\boldsymbol{\theta}^{(k)})$
- Then, update via gradient descent

$$\begin{aligned}\boldsymbol{\theta}^{(k+1)} &= \boldsymbol{\theta}^{(k)} - \alpha_k \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}^{(k)}) \\ &= \boldsymbol{\theta}^{(k)} - \alpha_k \sum_{j \in \mathcal{M}_k} \nabla_{\boldsymbol{\theta}} \left( -y_j g_{\boldsymbol{\theta}}(\mathbf{x}_j) \right).\end{aligned}$$

# Perceptron Algorithm

- We can show that

$$\begin{aligned}\nabla_{\theta} \left( -y_j g_{\theta}(\mathbf{x}_j) \right) &= \begin{cases} -y_j \nabla_{\theta} \left( \mathbf{w}^T \mathbf{x}_j + w_0 \right) & , \\ 0, & , \end{cases} \\ &= \begin{cases} = -y_j \begin{bmatrix} \mathbf{x}_j \\ 1 \end{bmatrix} & \text{if } j \in \mathcal{M}_k, \\ 0, & \text{if } j \notin \mathcal{M}_k. \end{cases}\end{aligned}$$

- Thus, the update is

$$\begin{bmatrix} \mathbf{w}^{(k+1)} \\ w_0^{(k+1)} \end{bmatrix} = \begin{bmatrix} \mathbf{w}^{(k)} \\ w_0^{(k)} \end{bmatrix} + \alpha_k \sum_{j \in \mathcal{M}_k} \begin{bmatrix} y_j \mathbf{x}_j \\ y_j \end{bmatrix}.$$

# Perceptron Algorithm

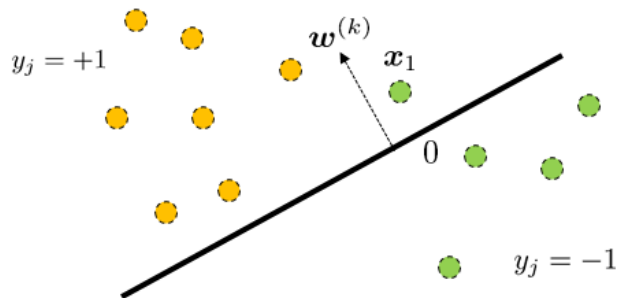
- The algorithm is
- For  $k = 1, 2, \dots$ ,
- Update  $\mathcal{M}_k = \{j \mid y_j g_{\theta}(\mathbf{x}_j) < 0\}$  for  $\theta = \theta^{(k)}$ .
- Gradient descent

$$\begin{bmatrix} \mathbf{w}^{(k+1)} \\ w_0^{(k+1)} \end{bmatrix} = \begin{bmatrix} \mathbf{w}^{(k)} \\ w_0^{(k)} \end{bmatrix} + \alpha_k \sum_{j \in \mathcal{M}_k} \begin{bmatrix} y_j \mathbf{x}_j \\ y_j \end{bmatrix}.$$

- End For
- The set  $\mathcal{M}_k$  can grow or can shrink from  $\mathcal{M}_{k-1}$ .
- If training samples are linearly separable, then converge. Zero training loss.
- If training samples are not linearly separable, then oscillates.



## Updating One Sample



# Outline

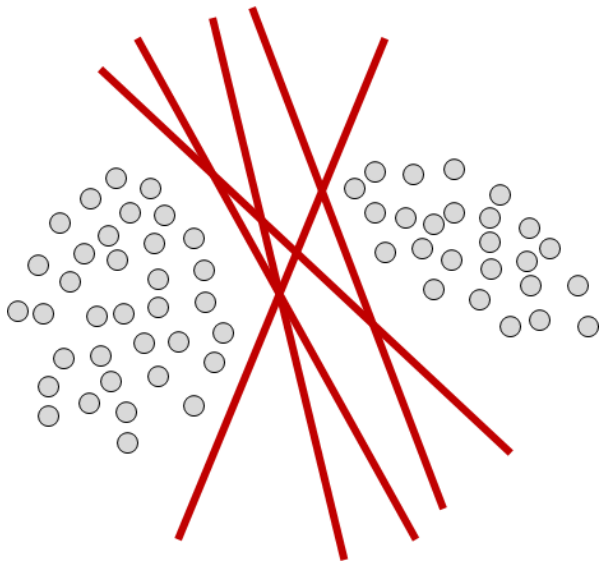
## Discriminative Approaches

- Lecture 16 Perceptron 1: Definition and Basic Concepts
- **Lecture 17 Perceptron 2: Algorithm and Property**
- Lecture 18 Multi-Layer Perceptron: Back Propagation

## This lecture: Perceptron 2

- Perceptron Algorithm
  - Loss Function
  - Algorithm
- Optimality
  - Uniqueness
  - Batch and Online Mode
- Convergence
  - Main Results
  - Implication

## Non-uniqueness of Global Minimizer



# Optimality of Perceptron Algorithm

- Let **perceptron algorithm** output

$$\theta_{\text{perceptron}}^* = \text{Perceptron Algorithm}(\{\mathbf{x}_1, \dots, \mathbf{x}_N\}).$$

- Let **ideal** solution

$$\theta_{\text{hard}}^* = \underset{\theta}{\operatorname{argmin}} J_{\text{hard}}(\theta).$$

That means

$$J_{\text{hard}}(\theta_{\text{hard}}^*) \leq J_{\text{hard}}(\theta), \quad \forall \theta.$$

- If the two classes are linearly separable, then  $\theta_{\text{perceptron}}^*$  is a global minimizer:

$$J_{\text{hard}}(\theta_{\text{perceptron}}^*) \leq J_{\text{hard}}(\theta), \quad \forall \theta.$$

and

$$J_{\text{hard}}(\theta_{\text{perceptron}}^*) = J_{\text{hard}}(\theta_{\text{hard}}^*) = 0.$$

## Uniqueness of Perceptron Solution

- If  $\theta^*$  minimizes  $J_{\text{hard}}(\theta^*)$ , then  $\alpha\theta^*$  for some constant  $\alpha > 0$  also minimizes  $J_{\text{hard}}(\theta^*)$ .
- This is because

$$\begin{aligned}g_{\alpha\theta}(\mathbf{x}) &= (\alpha\mathbf{w})^T \mathbf{x} + (\alpha w_0) \\ &= \alpha(\mathbf{w}^T \mathbf{x} + w_0).\end{aligned}$$

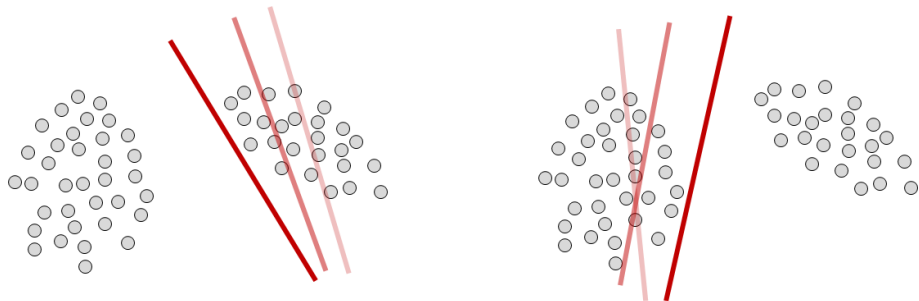
- If  $g_{\theta}(\mathbf{x}) > 0$ , then  $g_{\alpha\theta}(\mathbf{x}) > 0$ . So if  $h_{\theta}(\mathbf{x}) = +1$ , then  $h_{\alpha\theta}(\mathbf{x}) = +1$ .
- If  $g_{\theta}(\mathbf{x}) < 0$ , then  $g_{\alpha\theta}(\mathbf{x}) < 0$ . So if  $h_{\theta}(\mathbf{x}) = -1$ , then  $h_{\alpha\theta}(\mathbf{x}) = -1$ .
- The sign of  $\mathbf{w}^T \mathbf{x} + w_0$  is unchanged as long as  $\alpha > 0$ .

$$\begin{aligned}J_{\text{hard}}(\theta^*) &= \sum_{j=1}^N \max \left\{ -y_j h_{\theta^*}(\mathbf{x}_j), 0 \right\} \\ &= \sum_{j=1}^N \max \left\{ -y_j h_{\alpha\theta^*}(\mathbf{x}_j), 0 \right\} = J_{\text{hard}}(\alpha\theta^*)\end{aligned}$$

# Factors for Uniqueness

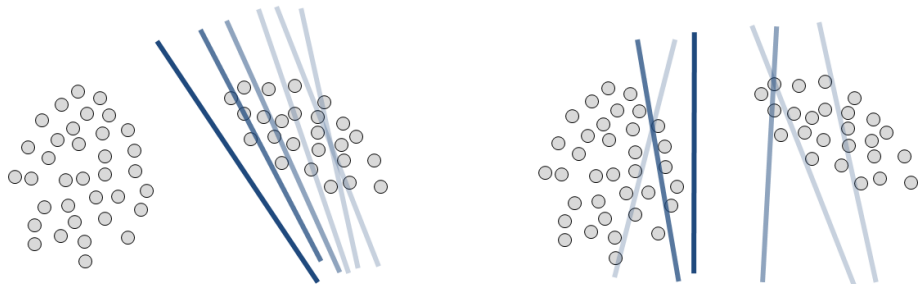
- **Initialization**

- Start at a different location, end on a different location
- You still converge, but no longer unique solution
- $\mathcal{M}_k$  changes



# Factors for Uniqueness

- **Step Size**
- Too large step: oscillate
- Too small step: slow movement
- Terminates as long as no misclassification



## Batch vs Online Mode

- **Batch mode**

$$\begin{bmatrix} \mathbf{w}^{(k+1)} \\ w_0^{(k+1)} \end{bmatrix} = \begin{bmatrix} \mathbf{w}^{(k)} \\ w_0^{(k)} \end{bmatrix} + \alpha_k \sum_{j \in \mathcal{M}_k} \begin{bmatrix} y_j \mathbf{x}_j \\ y_j \end{bmatrix}.$$

Update via the average of misclassified samples

- **Online mode**

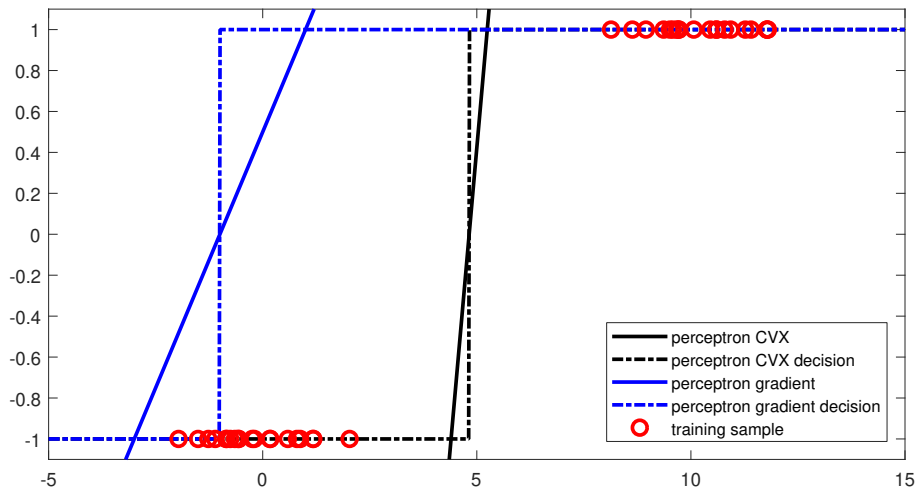
$$\begin{bmatrix} \mathbf{w}^{(k+1)} \\ w_0^{(k+1)} \end{bmatrix} = \begin{bmatrix} \mathbf{w}^{(k)} \\ w_0^{(k)} \end{bmatrix} + \alpha_k \begin{bmatrix} y_j \mathbf{x}_j \\ y_j \end{bmatrix},$$

Update via a single misclassified sample

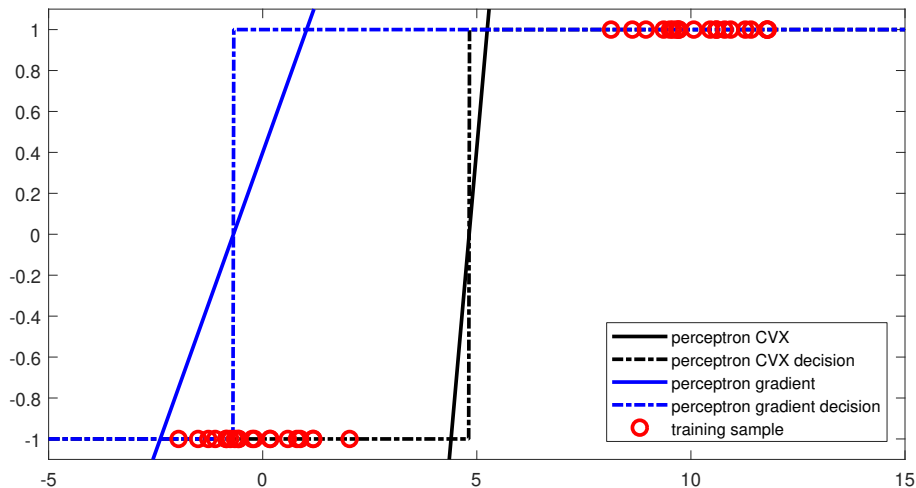
- $j$  is a sample randomly picked from  $\mathcal{M}_k$ .
- Stochastic gradient descent.



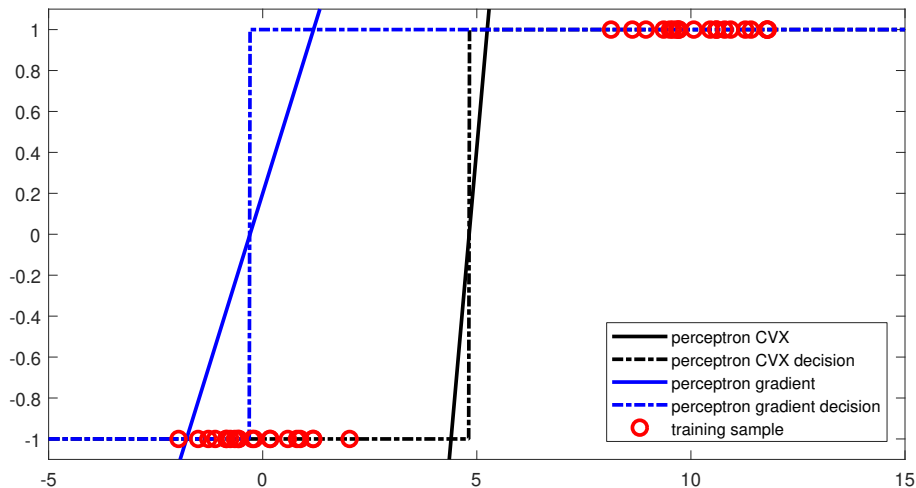
# Online Mode



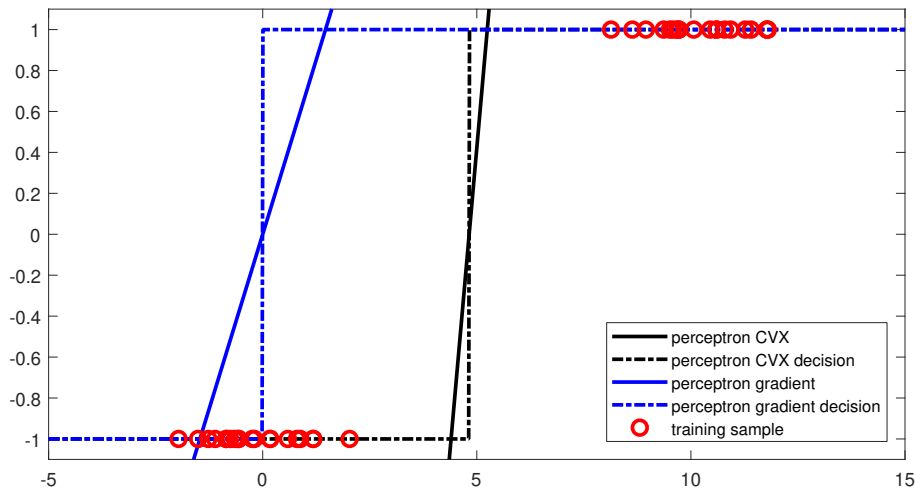
# Online Mode



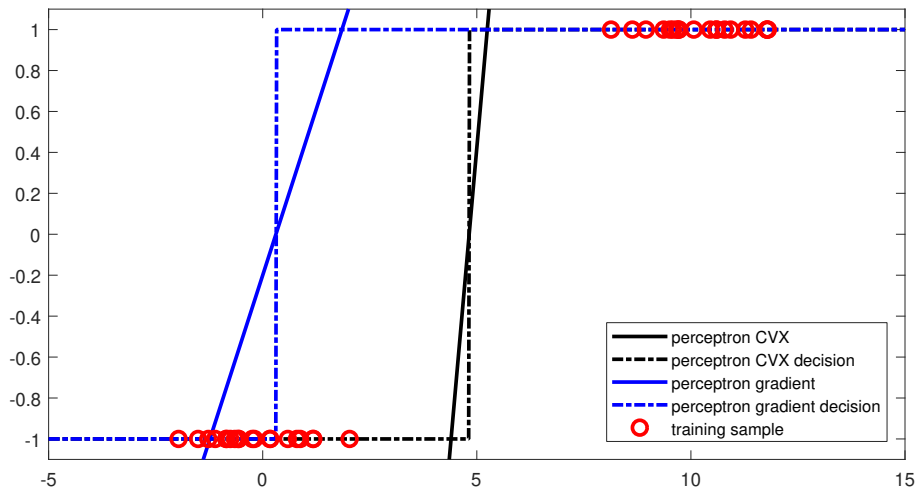
# Online Mode



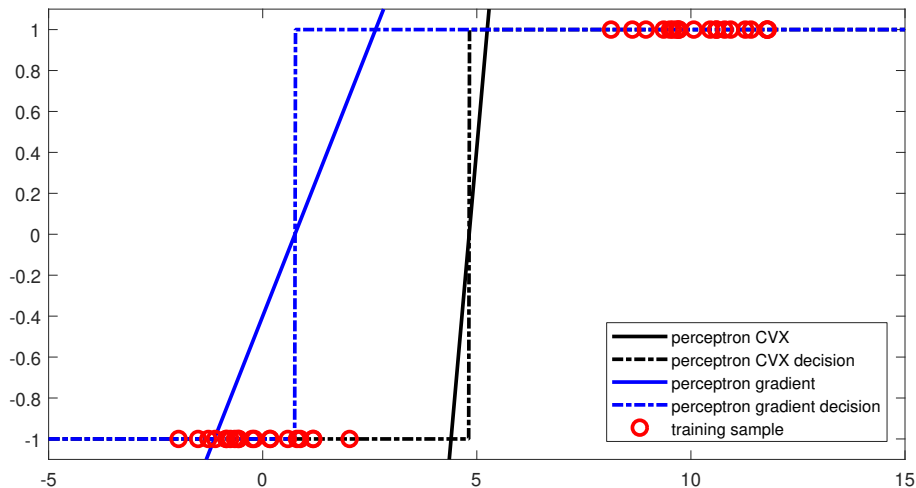
# Online Mode



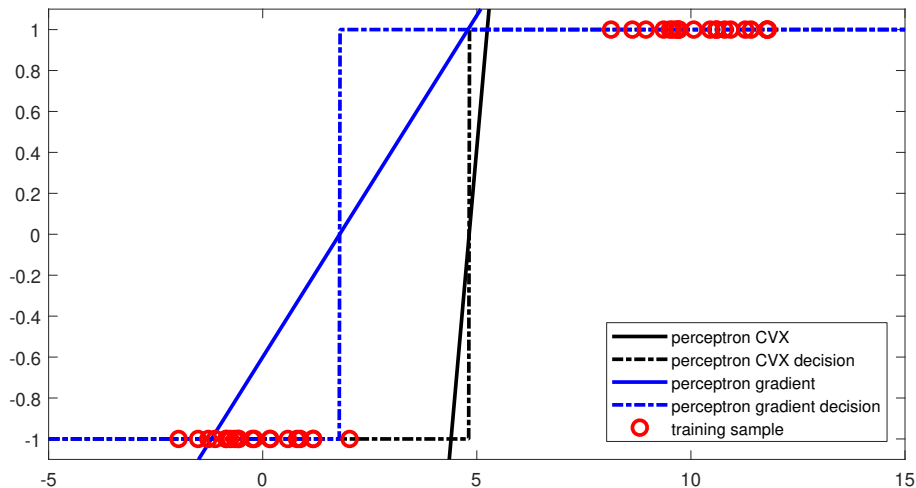
# Online Mode



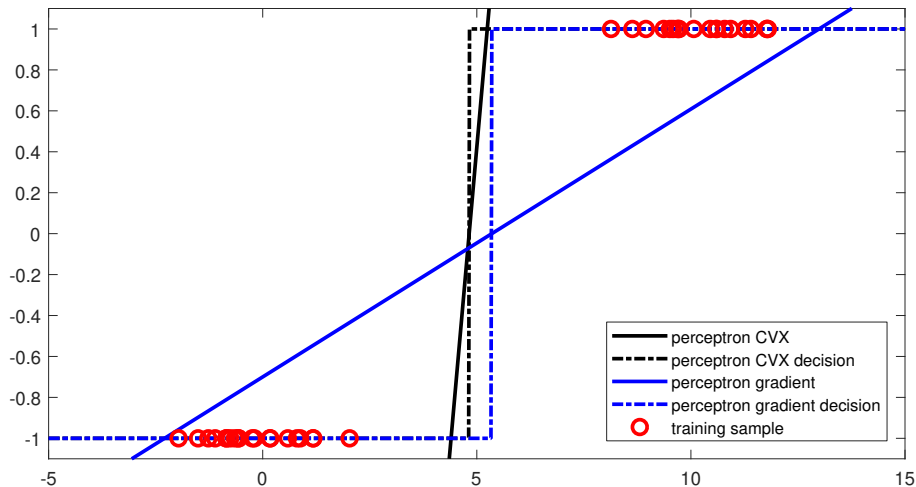
# Online Mode



# Online Mode

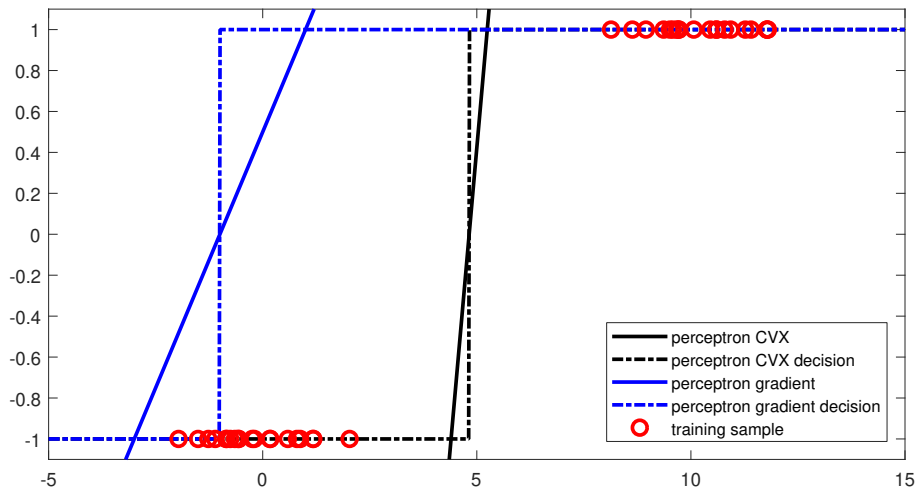


# Online Mode

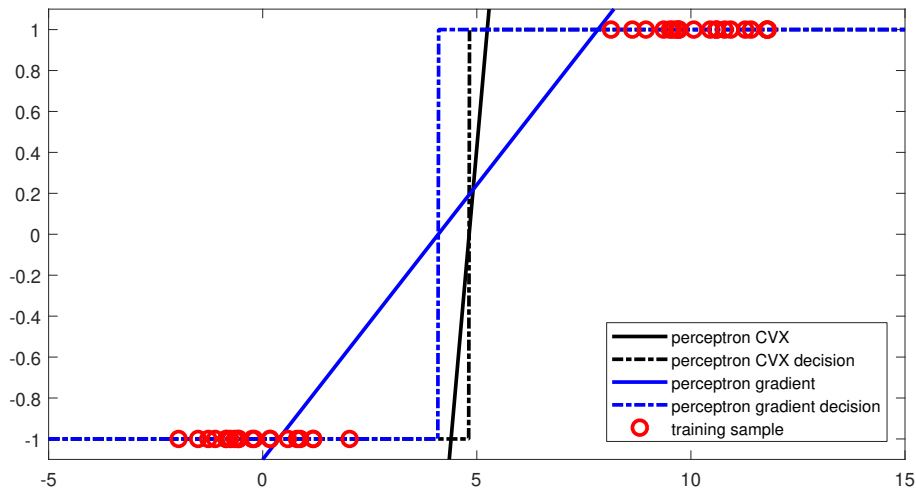




## Batch Mode

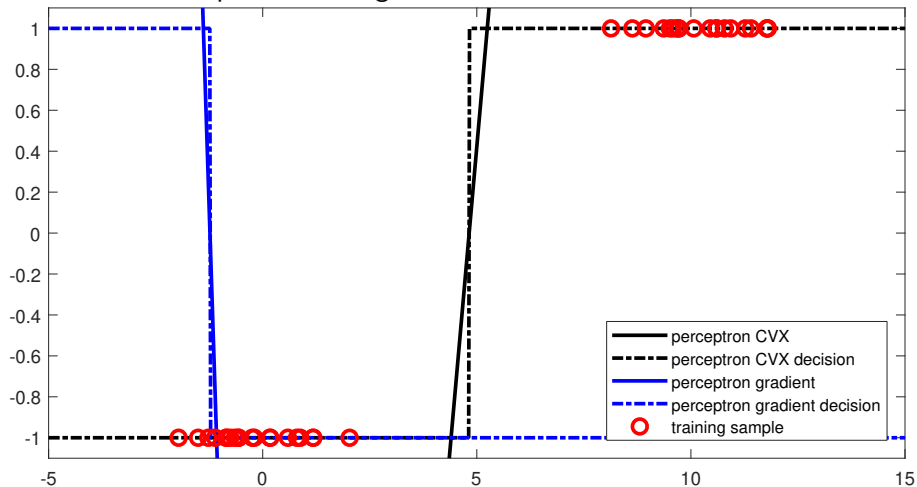


# Batch Mode



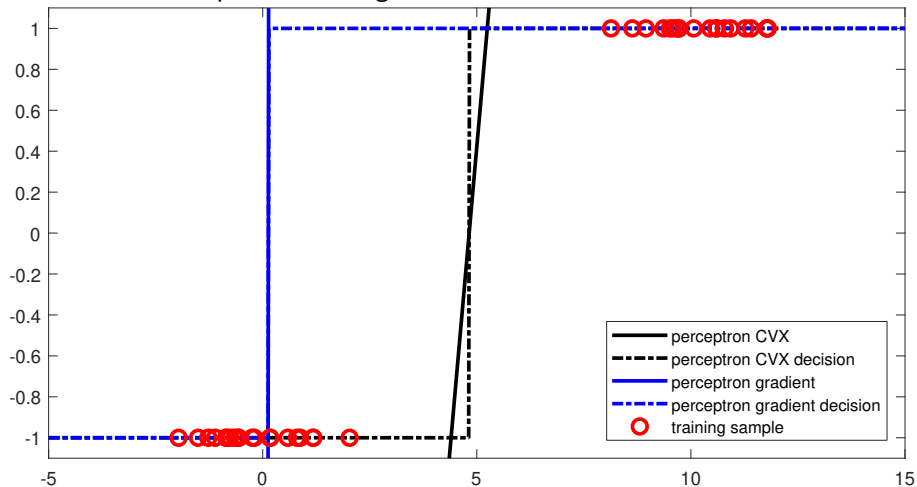
# Step Size

Batch mode: Step size too large.



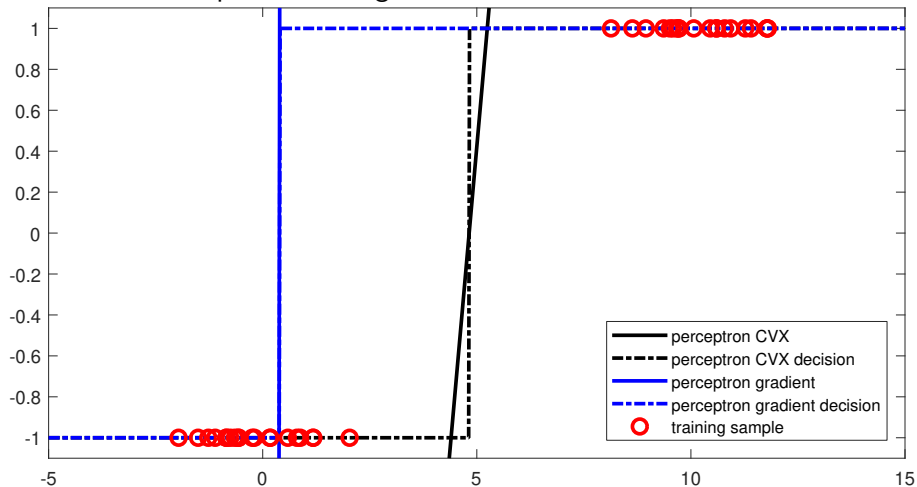
# Step Size

Batch mode: Step size too large.



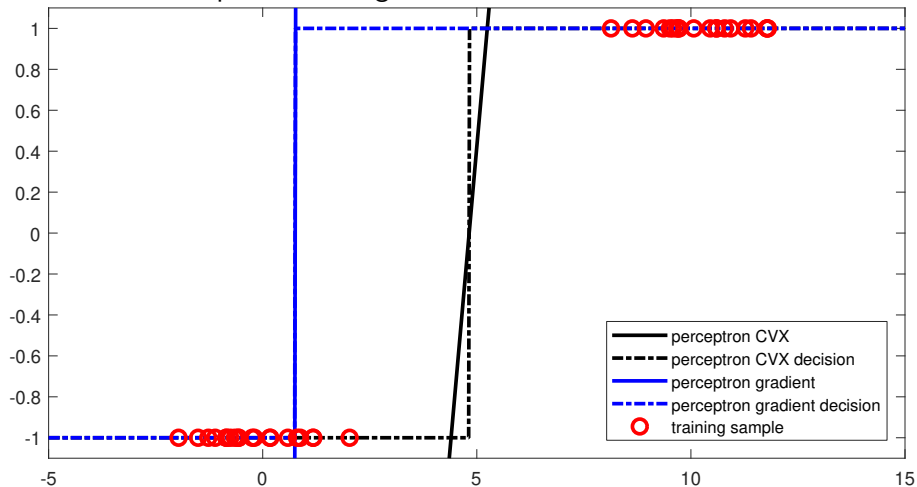
# Step Size

Batch mode: Step size too large.



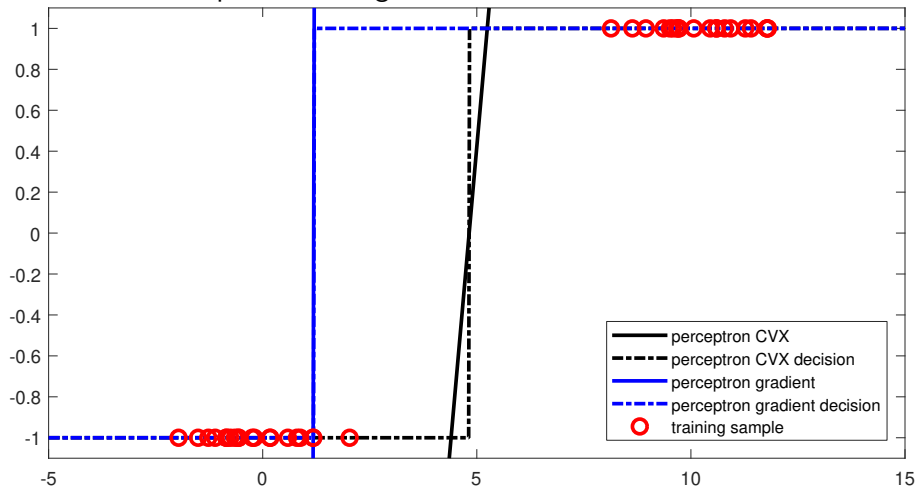
# Step Size

Batch mode: Step size too large.



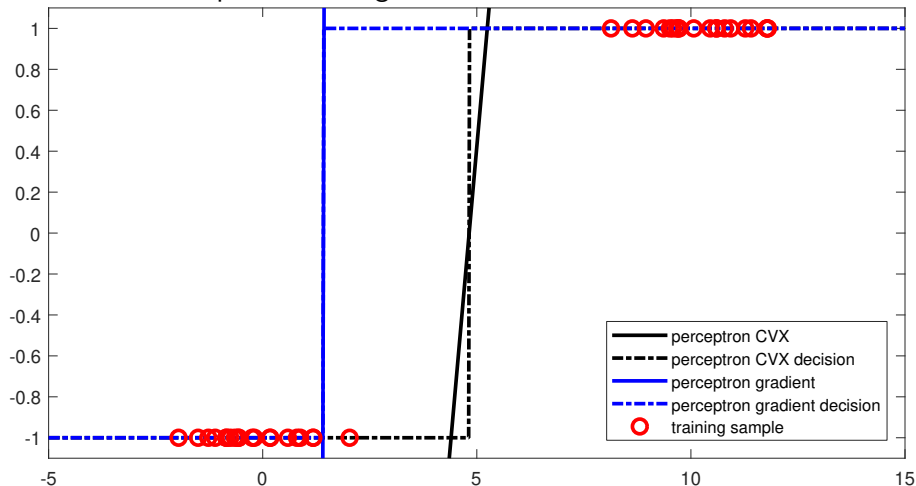
# Step Size

Batch mode: Step size too large.



# Step Size

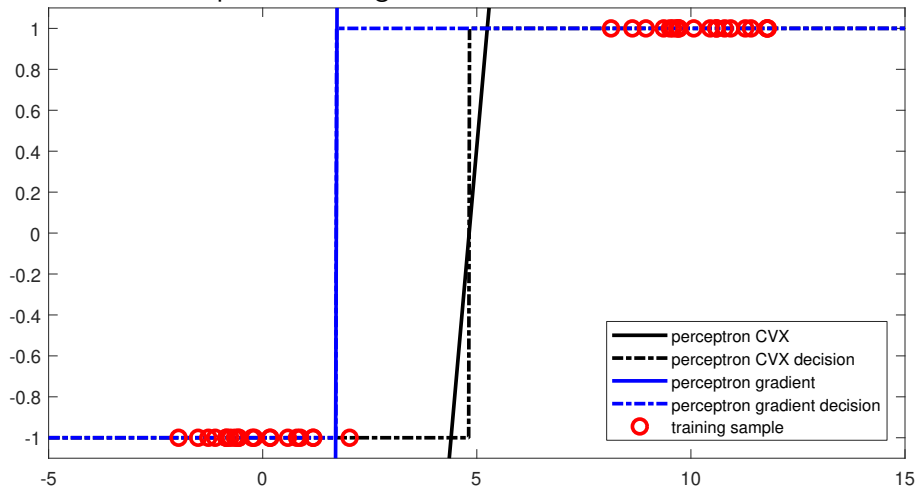
Batch mode: Step size too large.





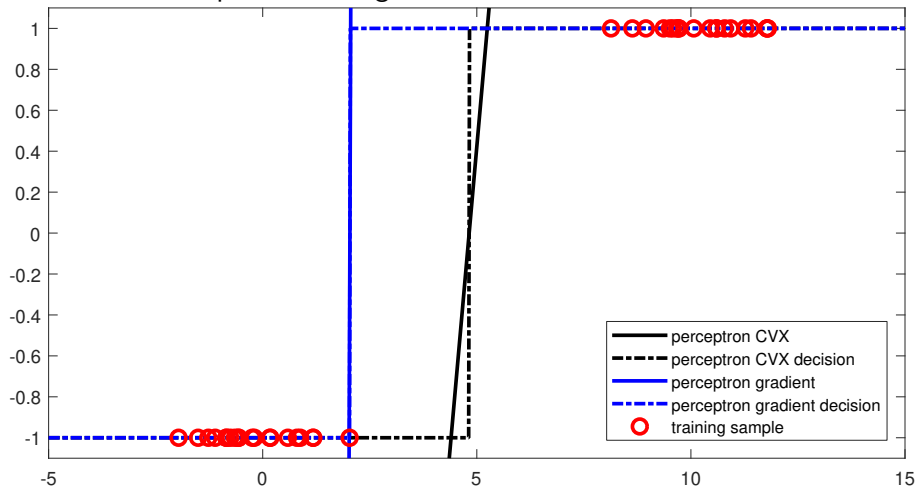
# Step Size

Batch mode: Step size too large.

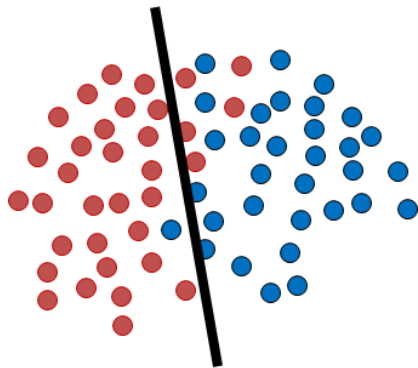


# Step Size

Batch mode: Step size too large.



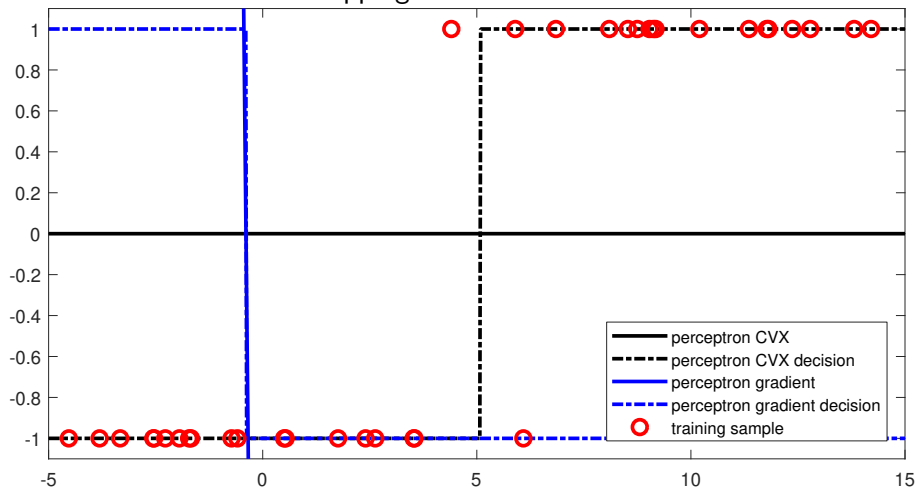
## Linearly Not Separable



- No separating hyperplane
- CVX will still find you a solution
- But loss is no longer zero
- Perceptron algorithm will not converge

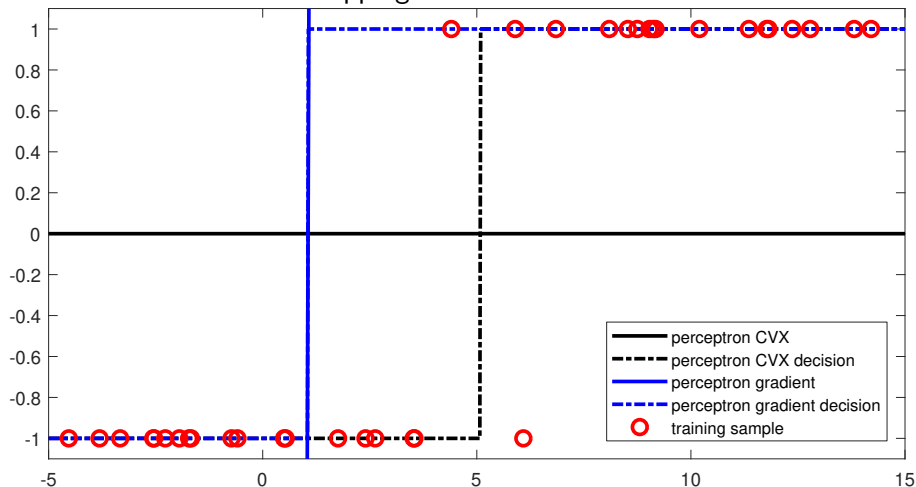
# Linearly Not Separable

If the two classes are overlapping



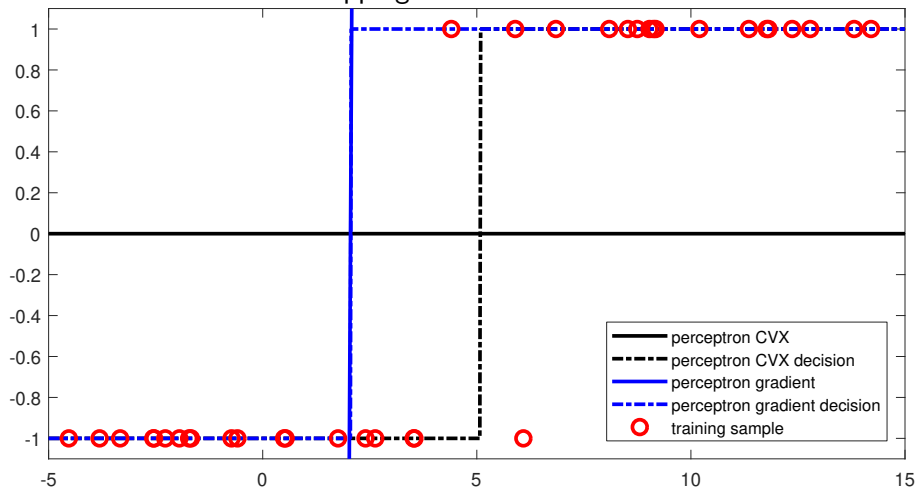
# Linearly Not Separable

If the two classes are overlapping



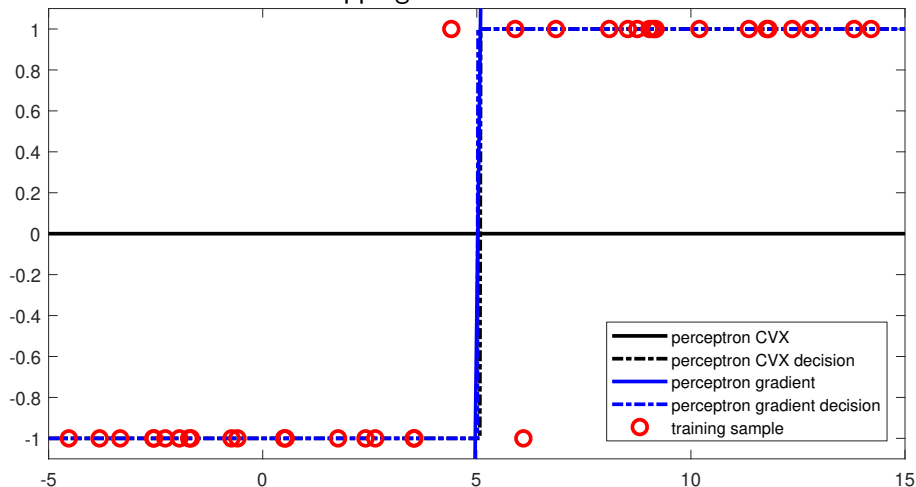
# Linearly Not Separable

If the two classes are overlapping



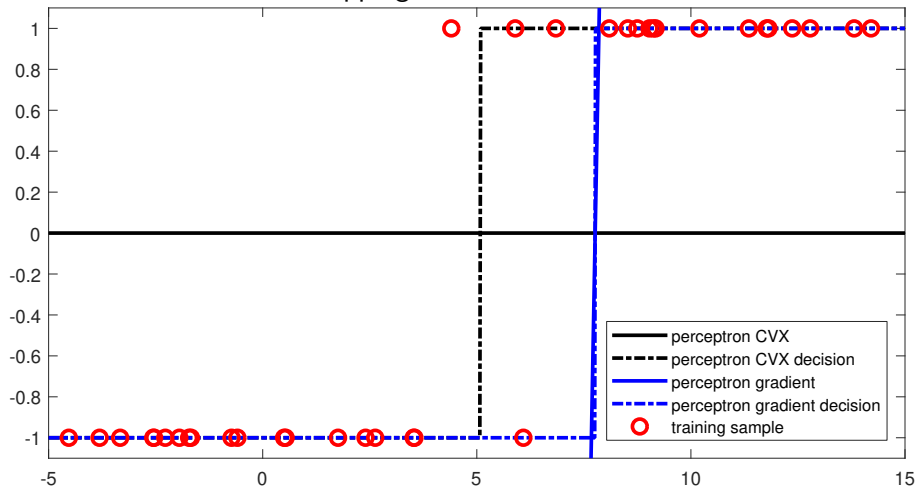
# Linearly Not Separable

If the two classes are overlapping



# Linearly Not Separable

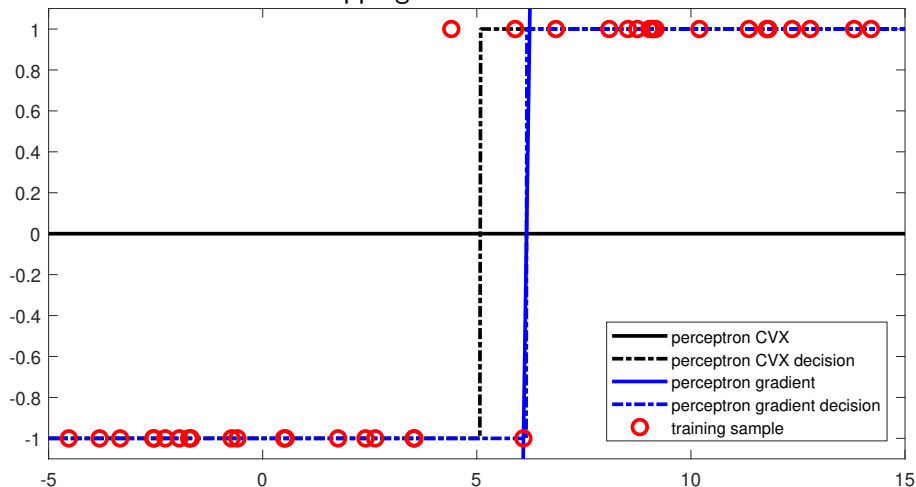
If the two classes are overlapping





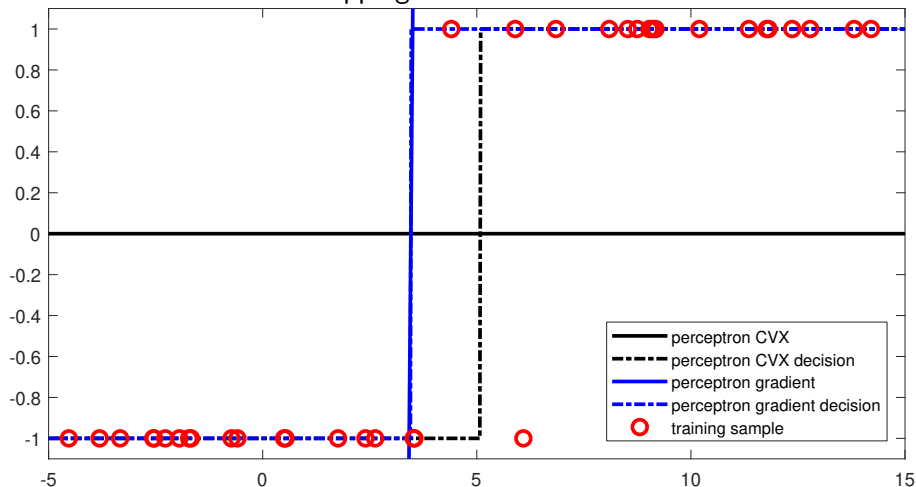
# Linearly Not Separable

If the two classes are overlapping



# Linearly Not Separable

If the two classes are overlapping



# Outline

## Discriminative Approaches

- Lecture 16 Perceptron 1: Definition and Basic Concepts
- **Lecture 17 Perceptron 2: Algorithm and Property**
- Lecture 18 Multi-Layer Perceptron: Back Propagation

## This lecture: Perceptron 2

- Perceptron Algorithm
  - Loss Function
  - Algorithm
- Optimality
  - Uniqueness
  - Batch and Online Mode
- Convergence
  - Main Results
  - Implication

# Convergence of Perceptron Algorithm

**Theorem.** Assume the following things:

- The two classes are linearly separable
- This means:  $(\boldsymbol{\theta}^*)^T (y_j \mathbf{x}_j) = y_j ((\mathbf{w}^*)^T \mathbf{x}_j + w_0^*) \geq \gamma$  for some  $\gamma > 0$
- $\|\mathbf{x}_j\|_2 \leq R$  for some constant
- Initialize  $\boldsymbol{\theta}^{(0)} = \mathbf{0}$

Then, batch mode perceptron algorithm converges to the true solution  $\boldsymbol{\theta}^*$

$$\|\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^*\|^2 = 0,$$

when the number of iterations  $k$  exceeds

$$k \geq \frac{\|\boldsymbol{\theta}^*\|^2 R^2}{\gamma^2}.$$

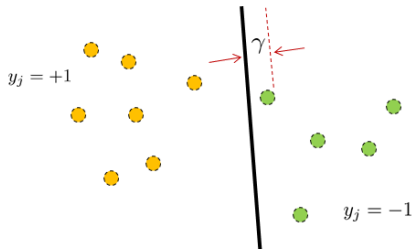
# Interpreting the Perceptron Convergence

**Theorem.** Assume the following things:

- The two classes are linearly separable
- This means:  $(\theta^*)^T (y_j \mathbf{x}_j) = y_j ((\mathbf{w}^*)^T \mathbf{x}_j + w_0^*) \geq \gamma$  for some  $\gamma > 0$
- $\|\mathbf{x}_j\|_2 \leq R$  for some constant
- Initialize  $\theta^{(0)} = \mathbf{0}$

**Comment.**

- $\gamma$  is the margin
- $\theta^*$  is ONE solution such that the margin is at least  $\gamma$



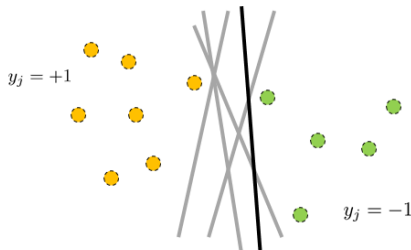
# Interpreting the Perceptron Convergence

**Theorem.** Assume the following things:

- The two classes are linearly separable
- This means:  $(\theta^*)^T (y_j \mathbf{x}_j) = y_j ((\mathbf{w}^*)^T \mathbf{x}_j + w_0^*) \geq \gamma$  for some  $\gamma > 0$
- $\|\mathbf{x}_j\|_2 \leq R$  for some constant
- Initialize  $\theta^{(0)} = \mathbf{0}$

**Comment.**

- If you do not initialize at  $\mathbf{0}$ , still converge.
- The solution  $\theta^*$  might be different.



## Interpreting the Perceptron Convergence

Then, **batch mode** perceptron algorithm converges to the true solution  $\theta^*$

$$\|\theta^{(k+1)} - \theta^*\|^2 = 0$$

when the number of iterations  $k$  exceeds

$$k \geq \frac{\|\theta^*\|^2 R^2}{\gamma^2}.$$

### Comment:

- You can turn batch mode to online mode by picking only one  $j \in \mathcal{M}_k$
- You will do slower, but you can still converge
- $\theta^*$  is the converging point of *this* particular sequence  $\{\theta^1, \theta^2, \dots, \theta^\infty\}$
- Not an arbitrary separating hyperplane

# Interpreting the Perceptron Convergence

Then, batch mode perceptron algorithm converges to the true solution  $\theta^*$

$$\|\theta^{(k+1)} - \theta^*\|^2 = 0,$$

when the number of iterations  $k$  exceeds

$$k \geq \frac{\|\theta^*\|^2 R^2}{\gamma^2}.$$

## Comment:

- $R$  controls the radius of the class.
- Large  $R$ : Wide spread. Difficult. Need large  $k$ .
- $\gamma$  controls the margin.
- Large  $\gamma$ : Big margin. Easy. Need small  $k$ .



## Summary of the Convergence Theorem

- **Algorithm:** You use gradient descent on  $J_{\text{soft}}(\theta)$
- **Solution:** You get a global minimizer for  $J_{\text{hard}}(\theta)$
- But this is just one of the many global minimizers
- **Assumption:** Linearly separable
- If not linearly separable, then will oscillate
- **Margin:** At optimal solution there is a margin because separable
- Applications: Not quite; There are many better methods
- Theoretical usage: Good for analyzing linear models. Very simple algorithm.

## Perceptron Algorithm

- Abu-Mostafa, Learning from Data, Chapter 1.2
- Duda, Hart, Stork, Pattern Classification, Chapter 5.5
- Cornell CS 4780 Lecture <https://www.cs.cornell.edu/courses/cs4780/2018fa/lectures/lecturenote03.html>
- UCSD ECE 271B Lecture <http://www.svcl.ucsd.edu/courses/ece271B-F09/handouts/perceptron.pdf>

# Appendix

## Proof Part 1

- Define

$$\bar{\mathbf{x}}^{(k)} = \sum_{j \in \mathcal{M}_k} y_j \mathbf{x}_j.$$

- Let  $\boldsymbol{\theta}^*$  be the minimizer. Then,

$$\begin{aligned} \|\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^*\|^2 &= \|\boldsymbol{\theta}^{(k)} + \alpha_k \bar{\mathbf{x}}^{(k)} - \boldsymbol{\theta}^*\|^2 \\ &= \|(\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*) + \alpha_k \bar{\mathbf{x}}^{(k)}\|^2 \\ &= \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|^2 + 2\alpha_k (\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)} + \alpha_k^2 \|\bar{\mathbf{x}}^{(k)}\|^2 \\ &= \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|^2 + 2\alpha_k (\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*)^T \left( \sum_{j \in \mathcal{M}_k} y_j \mathbf{x}_j \right) \\ &\quad + \alpha_k^2 \left\| \sum_{j \in \mathcal{M}_k} y_j \mathbf{x}_j \right\|^2. \end{aligned}$$

## Proof Part 2

- By construction,  $\boldsymbol{\theta}^{(k)}$  updates only the misclassified samples (during the  $k$ -th iteration)
- So for any  $j \in \mathcal{M}_k$  we must have  $(\boldsymbol{\theta}^{(k)})^T (y_j \mathbf{x}_j) \leq 0$ .
- This implies that

$$(\boldsymbol{\theta}^{(k)})^T \bar{\mathbf{x}}^{(k)} = \sum_{j \in \mathcal{M}_k} (\boldsymbol{\theta}^{(k)})^T y_j \mathbf{x}_j \leq 0.$$

## Proof Part 2

- By construction,  $\boldsymbol{\theta}^{(k)}$  updates only the misclassified samples (during the  $k$ -th iteration)
- So for any  $j \in \mathcal{M}_k$  we must have  $(\boldsymbol{\theta}^{(k)})^T (y_j \mathbf{x}_j) \leq 0$ .
- This implies that

$$(\boldsymbol{\theta}^{(k)})^T \bar{\mathbf{x}}^{(k)} = \sum_{j \in \mathcal{M}_k} (\boldsymbol{\theta}^{(k)})^T y_j \mathbf{x}_j \leq 0.$$

## Proof Part 2

- By construction,  $\boldsymbol{\theta}^{(k)}$  updates only the misclassified samples (during the  $k$ -th iteration)
- So for any  $j \in \mathcal{M}_k$  we must have  $(\boldsymbol{\theta}^{(k)})^T (y_j \mathbf{x}_j) \leq 0$ .
- This implies that

$$(\boldsymbol{\theta}^{(k)})^T \bar{\mathbf{x}}^{(k)} = \sum_{j \in \mathcal{M}_k} (\boldsymbol{\theta}^{(k)})^T y_j \mathbf{x}_j \leq 0.$$

## Proof Part 2

- By construction,  $\boldsymbol{\theta}^{(k)}$  updates only the misclassified samples (during the  $k$ -th iteration)
- So for any  $j \in \mathcal{M}_k$  we must have  $(\boldsymbol{\theta}^{(k)})^T (y_j \mathbf{x}_j) \leq 0$ .
- This implies that

$$(\boldsymbol{\theta}^{(k)})^T \bar{\mathbf{x}}^{(k)} = \sum_{j \in \mathcal{M}_k} (\boldsymbol{\theta}^{(k)})^T y_j \mathbf{x}_j \leq 0.$$



## Proof Part 2

- By construction,  $\boldsymbol{\theta}^{(k)}$  updates only the misclassified samples (during the  $k$ -th iteration)
- So for any  $j \in \mathcal{M}_k$  we must have  $(\boldsymbol{\theta}^{(k)})^T (y_j \mathbf{x}_j) \leq 0$ .
- This implies that

$$(\boldsymbol{\theta}^{(k)})^T \bar{\mathbf{x}}^{(k)} = \sum_{j \in \mathcal{M}_k} (\boldsymbol{\theta}^{(k)})^T y_j \mathbf{x}_j \leq 0.$$

- Therefore, we can show that

$$\begin{aligned} & \|\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^*\|^2 \\ & \leq \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|^2 + 2\alpha_k (\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)} + \alpha_k^2 \|\bar{\mathbf{x}}^{(k)}\|^2 \\ & = \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|^2 + \cancel{2\alpha_k (\boldsymbol{\theta}^{(k)})^T \bar{\mathbf{x}}^{(k)}} - 2\alpha_k (\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)} + \alpha_k^2 \|\bar{\mathbf{x}}^{(k)}\|^2 \\ & \leq \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|^2 - 2\alpha_k (\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)} + \alpha_k^2 \|\bar{\mathbf{x}}^{(k)}\|^2. \end{aligned}$$

## Proof Part 3

- So we have

$$\|\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^*\|^2 \leq \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|^2 \underbrace{-2\alpha_k(\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)} + \alpha_k^2 \|\bar{\mathbf{x}}^{(k)}\|^2}_{}$$

- The sum of the last two terms is

$$-2\alpha_k(\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)} + \alpha_k^2 \|\bar{\mathbf{x}}^{(k)}\|^2 = \alpha_k \left( -2(\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)} + \alpha_k \|\bar{\mathbf{x}}^{(k)}\|^2 \right),$$

- Negative if and only if  $\alpha_k < \frac{2(\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)}}{\|\bar{\mathbf{x}}^{(k)}\|^2}$
- Thus, we choose

$$\alpha_k = \frac{(\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)}}{\|\bar{\mathbf{x}}^{(k)}\|^2},$$

## Proof Part 4

- Then, we can have

$$-2\alpha_k(\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)} + \alpha_k^2 \|\bar{\mathbf{x}}^{(k)}\|^2 = -2\alpha_k(\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)} + \alpha_k^2 \|\bar{\mathbf{x}}^{(k)}\|^2$$

- By assumption  $\|\mathbf{x}_j\|^2 \leq R$  for any  $j$ , and  $y_j(\boldsymbol{\theta}^*)^T \mathbf{x}_j \geq \gamma$  for any  $j$
- So

$$\begin{aligned} \frac{\left((\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)}\right)^2}{\|\bar{\mathbf{x}}^{(k)}\|^2} &= \frac{\left(\sum_{j \in \mathcal{M}_k} y_j(\boldsymbol{\theta}^*)^T \mathbf{x}_j\right)^2}{\sum_{j \in \mathcal{M}_k} \|\mathbf{x}_j\|^2} \\ &\geq \frac{\left(\sum_{j \in \mathcal{M}_k} \gamma\right)^2}{\sum_{j \in \mathcal{M}_k} R^2} \end{aligned}$$

## Proof Part 4

- Then, we can have

$$-2\alpha_k(\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)} + \alpha_k^2 \|\bar{\mathbf{x}}^{(k)}\|^2 = -2\alpha_k(\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)} + \alpha_k^2 \|\bar{\mathbf{x}}^{(k)}\|^2$$

- By assumption  $\|\mathbf{x}_j\|^2 \leq R$  for any  $j$ , and  $y_j(\boldsymbol{\theta}^*)^T \mathbf{x}_j \geq \gamma$  for any  $j$
- So

$$\begin{aligned} \frac{\left((\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)}\right)^2}{\|\bar{\mathbf{x}}^{(k)}\|^2} &= \frac{\left(\sum_{j \in \mathcal{M}_k} y_j(\boldsymbol{\theta}^*)^T \mathbf{x}_j\right)^2}{\sum_{j \in \mathcal{M}_k} \|\mathbf{x}_j\|^2} \\ &\geq \frac{\left(\sum_{j \in \mathcal{M}_k} \gamma\right)^2}{\sum_{j \in \mathcal{M}_k} R^2} \end{aligned}$$

## Proof Part 4

- Then, we can have

$$\begin{aligned} -2\alpha_k(\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)} + \alpha_k^2 \|\bar{\mathbf{x}}^{(k)}\|^2 &= -2\alpha_k(\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)} + \alpha_k^2 \|\bar{\mathbf{x}}^{(k)}\|^2 \\ &= -\frac{\left((\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)}\right)^2}{\|\bar{\mathbf{x}}^{(k)}\|^2}. \end{aligned}$$

- By assumption  $\|\mathbf{x}_j\| \leq R$  for any  $j$ , and  $y_j(\boldsymbol{\theta}^*)^T \mathbf{x}_j \geq \gamma$  for any  $j$
- So

$$\begin{aligned} \frac{\left((\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)}\right)^2}{\|\bar{\mathbf{x}}^{(k)}\|^2} &= \frac{\left(\sum_{j \in \mathcal{M}_k} y_j(\boldsymbol{\theta}^*)^T \mathbf{x}_j\right)^2}{\sum_{j \in \mathcal{M}_k} \|\mathbf{x}_j\|^2} \\ &\geq \frac{\left(\sum_{j \in \mathcal{M}_k} \gamma\right)^2}{\sum_{j \in \mathcal{M}_k} R^2} \end{aligned}$$

## Proof Part 4

- Then, we can have

$$\begin{aligned} -2\alpha_k(\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)} + \alpha_k^2 \|\bar{\mathbf{x}}^{(k)}\|^2 &= -2\alpha_k(\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)} + \alpha_k^2 \|\bar{\mathbf{x}}^{(k)}\|^2 \\ &= -\frac{\left((\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)}\right)^2}{\|\bar{\mathbf{x}}^{(k)}\|^2}. \end{aligned}$$

- By assumption  $\|\mathbf{x}_j\| \leq R$  for any  $j$ , and  $y_j(\boldsymbol{\theta}^*)^T \mathbf{x}_j \geq \gamma$  for any  $j$
- So

$$\begin{aligned} \frac{\left((\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)}\right)^2}{\|\bar{\mathbf{x}}^{(k)}\|^2} &= \frac{\left(\sum_{j \in \mathcal{M}_k} y_j(\boldsymbol{\theta}^*)^T \mathbf{x}_j\right)^2}{\sum_{j \in \mathcal{M}_k} \|\mathbf{x}_j\|^2} \\ &\geq \frac{\left(\sum_{j \in \mathcal{M}_k} \gamma\right)^2}{\sum_{j \in \mathcal{M}_k} R^2} \end{aligned}$$

## Proof Part 4

- Then, we can have

$$\begin{aligned} -2\alpha_k(\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)} + \alpha_k^2 \|\bar{\mathbf{x}}^{(k)}\|^2 &= -2\alpha_k(\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)} + \alpha_k^2 \|\bar{\mathbf{x}}^{(k)}\|^2 \\ &= -\frac{\left((\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)}\right)^2}{\|\bar{\mathbf{x}}^{(k)}\|^2}. \end{aligned}$$

- By assumption  $\|\mathbf{x}_j\| \leq R$  for any  $j$ , and  $y_j(\boldsymbol{\theta}^*)^T \mathbf{x}_j \geq \gamma$  for any  $j$
- So

$$\begin{aligned} \frac{\left((\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)}\right)^2}{\|\bar{\mathbf{x}}^{(k)}\|^2} &= \frac{\left(\sum_{j \in \mathcal{M}_k} y_j(\boldsymbol{\theta}^*)^T \mathbf{x}_j\right)^2}{\sum_{j \in \mathcal{M}_k} \|\mathbf{x}_j\|^2} \\ &\geq \frac{\left(\sum_{j \in \mathcal{M}_k} \gamma\right)^2}{\sum_{j \in \mathcal{M}_k} R^2} \end{aligned}$$

## Proof Part 4

- Then, we can have

$$\begin{aligned} -2\alpha_k(\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)} + \alpha_k^2 \|\bar{\mathbf{x}}^{(k)}\|^2 &= -2\alpha_k(\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)} + \alpha_k^2 \|\bar{\mathbf{x}}^{(k)}\|^2 \\ &= -\frac{\left((\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)}\right)^2}{\|\bar{\mathbf{x}}^{(k)}\|^2}. \end{aligned}$$

- By assumption  $\|\mathbf{x}_j\| \leq R$  for any  $j$ , and  $y_j(\boldsymbol{\theta}^*)^T \mathbf{x}_j \geq \gamma$  for any  $j$
- So

$$\begin{aligned} \frac{\left((\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)}\right)^2}{\|\bar{\mathbf{x}}^{(k)}\|^2} &= \frac{\left(\sum_{j \in \mathcal{M}_k} y_j(\boldsymbol{\theta}^*)^T \mathbf{x}_j\right)^2}{\sum_{j \in \mathcal{M}_k} \|\mathbf{x}_j\|^2} \\ &\geq \frac{\left(\sum_{j \in \mathcal{M}_k} \gamma\right)^2}{\sum_{j \in \mathcal{M}_k} R^2} \end{aligned}$$



## Proof Part 4

- Then, we can have

$$\begin{aligned} -2\alpha_k(\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)} + \alpha_k^2 \|\bar{\mathbf{x}}^{(k)}\|^2 &= -2\alpha_k(\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)} + \alpha_k^2 \|\bar{\mathbf{x}}^{(k)}\|^2 \\ &= -\frac{\left((\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)}\right)^2}{\|\bar{\mathbf{x}}^{(k)}\|^2}. \end{aligned}$$

- By assumption  $\|\mathbf{x}_j\| \leq R$  for any  $j$ , and  $y_j(\boldsymbol{\theta}^*)^T \mathbf{x}_j \geq \gamma$  for any  $j$
- So

$$\begin{aligned} \frac{\left((\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)}\right)^2}{\|\bar{\mathbf{x}}^{(k)}\|^2} &= \frac{\left(\sum_{j \in \mathcal{M}_k} y_j(\boldsymbol{\theta}^*)^T \mathbf{x}_j\right)^2}{\sum_{j \in \mathcal{M}_k} \|\mathbf{x}_j\|^2} \\ &\geq \frac{\left(\sum_{j \in \mathcal{M}_k} \gamma\right)^2}{\sum_{j \in \mathcal{M}_k} R^2} \end{aligned}$$

## Proof Part 4

- Then, we can have

$$\begin{aligned} -2\alpha_k(\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)} + \alpha_k^2 \|\bar{\mathbf{x}}^{(k)}\|^2 &= -2\alpha_k(\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)} + \alpha_k^2 \|\bar{\mathbf{x}}^{(k)}\|^2 \\ &= -\frac{\left((\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)}\right)^2}{\|\bar{\mathbf{x}}^{(k)}\|^2}. \end{aligned}$$

- By assumption  $\|\mathbf{x}_j\|^2 \leq R$  for any  $j$ , and  $y_j(\boldsymbol{\theta}^*)^T \mathbf{x}_j \geq \gamma$  for any  $j$
- So

$$\begin{aligned} \frac{\left((\boldsymbol{\theta}^*)^T \bar{\mathbf{x}}^{(k)}\right)^2}{\|\bar{\mathbf{x}}^{(k)}\|^2} &= \frac{\left(\sum_{j \in \mathcal{M}_k} y_j(\boldsymbol{\theta}^*)^T \mathbf{x}_j\right)^2}{\sum_{j \in \mathcal{M}_k} \|\mathbf{x}_j\|^2} \\ &\geq \frac{\left(\sum_{j \in \mathcal{M}_k} \gamma\right)^2}{\sum_{j \in \mathcal{M}_k} R^2} = |\mathcal{M}_k| \frac{\gamma^2}{R^2} \end{aligned}$$

## Proof Part 5

- Then by induction we can show that

$$\|\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^*\|^2 < \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|^2 - \sum_{i=1}^k |\mathcal{M}_i| \frac{\gamma^2}{R^2}.$$

- We can conclude that

$$\sum_{i=1}^k |\mathcal{M}_i| \frac{\gamma^2}{R^2} < \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|^2 = \|\boldsymbol{\theta}^*\|^2,$$

- Therefore,

$$\underbrace{\sum_{i=1}^k |\mathcal{M}_i|}_{k \leq (\cdot)} < \frac{\|\boldsymbol{\theta}^*\|^2 R^2}{\gamma^2}$$

## Proof Part 5

- Then by induction we can show that

$$\|\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^*\|^2 < \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|^2 - \sum_{i=1}^k |\mathcal{M}_i| \frac{\gamma^2}{R^2}.$$

- We can conclude that

$$\sum_{i=1}^k |\mathcal{M}_i| \frac{\gamma^2}{R^2} < \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|^2 = \|\boldsymbol{\theta}^*\|^2,$$

- Therefore,

$$\underbrace{\sum_{i=1}^k |\mathcal{M}_i|}_{k \leq (\cdot)} < \frac{\|\boldsymbol{\theta}^*\|^2 R^2}{\gamma^2}$$

## Proof Part 5

- Then by induction we can show that

$$\|\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^*\|^2 < \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|^2 - \sum_{i=1}^k |\mathcal{M}_i| \frac{\gamma^2}{R^2}.$$

- We can conclude that

$$\sum_{i=1}^k |\mathcal{M}_i| \frac{\gamma^2}{R^2} < \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|^2 = \|\boldsymbol{\theta}^*\|^2,$$

- Therefore,

$$\underbrace{\sum_{i=1}^k |\mathcal{M}_i|}_{k \leq (\cdot)} < \frac{\|\boldsymbol{\theta}^*\|^2 R^2}{\gamma^2}$$

## Proof Part 5

- Then by induction we can show that

$$\|\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^*\|^2 < \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|^2 - \sum_{i=1}^k |\mathcal{M}_i| \frac{\gamma^2}{R^2}.$$

- We can conclude that

$$\sum_{i=1}^k |\mathcal{M}_i| \frac{\gamma^2}{R^2} < \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|^2 = \|\boldsymbol{\theta}^*\|^2,$$

- Therefore,

$$\underbrace{\sum_{i=1}^k |\mathcal{M}_i|}_{k \leq (\cdot)} < \frac{\|\boldsymbol{\theta}^*\|^2 R^2}{\gamma^2}$$

## Proof Part 5

- Then by induction we can show that

$$\|\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^*\|^2 < \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|^2 - \sum_{i=1}^k |\mathcal{M}_i| \frac{\gamma^2}{R^2}.$$

- We can conclude that

$$\sum_{i=1}^k |\mathcal{M}_i| \frac{\gamma^2}{R^2} < \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|^2 = \|\boldsymbol{\theta}^*\|^2,$$

- Therefore,

$$\underbrace{\sum_{i=1}^k |\mathcal{M}_i|}_{k \leq (\cdot)} < \frac{\|\boldsymbol{\theta}^*\|^2 R^2}{\gamma^2}$$

## Proof Part 5

- Then by induction we can show that

$$\|\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^*\|^2 < \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|^2 - \sum_{i=1}^k |\mathcal{M}_i| \frac{\gamma^2}{R^2}.$$

- We can conclude that

$$\sum_{i=1}^k |\mathcal{M}_i| \frac{\gamma^2}{R^2} < \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|^2 = \|\boldsymbol{\theta}^*\|^2,$$

- Therefore,

$$\underbrace{\sum_{i=1}^k |\mathcal{M}_i|}_{k \leq (\cdot)} < \frac{\|\boldsymbol{\theta}^*\|^2 R^2}{\gamma^2}$$



## Proof Part 5

- Then by induction we can show that

$$\|\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^*\|^2 < \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|^2 - \sum_{i=1}^k |\mathcal{M}_i| \frac{\gamma^2}{R^2}.$$

- We can conclude that

$$\sum_{i=1}^k |\mathcal{M}_i| \frac{\gamma^2}{R^2} < \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|^2 = \|\boldsymbol{\theta}^*\|^2,$$

- Therefore,

$$\underbrace{\sum_{i=1}^k |\mathcal{M}_i|}_{k \leq (\cdot)} < \frac{\|\boldsymbol{\theta}^*\|^2 R^2}{\gamma^2}$$

## Proof Part 5

- Then by induction we can show that

$$\|\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^*\|^2 < \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|^2 - \sum_{i=1}^k |\mathcal{M}_i| \frac{\gamma^2}{R^2}.$$

- We can conclude that

$$\sum_{i=1}^k |\mathcal{M}_i| \frac{\gamma^2}{R^2} < \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|^2 = \|\boldsymbol{\theta}^*\|^2,$$

- Therefore,

$$\underbrace{\sum_{i=1}^k |\mathcal{M}_i|}_{k \leq (\cdot)} < \frac{\|\boldsymbol{\theta}^*\|^2 R^2}{\gamma^2} = \frac{\max_j \|\boldsymbol{\theta}^*\|^2 \|\mathbf{x}_j\|^2}{(\min_j (\boldsymbol{\theta}^*)^T \mathbf{x}_j)^2}$$