

ECE595 / STAT598: Machine Learning I

Lecture 17.2: Perceptron 2 - Optimality

Spring 2020

Stanley Chan

School of Electrical and Computer Engineering
Purdue University



Outline

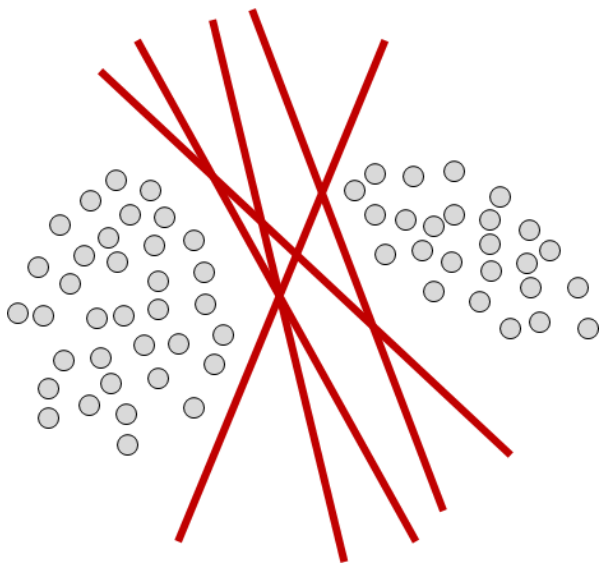
Discriminative Approaches

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

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

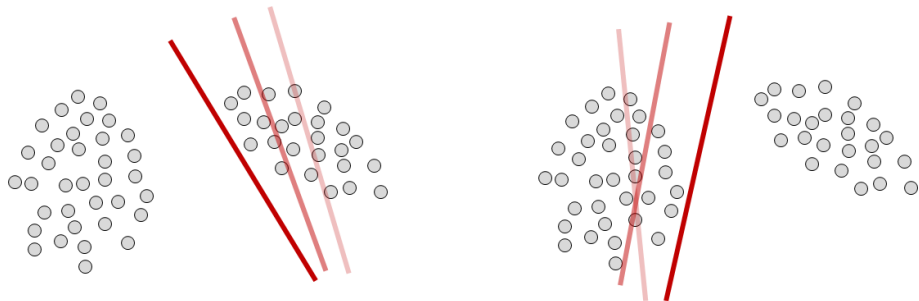
- If $g_{\theta}(\mathbf{x}) > 0$, then $g_{\alpha\theta}(\mathbf{x}) > 0$. So if $h_{\theta}(\mathbf{x}) = +1$, then $h_{\alpha\theta}(\mathbf{x}) = +1$.
- If $g_{\theta}(\mathbf{x}) < 0$, then $g_{\alpha\theta}(\mathbf{x}) < 0$. So if $h_{\theta}(\mathbf{x}) = -1$, then $h_{\alpha\theta}(\mathbf{x}) = -1$.
- 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

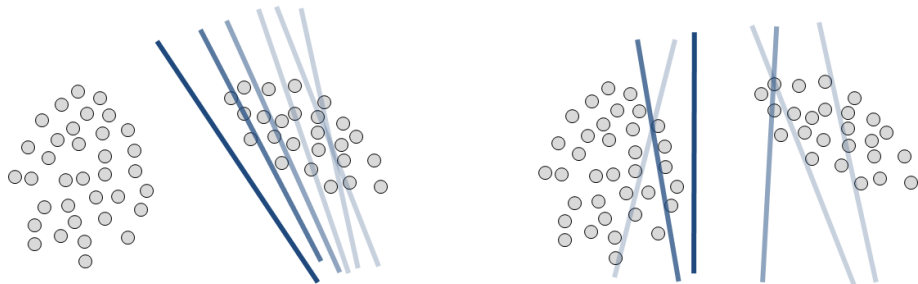
- **Initialization**

- 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

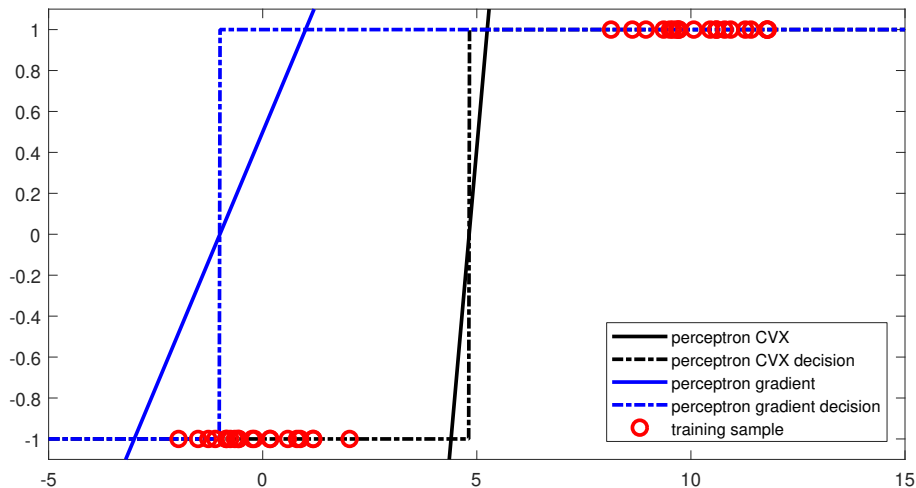
- **Online mode**

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

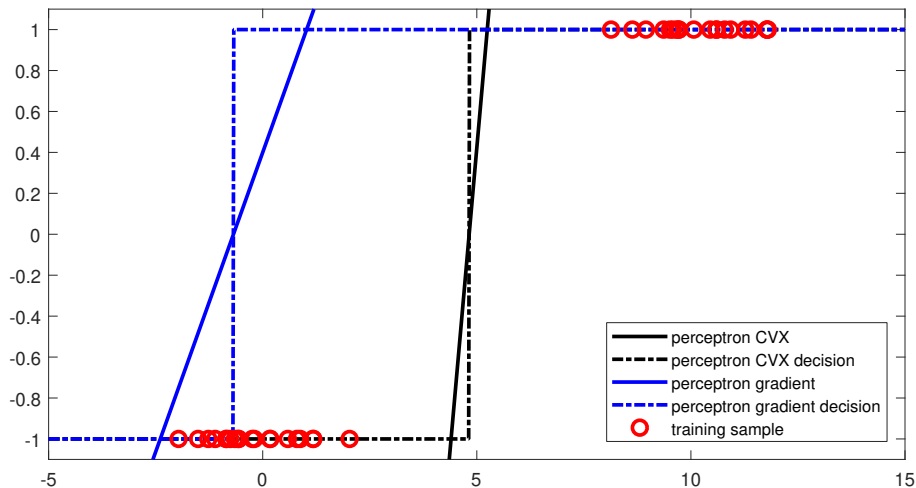
Update via a single misclassified sample

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

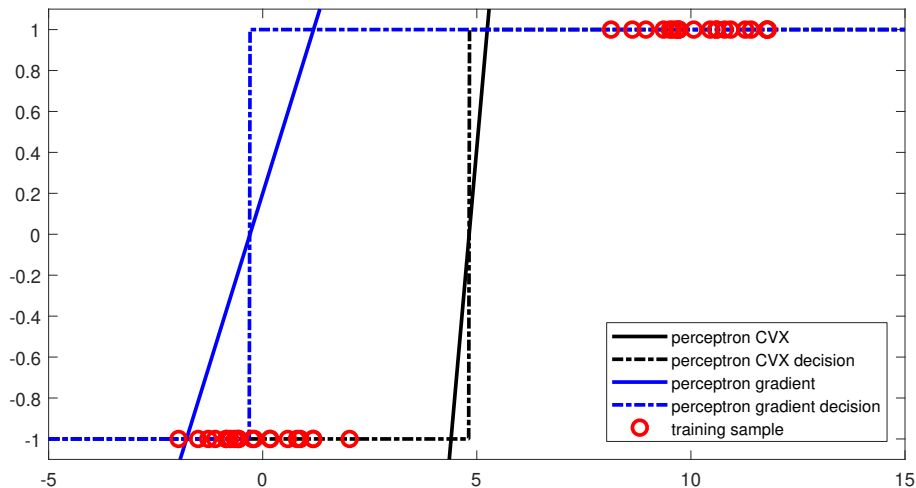
Online Mode



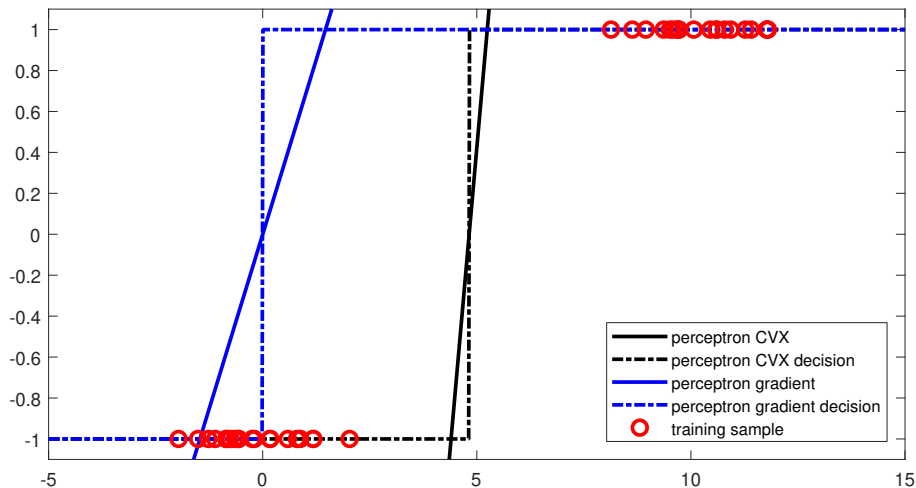
Online Mode



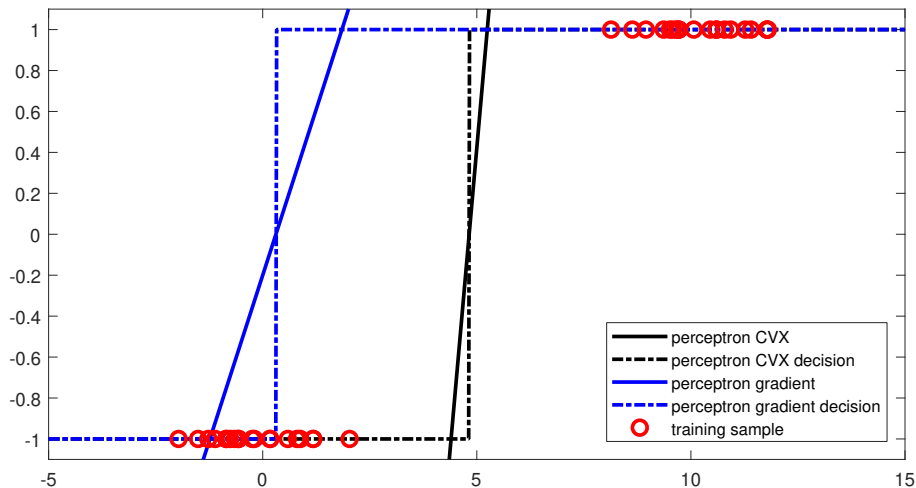
Online Mode



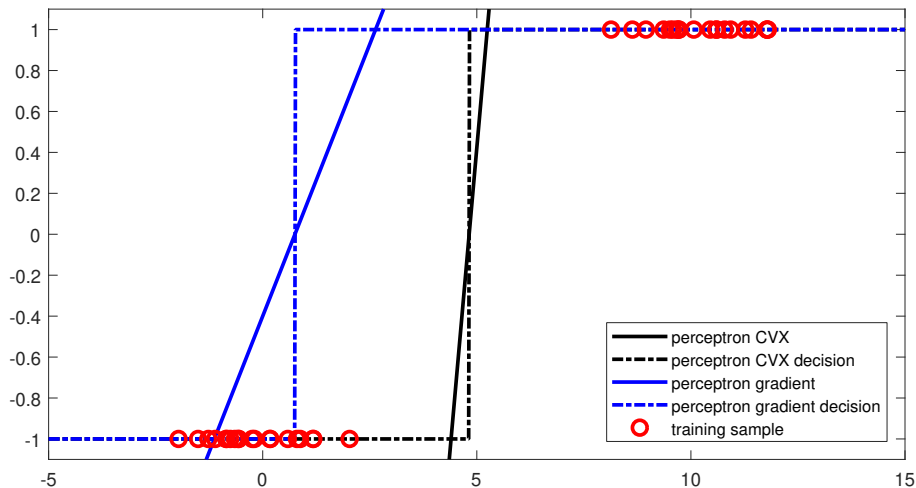
Online Mode



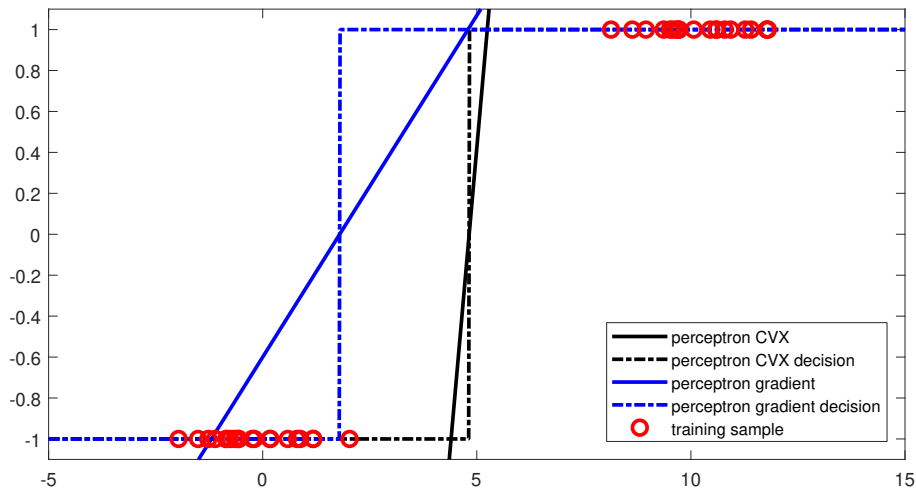
Online Mode



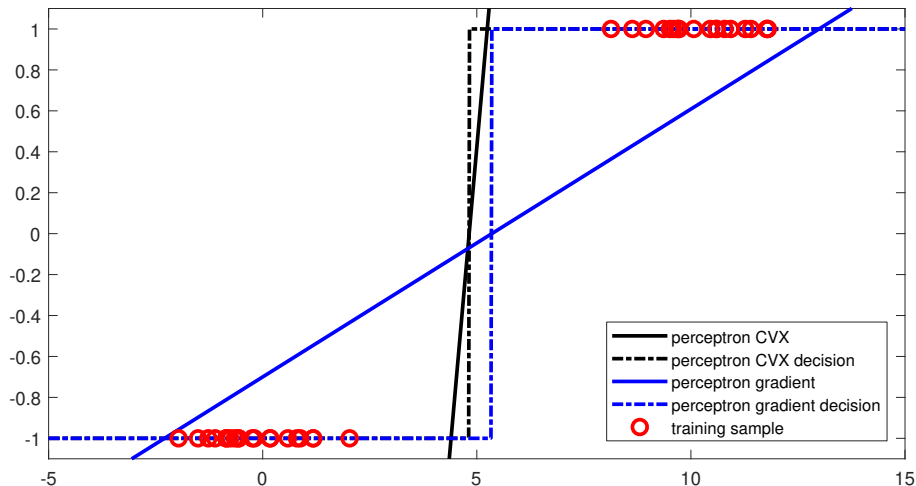
Online Mode



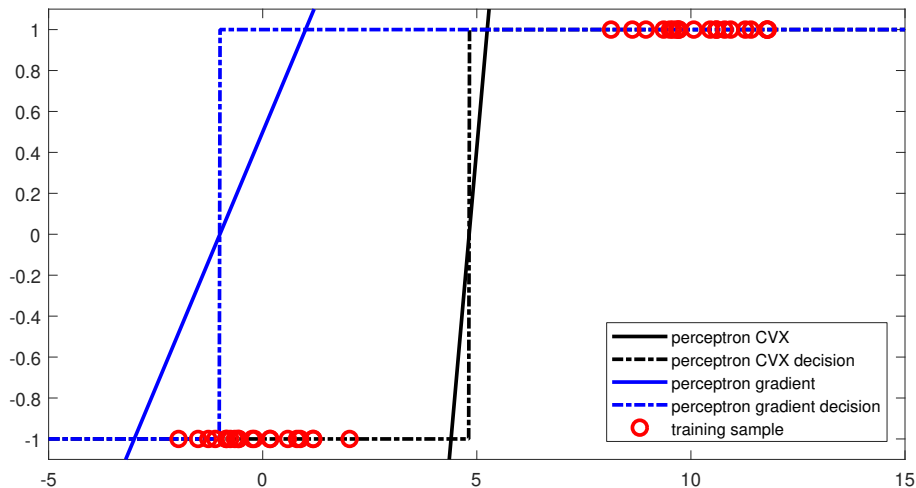
Online Mode



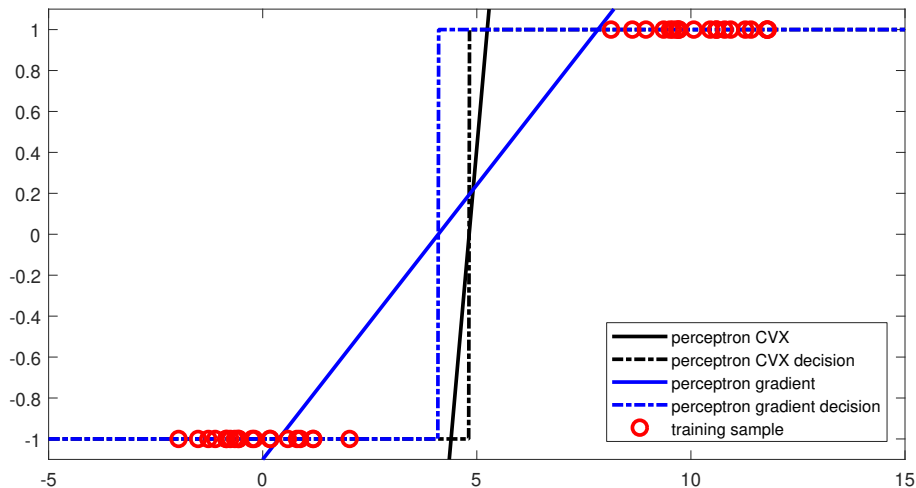
Online Mode



Batch Mode

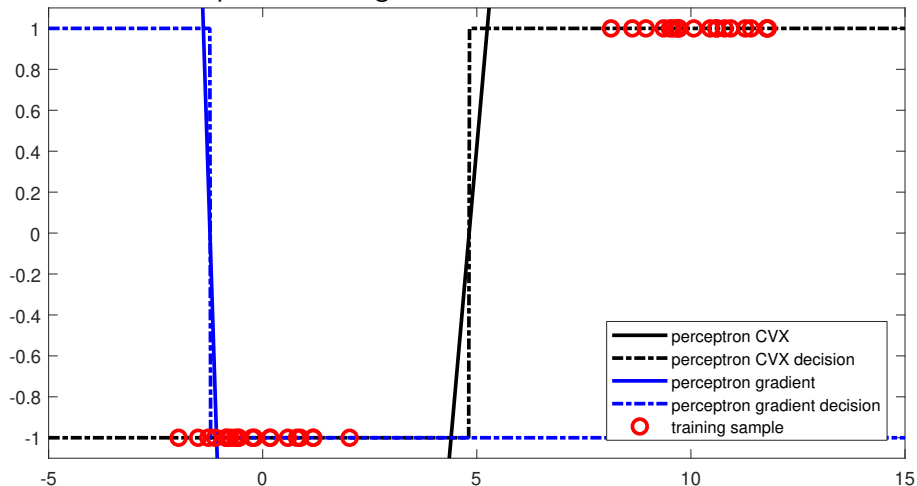


Batch Mode



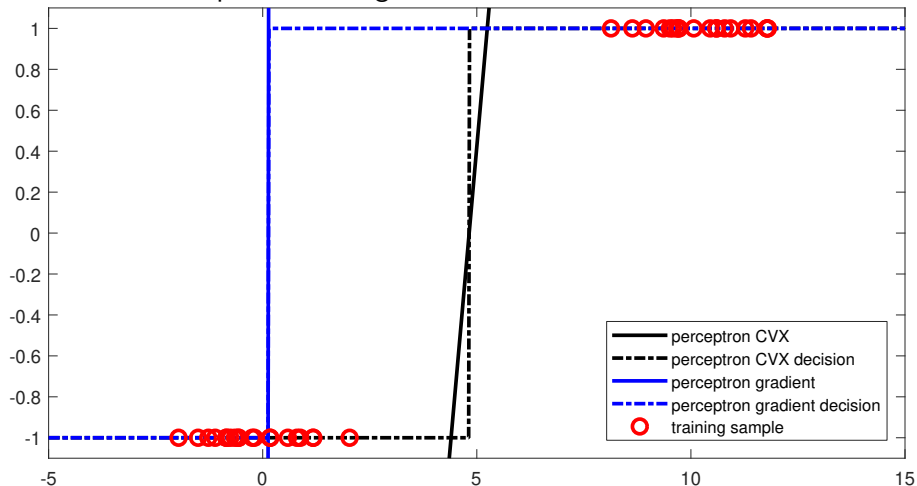
Step Size

Batch mode: Step size too large.



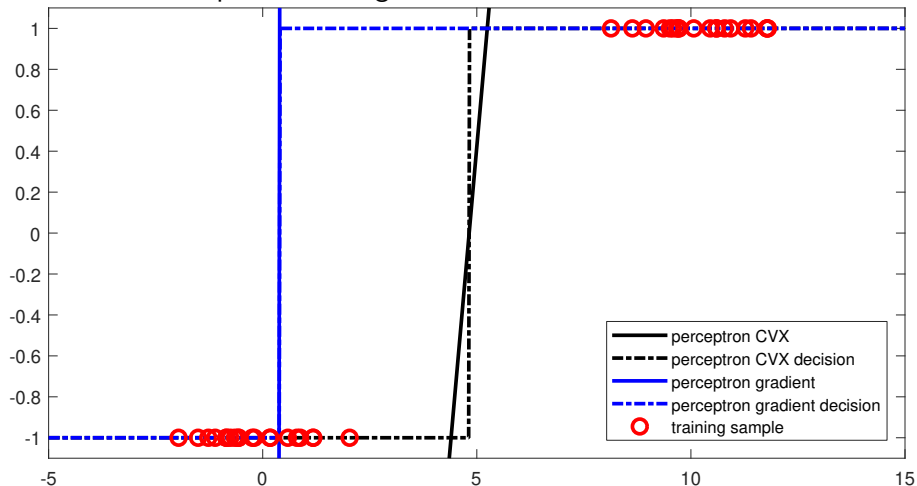
Step Size

Batch mode: Step size too large.



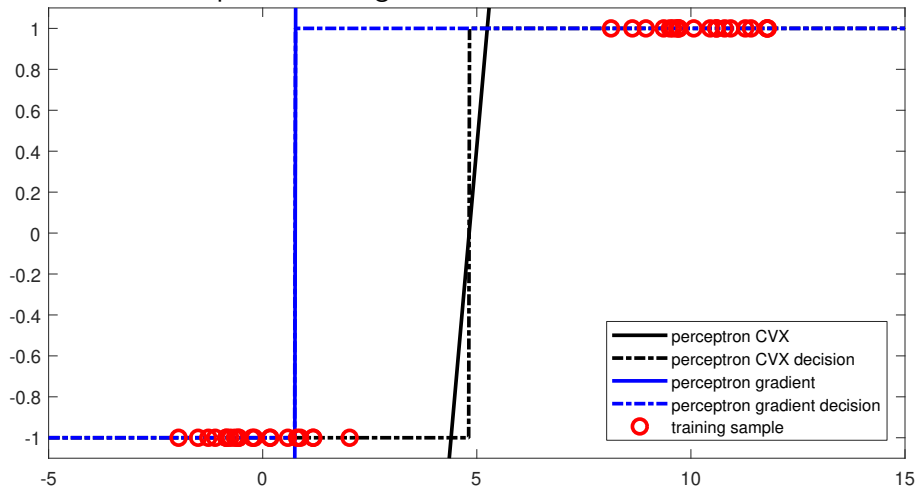
Step Size

Batch mode: Step size too large.



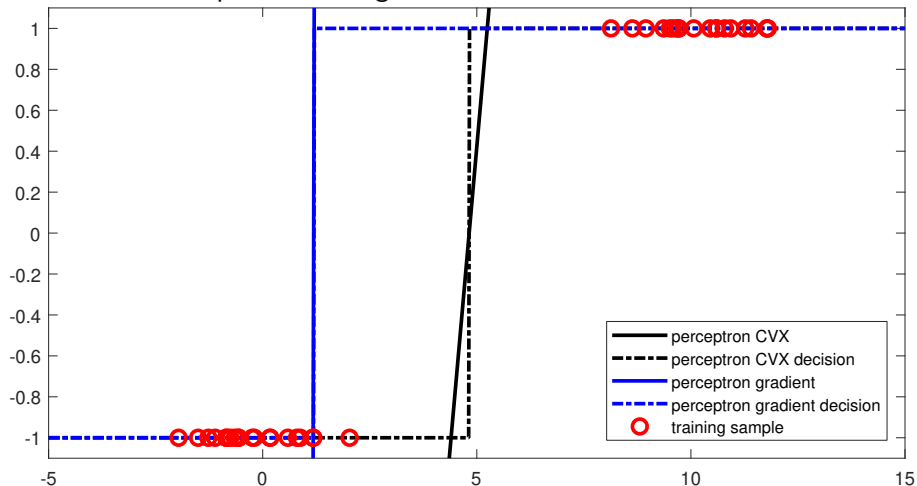
Step Size

Batch mode: Step size too large.



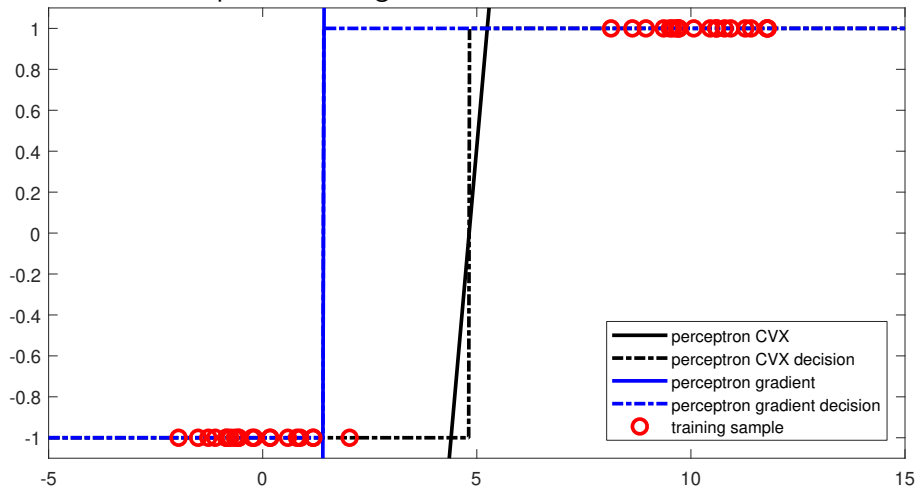
Step Size

Batch mode: Step size too large.



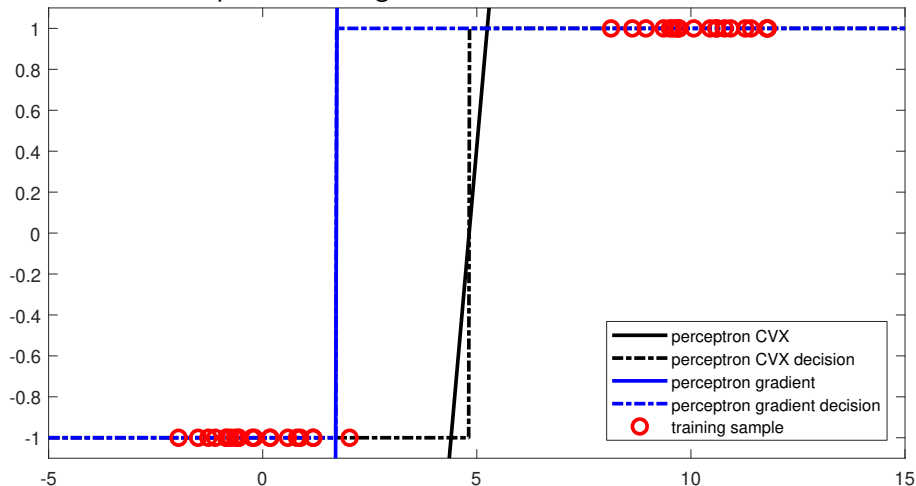
Step Size

Batch mode: Step size too large.



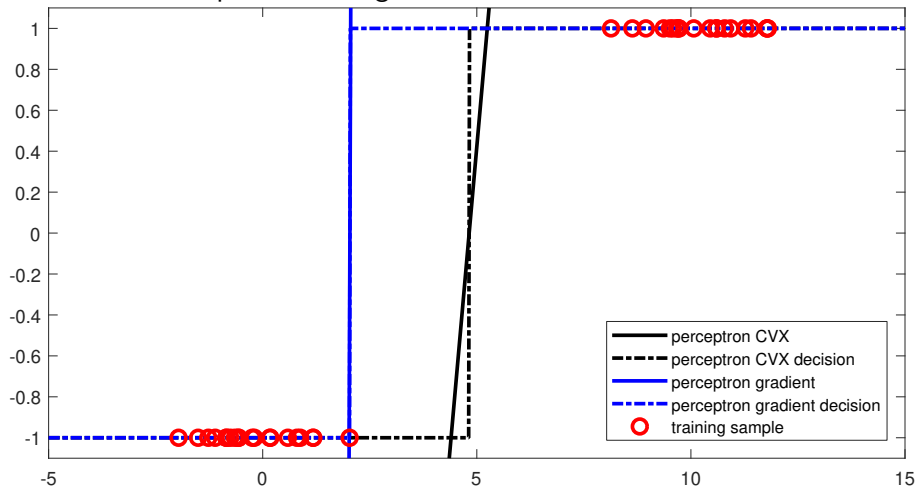
Step Size

Batch mode: Step size too large.

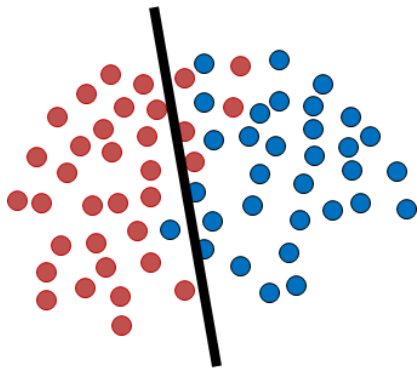


Step Size

Batch mode: Step size too large.



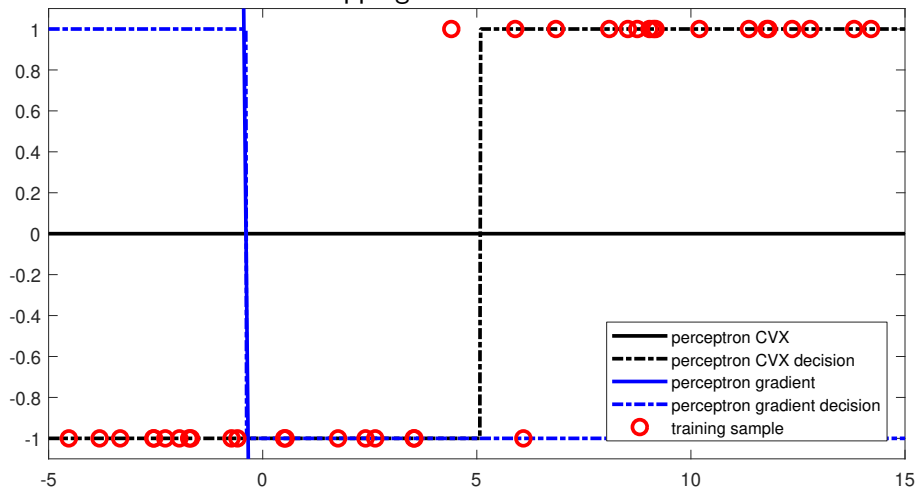
Linearly Not Separable



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

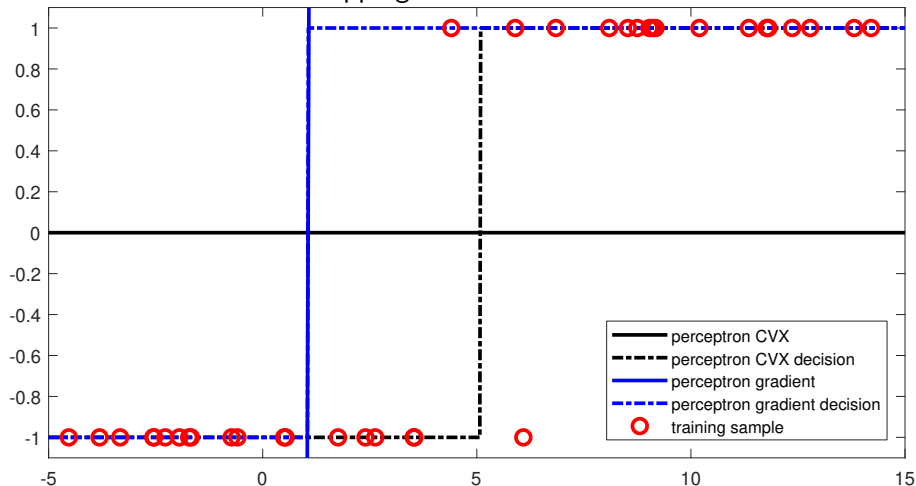
Linearly Not Separable

If the two classes are overlapping



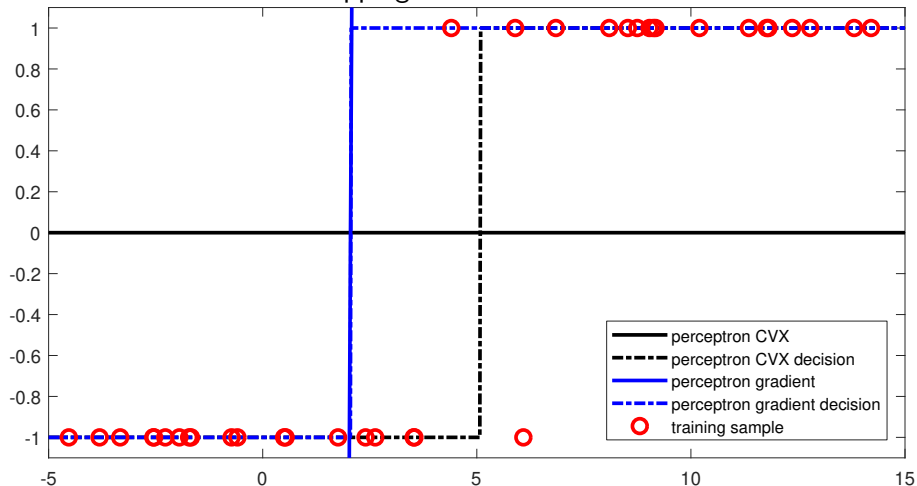
Linearly Not Separable

If the two classes are overlapping



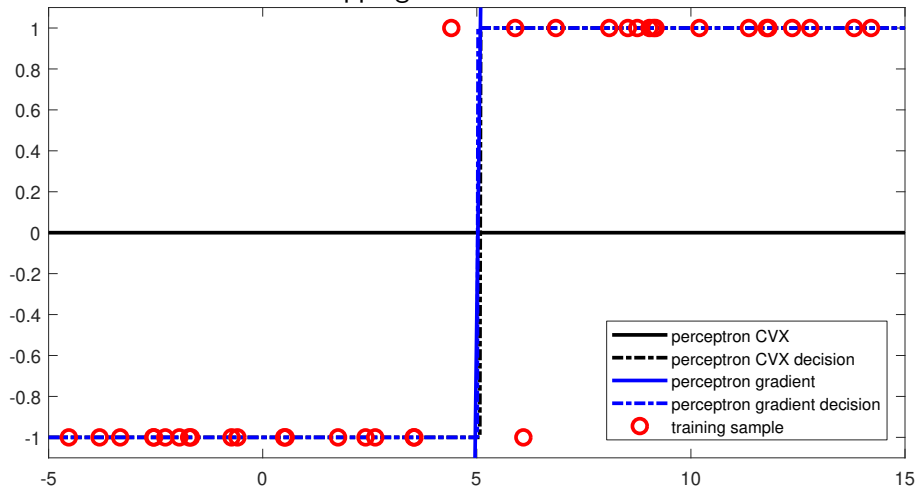
Linearly Not Separable

If the two classes are overlapping



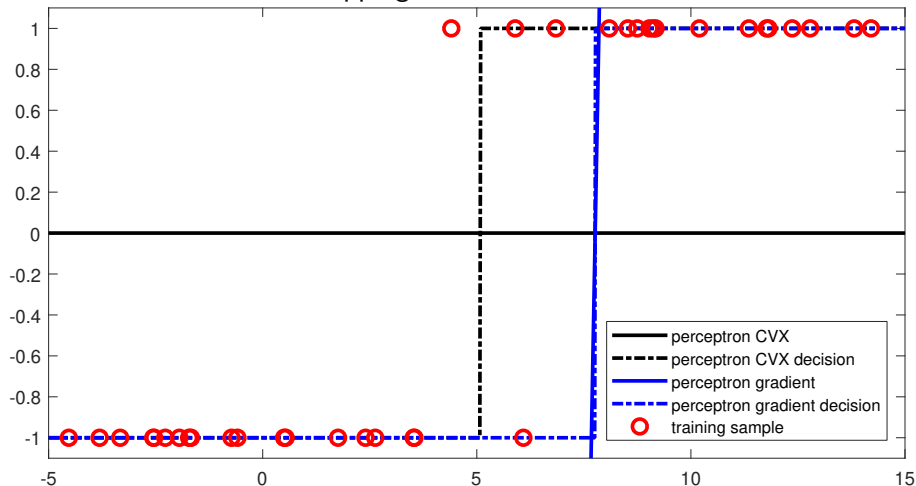
Linearly Not Separable

If the two classes are overlapping



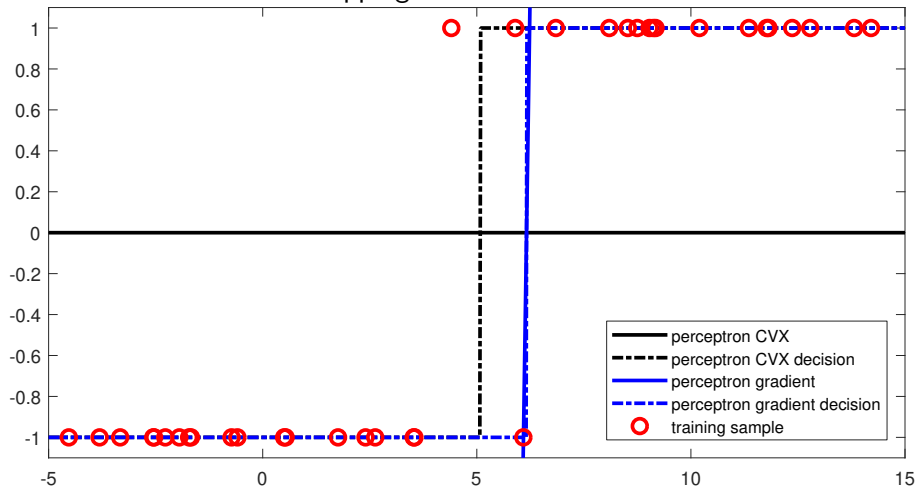
Linearly Not Separable

If the two classes are overlapping



Linearly Not Separable

If the two classes are overlapping



Linearly Not Separable

If the two classes are overlapping

