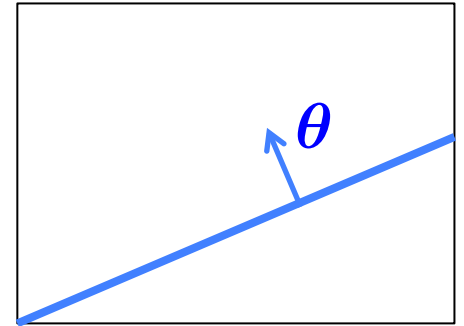




Linear Classification: The Perceptron

Linear Classifiers

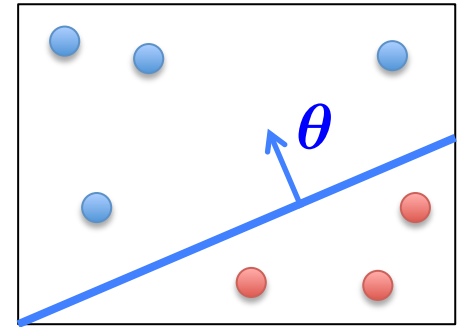
- A **hyperplane** partitions \mathbb{R}^d into two half-spaces
 - Defined by the normal vector $\theta \in \mathbb{R}^d$
 - θ is orthogonal to any vector lying on the hyperplane
 - Assumed to pass through the origin
 - This is because we incorporated bias term θ_0 into it by $x_0 = 1$
- Consider classification with +1, -1 labels ...



Linear Classifiers

- **Linear classifiers:** represent decision boundary by hyperplane

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \quad \mathbf{x}^\top = \begin{bmatrix} 1 & x_1 & \dots & x_d \end{bmatrix}$$



$$h(\mathbf{x}) = \text{sign}(\boldsymbol{\theta}^\top \mathbf{x}) \quad \text{where} \quad \text{sign}(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ -1 & \text{if } z < 0 \end{cases}$$

– Note that: $\boldsymbol{\theta}^\top \mathbf{x} > 0 \implies y = +1$

$\boldsymbol{\theta}^\top \mathbf{x} < 0 \implies y = -1$

The Perceptron

$$h(\mathbf{x}) = \text{sign}(\boldsymbol{\theta}^\top \mathbf{x}) \quad \text{where} \quad \text{sign}(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ -1 & \text{if } z < 0 \end{cases}$$

- The perceptron uses the following update rule each time it receives a new training instance $(\mathbf{x}^{(i)}, y^{(i)})$

$$\theta_j \leftarrow \theta_j - \frac{\alpha}{2} \underbrace{\left(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)}_{\text{either 2 or -2}} x_j^{(i)}$$

- If the prediction matches the label, make no change
- Otherwise, adjust $\boldsymbol{\theta}$

The Perceptron

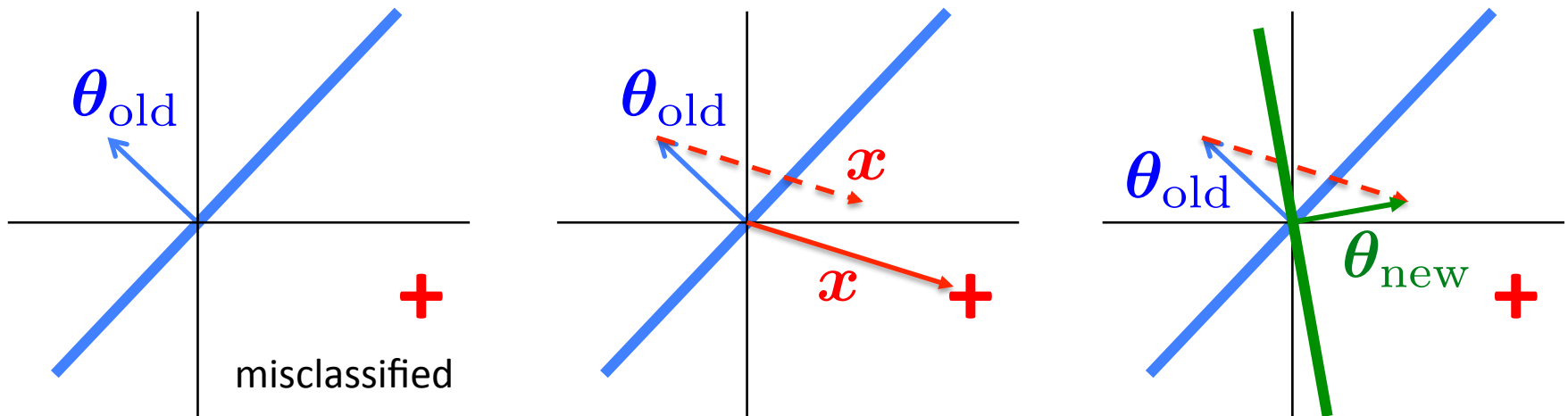
- The perceptron uses the following update rule each time it receives a new training instance $(\mathbf{x}^{(i)}, y^{(i)})$

$$\theta_j \leftarrow \theta_j - \frac{\alpha}{2} \underbrace{\left(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right)}_{\text{either 2 or -2}} x_j^{(i)}$$

- Re-write as $\theta_j \leftarrow \theta_j + \alpha y^{(i)} x_j^{(i)}$ (only upon misclassification)
 - Can eliminate α in this case, since its only effect is to scale θ by a constant, which doesn't affect performance

Perceptron Rule: If $\mathbf{x}^{(i)}$ is misclassified, do $\theta \leftarrow \theta + y^{(i)} \mathbf{x}^{(i)}$

Why the Perceptron Update Works




Why the Perceptron Update Works

- Consider the misclassified example ($y = +1$)
 - Perceptron wrongly thinks that $\theta_{\text{old}}^T \mathbf{x} < 0$
- Update:
$$\theta_{\text{new}} = \theta_{\text{old}} + y\mathbf{x} = \theta_{\text{old}} + \mathbf{x} \quad (\text{since } y = +1)$$

- Note that

$$\begin{aligned}\theta_{\text{new}}^T \mathbf{x} &= (\theta_{\text{old}} + \mathbf{x})^T \mathbf{x} \\ &= \theta_{\text{old}}^T \mathbf{x} + \mathbf{x}^T \mathbf{x}\end{aligned}$$

 $\|\mathbf{x}\|_2^2 > 0$

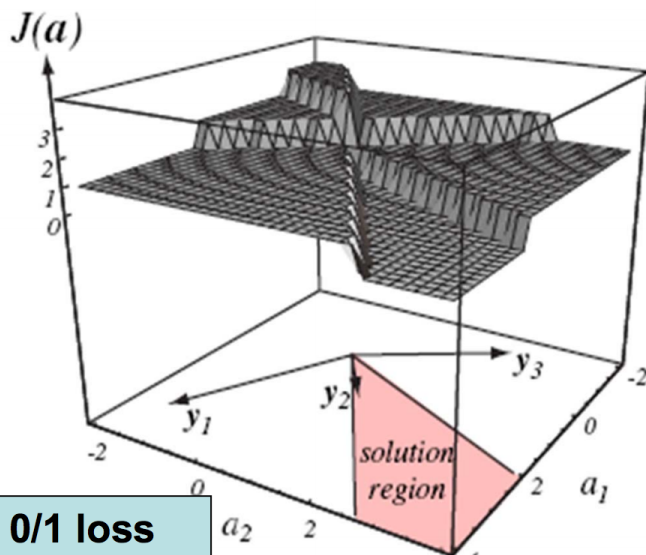
- Therefore, $\theta_{\text{new}}^T \mathbf{x}$ is less negative than $\theta_{\text{old}}^T \mathbf{x}$
 - So, we are making ourselves more correct on this example!

The Perceptron Cost Function

- Prediction is correct if $y^{(i)} \mathbf{x}^{(i)} \boldsymbol{\theta} > 0$
- Could have used 0/1 loss

$$J_{0/1}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \ell(\text{sign}(x^{(i)} \boldsymbol{\theta}), y^{(i)})$$

where $\ell()$ is 0 if the prediction is correct, 1 otherwise



0/1 loss

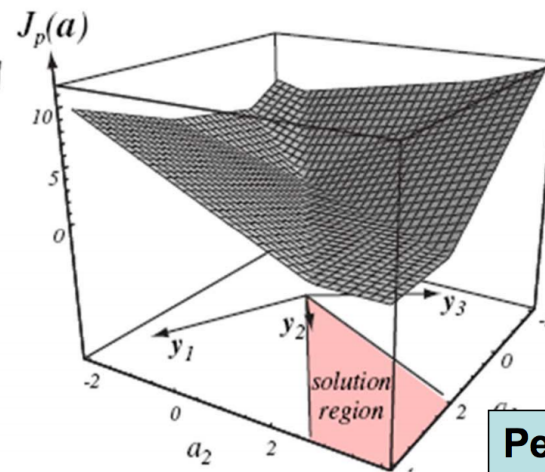
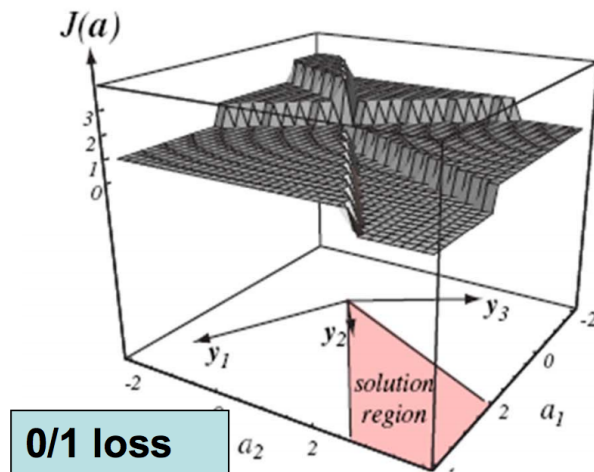
Doesn't produce a useful gradient

The Perceptron Cost Function

- The perceptron uses the following cost function

$$J_p(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \max(0, -y^{(i)} x^{(i)} \boldsymbol{\theta})$$

- $\max(0, -y^{(i)} x^{(i)} \boldsymbol{\theta})$ is 0 if the prediction is correct
- Otherwise, it is the confidence in the misprediction



Nice gradient

Online Perceptron Algorithm

Let $\theta \leftarrow [0, 0, \dots, 0]$

Repeat:

 Receive training example $(\mathbf{x}^{(i)}, y^{(i)})$

 if $y^{(i)} \mathbf{x}^{(i)} \theta \leq 0$ // prediction is incorrect

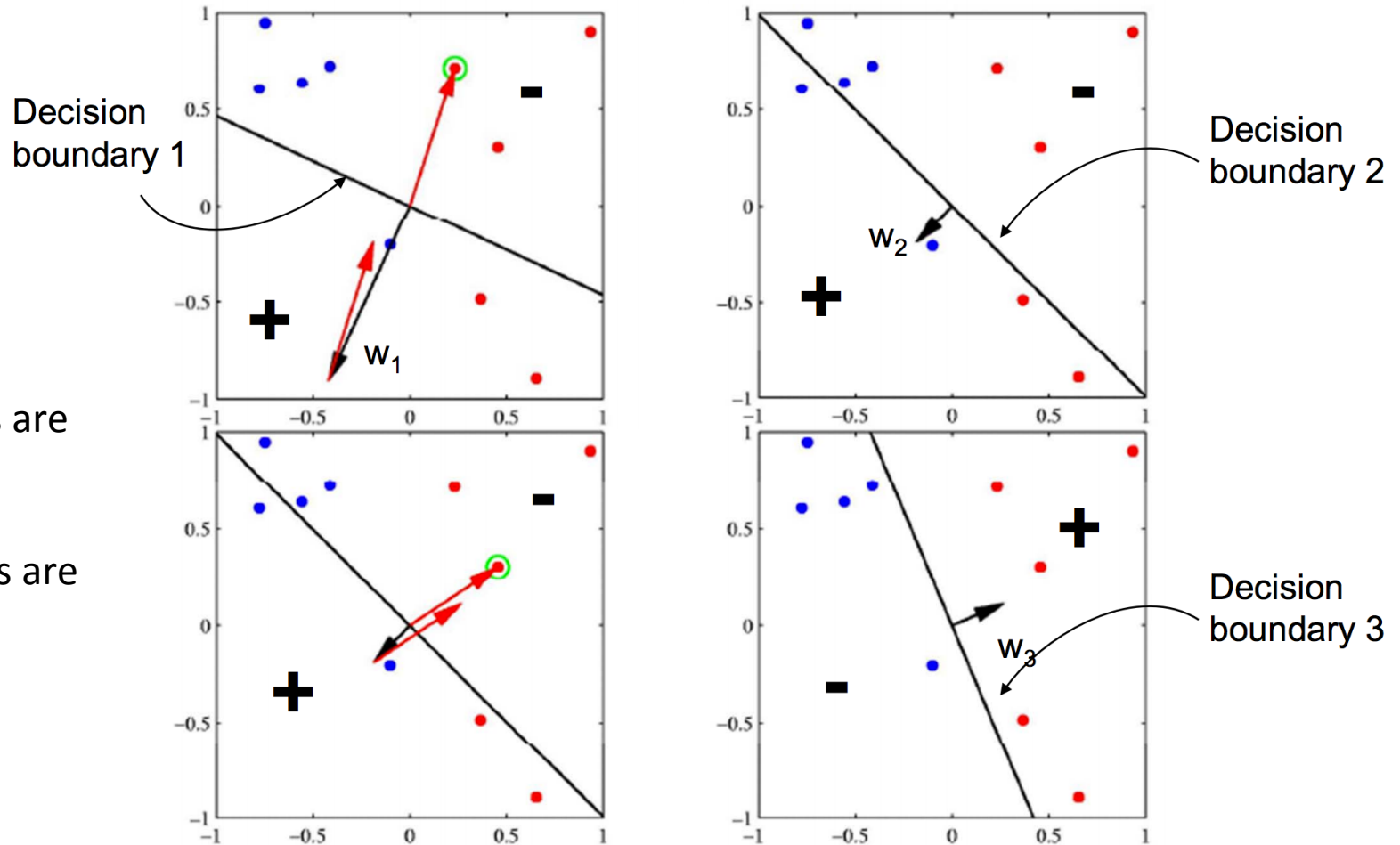
$\theta \leftarrow \theta + y^{(i)} \mathbf{x}^{(i)}$

Online learning – the learning mode where the model update is performed each time a single observation is received

Batch learning – the learning mode where the model update is performed after observing the entire training set

Online Perceptron Algorithm

When an error is made, moves the weight in a direction that corrects the error



See the perceptron in action: www.youtube.com/watch?v=vGwemZhPIsA

Batch Perceptron

Given training data $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$

Let $\boldsymbol{\theta} \leftarrow [0, 0, \dots, 0]$

Repeat:

Let $\boldsymbol{\Delta} \leftarrow [0, 0, \dots, 0]$

for $i = 1 \dots n$, do

if $y^{(i)} \mathbf{x}^{(i)} \boldsymbol{\theta} \leq 0$

// prediction for i^{th} instance is incorrect

$\boldsymbol{\Delta} \leftarrow \boldsymbol{\Delta} + y^{(i)} \mathbf{x}^{(i)}$

$\boldsymbol{\Delta} \leftarrow \boldsymbol{\Delta} / n$

// compute average update

$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \boldsymbol{\Delta}$

Until $\|\boldsymbol{\Delta}\|_2 < \epsilon$

- Simplest case: $\alpha = 1$ and don't normalize, yields the fixed increment perceptron
- Guaranteed to find a separating hyperplane if one exists

Improving the Perceptron

- The Perceptron produces many θ 's during training
- The standard Perceptron simply uses the final θ at test time
 - This may sometimes not be a good idea!
 - Some other θ may be correct on 1,000 consecutive examples, but one mistake ruins it!
- **Idea:** Use a combination of multiple perceptrons
 - (i.e., neural networks!)
- **Idea:** Use the intermediate θ 's
 - **Voted Perceptron:** vote on predictions of the intermediate θ 's
 - **Averaged Perceptron:** average the intermediate θ 's