

COMS 4030A/7047A

Adaptive Computation and Machine Learning

Hima Vadapalli

Semester I, 2022

Perceptron

Slides based heavily on course material by Eric Eton, Alan Fern, Piyush Rai and Andrew Ng

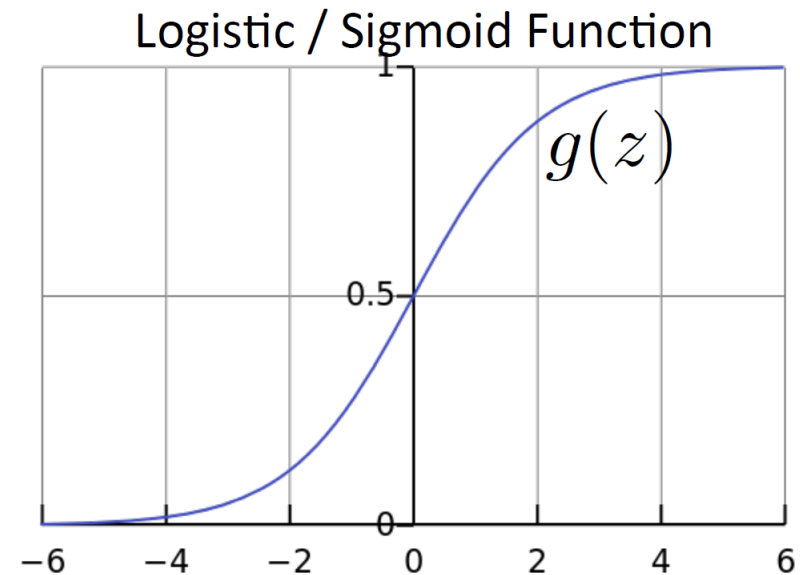
Logistic Regression

- Takes a probabilistic approach to learning discriminative functions (i.e., a classifier)
- $h_{\theta}(\mathbf{x})$ should give $p(y = 1 \mid \mathbf{x}; \theta)$
 - Want $0 \leq h_{\theta}(\mathbf{x}) \leq 1$
- Logistic regression model:

$$h_{\theta}(\mathbf{x}) = g(\theta^{\top} \mathbf{x})$$

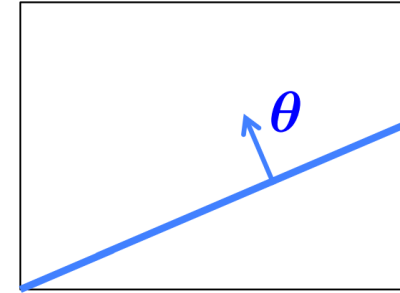
$$g(z) = \frac{1}{1 + e^{-z}}$$

$$h_{\theta}(\mathbf{x}) = \frac{1}{1 + e^{-\theta^{\top} \mathbf{x}}}$$



Linear Classifiers

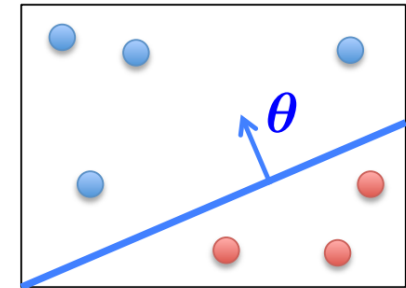
- A **hyperplane** partitions \mathbb{R}^d into two half-spaces
 - Defined by the normal vector $\theta \in \mathbb{R}^d$
 - θ is orthogonal to any vector lying on the hyperplane



- Consider classification with +1, -1 labels ...

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

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



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 $\boldsymbol{\theta}$

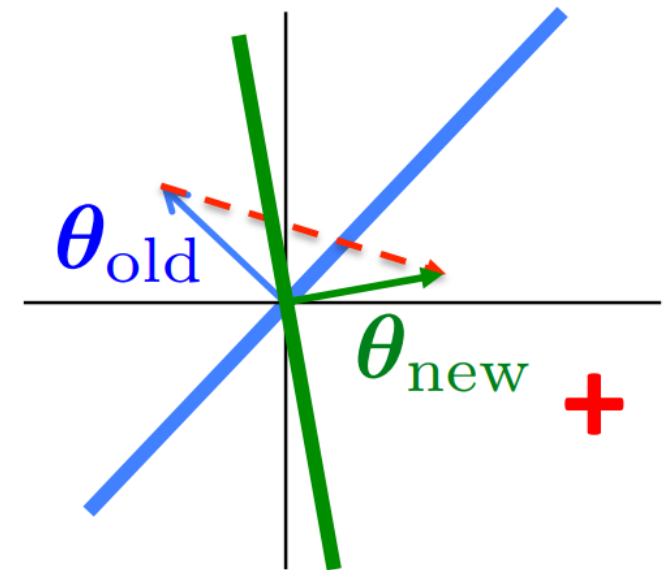
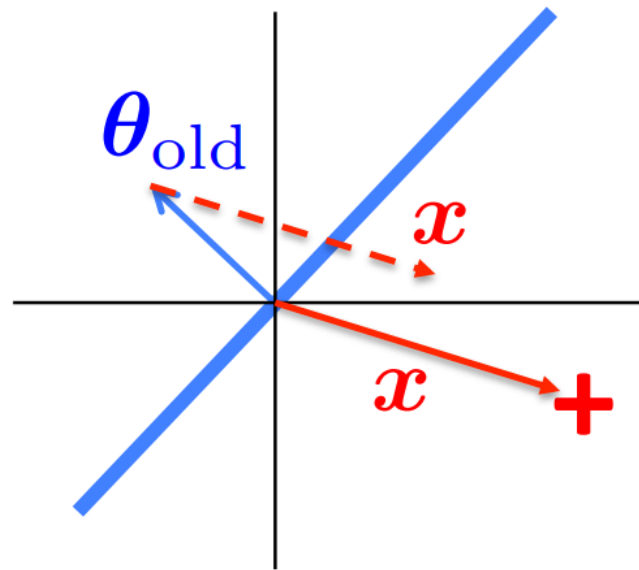
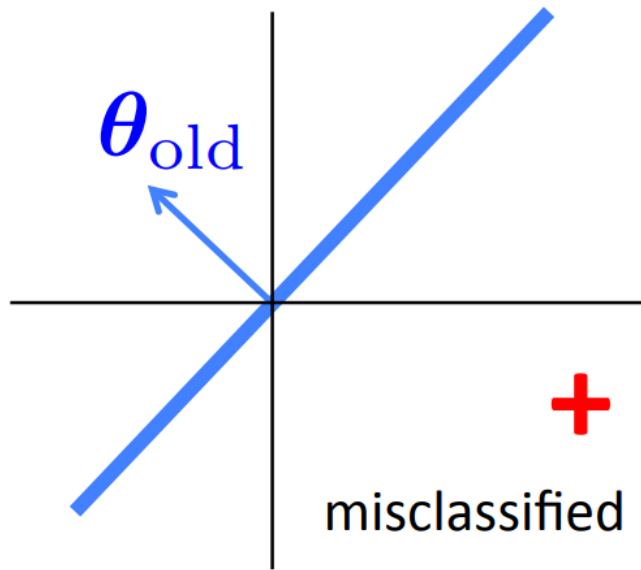
The Perceptron

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 $x^{(i)}$ is misclassified, do $\theta \leftarrow \theta + y^{(i)} x^{(i)}$

Why the Perceptron Update Works



Why the Perceptron Update Works

- Consider the misclassified example ($y = +1$)
 - Perceptron wrongly thinks that $\theta_{\text{old}}^\top \mathbf{x} < 0$

- Update:

$$\theta_{\text{new}} = \theta_{\text{old}} + y\mathbf{x} = \theta_{\text{old}} + \mathbf{x} \quad (\text{since } y = +1)$$

- Note that

$$\begin{aligned}\theta_{\text{new}}^\top \mathbf{x} &= (\theta_{\text{old}} + \mathbf{x})^\top \mathbf{x} \\ &= \theta_{\text{old}}^\top \mathbf{x} + \underbrace{\mathbf{x}^\top \mathbf{x}}_{\|\mathbf{x}\|_2^2 > 0}\end{aligned}$$

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

The Perceptron Cost Function

- The perceptron uses the following cost function

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

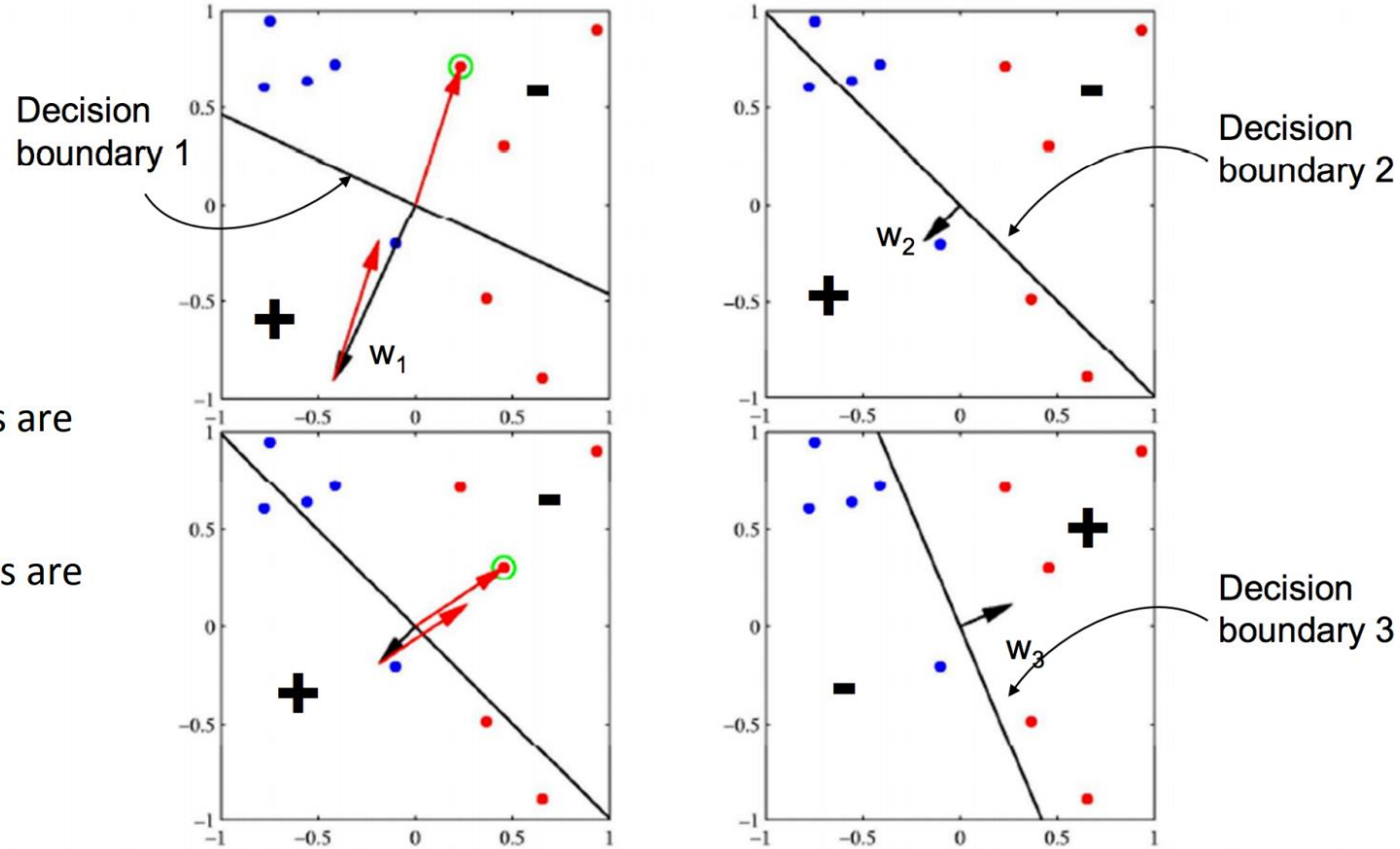
Where $\max(0, -y^{(i)} \theta^T x^{(i)})$
= 0 if prediction is correct
otherwise it is the confidence in the misprediction

Online Perceptron Algorithm

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

Red points are
labeled +

Blue points are
labeled -



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