

Logistic Regression with Neural Networks

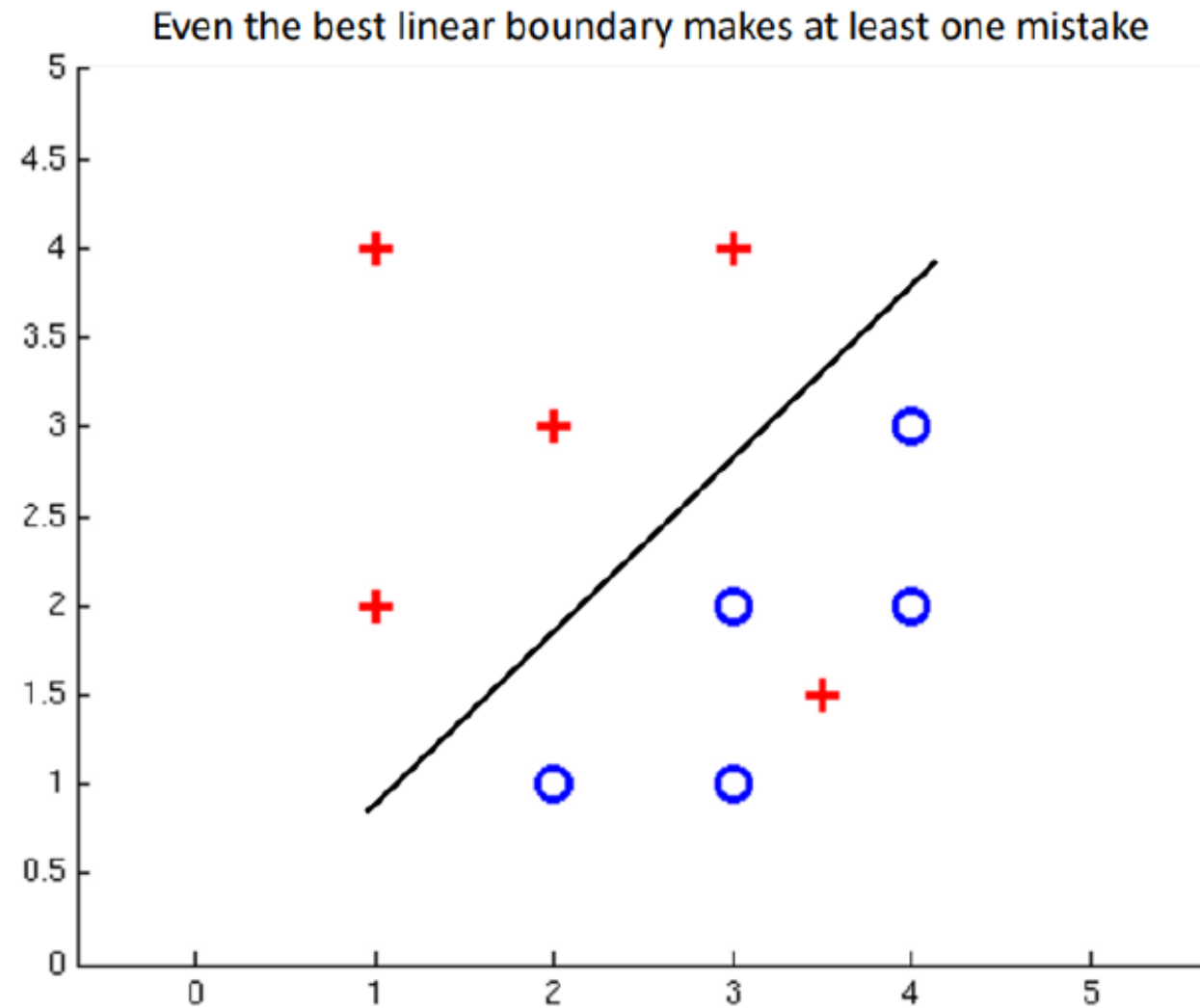
Fatemeh Mansoori
University of Isfahan

Perceptron Conlcultion

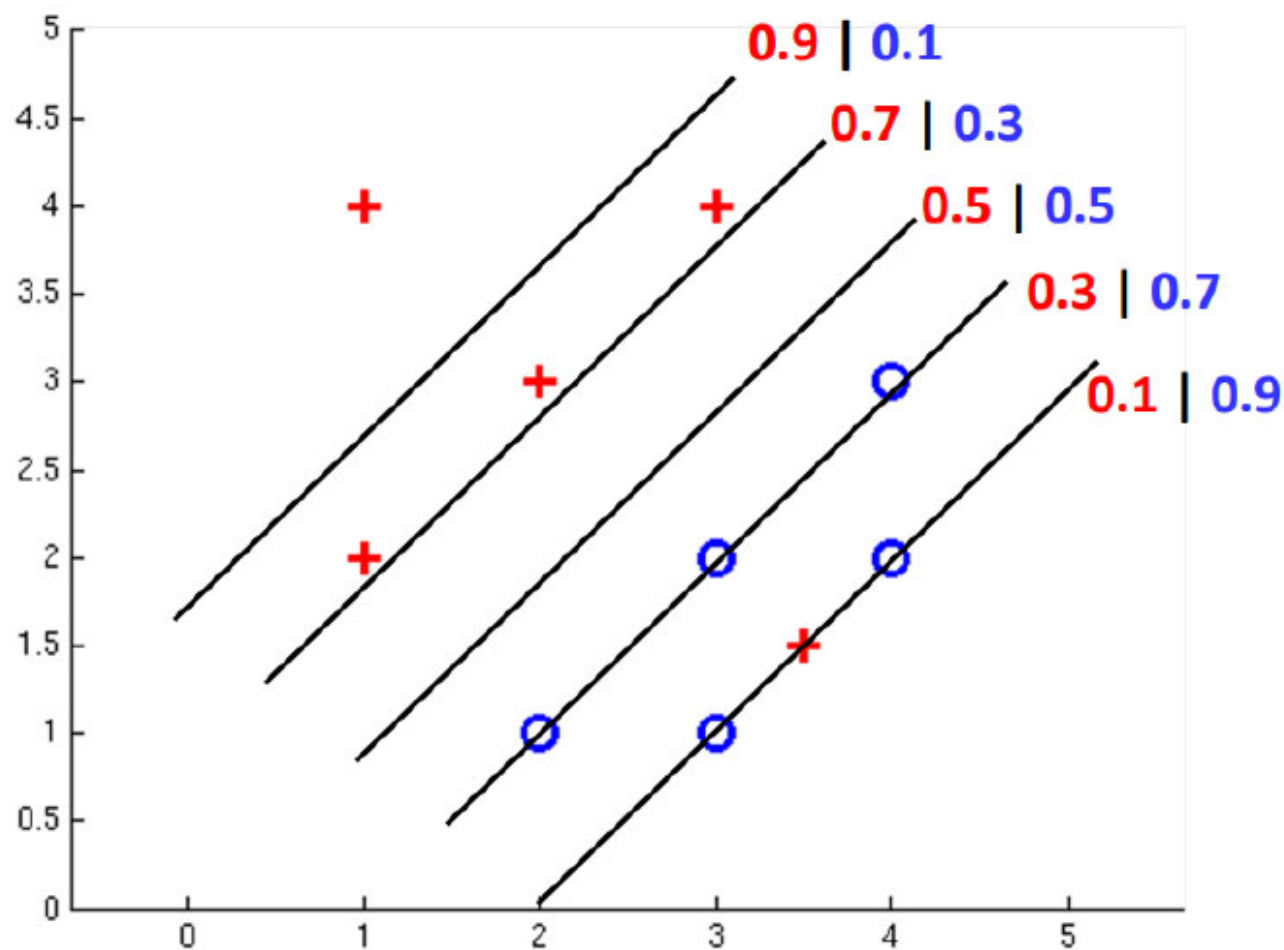
The (classic) Perceptron has many problems
(as discussed in the previous lecture)

- Linear classifier, no non-linear boundaries possible
- Binary classifier
- Does not converge if classes are not linearly separable
- Many "optimal" solutions in terms of 0/1 loss on the training data, most will not be optimal in terms of generalization performance

Non-Separable Case



Non-Separable Case: Probabilistic Decision

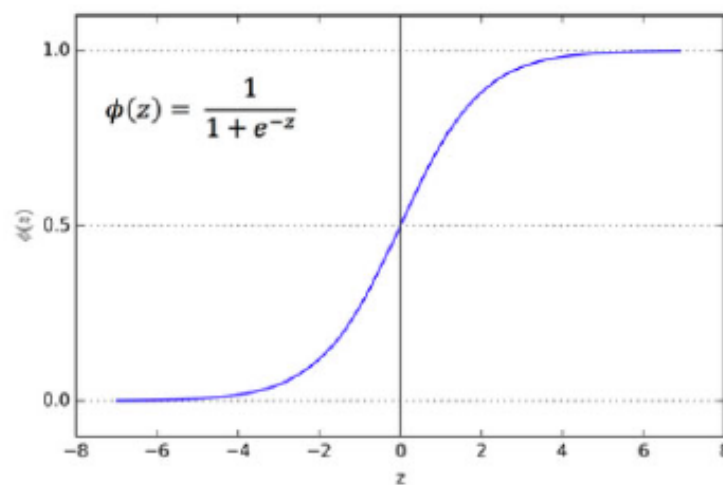


How to get probabilistic decisions?

- Perceptron scoring: $z = w \cdot f(x)$
- If $z = w \cdot f(x)$ very positive \rightarrow want probability of + going to 1
- If $z = w \cdot f(x)$ very negative \rightarrow want probability of + going to 0

- Sigmoid function

$$\begin{aligned}\phi(z) &= \frac{1}{1 + e^{-z}} \\ &= \frac{e^z}{e^z + 1}\end{aligned}$$



How to get probabilistic decisions?

- Perceptron scoring: $z = w \cdot f(x)$
- If $z = w \cdot f(x)$ very positive \rightarrow want probability of + going to 1
- If $z = w \cdot f(x)$ very negative \rightarrow want probability of + going to 0

- Sigmoid function

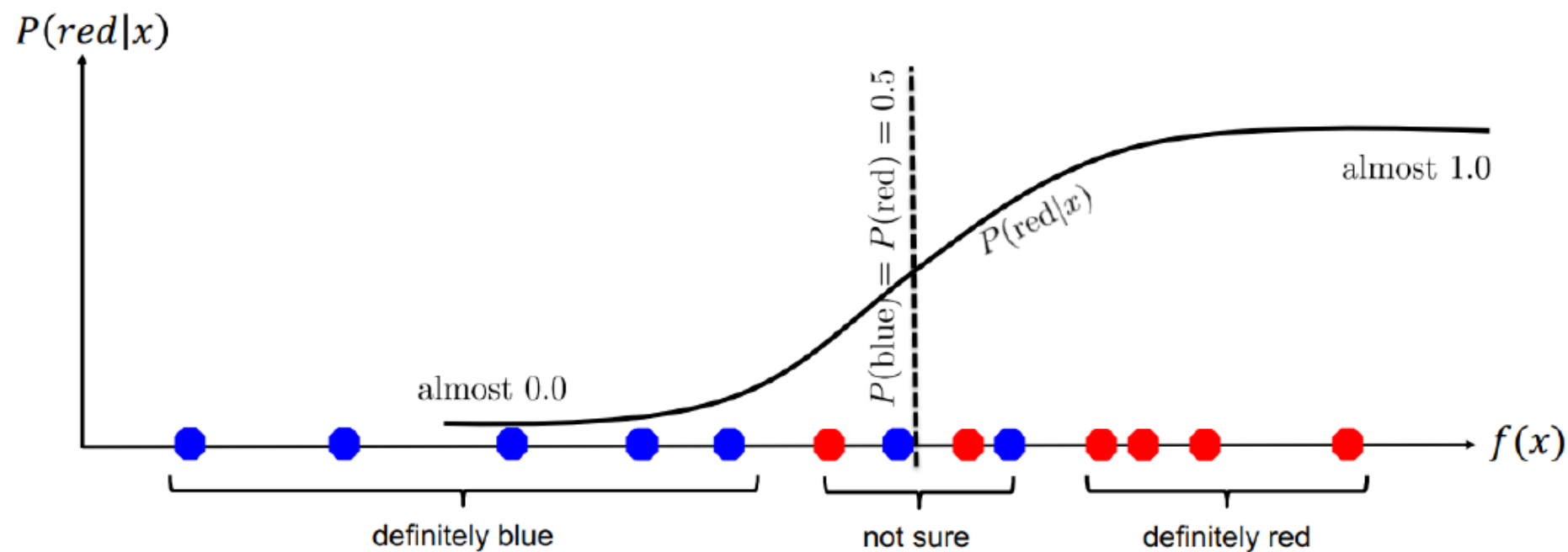
$$\phi(z) = \frac{1}{1 + e^{-z}}$$

$$P(y = +1 \mid x; w) = \frac{1}{1 + e^{-w \cdot f(x)}}$$

$$P(y = -1 \mid x; w) = 1 - \frac{1}{1 + e^{-w \cdot f(x)}}$$

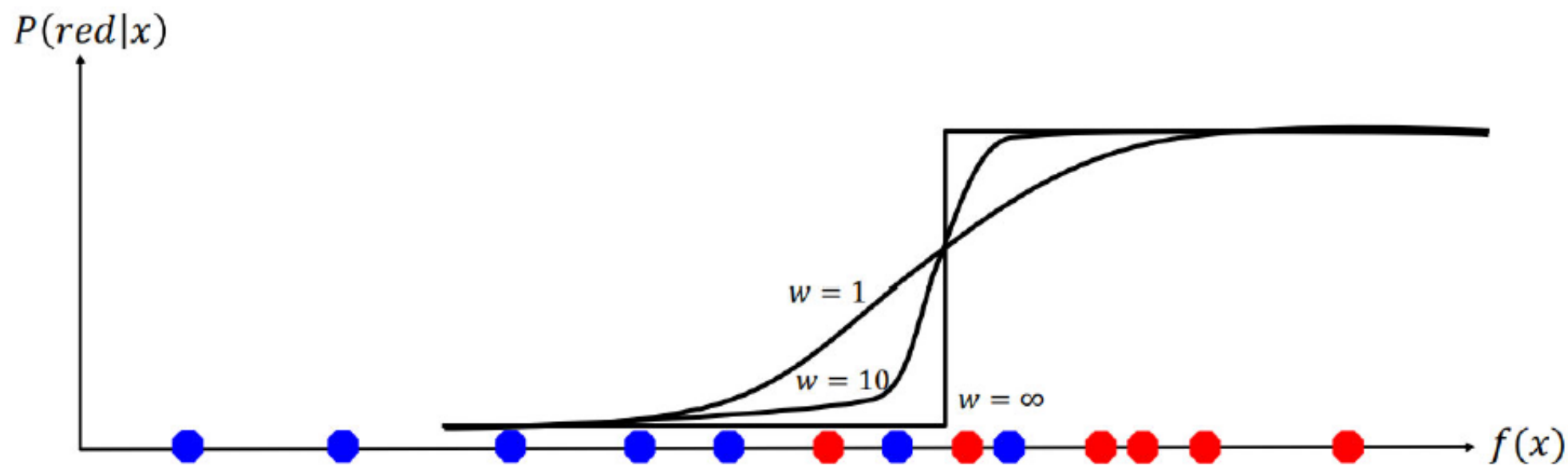
= **Logistic Regression**

A 1D Example



$$P(\text{red}|x; w) = \phi(w \cdot f(x)) = \frac{1}{1 + e^{-w \cdot f(x)}}$$

A 1D Example: varying w



$$P(\text{red}|x; w) = \phi(w \cdot f(x)) = \frac{1}{1 + e^{-w \cdot f(x)}}$$

Best w ?

$$\text{Likelihood} = P(\text{training data} | w)$$

$$= \prod_i P(\text{training datapoint } i \mid w)$$

$$= \prod_i P(\text{point } x^{(i)} \text{ has label } y^{(i)} | w)$$

$$= \prod_i P(y^{(i)} | x^{(i)}; w)$$

$$\text{Log Likelihood} = \sum_i \log P(y^{(i)} | x^{(i)}; w)$$

Best w ?

$$\max_w ll(w) = \max_w \sum_i \log P(y^{(i)} | x^{(i)}; w)$$

Logistic regression cost function

If $y = 1$: $p(y|x) = \hat{y}$

If $y = 0$: $p(y|x) = 1 - \hat{y}$

$$p(y|x) = \hat{y}^y (1-\hat{y})^{(1-y)}$$

If $y=1$: $p(y|x) = \hat{y} \underbrace{(1-\hat{y})^0}_{=1}$

If $y=0$: $p(y|x) = \hat{y}^0 \underbrace{(1-\hat{y})^{(1-0)}}_{=1} = 1 \times (1-\hat{y}) = 1-\hat{y}$

$$\begin{aligned} \uparrow \log p(y|x) &= \log \hat{y}^y (1-\hat{y})^{(1-y)} = y \log \hat{y} + (1-y) \log (1-\hat{y}) \\ &= -\frac{1}{\epsilon} \ell(\hat{y}, y) \downarrow \end{aligned}$$

Loss Optimization

We want to find the network weights that **achieve the lowest loss**

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \mathbf{W}), y^{(i)})$$

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} J(\mathbf{W})$$

Cross Entropy Loss Optimization

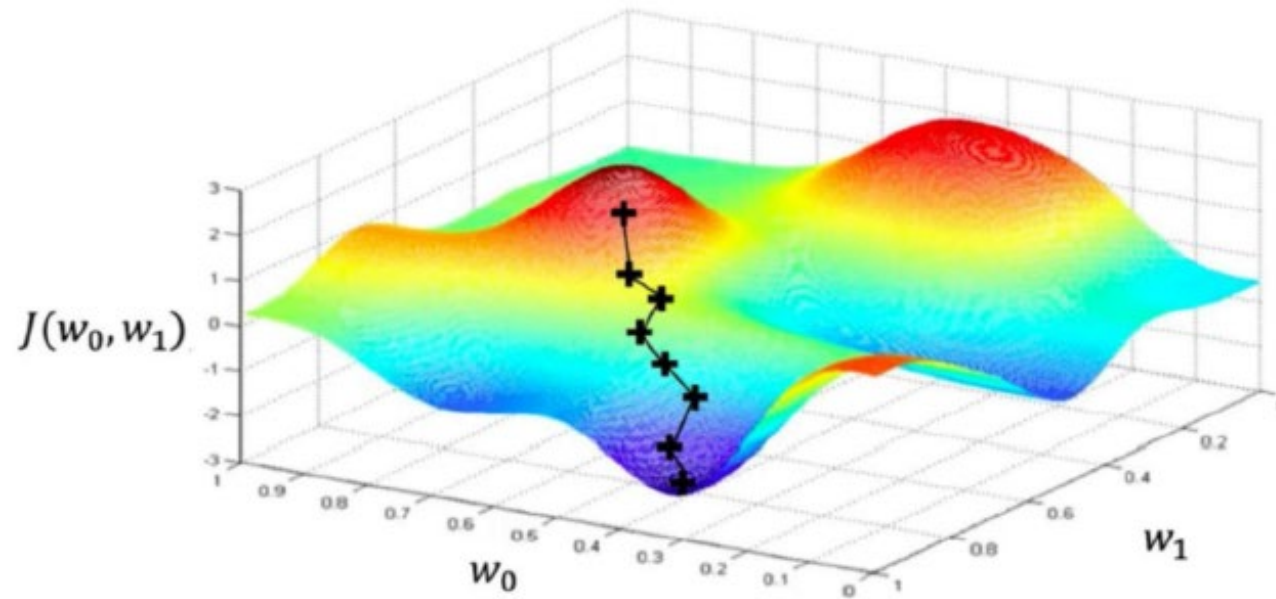
$$\text{Recap: } \hat{y} = \sigma(w^T x + b), \quad \sigma(z) = \frac{1}{1+e^{-z}}$$

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

Want to find w, b that minimize $J(w, b)$

Loss Optimization

Repeat until convergence



Gradient Descent

“Walking downhill and always taking a step in the direction that goes down the most.”

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}}$
4. Update weights, $\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}}$
5. Return weights

Logistic regression derivatives

