

Deep learning with Neural Networks

Logistic Regression, Stochastic Optimization & Regularization

Pablo Martínez Olmos, pamartin@ing.uc3m.es

Binary Logistic Regression

- Binary classification method

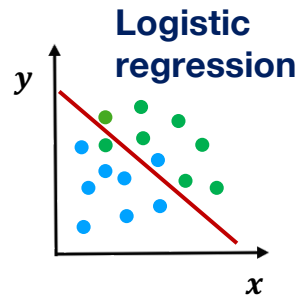


x

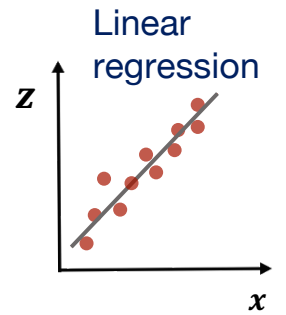


1: CAT
0: NOT CAT

y

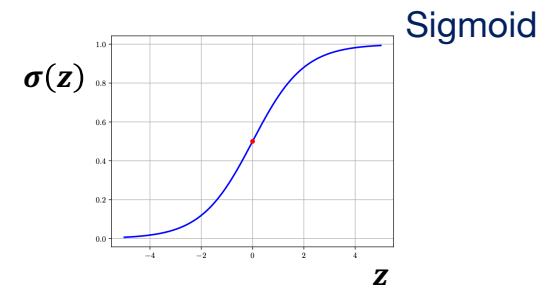


=



$$z = w^T x + w_0$$

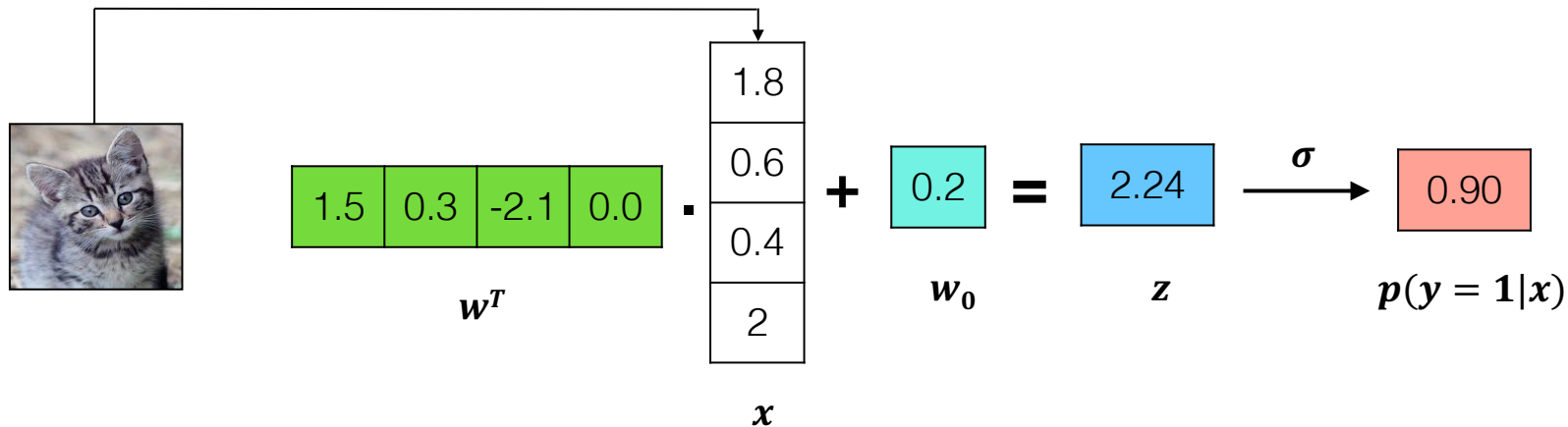
+



$$p(y=1|x) = \sigma(w^T x + w_0) = \frac{1}{1 + e^{-(w^T x + w_0)}}$$

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

- Example



Binary Logistic Regression

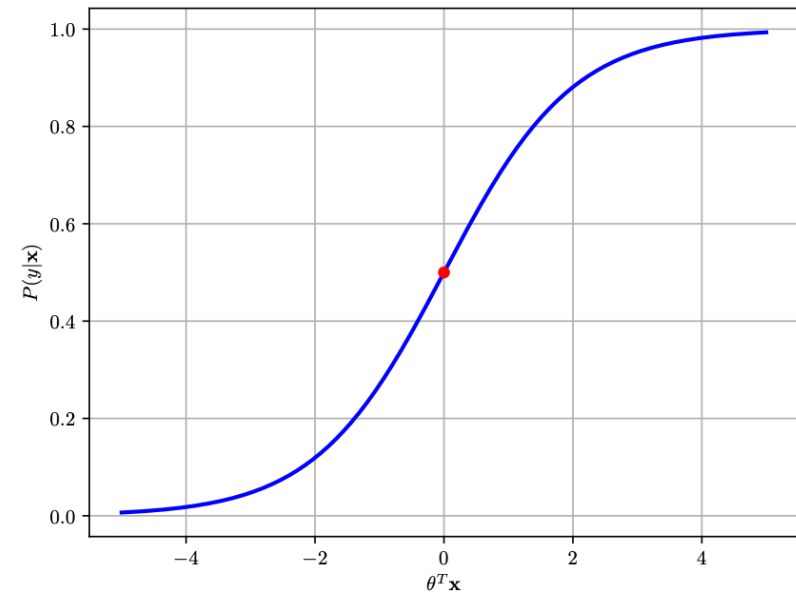
- Training database: $\mathcal{D} \doteq (\mathbf{x}^{(i)}, y^{(i)})_{i=1}^N$ $\mathbf{x}^{(i)} \in \mathbb{R}^m$ $y^{(i)} \in \{0, 1\}$

$$P(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + w_0)}} \doteq \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

- Decision boundary:

$$\{\mathbf{x} \in \mathbb{R}^m : \mathbf{w}^T \mathbf{x} + w_0 = 0\} \leftarrow$$

Hyperplane! Linear classifier



Logistic Regression: training

Binary Logistic Regression. Binary Cross Entropy Function

$$P(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + w_0)}} \doteq \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

Binary Cross Entropy Function

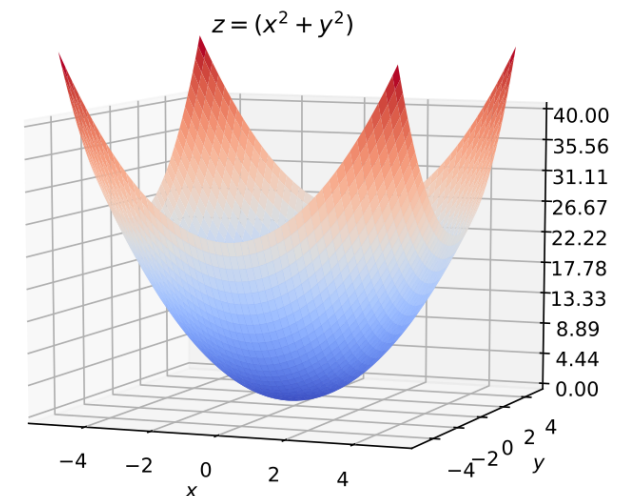
$$\mathcal{L} = -\log P(\mathbf{y}|\mathbf{X})$$

$$= -\sum_{i=1}^N \mathbb{1}[y^{(i)} = 1] \log P(y^{(i)} = 1|\mathbf{x}^{(i)}) + \mathbb{1}[y^{(i)} = 0] \log P(y^{(i)} = 0|\mathbf{x}^{(i)})$$

\mathcal{L} is convex!!

Logistic Regression Training

$$\hat{\mathbf{w}}, \hat{w}_0 = \arg \min_{\mathbf{w}, w_0} \mathcal{L}$$



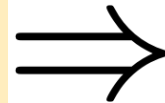
Logistic Regression Training

Binary Cross Entropy Function

$$\begin{aligned}\mathcal{L} &= -\log P(\mathbf{y}|\mathbf{X}) \\ &= -\sum_{i=1}^N \mathbb{1}[y^{(i)} = 1] \log P(y^{(i)} = 1|\mathbf{x}^{(i)}) + \mathbb{1}[y^{(i)} = 0] \log P(y^{(i)} = 0|\mathbf{x}^{(i)})\end{aligned}$$

Logistic Regression Training

$$\hat{\mathbf{w}}, \hat{w}_0 = \arg \min_{\mathbf{w}, w_0} \mathcal{L}$$



$$\nabla_{\mathbf{w}, w_0} \mathcal{L}(\hat{\mathbf{w}}, \hat{w}_0) = \mathbf{0}$$

No closed-form solution!

Gradient Descent (GD)

$$\mathbf{w}_+ \doteq \begin{bmatrix} w_0 & \mathbf{w} \end{bmatrix} \Rightarrow \nabla_{\mathbf{w}_+} \mathcal{L} = \mathbf{0}$$

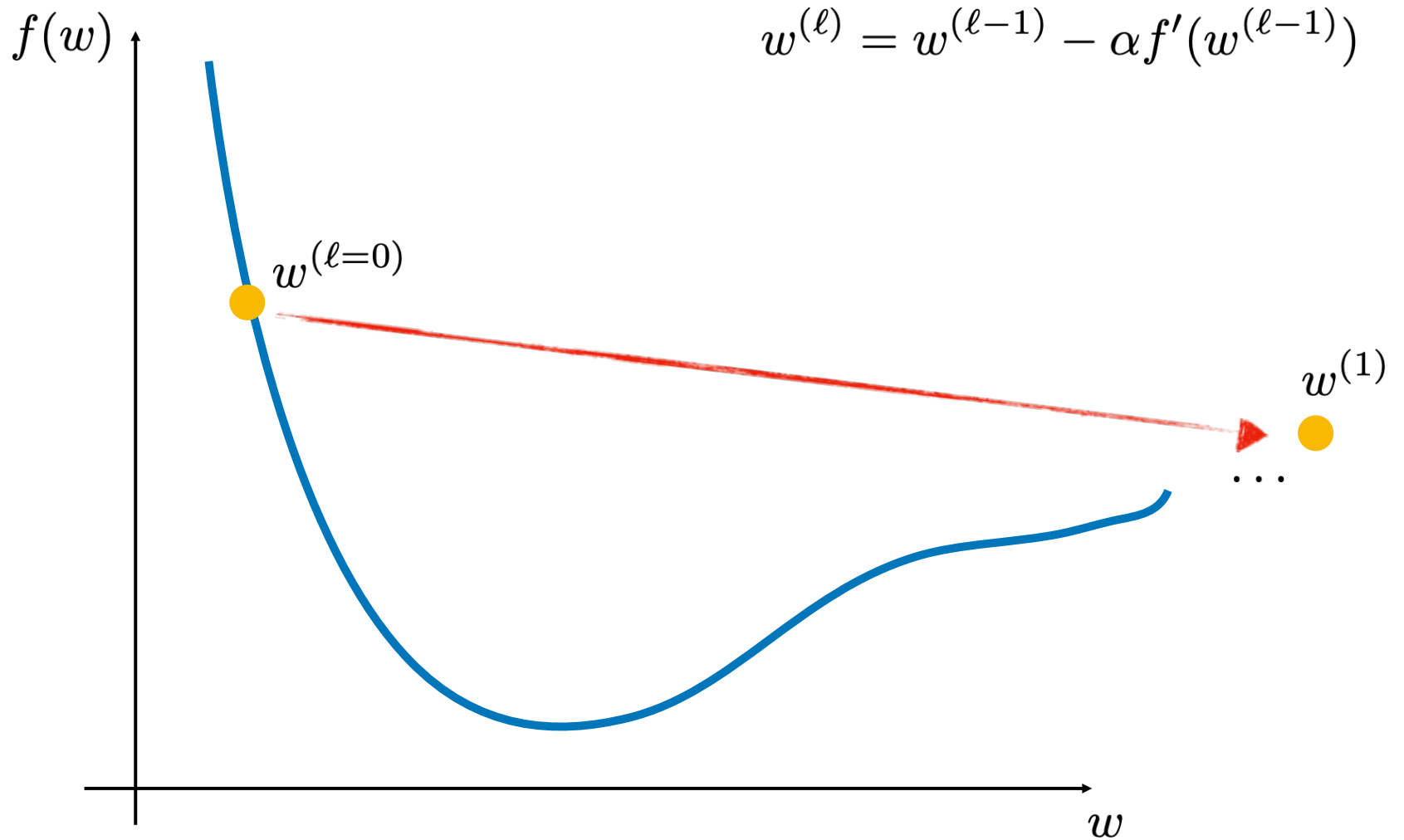
$$\mathbf{w}_+^{(\ell)} = \mathbf{w}_+^{(\ell-1)} - (\alpha \nabla \mathcal{L})|_{\mathbf{w}_+^{(\ell-1)}}$$



the step size

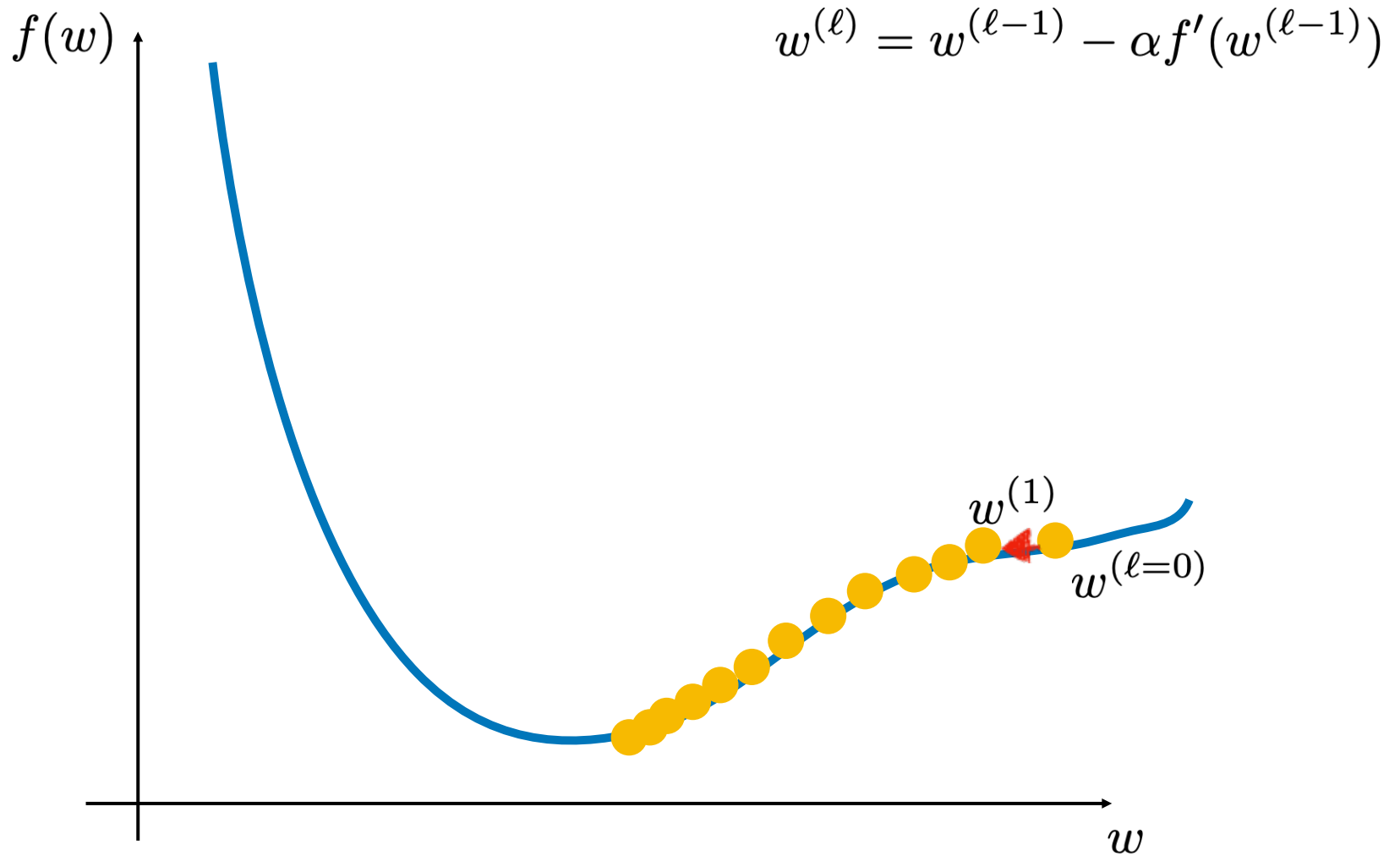
Choosing the step size manually is a **hard problem!**

Gradient Descent



If α is too large

Gradient Descent



If α is too small

Line Search methods: automatically running the step size

$$\mathbf{w}_+^{(\ell)} = \mathbf{w}_+^{(\ell-1)} - \left(\alpha^{(\ell)} \nabla \mathcal{L} \right) |_{\mathbf{w}_+^{(\ell-1)}}$$

- Estimate the largest possible step size that guarantees that the function decreases
- Every state-of-the-art Deep Learning library contains implementations of various algorithms to optimize gradient descent
- Momentum, Adam, Adagrad, ...

Check out in Aula Global two excellent posts on SGD methods for Deep Learning

[\(Link1, Link2\)](#)

Logistic Regression: Stochastic gradient descent

Stochastic Optimization (Mini-batch Optimization)

Binary Cross Entropy Function

$$\begin{aligned}\mathcal{L} &= -\log P(\mathbf{y}|\mathbf{X}) \\ &= -\sum_{i=1}^N \mathbb{1}[y^{(i)} = 1] \log P(y^{(i)} = 1|\mathbf{x}^{(i)}) + \mathbb{1}[y^{(i)} = 0] \log P(y^{(i)} = 0|\mathbf{x}^{(i)})\end{aligned}$$

Expensive for very large
databases!!

$$\nabla_{\mathbf{w}_+} \mathcal{L} = \sum_{i=1}^N \left(\sigma \left(\mathbf{w}_+ \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix} \right) - y^{(i)} \right) \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix}$$

Mini-batch optimization

Select at random a mini-batch \mathcal{B} of data at every SGD iteration

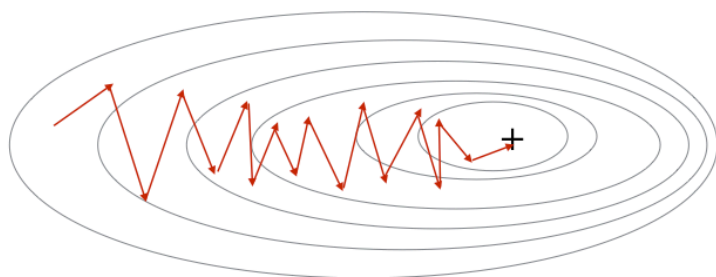
Stochastic Optimization (Mini-batch Optimization)

$$\mathbf{w}_+ \doteq \begin{bmatrix} w_0 & \mathbf{w} \end{bmatrix} \Rightarrow \nabla_{\mathbf{w}_+} \mathcal{L} = \mathbf{0}$$

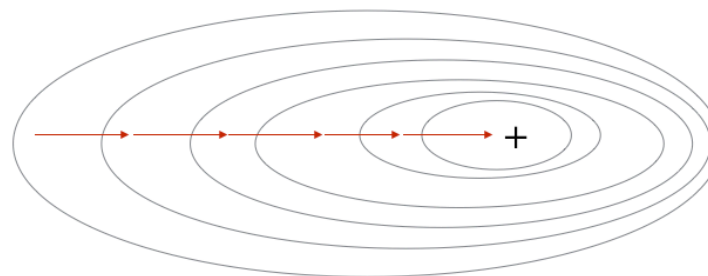
$$\mathbf{w}_+^{(\ell)} = \mathbf{w}_+^{(\ell-1)} - (\alpha \nabla \mathcal{L})|_{\mathbf{w}_+^{(\ell-1)}}$$

$$\nabla_{\mathbf{w}_+} \mathcal{L} \approx \sum_{i \in \mathcal{B}} \left(\sigma \left(\mathbf{w}_+ \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix} \right) - y^{(i)} \right) \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix}$$

Stochastic Gradient Descent



Gradient Descent



Source: this [post](#)

Logistic Regression: multiple classes

Binary Logistic Regression

- Discriminative classifier:

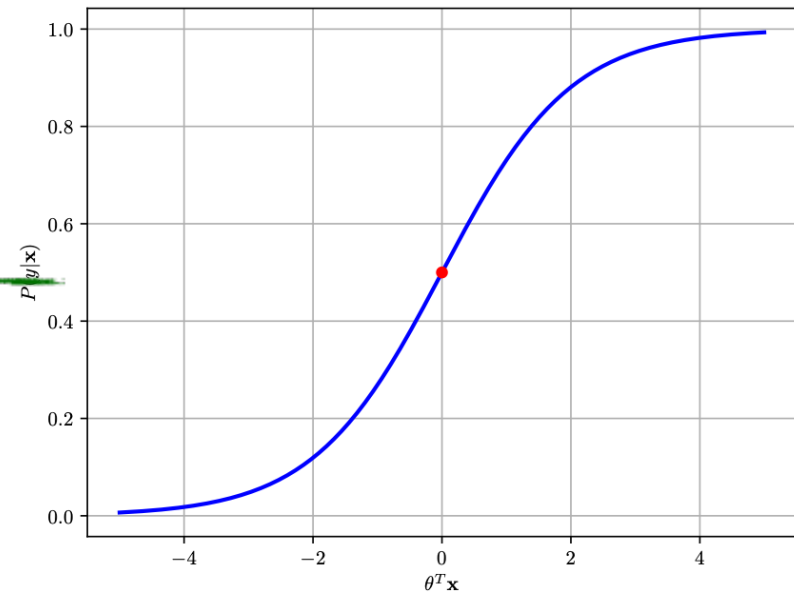
$$P(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + w_0)}} \doteq \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

Generalized linear model

- Decision boundary:

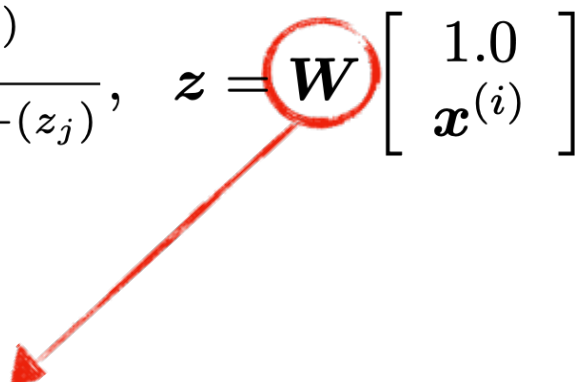
$$\{\mathbf{x} \in \mathbb{R}^m : \mathbf{w}^T \mathbf{x} + w_0 = 0\}$$

Hyperplane! Linear classifier



Multiclass Logistic Regression

- Training database: $\mathcal{D} \doteq (\mathbf{x}^{(i)}, y^{(i)})_{i=1}^N$ $\mathbf{x}^{(i)} \in \mathbb{R}^m$ $y^{(i)} \in \{1, \dots, K\}$
- Discriminative classifier based on the **Softmax function**:

$$P(y = k | \mathbf{x}) = \frac{e^{-(z_k)}}{\sum_{j=1}^K e^{-(z_j)}}, \quad \mathbf{z} = \mathbf{W} \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix}$$


$K \times (m + 1)$

- Loss function (Cross entropy):

$$\mathcal{L} = -\log P(\mathbf{y} | \mathbf{X}) = -\sum_{n=1}^N \sum_{k=1}^K \mathbb{1}[y^{(n)} == k] \log P(y = k | \mathbf{x})$$

- Similar gradient expression