

# Outline

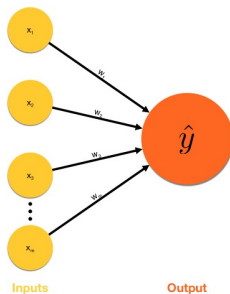
Introduction to Neural nets

Training Deep Neural Networks

Introduction to Statistical Decision Theory

## The Formal Neuron: 1943 [?]

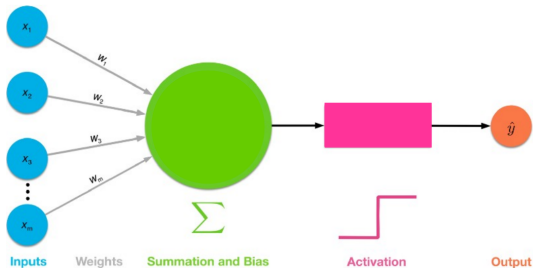
- ▶ Basis of Neural Networks
- ▶ Input: vector  $x \in \mathbb{R}^m$ , i.e.  $x = \{x_i\}_{i \in \{1,2,\dots,m\}}$
- ▶ Neuron output  $\hat{y} \in \mathbb{R}$ : scalar



# The Formal Neuron: 1943 [?]

► Mapping from  $x$  to  $\hat{y}$ :

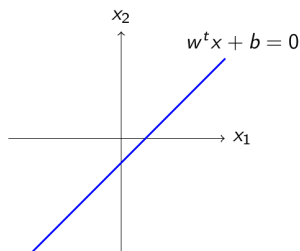
1. Linear (affine) mapping:  $s = w^T x + b$
2. Non-linear activation function:  $f: \hat{y} = f(s)$



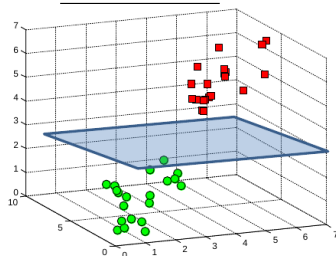
# The Formal Neuron: Linear Mapping

- ▶ Linear (affine) mapping:  $s = w^T x + b = \sum_{i=1}^m w_i x_i + b$ 
  - ▶  $w$ : normal vector to an hyperplane in  $\mathbb{R}^m \Rightarrow$  **linear boundary**
  - ▶  $b$  bias, shift the hyperplane position

## 2D hyperplane: line

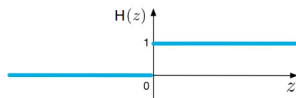


## 3D hyperplane: plane

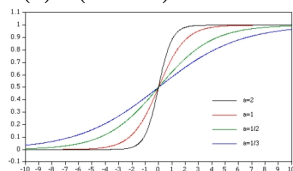


# The Formal Neuron: Activation Function

- ▶  $\hat{y} = f(\mathbf{w}^T \mathbf{x} + b)$ ,
- ▶  $f$ : activation function
  - ▶ Bio-inspired choice: Step (Heaviside) function:  $H(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{otherwise} \end{cases}$

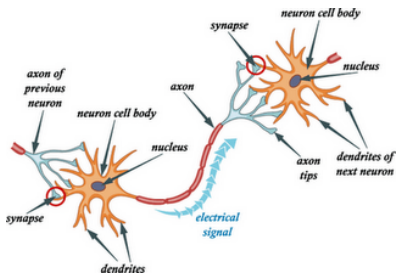
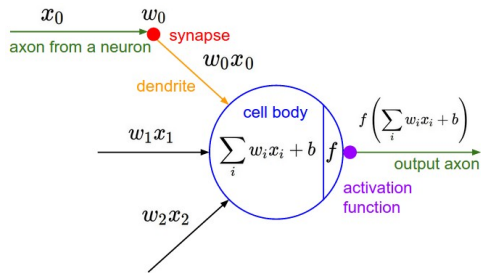


- ▶ Popular  $f$  choices: sigmoid, tanh, ReLU, GELU, ...
- ▶ Sigmoid:  $\sigma(z) = (1 + e^{-az})^{-1}$



- ▶  $a \uparrow$ : more similar to step function (step:  $a \rightarrow \infty$ )
- ▶ Sigmoid: linear and saturating regimes

## Step function: Connection to Biological Neurons



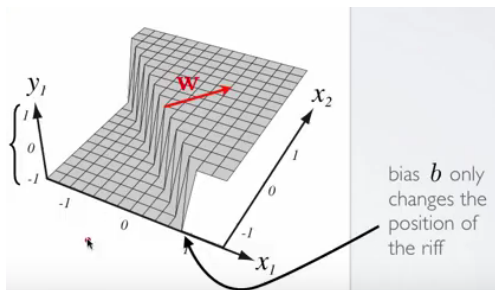
- ▶ Formal neuron, step activation  $H$ :  $\hat{y} = H(w^T x + b)$ 
  - ▶  $\hat{y} = 1$  (activated)  $\Leftrightarrow w^T x \geq -b$
  - ▶  $\hat{y} = 0$  (unactivated)  $\Leftrightarrow w^T x < -b$
- ▶ **Biological Neurons:** output activated  $\Leftrightarrow$  input weighted by synaptic weight  $\geq$  threshold

## The Formal neuron: Application to Binary Classification

- ▶ Binary Classification: label input  $x$  as belonging to class 1 or 0
- ▶ Neuron output with sigmoid:

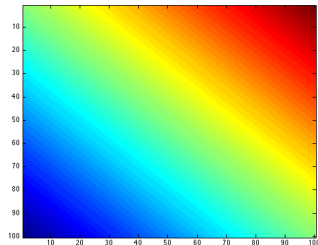
$$\hat{y} = \frac{1}{1 + e^{-a(w^T x + b)}}$$

- ▶ Sigmoid: probabilistic interpretation  $\Rightarrow \hat{y} \sim P(1|x)$ 
  - ▶ Input  $x$  classified as 1 if  $P(1|x) > 0.5 \Leftrightarrow w^T x + b > 0$
  - ▶ Input  $x$  classified as 0 if  $P(1|x) < 0.5 \Leftrightarrow w^T x + b < 0$   
 $\Rightarrow \text{sign}(w^T x + b)$ : linear boundary decision in input space !



## The Formal neuron: Toy Example for Binary Classification

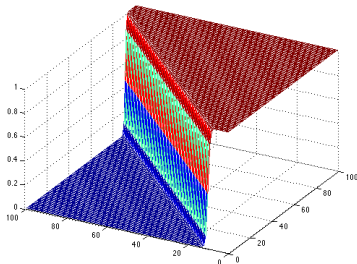
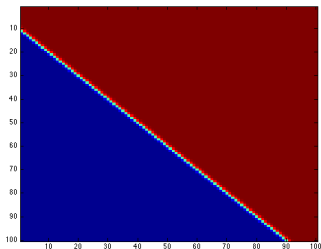
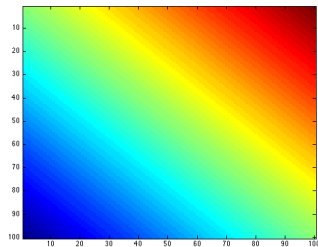
- ▶ 2d example:  $m = 2$ ,  $\mathbf{x} = \{x_1, x_2\} \in [-5; 5] \times [-5; 5]$
- ▶ Linear mapping:  $\mathbf{w} = [1; 1]$  and  $b = -2$
- ▶ Result of linear mapping :  $s = \mathbf{w}^\top \mathbf{x} + b$





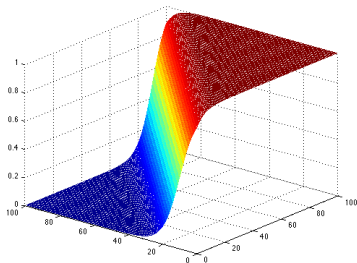
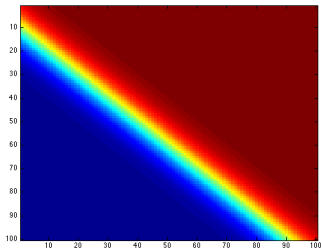
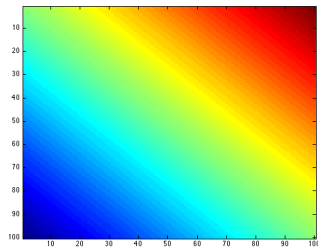
# The Formal neuron: Toy Example for Binary Classification

- ▶ 2d example:  $m = 2$ ,  $\mathbf{x} = \{x_1, x_2\} \in [-5; 5] \times [-5; 5]$
- ▶ Linear mapping:  $\mathbf{w} = [1; 1]$  and  $b = -2$
- ▶ Result of linear mapping :  $s = \mathbf{w}^\top \mathbf{x} + b$
- ▶ Sigmoid activation function:  $\hat{y} = \left(1 + e^{-a(\mathbf{w}^\top \mathbf{x} + b)}\right)^{-1}$ ,  
 $a = 10$



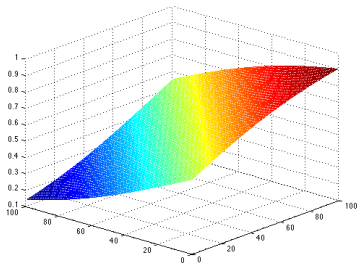
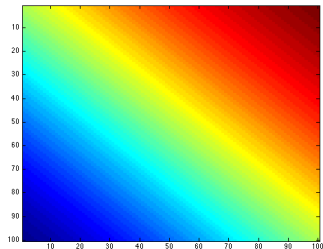
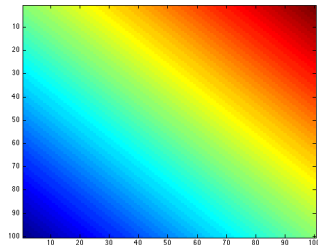
# The Formal neuron: Toy Example for Binary Classification

- ▶ 2d example:  $m = 2$ ,  $\mathbf{x} = \{x_1, x_2\} \in [-5; 5] \times [-5; 5]$
- ▶ Linear mapping:  $\mathbf{w} = [1; 1]$  and  $b = -2$
- ▶ Result of linear mapping :  $s = \mathbf{w}^\top \mathbf{x} + b$
- ▶ Sigmoid activation function:  $\hat{y} = \left(1 + e^{-a(\mathbf{w}^\top \mathbf{x} + b)}\right)^{-1}$ ,  
 $a = 1$



# The Formal neuron: Toy Example for Binary Classification

- ▶ 2d example:  $m = 2$ ,  $x = \{x_1, x_2\} \in [-5; 5] \times [-5; 5]$
- ▶ Linear mapping:  $w = [1; 1]$  and  $b = -2$
- ▶ Result of linear mapping :  $s = w^\top x + b$
- ▶ Sigmoid activation function:  $\hat{y} = \left(1 + e^{-a(w^\top x + b)}\right)^{-1}$ ,  
 $a = 0.1$



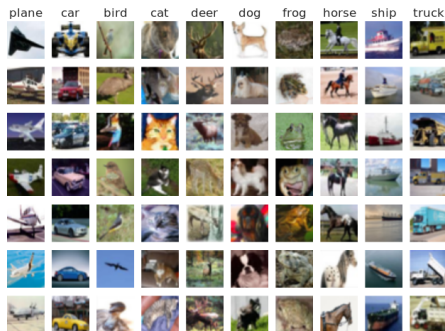
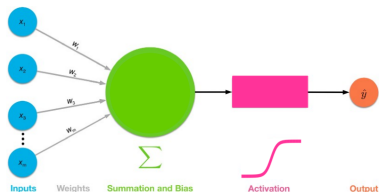
# From Formal Neuron to Neural Networks

## ▶ Formal Neuron:

1. A single scalar output
2. Linear decision boundary for binary classification

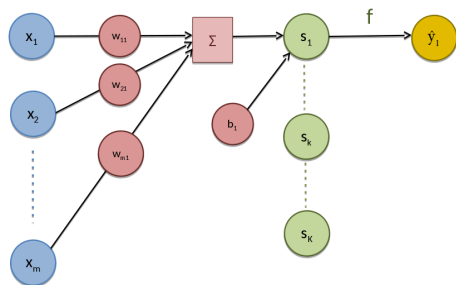
## ▶ Single scalar output: limited for several tasks

- ▶ Ex: multi-class classification, e.g. MNIST or CIFAR



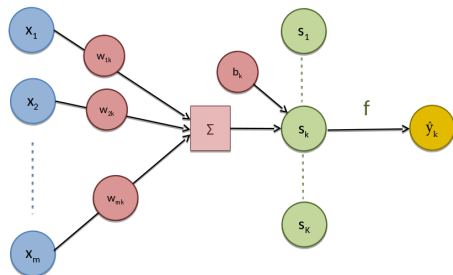
# Perceptron and Multi-Class Classification

- ▶ Formal Neuron: limited to binary classification
- ▶ **Multi-Class Classification:** use several output neurons instead of a single one!  
⇒ **Perceptron**
- ▶ Input  $x$  in  $\mathbb{R}^m$
- ▶ Output neuron  $\hat{y}_1$  is a formal neuron:
  - ▶ Linear (affine) mapping:  $s_1 = w_1^T x + b_1$
  - ▶ Non-linear activation function:  $f$ :  
 $\hat{y}_1 = f(s_1)$
- ▶ Linear mapping parameters:
  - ▶  $w_1 = \{w_{11}, \dots, w_{m1}\} \in \mathbb{R}^m$
  - ▶  $b_1 \in \mathbb{R}$



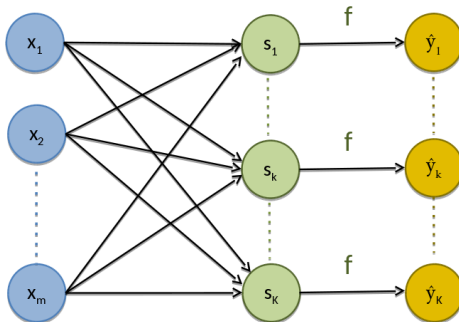
# Perceptron and Multi-Class Classification

- ▶ Input  $x$  in  $\mathbb{R}^m$
- ▶ Output neuron  $\hat{y}_k$  is a formal neuron:
  - ▶ Linear (affine) mapping:  $s_k = w_k^\top x + b_k$
  - ▶ Non-linear activation function:  $f$ :  
 $\hat{y}_k = f(s_k)$
- ▶ Linear mapping parameters:
  - ▶  $w_k = \{w_{1k}, \dots, w_{mk}\} \in \mathbb{R}^m$
  - ▶  $b_k \in \mathbb{R}$



# Perceptron and Multi-Class Classification

- ▶ Input  $x$  in  $\mathbb{R}^m$  ( $1 \times m$ ), output  $\hat{y}$ : concatenation of  $K$  formal neurons
- ▶ Linear (affine) mapping  $\sim$  matrix multiplication:  $s = xW + b$ 
  - ▶  $W$  matrix of size  $m \times K$  - columns are  $w_k$
  - ▶  $b$ : bias vector - size  $1 \times K$
- ▶ Element-wise non-linear activation:  $\hat{y} = f(s)$



# Perceptron and Multi-Class Classification

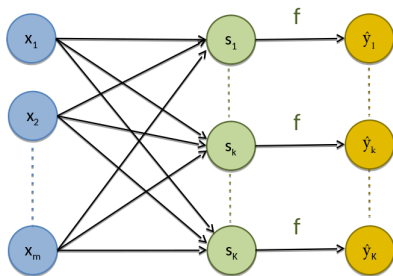
- ▶ **Soft-max Activation:**

$$\hat{y}_k = f(s_k) = \frac{e^{s_k}}{\sum_{k'=1}^K e^{s_{k'}}$$

- ▶ Note that  $f(s_k)$  depends on the other  $s_{k'}$ , the arrow is a functional link
- ▶ **Probabilistic interpretation for multi-class classification:**

- ▶ Each output neuron  $\leftrightarrow$  class
- ▶  $\hat{y}_k \sim P(k|x, W)$

⇒ **Logistic Regression (LR) Model!**





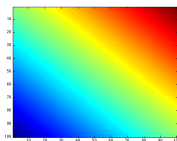
## 2d Toy Example for Multi-Class Classification

- $x = \{x_1, x_2\} \in [-5; 5] \times [-5; 5]$ ,  $\hat{y}$ : 3 outputs (classes)

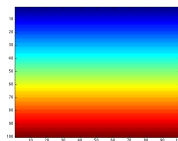
Linear mapping for each class:

$$s_k = w_k^T x + b_k$$

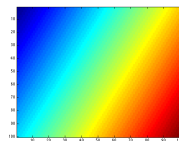
$$w_1 = [1; 1], b_1 = -2$$



$$w_2 = [0; -1], b_2 = 1$$

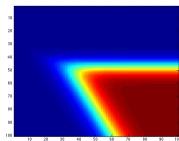
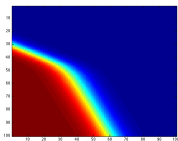
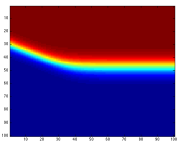


$$w_3 = [1; -0.5], b_3 = 10$$



Soft-max output:

$$P(k|x, W)$$

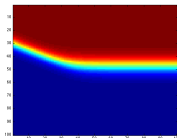


## 2d Toy Example for Multi-Class Classification

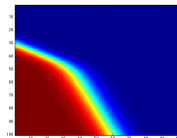
- $x = \{x_1, x_2\} \in [-5; 5] \times [-5; 5]$ ,  $\hat{y}$ : 3 outputs (classes)

Soft-max output:  
 $P(k|x, W)$

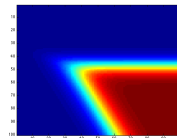
$$w_1 = [1; 1], b_1 = -2$$



$$w_2 = [0; -1], b_2 = 1$$

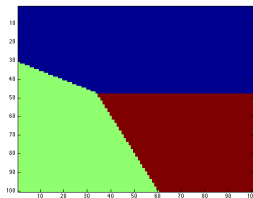


$$w_3 = [1; -0.5], b_3 = 10$$



Class Prediction:

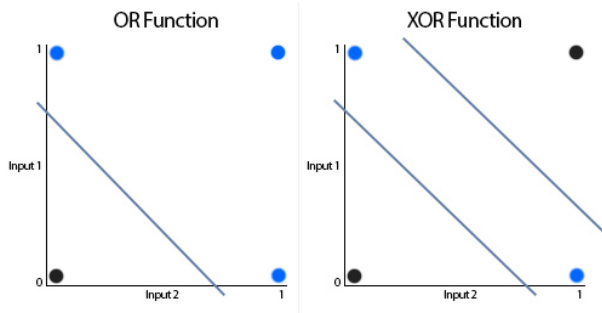
$$k^* = \arg \max_k P(k|x, W)$$



# Beyond Linear Classification

## X-OR Problem

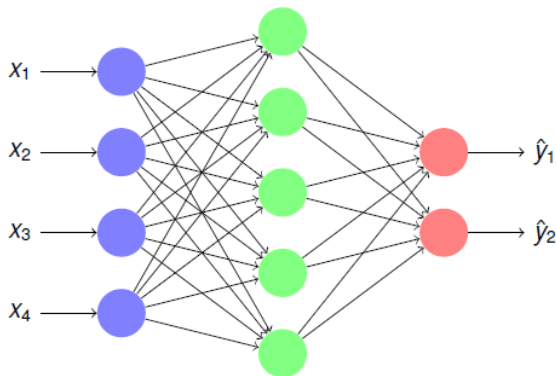
- ▶ Logistic Regression (LR): NN with 1 input layer & 1 output layer
- ▶ LR: limited to linear decision boundaries
- ▶ **X-OR**: NOT 1 and 2 OR NOT 2 AND 1
  - ▶ **X-OR**: Non linear decision function



## Beyond Linear Classification

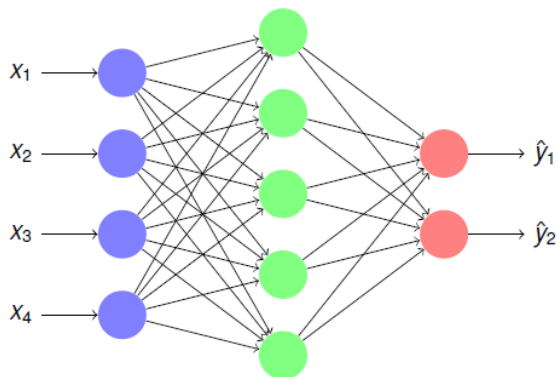
- ▶ LR: limited to linear boundaries
- ▶ **Solution:** add a layer!

- ▶ Input  $x$  in  $\mathbb{R}^m$ , e.g.  $m = 4$
- ▶ Output  $\hat{y}$  in  $\mathbb{R}^K$  ( $K$  # classes), e.g.  $K = 2$
- ▶ **Hidden layer  $h$  in  $\mathbb{R}^L$**



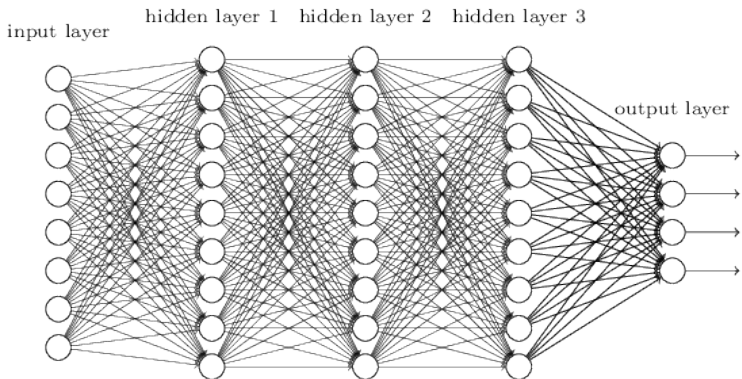
## Multi-Layer Perceptron

- ▶ **Hidden layer  $h$ :**  $x$  projection to a new space  $\mathbb{R}^L$
  - ▶ Neural Net with  $\geq 1$  hidden layer:  
**Multi-Layer Perceptron (MLP)**
  - ▶  $h$ : intermediate representations of  $x$  for classification  $\hat{y}$ :
- ▶  $h = f(xW_1 + b_1)$   
 $f$  non-linear activation,  
 $s = hW_2 + b_2$   
 $\hat{y} = \text{SoftMax}(s)$
  - ▶ **Mapping from  $x$  to  $\hat{y}$ : non-linear boundary!**  
 $\Rightarrow$  **non-linear activation  $f$  crucial!**



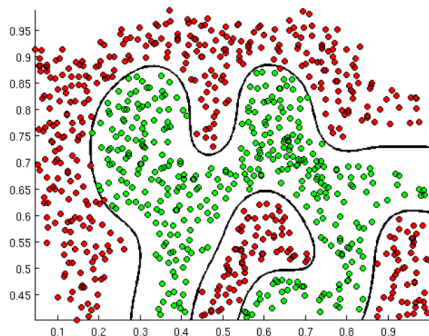
# Deep Neural Networks

- ▶ Adding more hidden layers: Deep Neural Networks (DNN)  $\Rightarrow$  **Basis of Deep Learning**
- ▶ Each layer  $h^l$  projects layer  $h^{l-1}$  into a new space
- ▶ Gradually learning intermediate representations useful for the task



## Conclusion

- ▶ Deep Neural Networks: applicable to classification problems with non-linear decision boundaries



- ▶ Visualize prediction from fixed model parameters
- ▶ Reverse problem: **Supervised Learning**

# Outline

Introduction to Neural nets

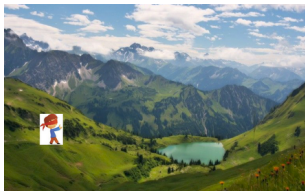
**Training Deep Neural Networks**

Introduction to Statistical Decision Theory



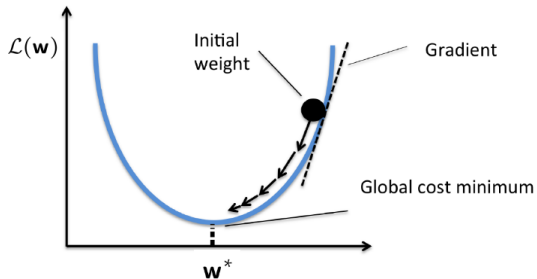
# Training Multi-Layer Perceptron (MLP)

- ▶ Input  $x$ , output  $y$
- ▶ A parametrized ( $w$ ) model  $x \Rightarrow y$ :  $f_w(x_i) = \hat{y}_i$
- ▶ Supervised context:
  - ▶ Training set  $\mathcal{A} = \{(x_i, y_i^*)\}_{i \in \{1, 2, \dots, N\}}$
  - ▶ Loss function  $\ell(\hat{y}_i, y_i^*)$  for each annotated pair  $(x_i, y_i^*)$
  - ▶ Goal: Minimizing average loss  $\mathcal{L}$  over training set:  $\mathcal{L}(w) = \frac{1}{N} \sum_{i=1}^N \ell(\hat{y}_i, y_i^*)$
- ▶ Assumptions: parameters  $w \in \mathbb{R}^d$  continuous,  $\mathcal{L}$  differentiable
- ▶ Gradient  $\nabla_w = \frac{\partial \mathcal{L}}{\partial w}$ : steepest direction to decrease loss  $\mathcal{L}(w)$



# MLP Training

- ▶ Gradient descent algorithm:
  - ▶ Initialize parameters  $w$
  - ▶ Update:  $w^{(t+1)} = w^{(t)} - \eta \frac{\partial \mathcal{L}}{\partial w}$
  - ▶ Until convergence, e.g.  $\|\nabla_w\|^2 \approx 0$



# Supervised Learning: Multi-Class Classification

- ▶ Logistic Regression for multi-class classification

- ▶  $\mathbf{s}_i = \mathbf{x}_i \mathbf{W} + \mathbf{b}$

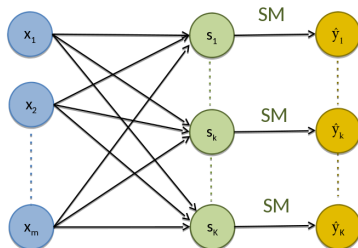
- ▶ Soft-Max (SM):  $\hat{y}_k \sim P(k/x_i, \mathbf{W}, \mathbf{b}) = \frac{e^{\mathbf{s}_k}}{\sum_{k'=1}^K e^{\mathbf{s}_{k'}}$

- ▶ Supervised loss function:  $\mathcal{L}(\mathbf{W}, \mathbf{b}) = \frac{1}{N} \sum_{i=1}^N \ell(\hat{y}_i, y_i^*)$

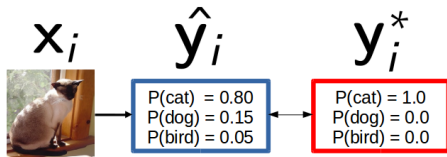
1.  $y \in \{1; 2; \dots; K\}$

2.  $\hat{y}_i = \arg \max_k P(k/x_i, \mathbf{W}, \mathbf{b})$

3.  $\ell_{0/1}(\hat{y}_i, y_i^*) = \begin{cases} 1 & \text{if } \hat{y}_i \neq y_i^* \\ 0 & \text{otherwise} \end{cases} : \mathbf{0/1 loss}$



## Logistic Regression Training Formulation



- ▶ Input  $x_i$ , ground truth output supervision  $y_i^*$
- ▶ One hot-encoding for  $y_i^*$ :

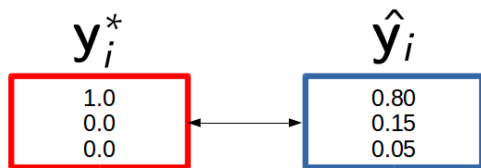
$$y_{c,i}^* = \begin{cases} 1 & \text{if } c \text{ is the ground truth class for } x_i \\ 0 & \text{otherwise} \end{cases}$$

## Logistic Regression Training Formulation

- ▶ Loss function: multi-class Cross-Entropy (CE)  $\ell_{CE}$
- ▶  $\ell_{CE}$ : Kullback-Leiber divergence between  $y_i^*$  and  $\hat{y}_i$

$$\ell_{CE}(\hat{y}_i, y_i^*) = KL(y_i^*, \hat{y}_i) = - \sum_{c=1}^K y_{c,i}^* \log(\hat{y}_{c,i}) = -\log(\hat{y}_{c^*,i})$$

- ▶  $\triangleleft$  KL asymmetric:  $KL(\hat{y}_i, y_i^*) \neq KL(y_i^*, \hat{y}_i)$   $\triangleleft$

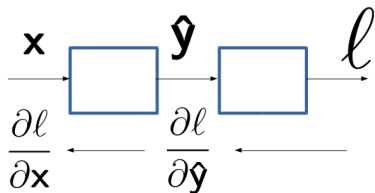


$$KL(y_i^*, \hat{y}_i) = -\log(\hat{y}_{c^*,i}) = -\log(0.8) \approx 0.22$$

## Logistic Regression Training

- ▶  $\mathcal{L}_{CE}(\mathbf{W}, \mathbf{b}) = \frac{1}{N} \sum_{i=1}^N \ell_{CE}(\hat{y}_i, y_i^*) = -\frac{1}{N} \sum_{i=1}^N \log(\hat{y}_{c^*, i})$
- ▶  $\ell_{CE}$  smooth convex upper bound of  $\ell_{0/1}$   
⇒ **gradient descent optimization**
- ▶ Gradient descent:  $\mathbf{W}^{(t+1)} = \mathbf{W}^{(t)} - \eta \frac{\partial \mathcal{L}_{CE}}{\partial \mathbf{W}}$      $(\mathbf{b}^{(t+1)} = \mathbf{b}^{(t)} - \eta \frac{\partial \mathcal{L}_{CE}}{\partial \mathbf{b}})$
- ▶ **MAIN CHALLENGE: computing**  $\frac{\partial \mathcal{L}_{CE}}{\partial \mathbf{W}} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \ell_{CE}}{\partial \mathbf{W}}$  ?  
  
⇒ Key Property: chain rule  $\frac{\partial x}{\partial z} = \frac{\partial x}{\partial y} \frac{\partial y}{\partial z}$   
⇒ **Backpropagation of gradient error!**

## Chain Rule



$$\frac{\partial \ell}{\partial \mathbf{x}} = \frac{\partial \ell}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{x}}$$

- Logistic regression:  $\frac{\partial \ell_{CE}}{\partial \mathbf{W}} = \frac{\partial \ell_{CE}}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial s_i} \frac{\partial s_i}{\partial \mathbf{W}}$

