

Deep Learning for Computer Vision

113-1/Fall 2024

<https://cool.ntu.edu.tw/courses/41702> (NTU COOL)

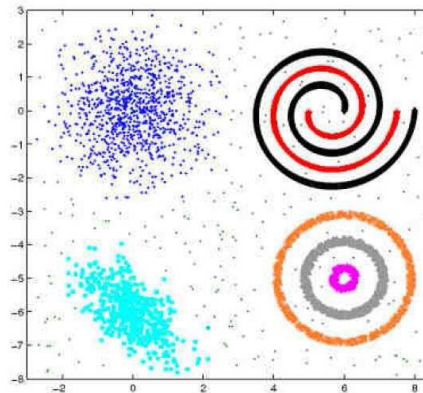
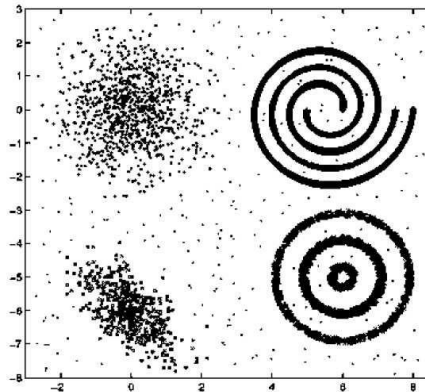
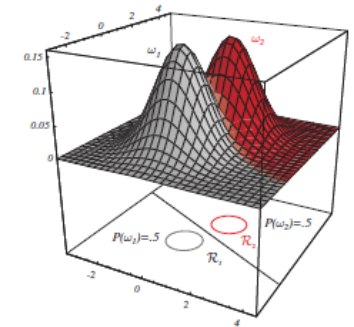
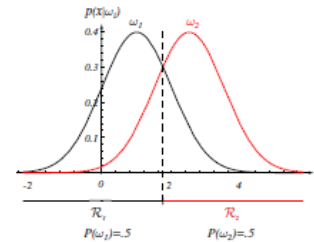
<http://vllab.ee.ntu.edu.tw/dlcv.html> (Public website)

Yu-Chiang Frank Wang 王鈺強, Professor

Dept. Electrical Engineering, National Taiwan University

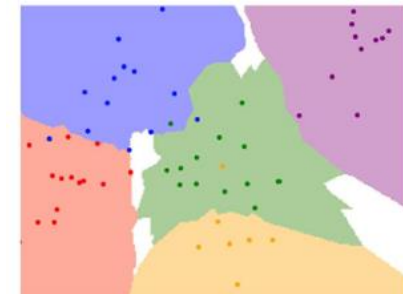
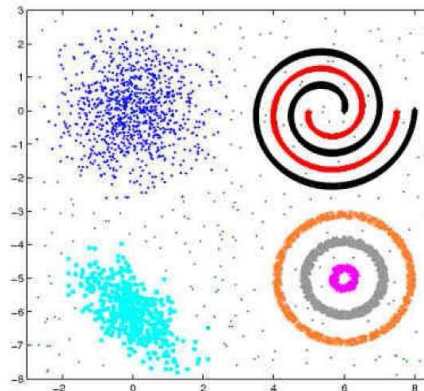
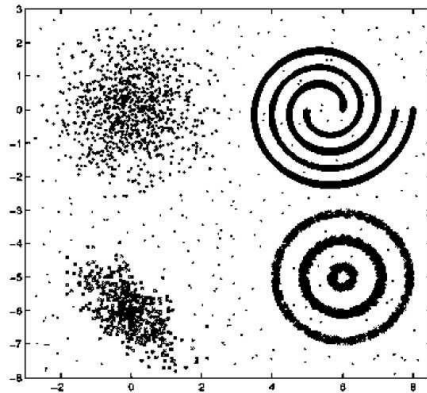
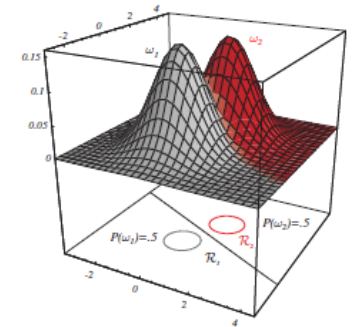
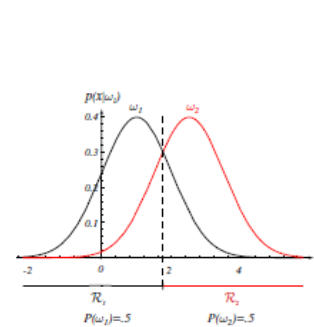
What's to Be Covered in This Lecture...

- Unsupervised vs. Supervised Learning
 - Clustering & Dimension Reduction
 - Training, testing, & validation
 - Linear Classification
 - From Linear Classifier to Neural Nets



What's to Be Covered Today...

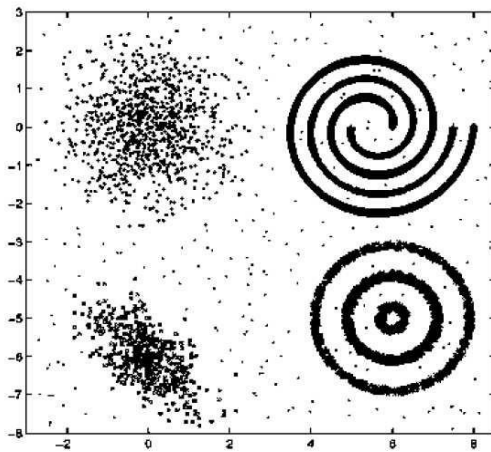
- Unsupervised vs. Supervised Learning
 - Clustering & Dimension Reduction
 - Training, testing, & validation
 - Linear Classification
 - From Linear Classifier to Neural Nets



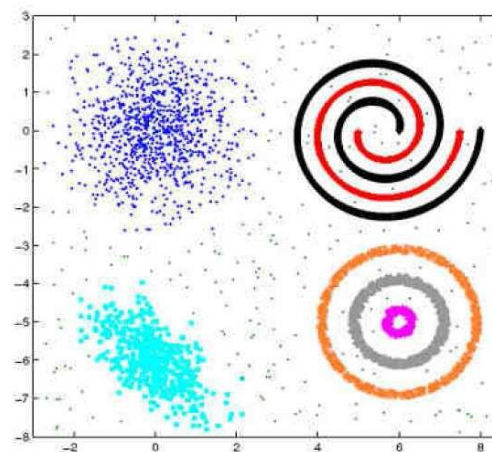
Clustering



- Clustering is an unsupervised algorithm.
 - Given:
 - a set of N unlabeled instances $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$; # of clusters K
 - Goal: group the samples into K partitions
- Remarks:
 - High within-cluster (intra-cluster) similarity
 - Low between-cluster (inter-cluster) similarity
 - But...how to determine a proper similarity measure?



(a) Input data

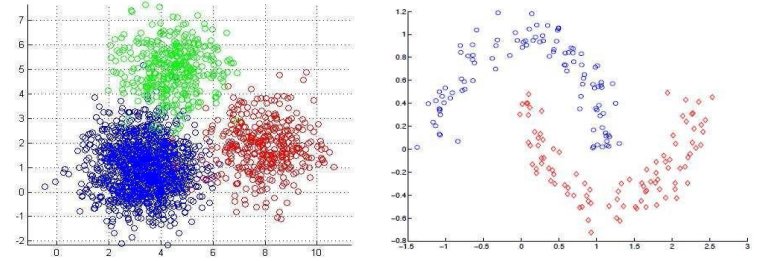


(b) Desired clustering

Similarity is NOT Always Objective...



Clustering (cont'd)



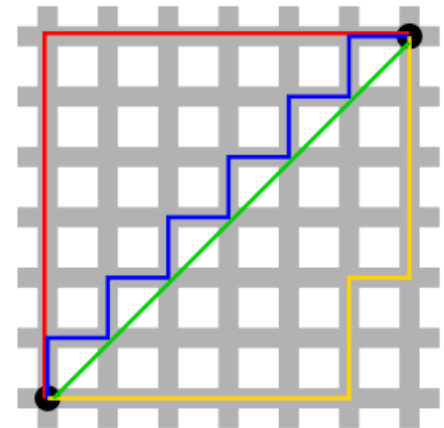
- Similarity:

- A key component/measure to perform data clustering
- **Inversely proportional** to distance
- Example distance metrics:

- Euclidean distance (L2 norm): $d(x, z) = \|x - z\|_2 = \sqrt{\sum_{i=1}^D (x_i - z_i)^2}$

- Manhattan distance (L1 norm): $d(x, z) = \|x - z\|_1 = \sum_{i=1}^D |x_i - z_i|$

- Note that p -norm of x is denoted as:

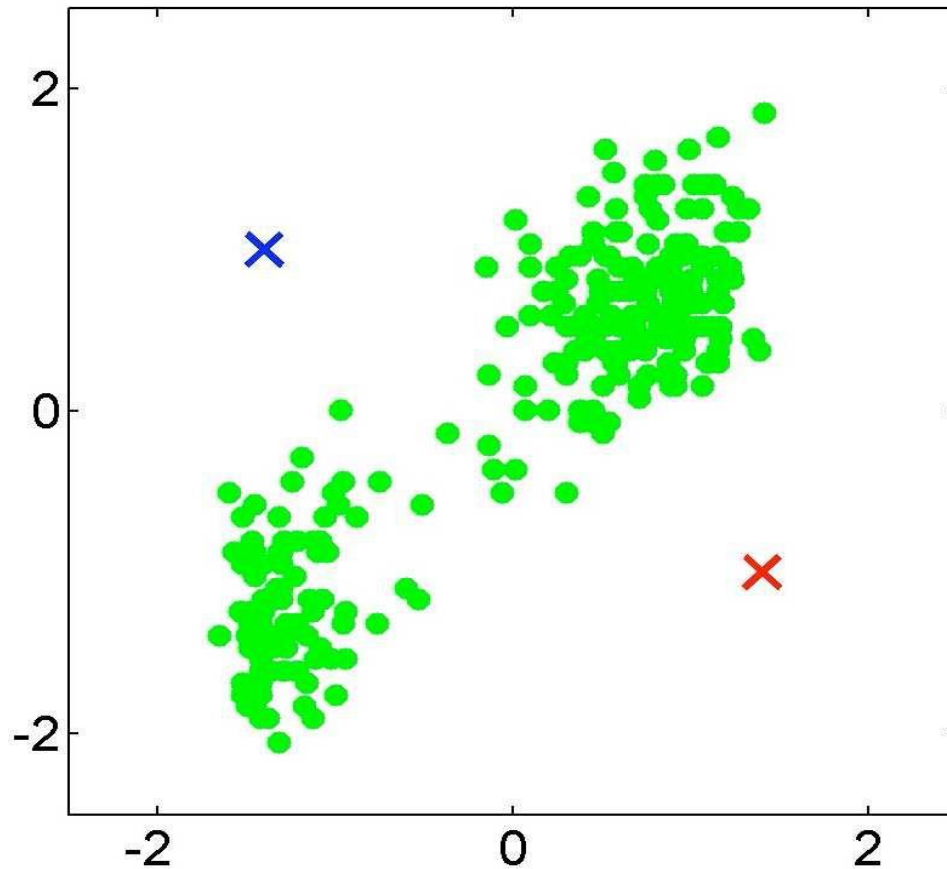


K-Means Clustering

- **Input:** N examples $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ ($\mathbf{x}_n \in \mathbb{R}^D$); number of partitions K
- **Initialize:** K cluster centers μ_1, \dots, μ_K . Several initialization options:
 - Randomly initialize μ_1, \dots, μ_K anywhere in \mathbb{R}^D
 - Or, simply choose any K examples as the cluster centers
- **Iterate:**
 - Assign each of example \mathbf{x}_n to its closest cluster center
 - Recompute the new cluster centers μ_k (mean/centroid of the set C_k)
 - Repeat while not converge
- **Possible convergence criteria:**
 - Cluster centers do not change anymore
 - Max. number of iterations reached
- **Output:**
 - K clusters (with centers/means of each cluster)

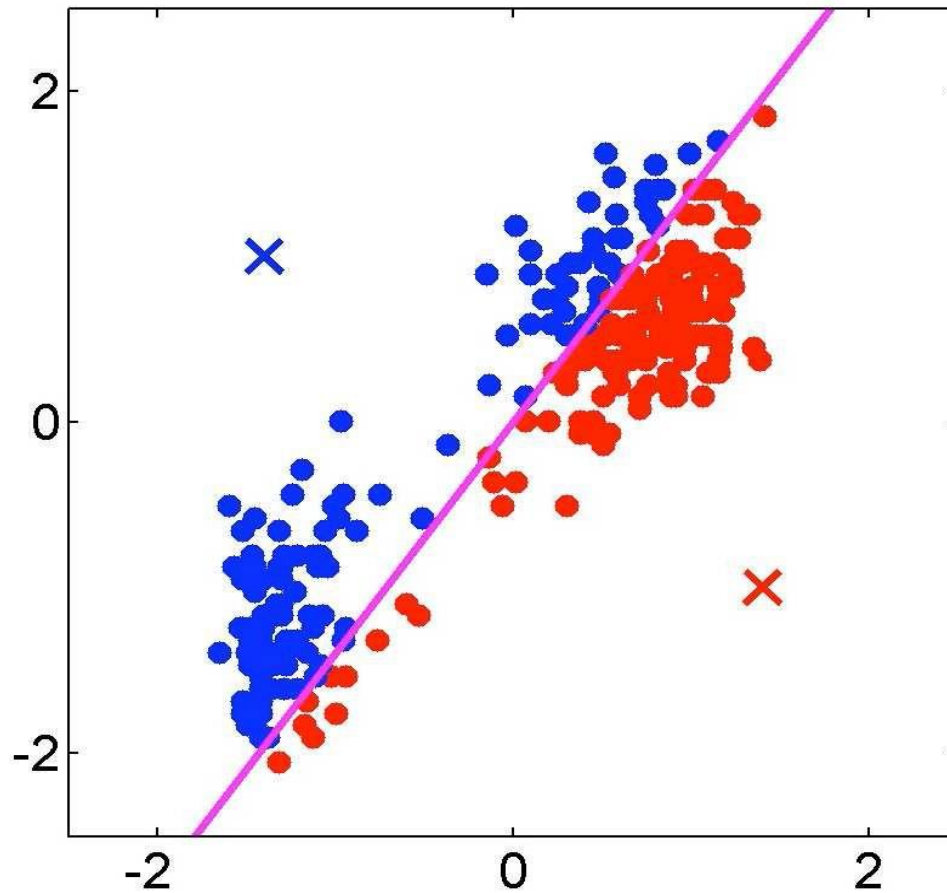
K-Means Clustering

- Example ($K = 2$): Initialization, iteration #1: pick cluster centers



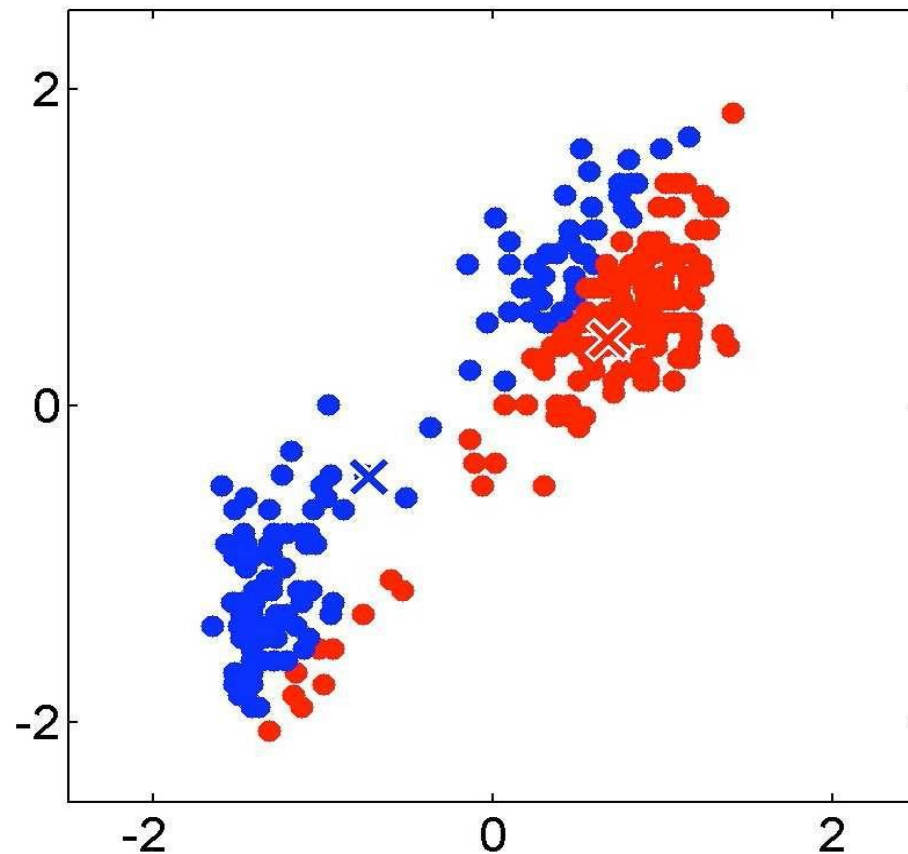
K-Means Clustering

- Example ($K = 2$): iteration #1-2, assign data to each cluster



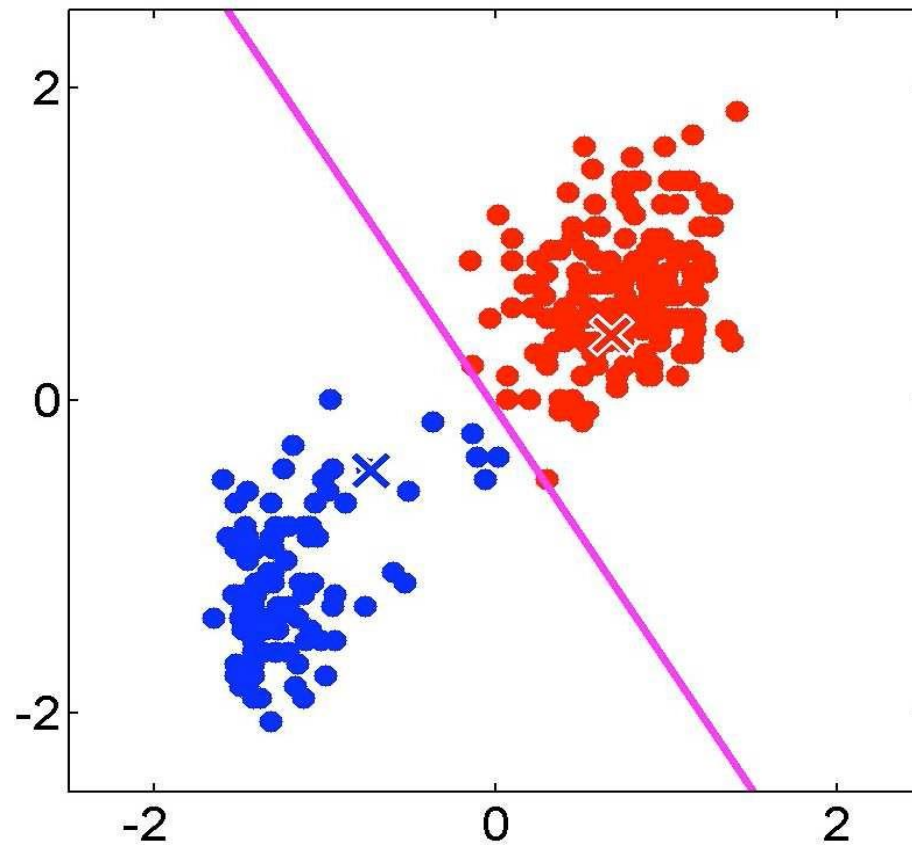
K-Means Clustering

- Example ($K = 2$): iteration #2-1, update cluster centers



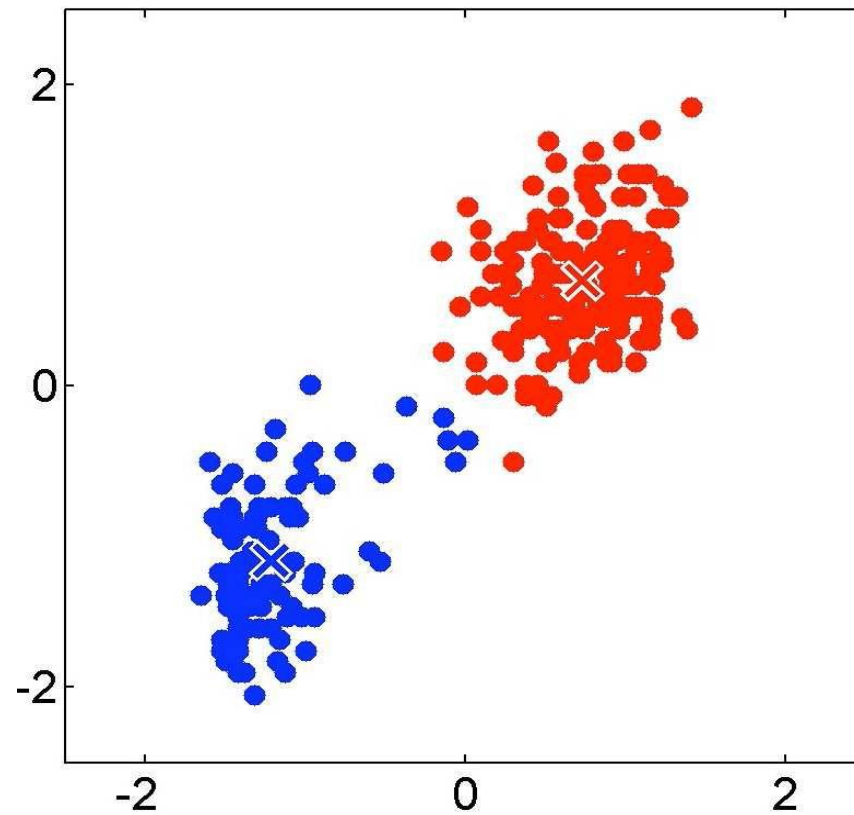
K-Means Clustering

- Example ($K = 2$): iteration #2, assign data to each cluster



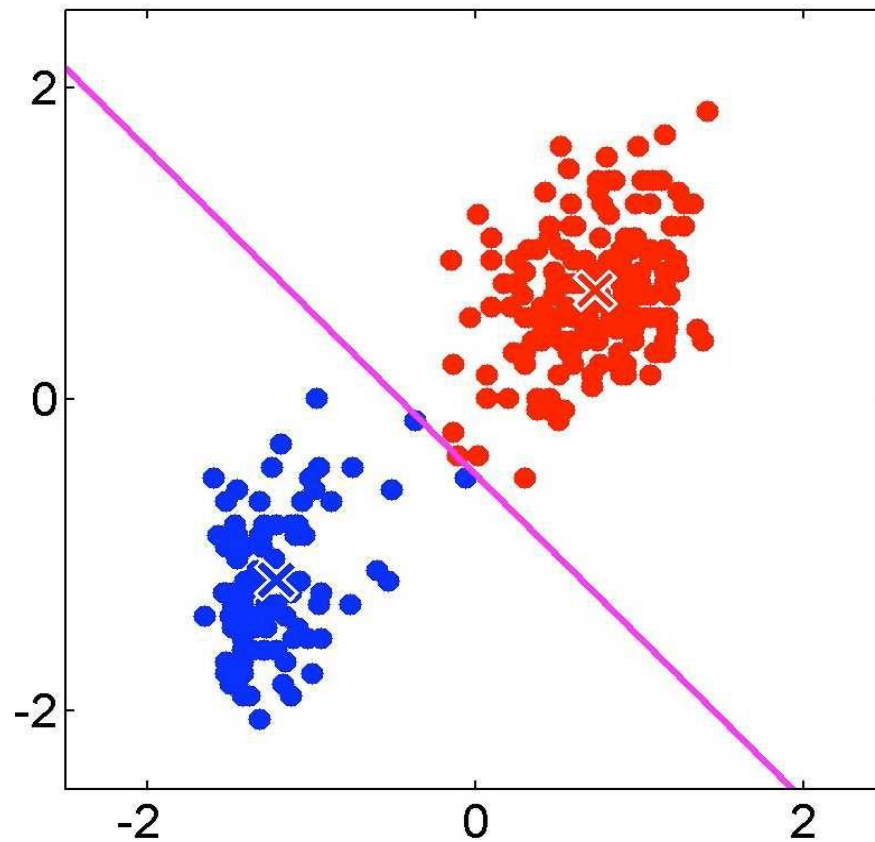
K-Means Clustering

- Example ($K = 2$): iteration #3-1



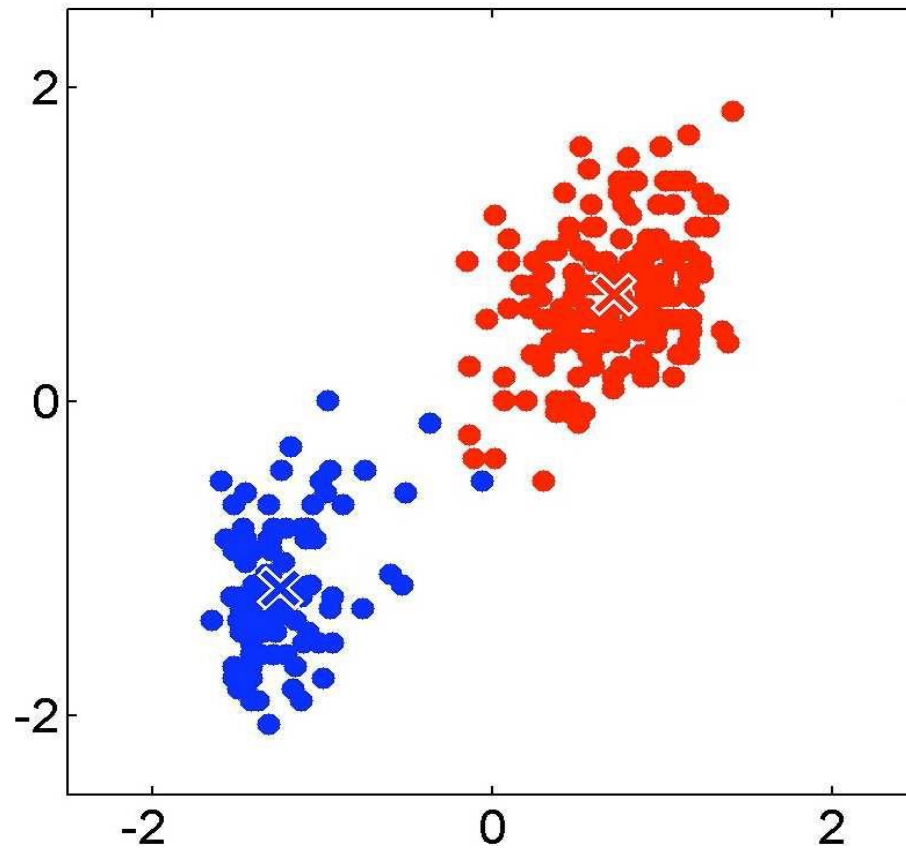
K-Means Clustering

- Example (K = 2): iteration #3-2



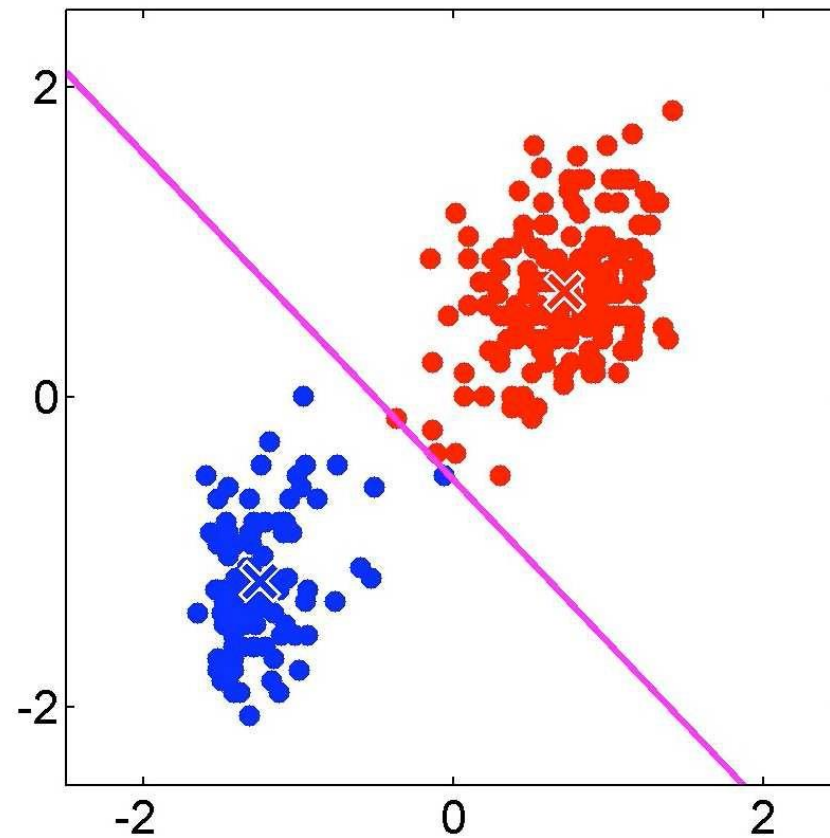
K-Means Clustering

- Example ($K = 2$): iteration #4-1



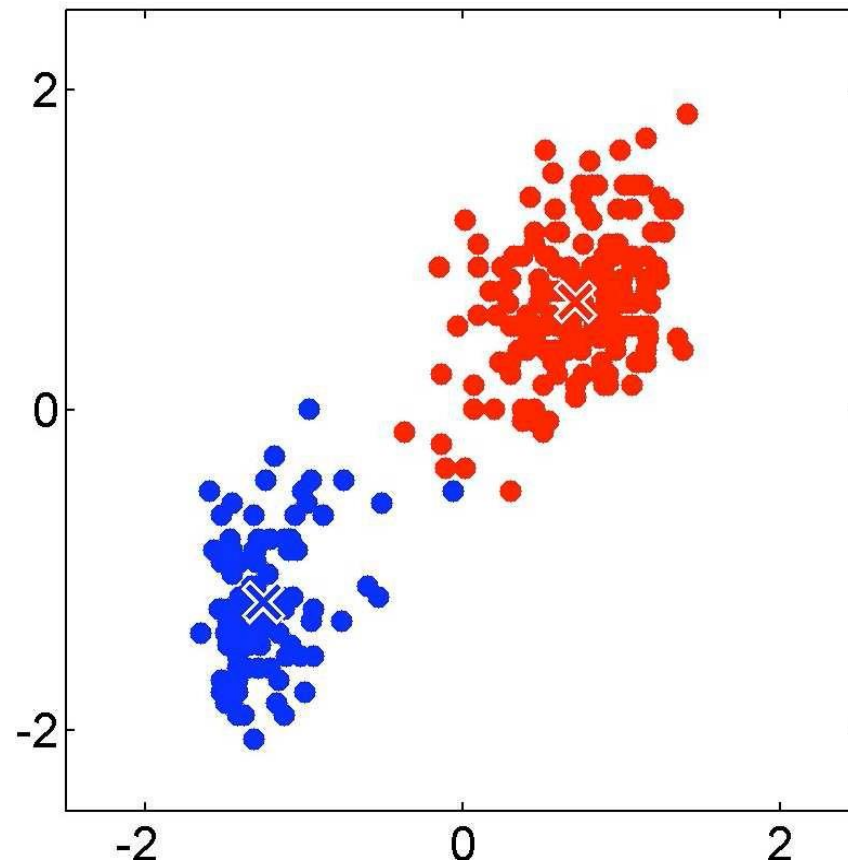
K-Means Clustering

- Example (K = 2): iteration #4-2



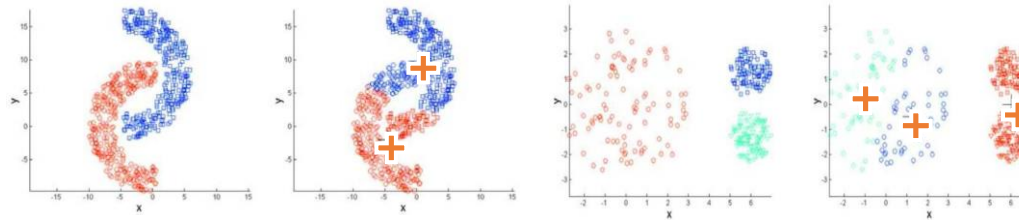
K-Means Clustering

- Example ($K = 2$): iteration #5, cluster means are not changed.



K-Means Clustering (cont'd)

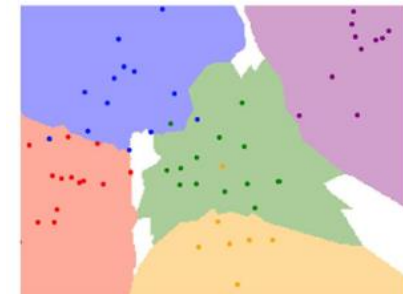
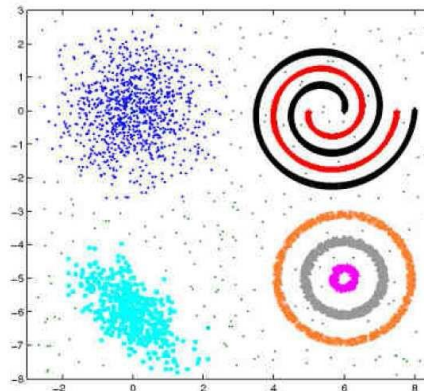
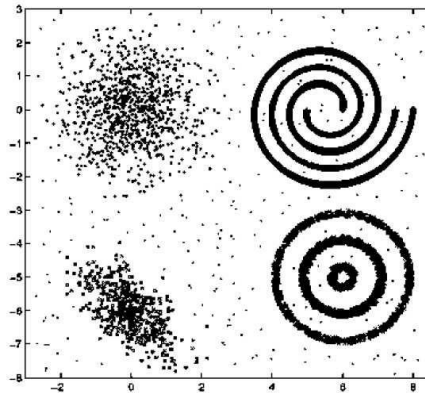
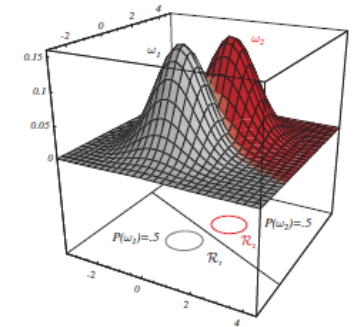
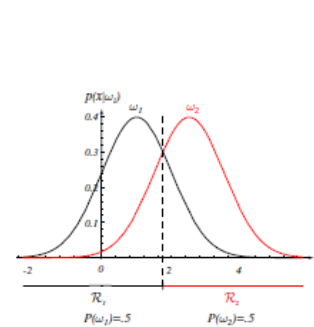
- Easy to implement, but...
 - Preferable for round shaped clusters with similar sizes



- Limitations
 - Sensitive to initialization →
 - Sensitive to outliers →
 - Hard assignment only →
- Remarks
 - Expectation-maximization (EM) algorithm
 - Speed-up possible by hierarchical clustering (e.g., $100 = 10^2$ clusters)

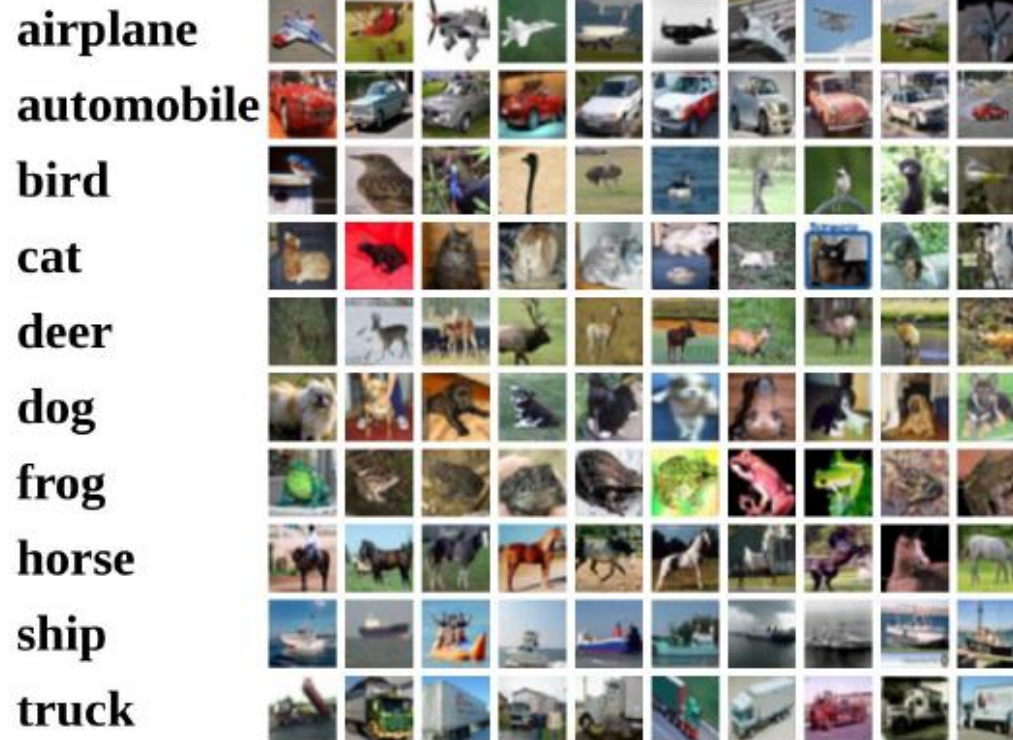
What's to Be Covered in This Lecture...

- Unsupervised vs. Supervised Learning
 - Clustering & Dimension Reduction
 - Training, testing, & validation
- Linear Classification
 - From Linear Classifier to Neural Nets



Linear Classification

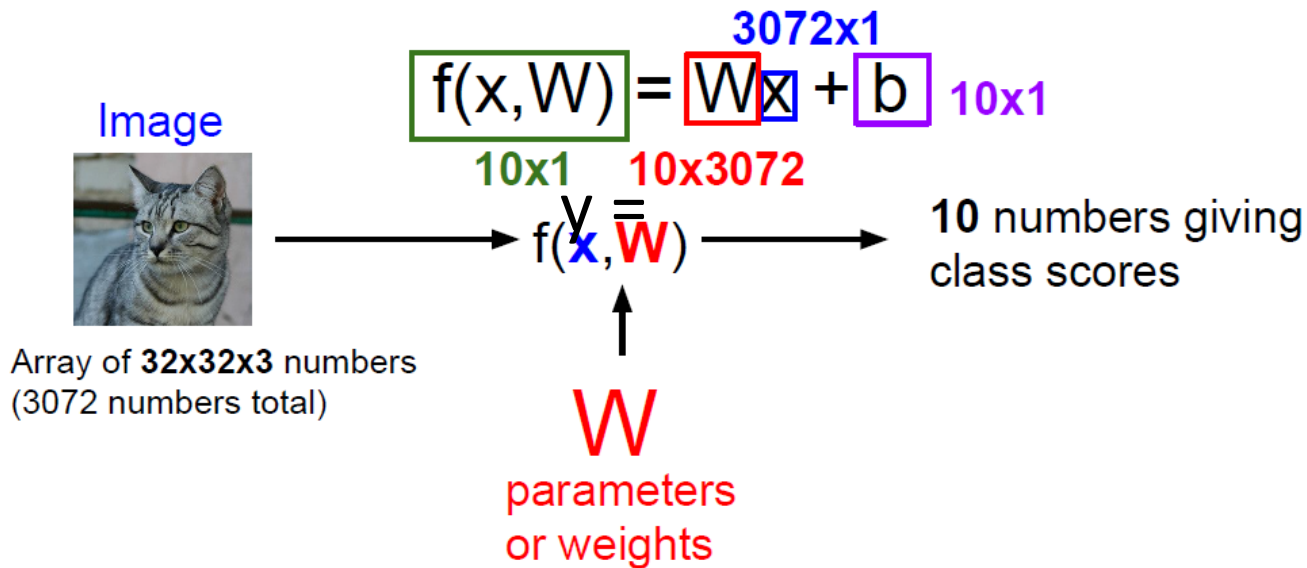
- Linear Classifier
 - Consider that we have 10 object categories of interest
 - E.g., CIFAR10 with 50K training & 10K test images of 10 categories.
And, each image is of size 32 x 32 x 3 pixels.



Linear Classification (cont'd)

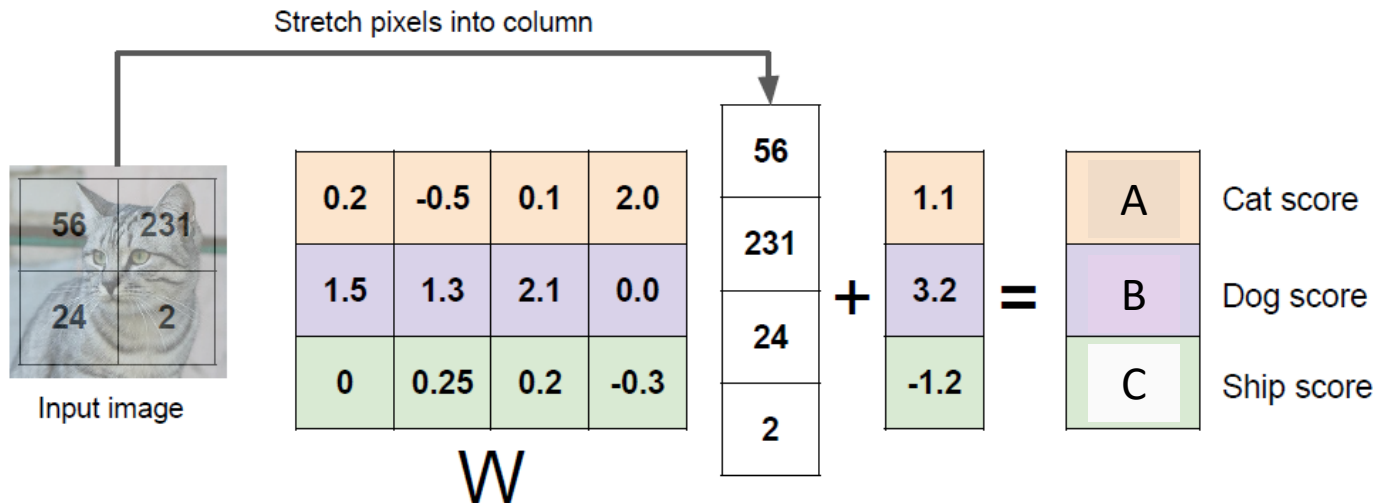


- Linear Classifier
 - Let's take the input image as x , and the linear classifier as W . We need $y = Wx + b$ as a 10-dimensional output vector, indicating the score for each class.



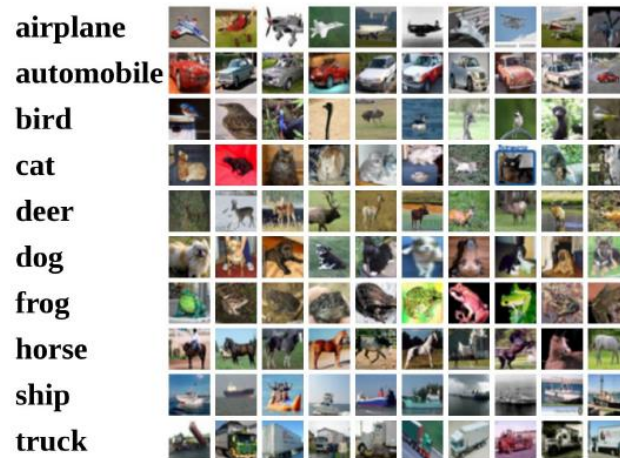
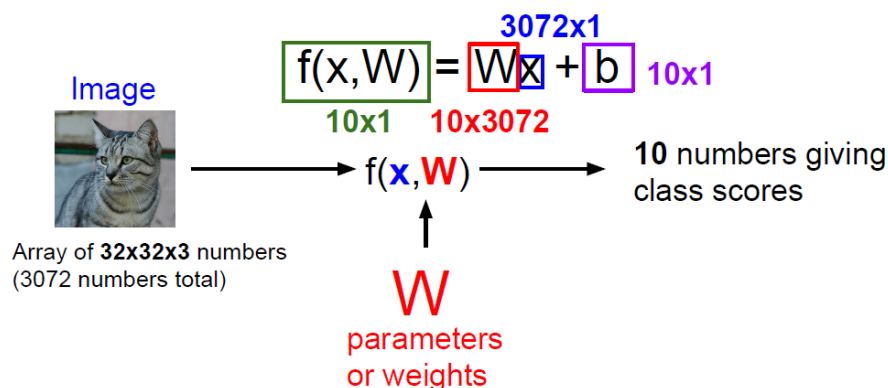
Linear Classification (cont'd)

- Linear Classifier
 - For example, an image with 2 x 2 pixels & 3 classes of interest we need to learn a linear classifier \mathbf{W} (plus a bias \mathbf{b}), so that desirable outputs $\mathbf{y} = \mathbf{W}\mathbf{x} + \mathbf{b}$ can be expected.



Final Remarks

- Interpreting \mathbf{W} in $\mathbf{y} = \mathbf{W}\mathbf{x} + \mathbf{b}$
 - Once \mathbf{W} is learned for predicting \mathbf{y} , each row in \mathbf{W} can be viewed as an **exemplar** of the corresponding class.
 - Recall that, $\mathbf{W}\mathbf{x}$ basically performs **inner product** (or **correlation**) between the input \mathbf{x} and the exemplar of each class -> **similarity** matters!
 - Any potential problem or limitation?



Loss Function (or Cost/Objective Function)

- **Loss is a function of model parameter W**

- Tells us how **good/bad** our learned model W in $y = Wx + b$ is.
If loss function, the lower, the better!
- Given a labeled dataset $\{(x_i, y_i)\}_{i=1}^N$
where x and y indicate the input instance and its label, respectively,

Loss of a single input instance is denoted as $L_i(f(x_i, W), y_i)$,

and that for the entire dataset is the sum or average of per-instance losses:

$$L = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i) .$$

- In practice, calculating full sum for L is expensive.
 - Approximate sum using a **minibatch** of instances (e.g., 32, 64, 128 samples, etc.)

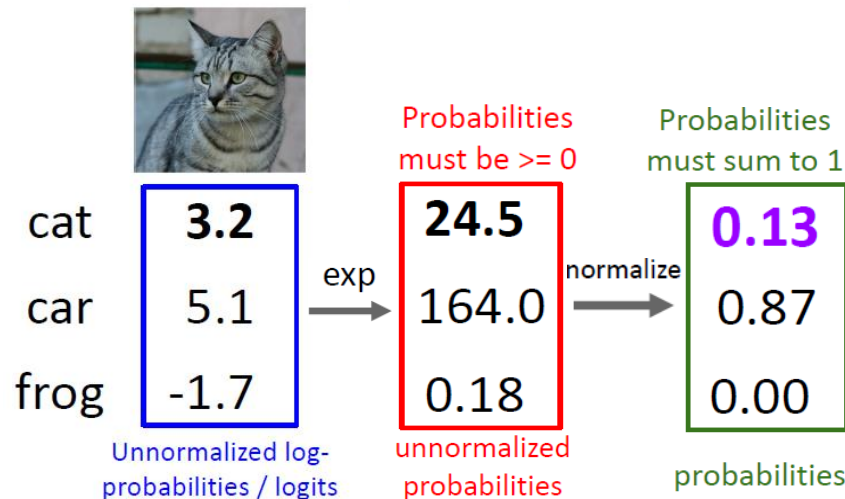
Loss Function (cont'd)

- **Cross-Entropy Loss (Multinomial Logistic Regression)**

- Interpret the classifier scores as **probabilities**
- **Softmax function:**

$$P(Y = k | X = x_i) = \frac{\exp(s_k)}{\sum_j \exp(s_j)} \quad \text{with } s = f(x_i; W) \text{ as the classifier output for input } \mathbf{x}_i$$

- See example below:



$$L_i = -\log P(Y = y_i | X = x_i)$$

$$L_{\text{cat}} = -\log(0.13) = 2.04$$

What about this L?

What are its possibly **min/max** value??

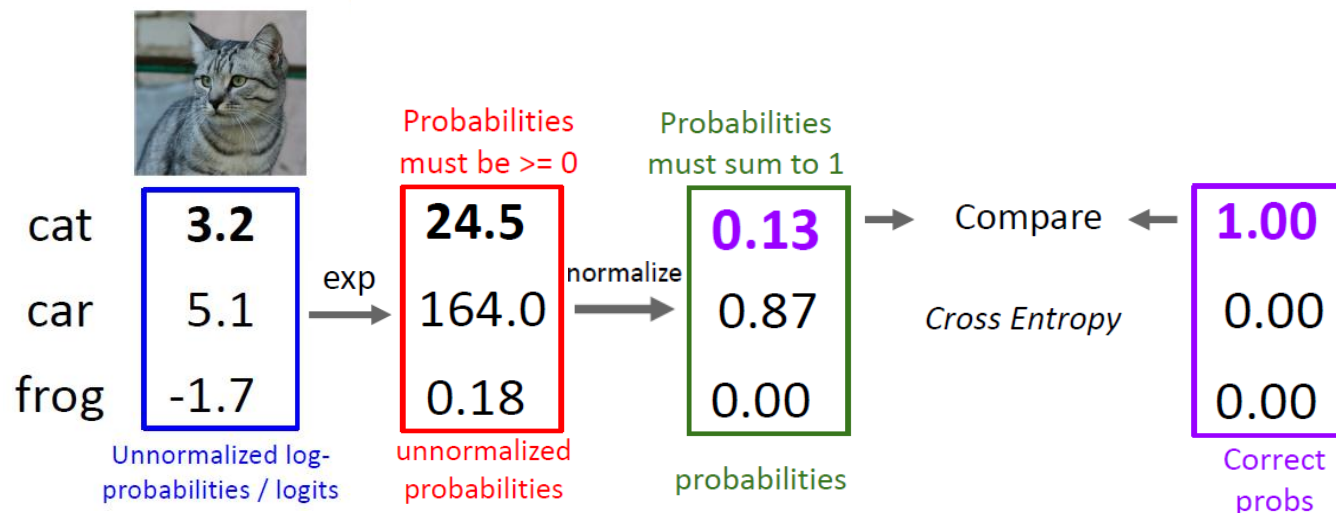
- **Cross-Entropy Loss (cont'd)**

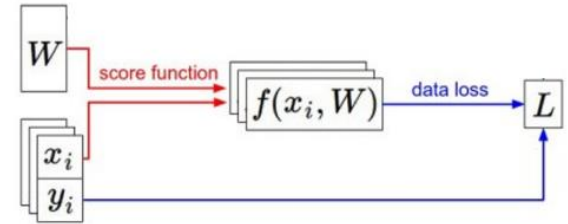
- **Softmax function:**

$$P(Y = k | X = x_i) = \frac{\exp(s_k)}{\sum_j \exp(s_j)} \quad \text{with } s = f(x_i; W) \text{ as the classifier output for input } \mathbf{x}_i$$

→ $L_i = -\log P(Y = y_i | X = x_i)$ or $L_i = -\log \left(\frac{\exp(s_{y_i})}{\sum_j \exp(s_j)} \right)$

- **(Binary) Cross Entropy Loss (or L_{BCE} ; see example below):**

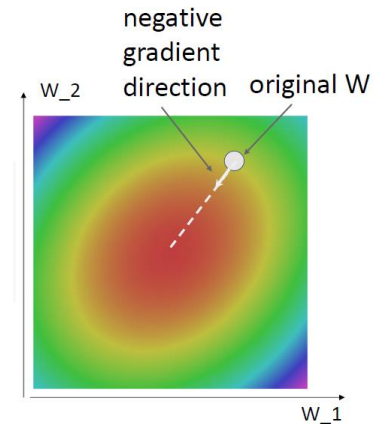




- **Learning W from L_{BCE}**

- Computing **gradients**:
Following the slope to reach the (hopefully global) minimum for W .
- **Gradient Descent** via numeric or analytic gradients:
 - Iteratively step in the direction of the **negative gradient** & search for W
 - Hyperparameters: weight initialization, # of steps, learning rate, etc.

```
# Vanilla gradient descent
w = initialize_weights()
for t in range(num_steps):
    dw = compute_gradient(loss_fn, data, w)
    w -= learning_rate * dw
```



- **Stochastic Gradient Descent**

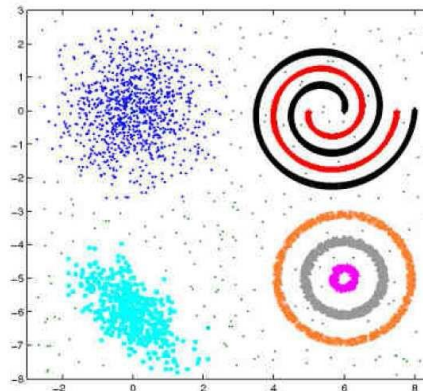
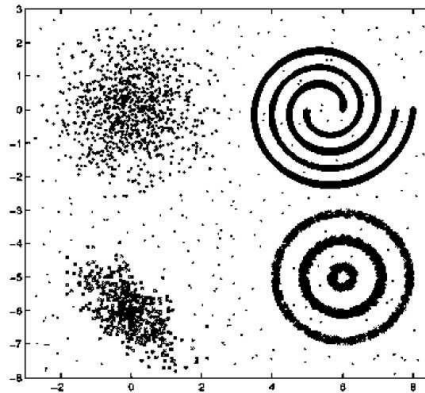
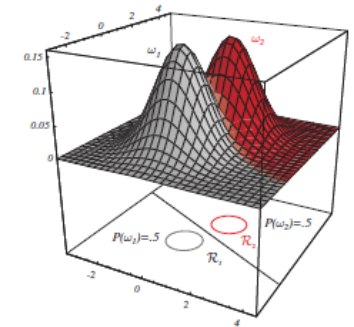
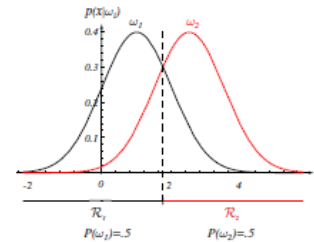
- Full sum in L is **expensive when large N**
- *Approximate* sum using a minibatch of instances (e.g., 32, 64, 128, etc.)
- Additional hyperparameters of batch size and data sampling

```
# Stochastic gradient descent
w = initialize_weights()
for t in range(num_steps):
    minibatch = sample_data(data, batch_size)
    dw = compute_gradient(loss_fn, minibatch, w)
    w -= learning_rate * dw
```



What's to Be Covered in This Lecture...

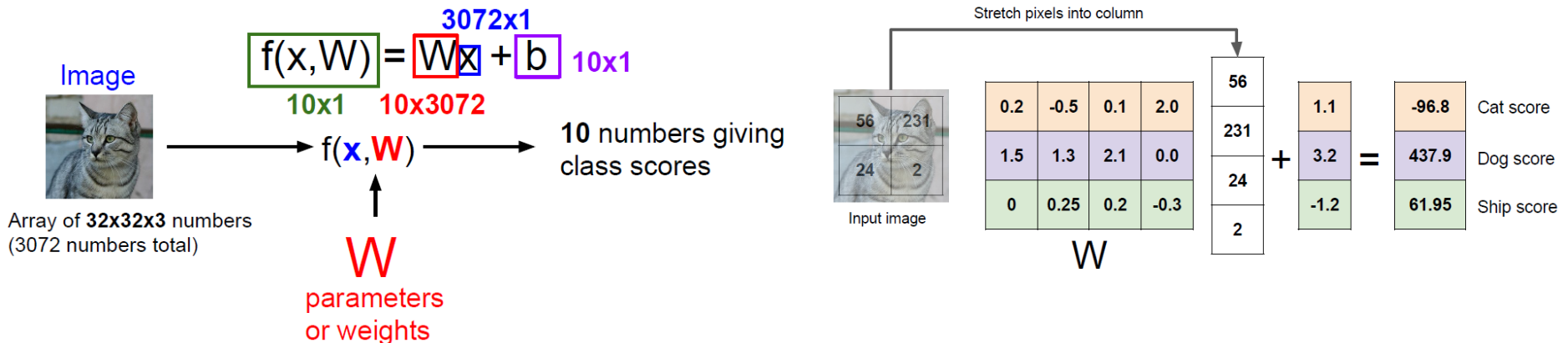
- Unsupervised vs. Supervised Learning
 - Clustering & Dimension Reduction
 - Training, testing, & validation
 - Linear Classification
 - From Linear Classifier to Neural Nets



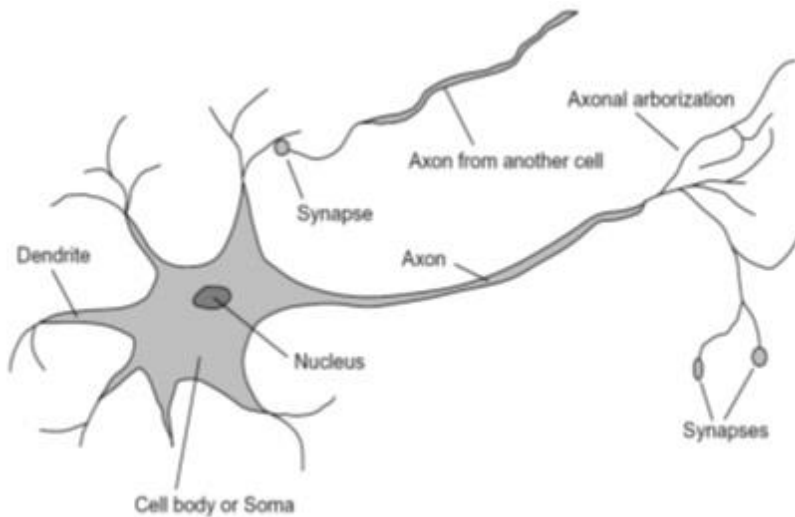
Recap: Linear Classification

- Linear Classifier

- Let's take the input image as x , and the linear classifier as W . We need $y = Wx + b$ as a 10-dimensional output vector, indicating the score for each class.
- For example, an image with 2 x 2 pixels & 3 classes of interest we need to learn a linear classifier W (plus a bias b), so that desirable outputs $y = Wx + b$ can be expected.

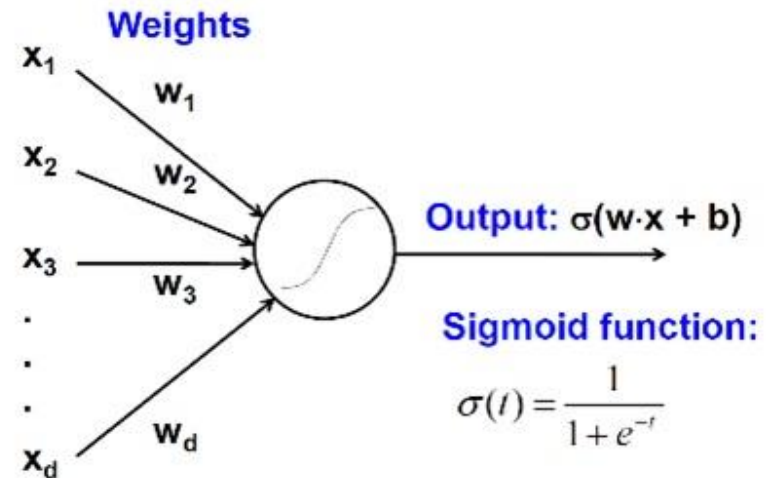


Biological neuron and Perceptrons



A biological neuron

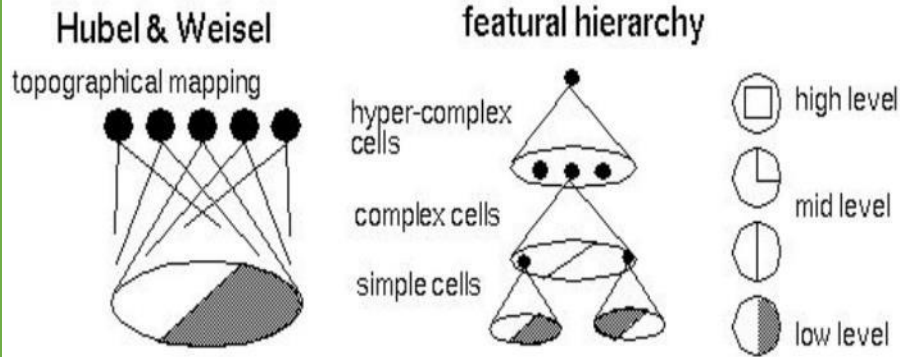
Input



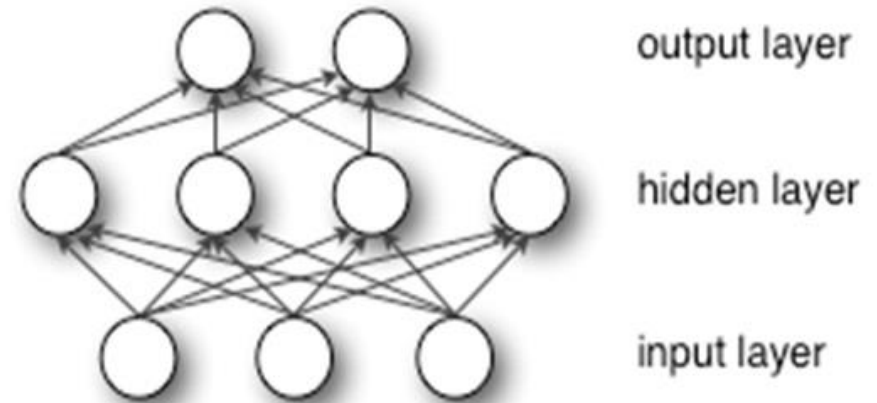
An artificial neuron (Perceptron)
- a linear classifier



Hubel/Wiesel Architecture and Multi-layer Neural Network



Hubel and Wiesel's architecture

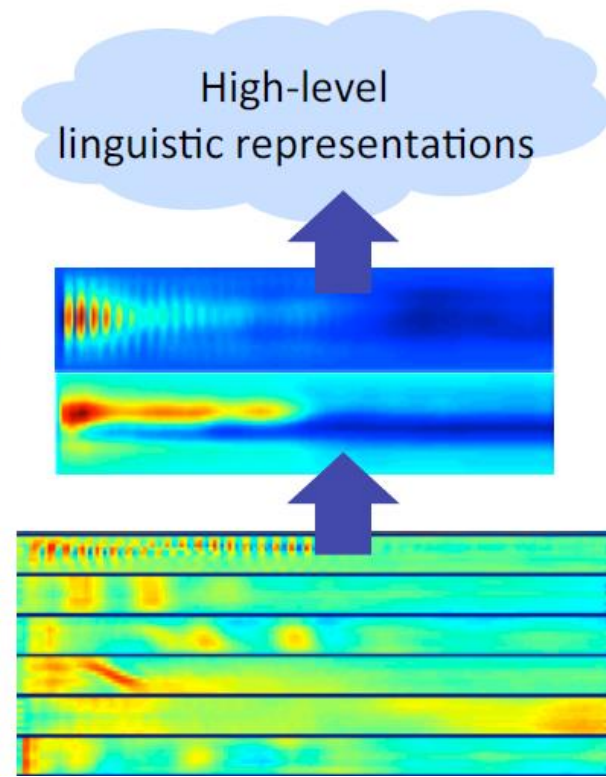
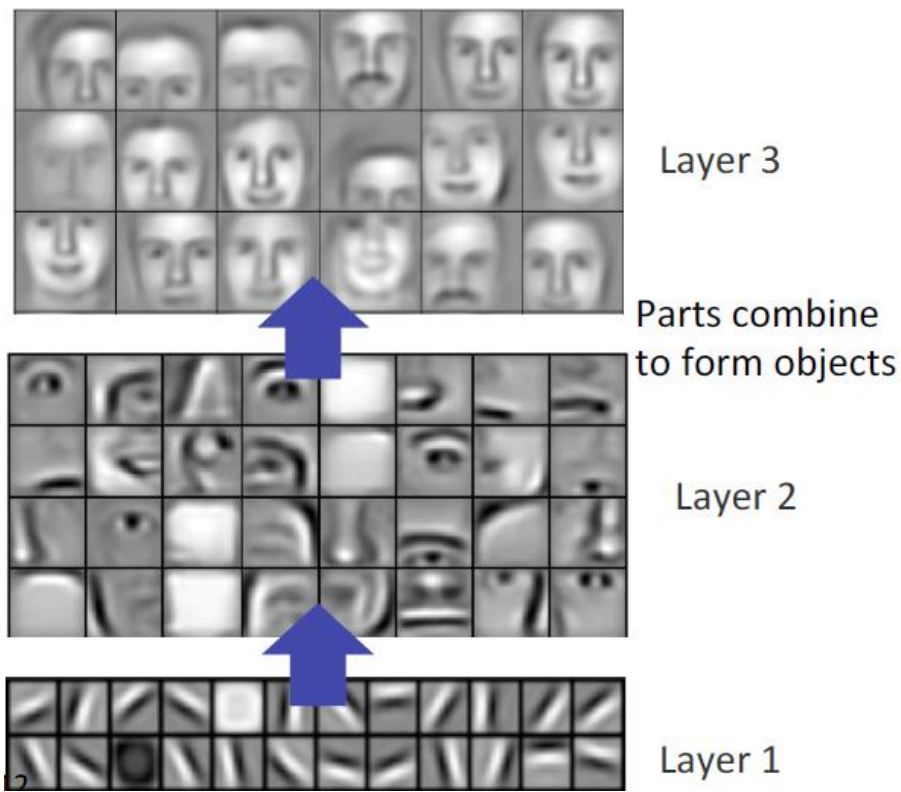


Multi-Layer Perception or Neural Network
- A *non-linear* classifier

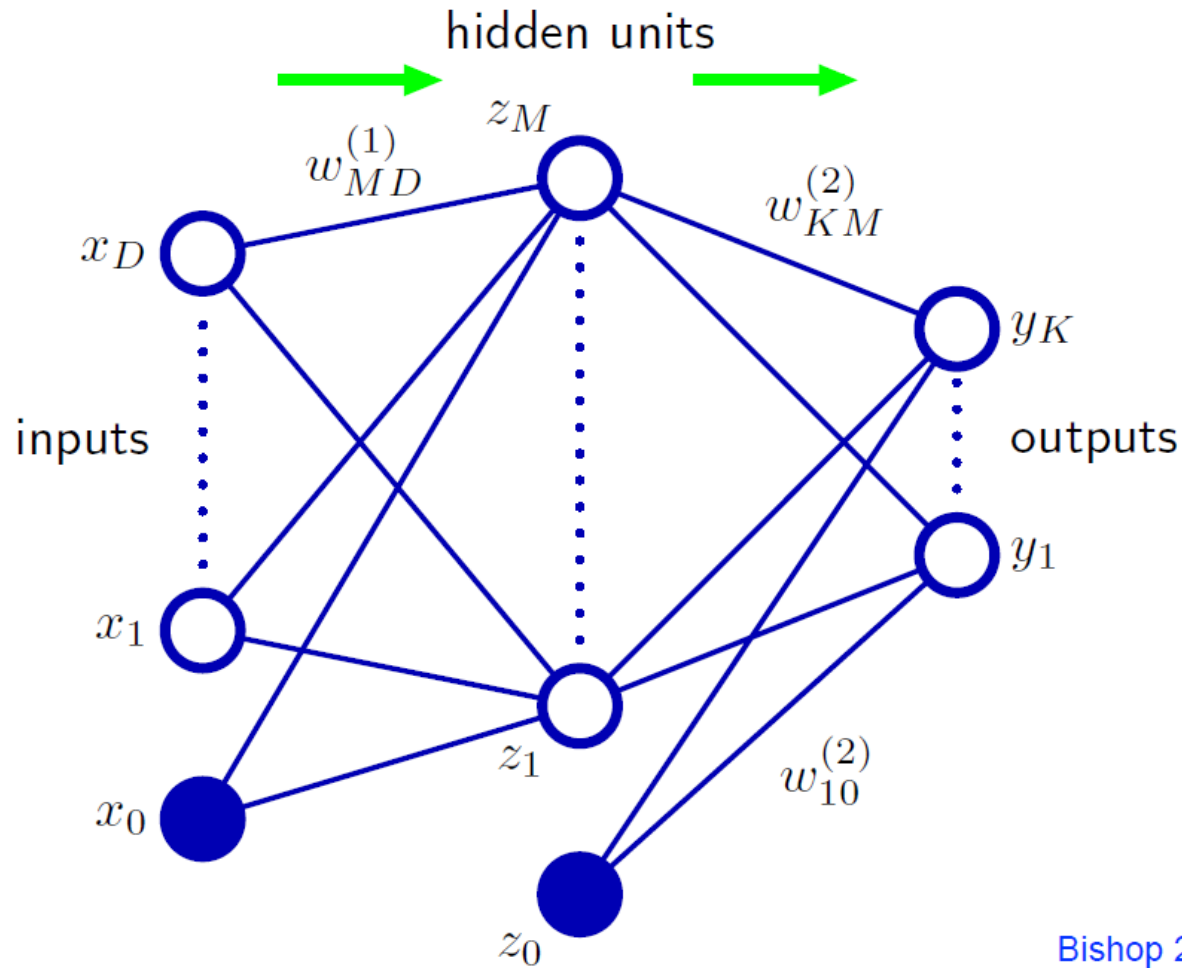


Hierarchical Representation Learning

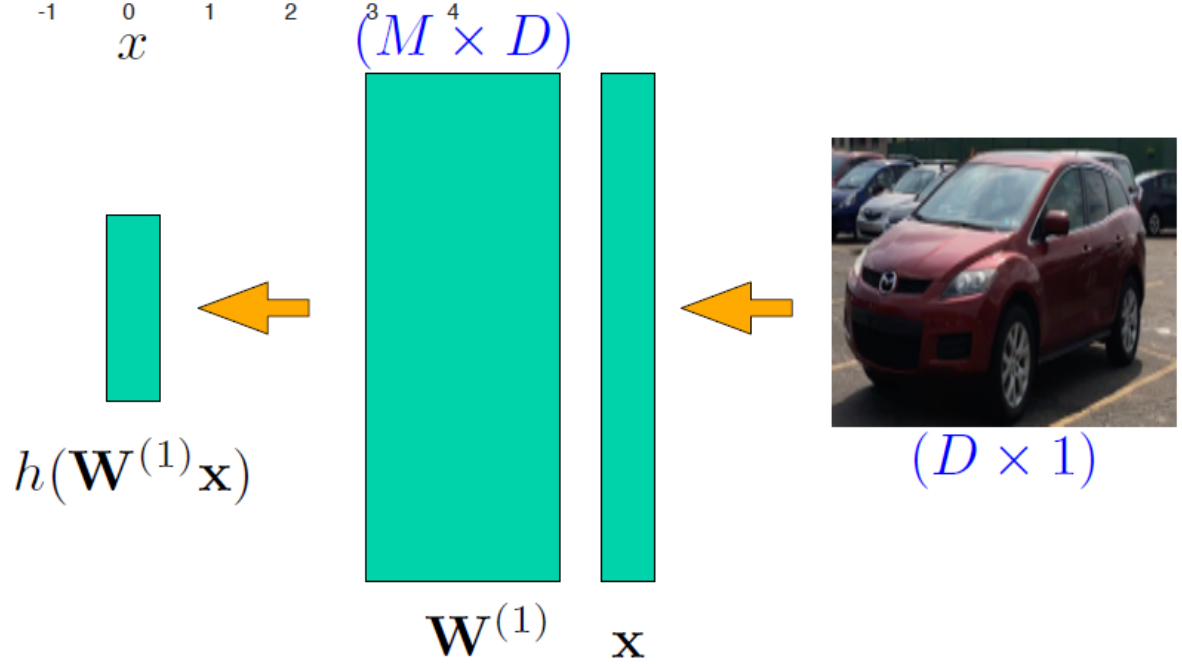
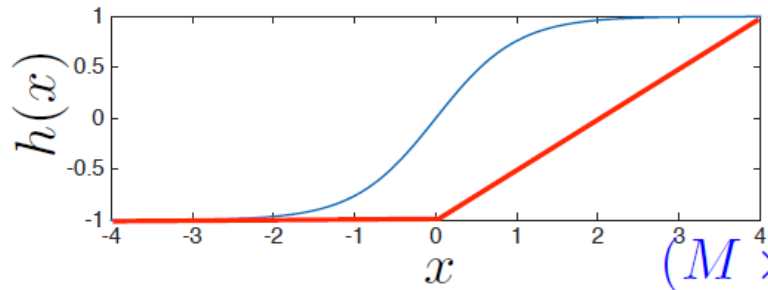
- Successive model layers learn deeper intermediate representations.



Multi-Layer Perceptron: A Nonlinear Classifier (cont'd)

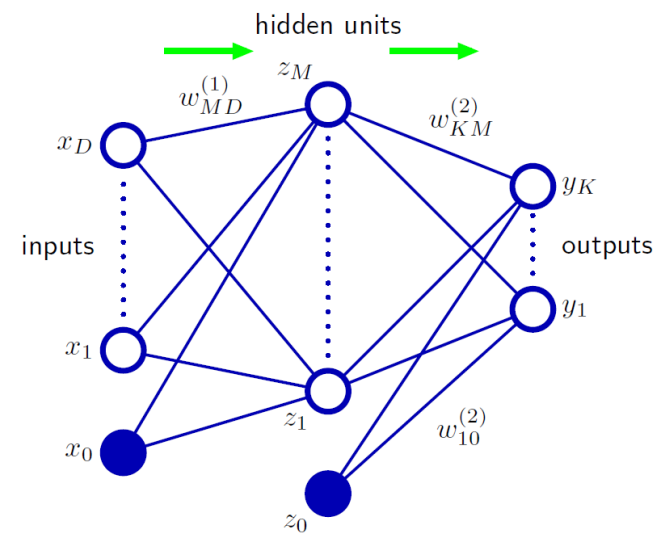


Multi-Layer Perceptron: A Nonlinear Classifier



Layer 1 in MLP

$$\mathbf{z} = \begin{bmatrix} z_1 \\ \vdots \\ z_M \end{bmatrix} \leftarrow \begin{bmatrix} h[\mathbf{x}^T \mathbf{w}_1^{(1)}] \\ \vdots \\ h[\mathbf{x}^T \mathbf{w}_M^{(1)}] \end{bmatrix}$$

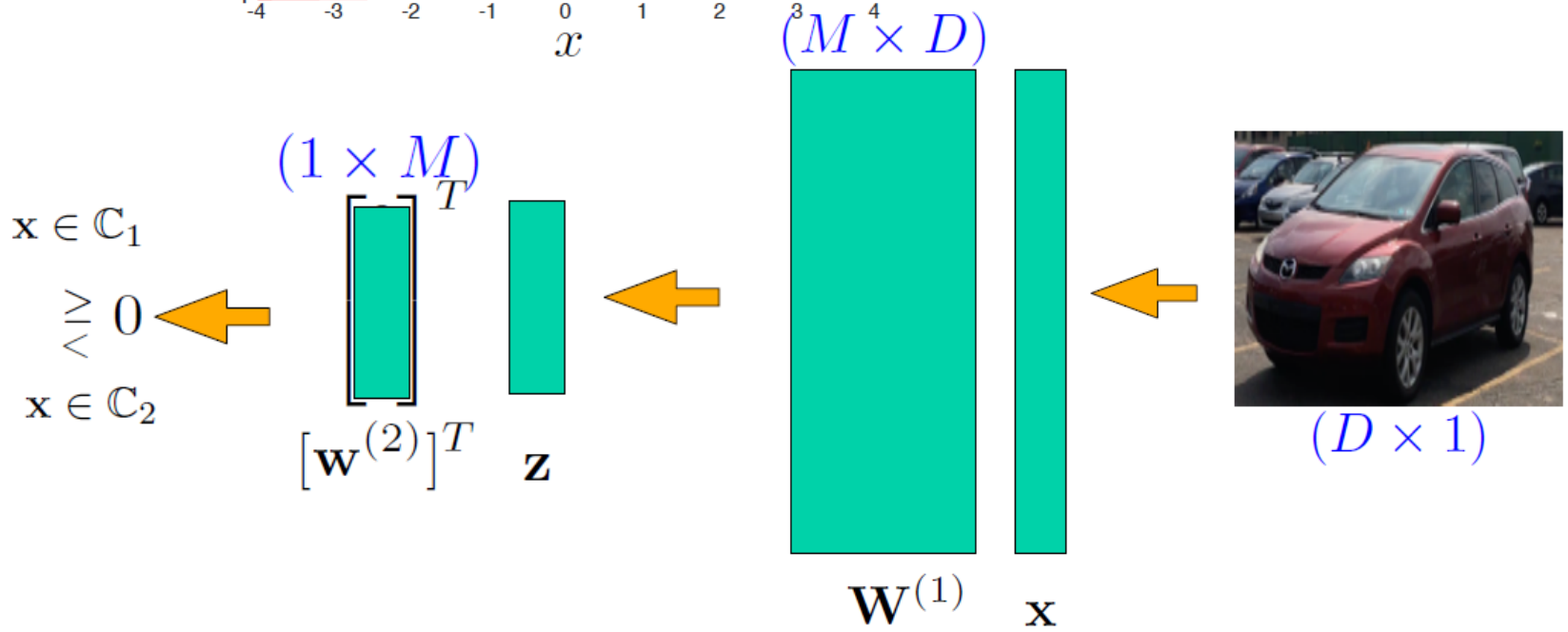
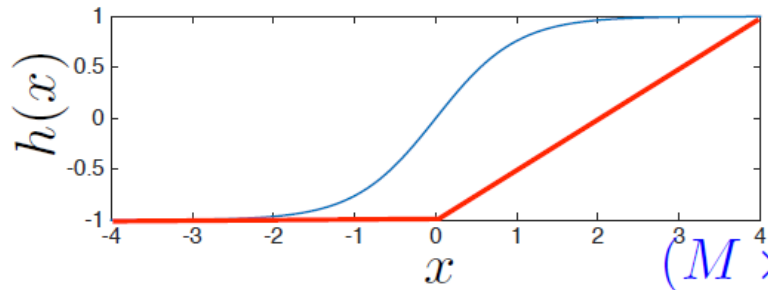


$h()$ = non-linear function

$[\mathbf{w}_1^{(1)}, \dots, \mathbf{w}_M^{(1)}]$ = 1st layer's $D \times M$ weights

\mathbf{x} = $D \times 1$ row input

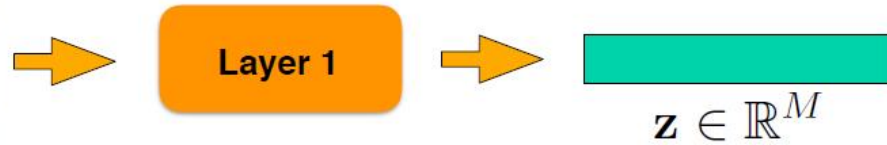
Multi-Layer Perceptron: A Nonlinear Classifier (cont'd)



Layer 2 in MLP



$$\mathbf{x} \in \mathbb{R}^D$$



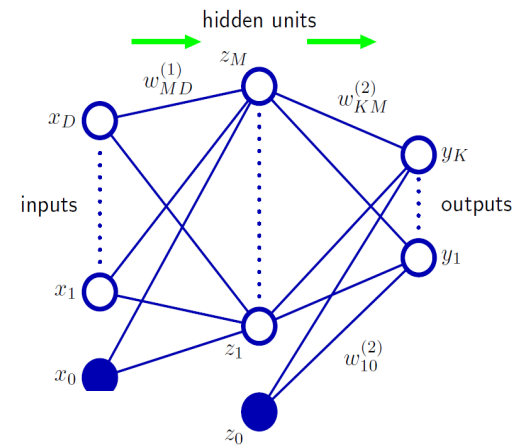
$$\mathbf{z} \in \mathbb{R}^M$$

$$\mathbf{z} \in \mathbb{C}_1$$

$$\mathbf{z}^T \mathbf{w}^{(2)} \geq 0$$

$$\mathbf{z}^T \mathbf{w}^{(2)} < 0$$

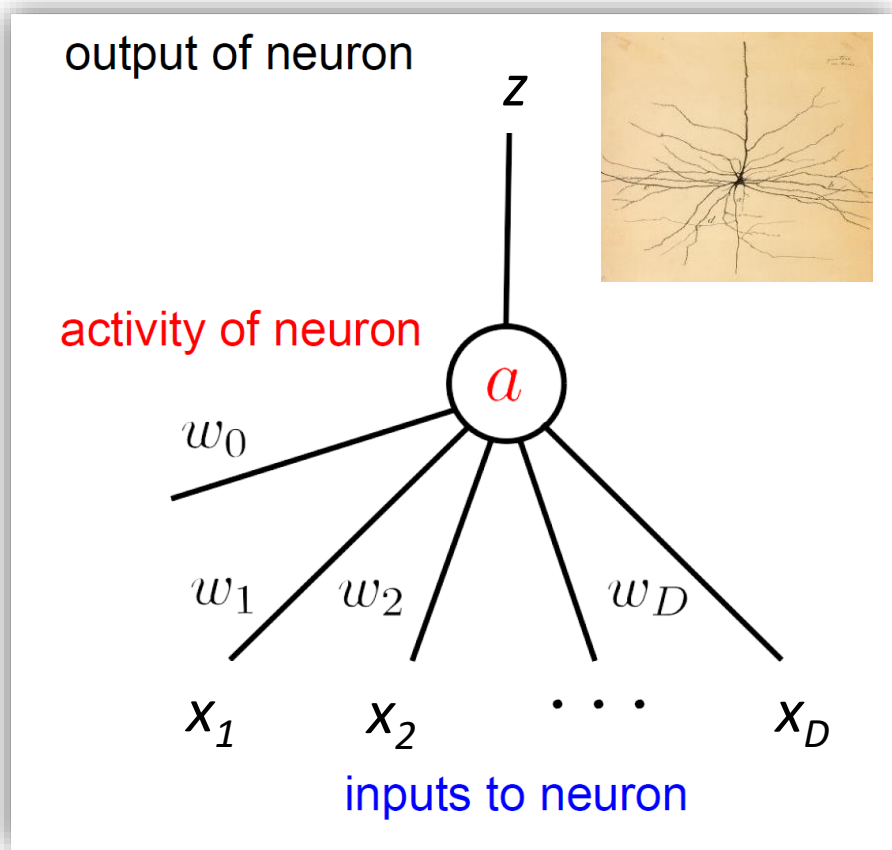
$$\mathbf{z} \in \mathbb{C}_2$$



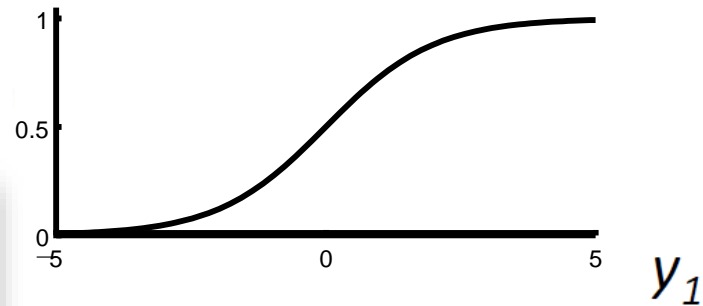
$\mathbf{z} = M \times 1$ output of layer 1
 $\mathbf{w}^{(2)} = 2\text{nd layer's } M \times 1$ weight vector

Let's Take a Closer Look... Output $y = Wx + b$

- A single neuron



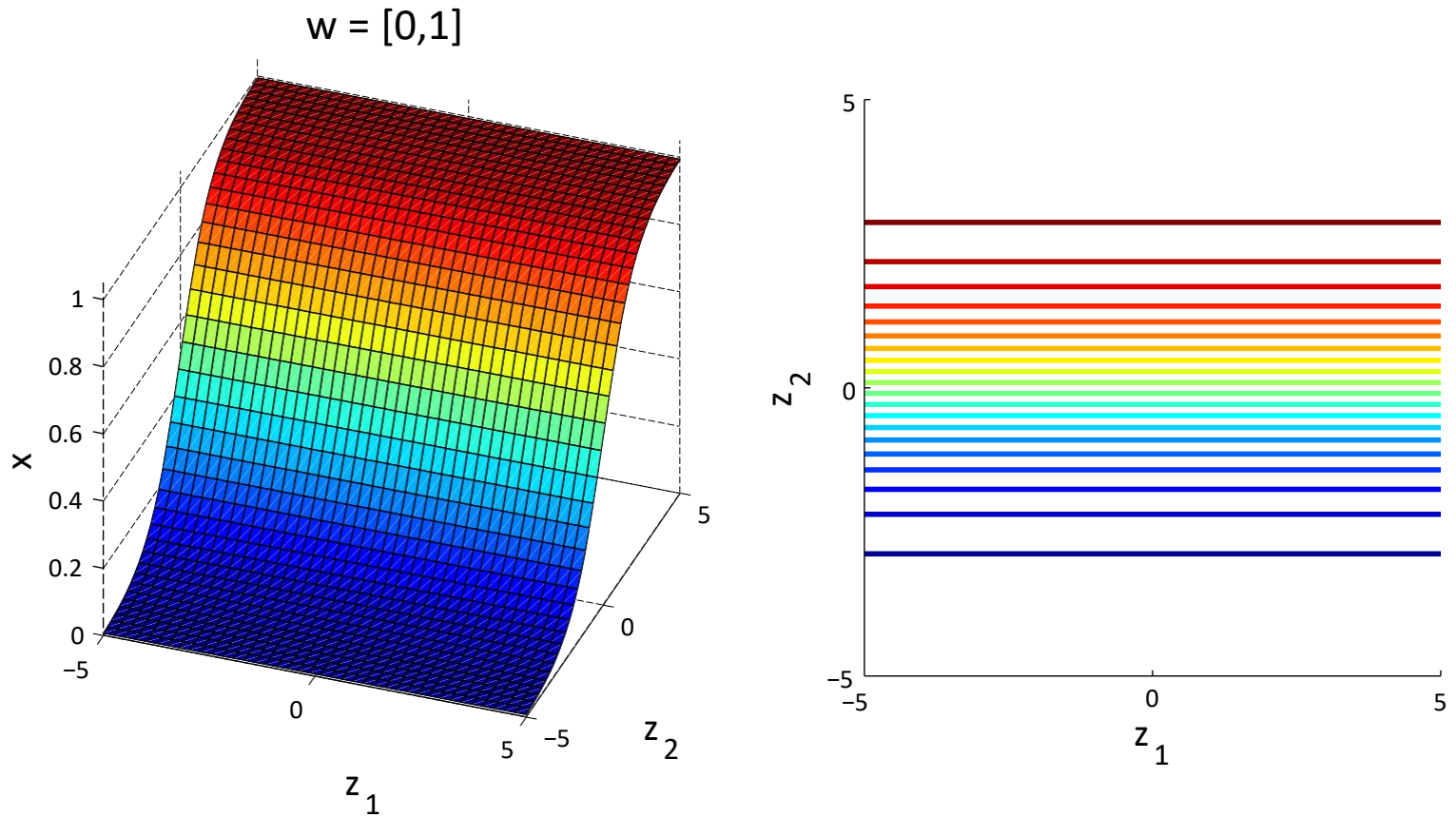
z



$$z(y_1) = \frac{1}{1 + \exp(-y_1)}, z \in (0,1)$$

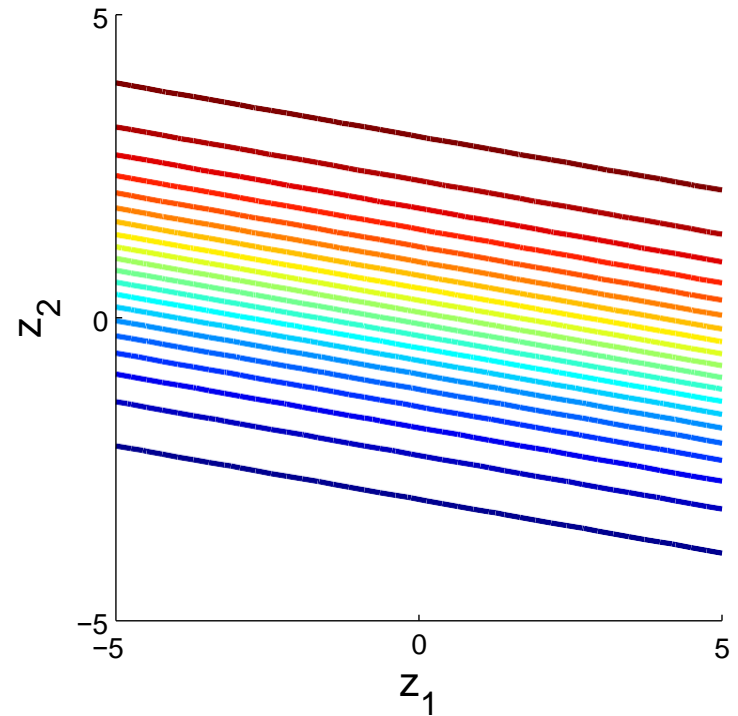
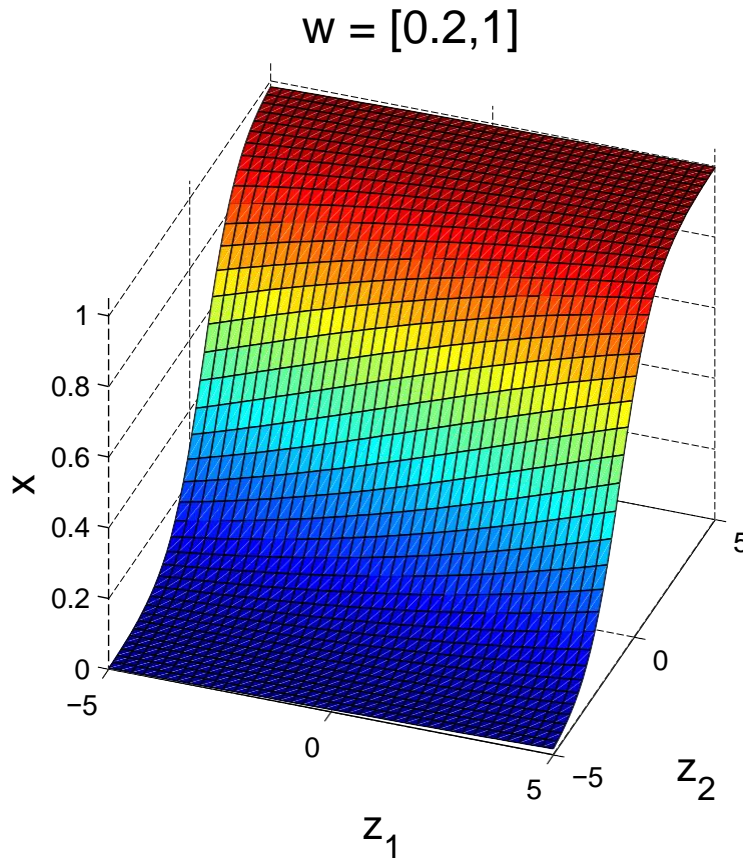
$$\begin{aligned} y_1 &= w_0 + \sum_{d=1}^D w_d z_d \\ &= \sum_{d=0}^D w_d z_d \end{aligned}$$

Input-Output Function of a Single Neuron



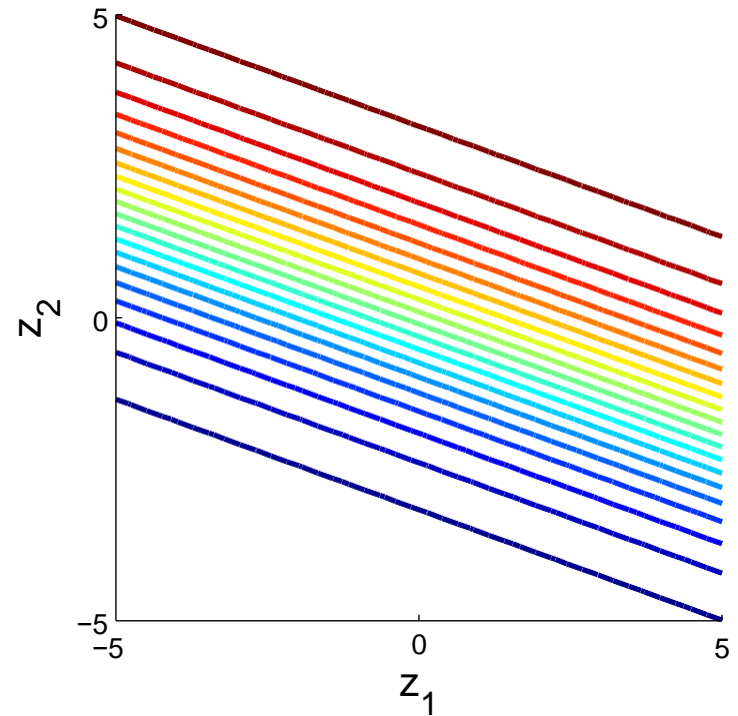
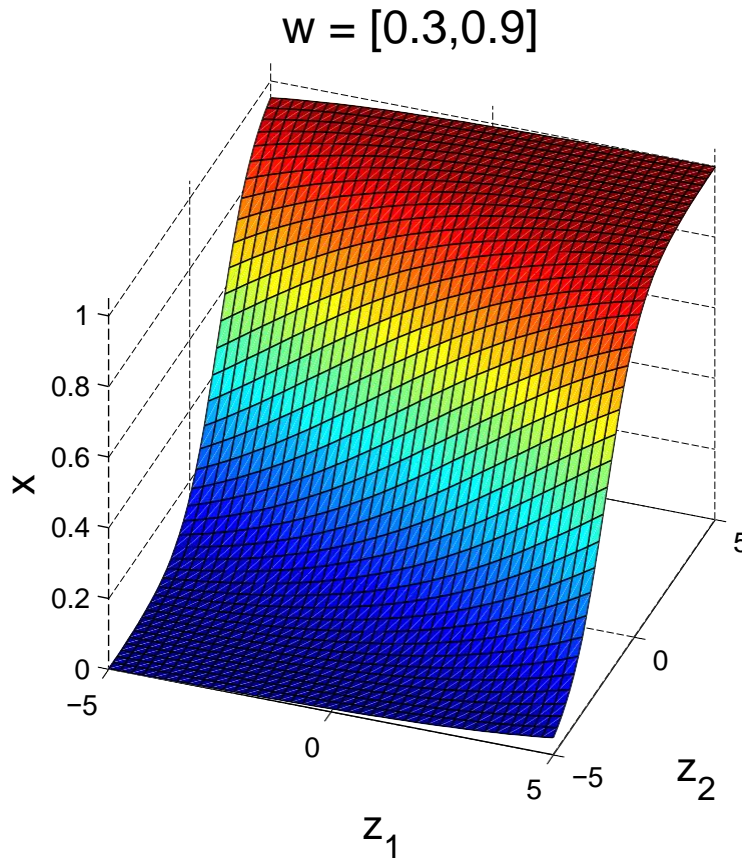
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron



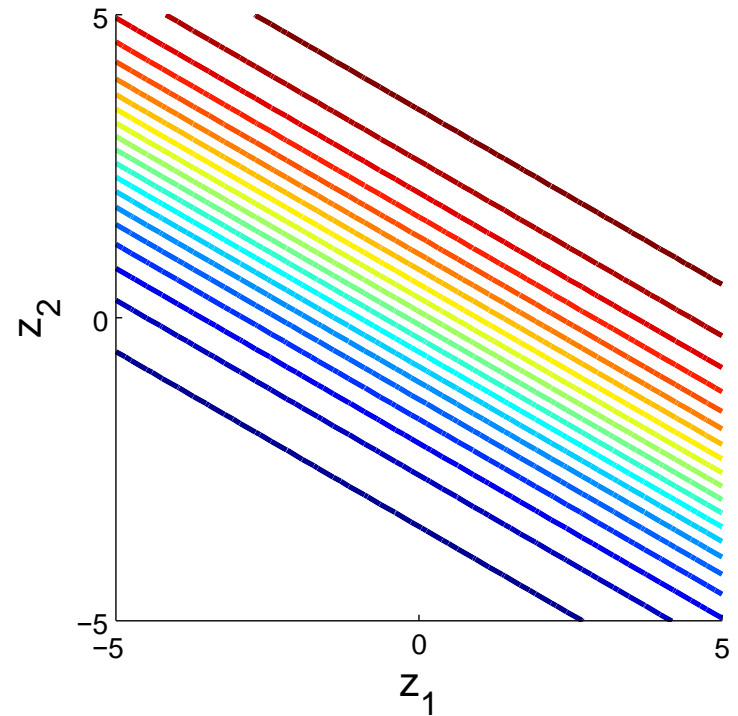
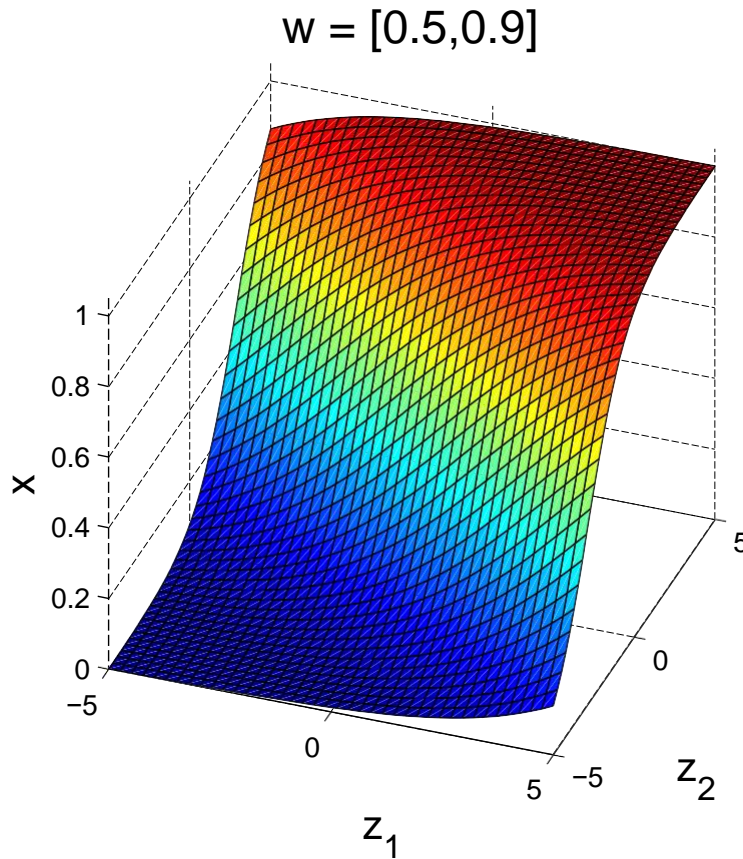
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron



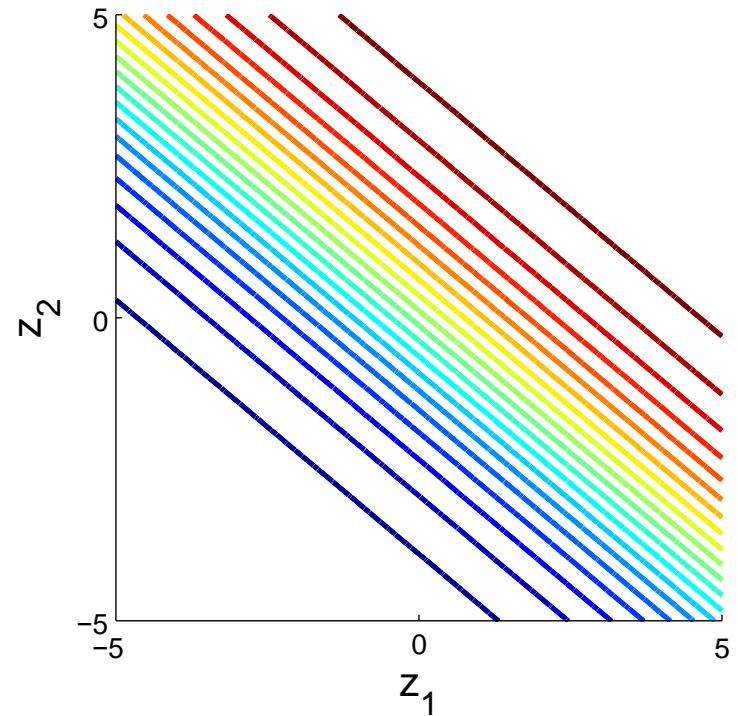
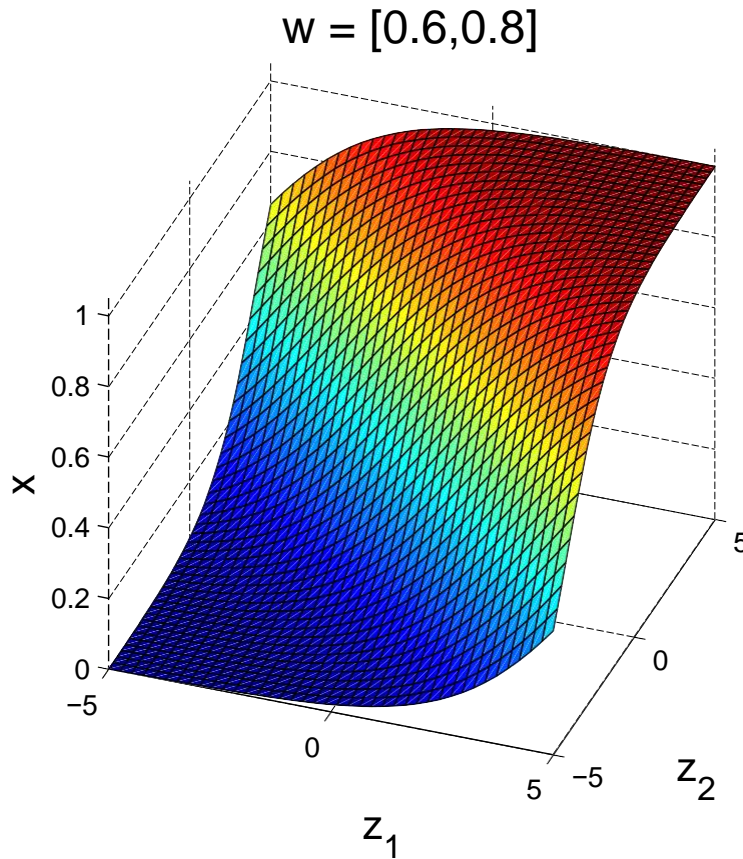
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron



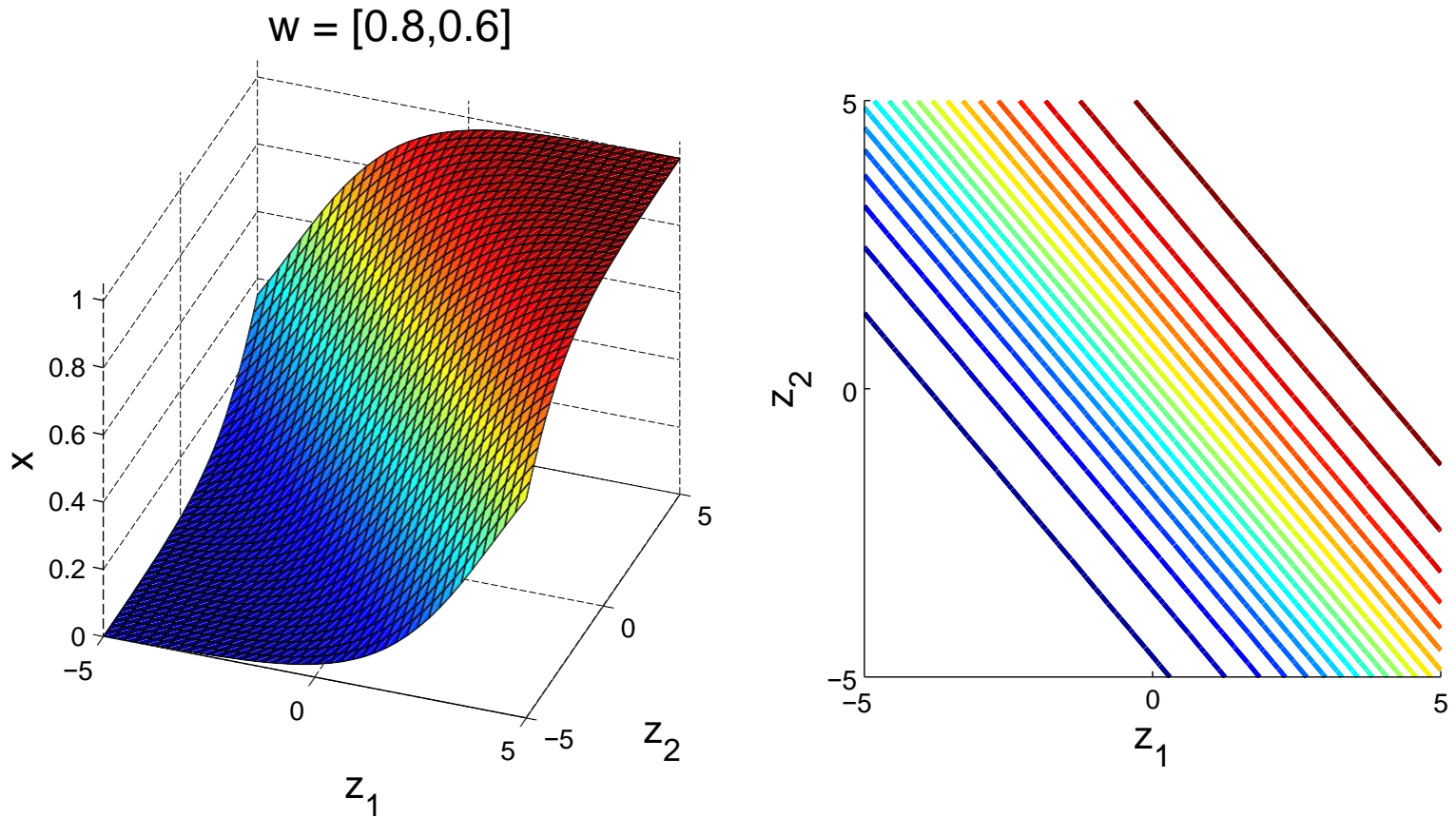
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron



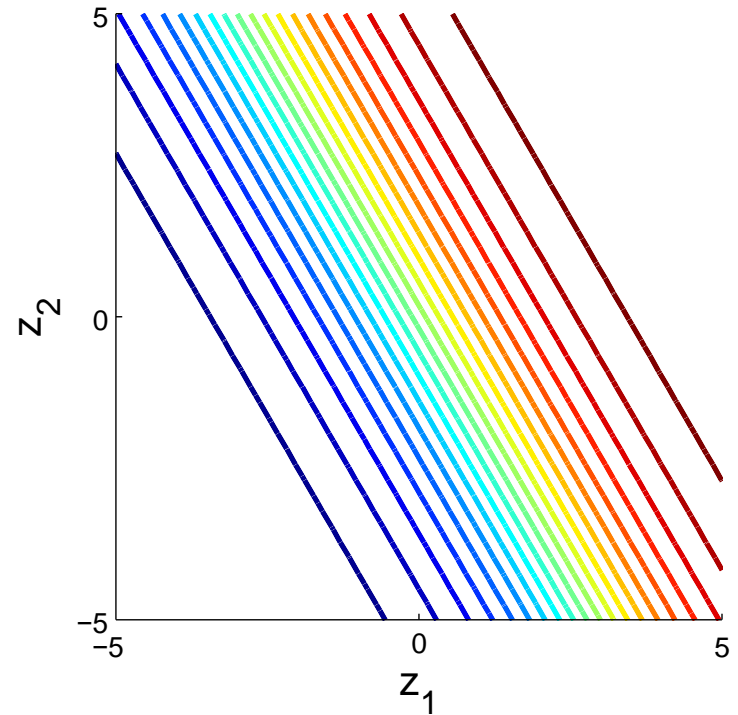
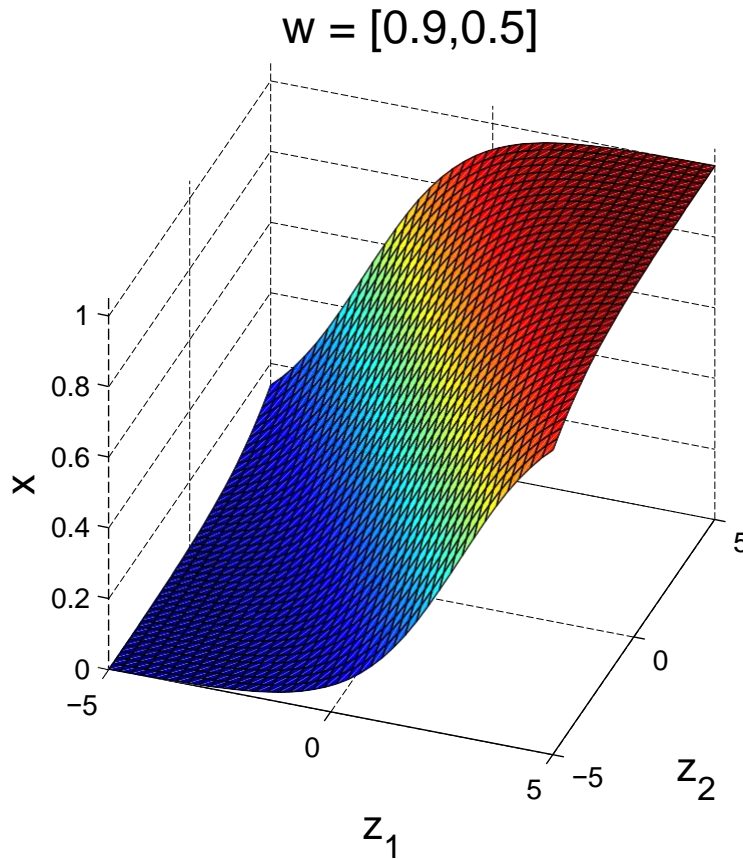
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron



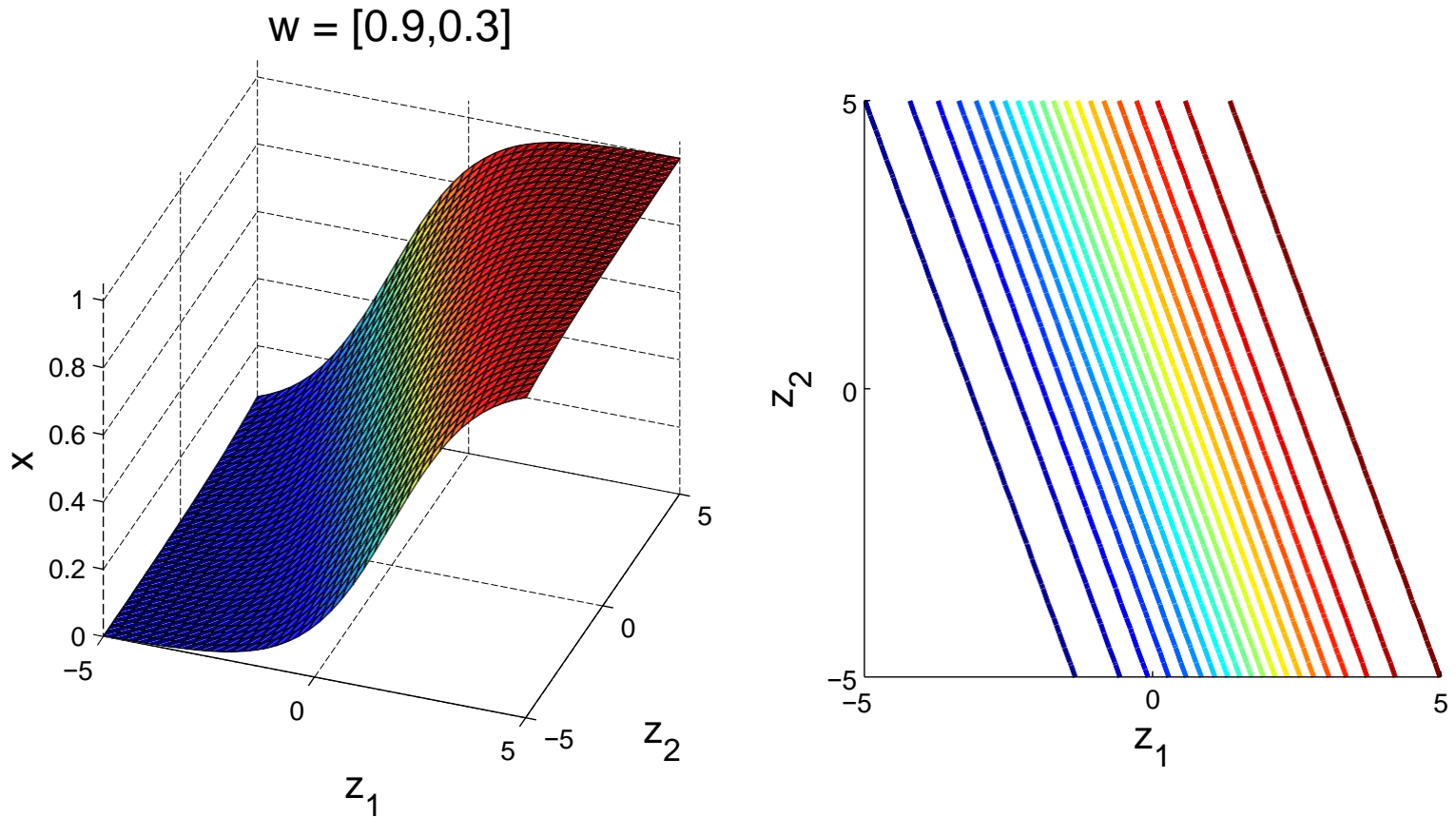
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron



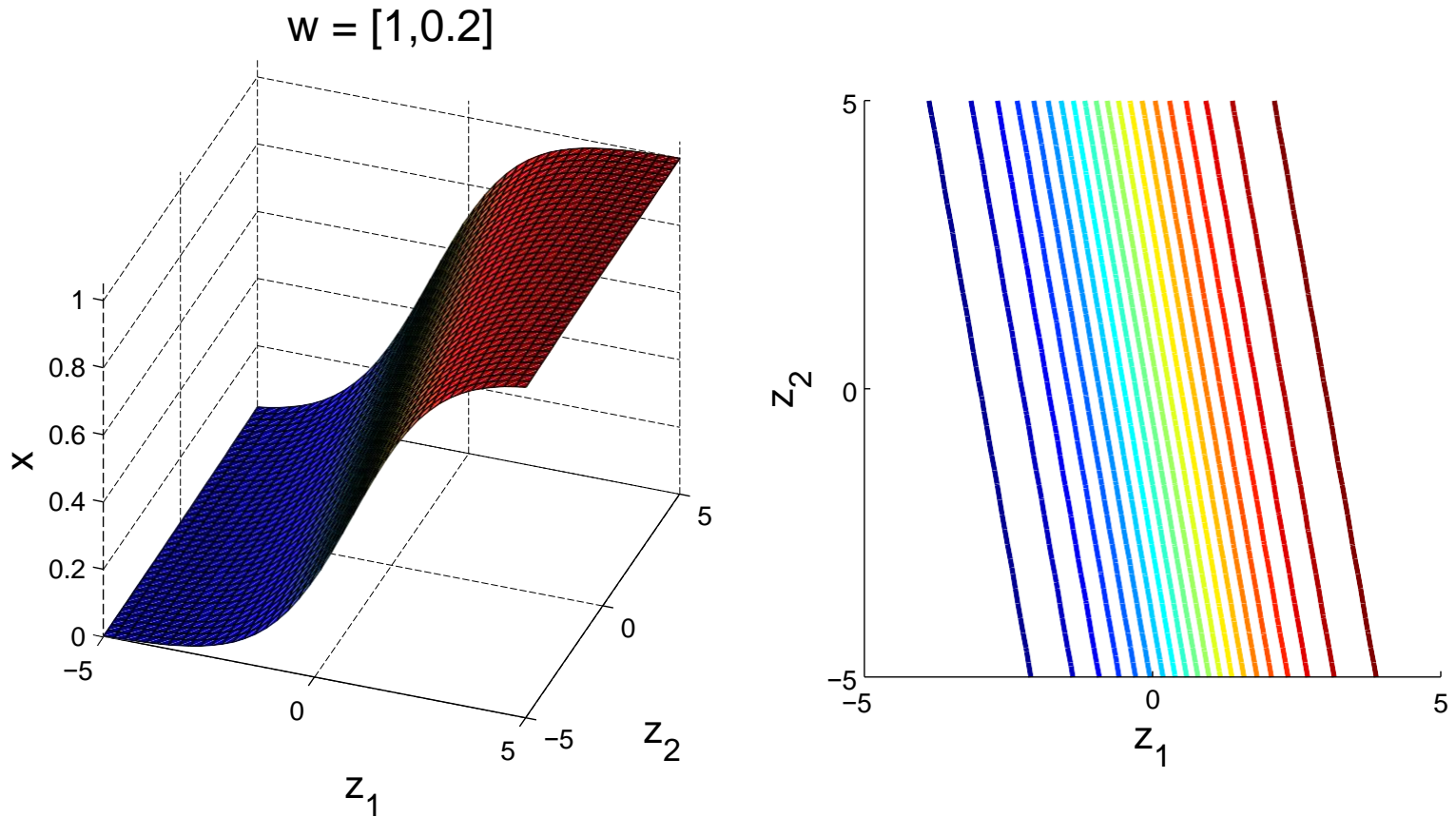
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron



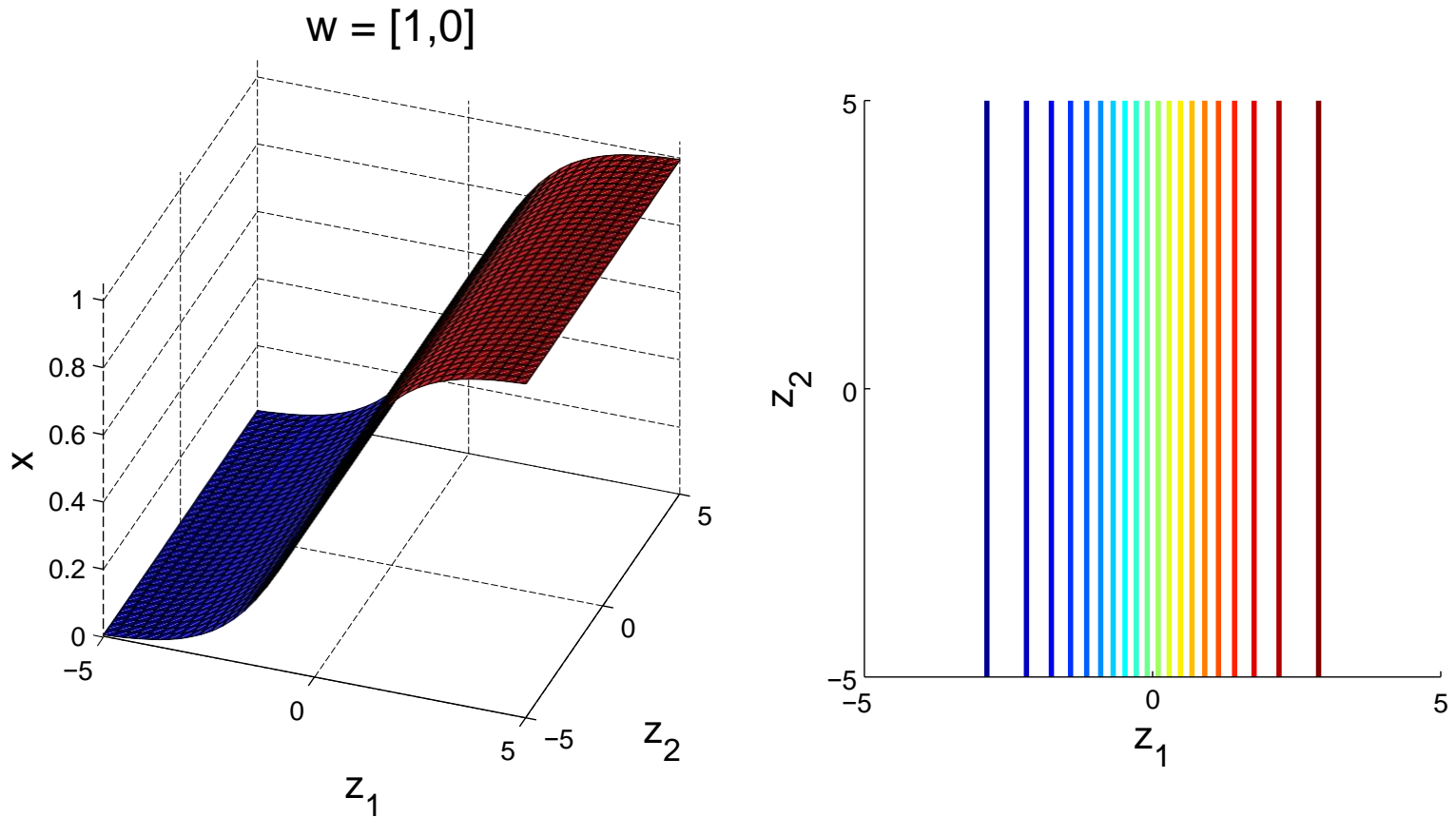
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron



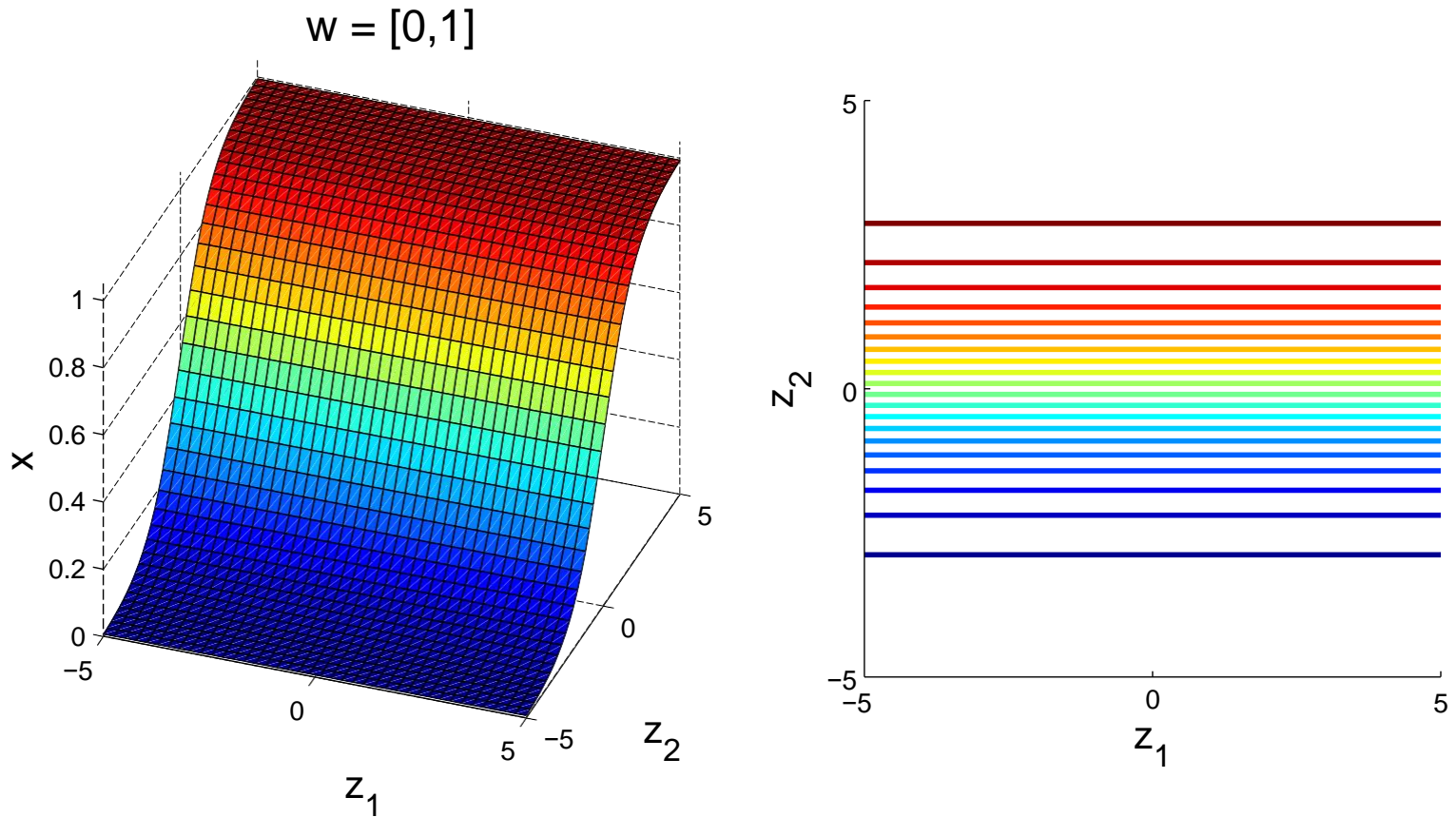
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron



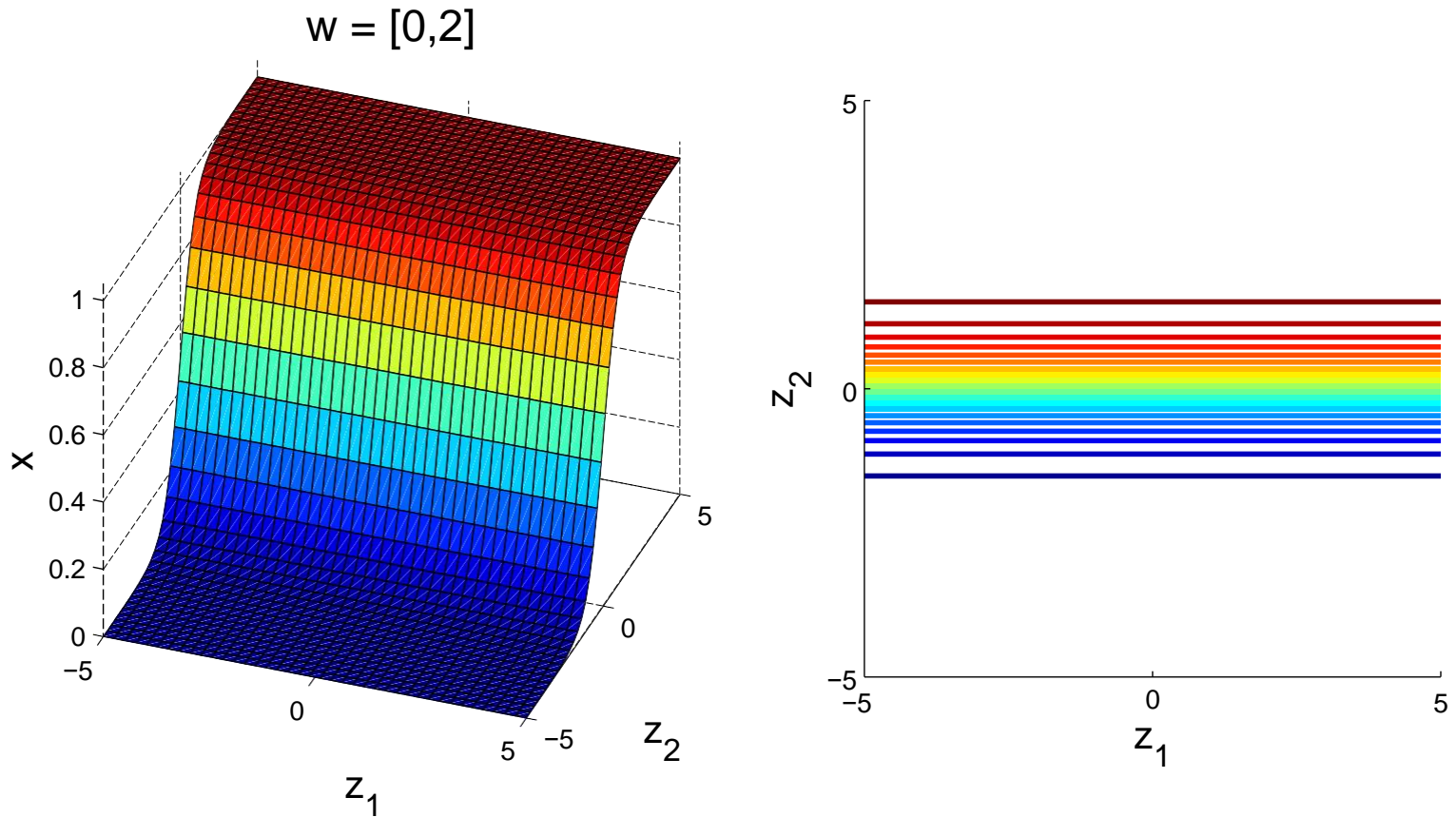
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron (cont'd)



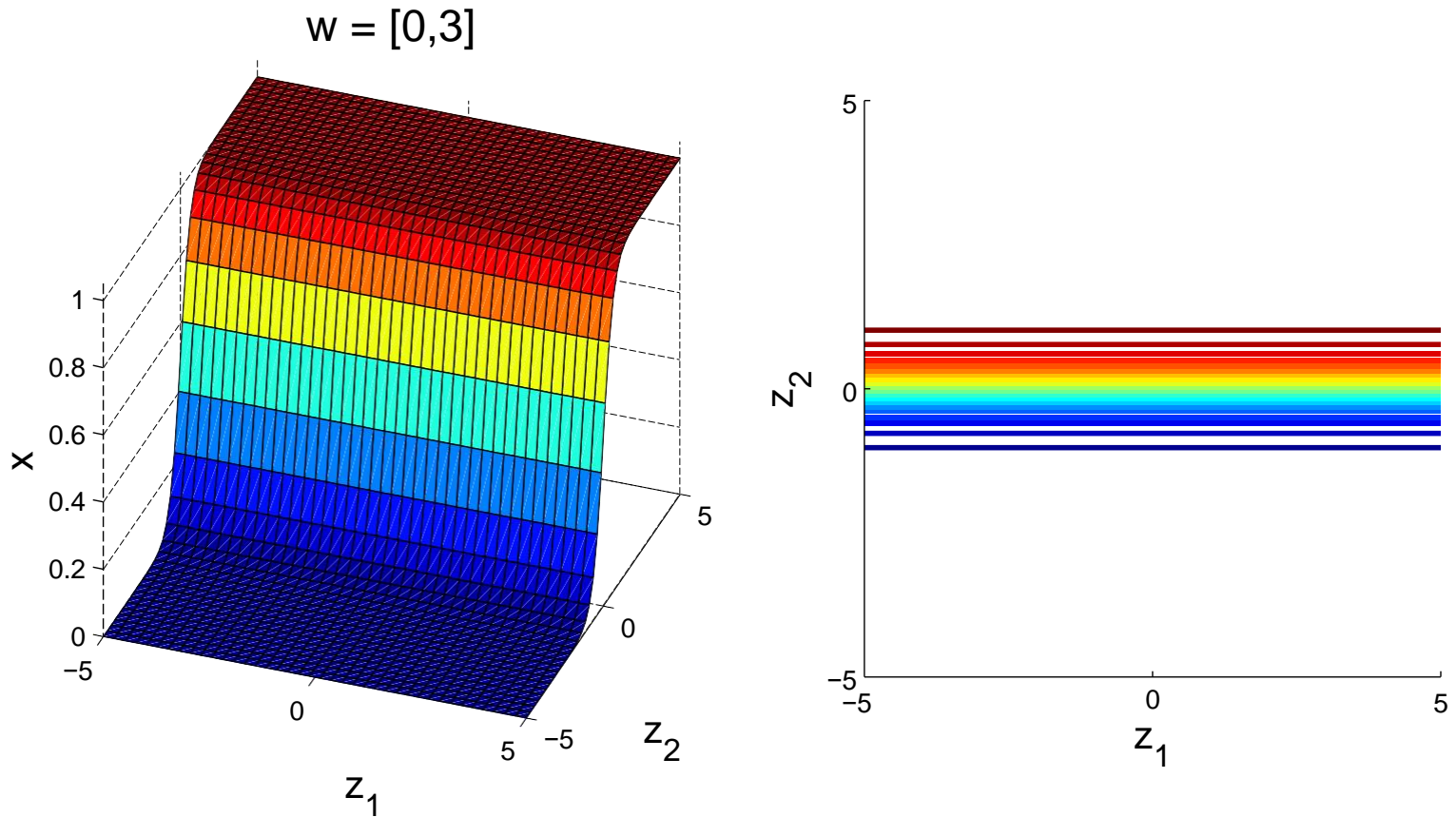
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron (cont'd)



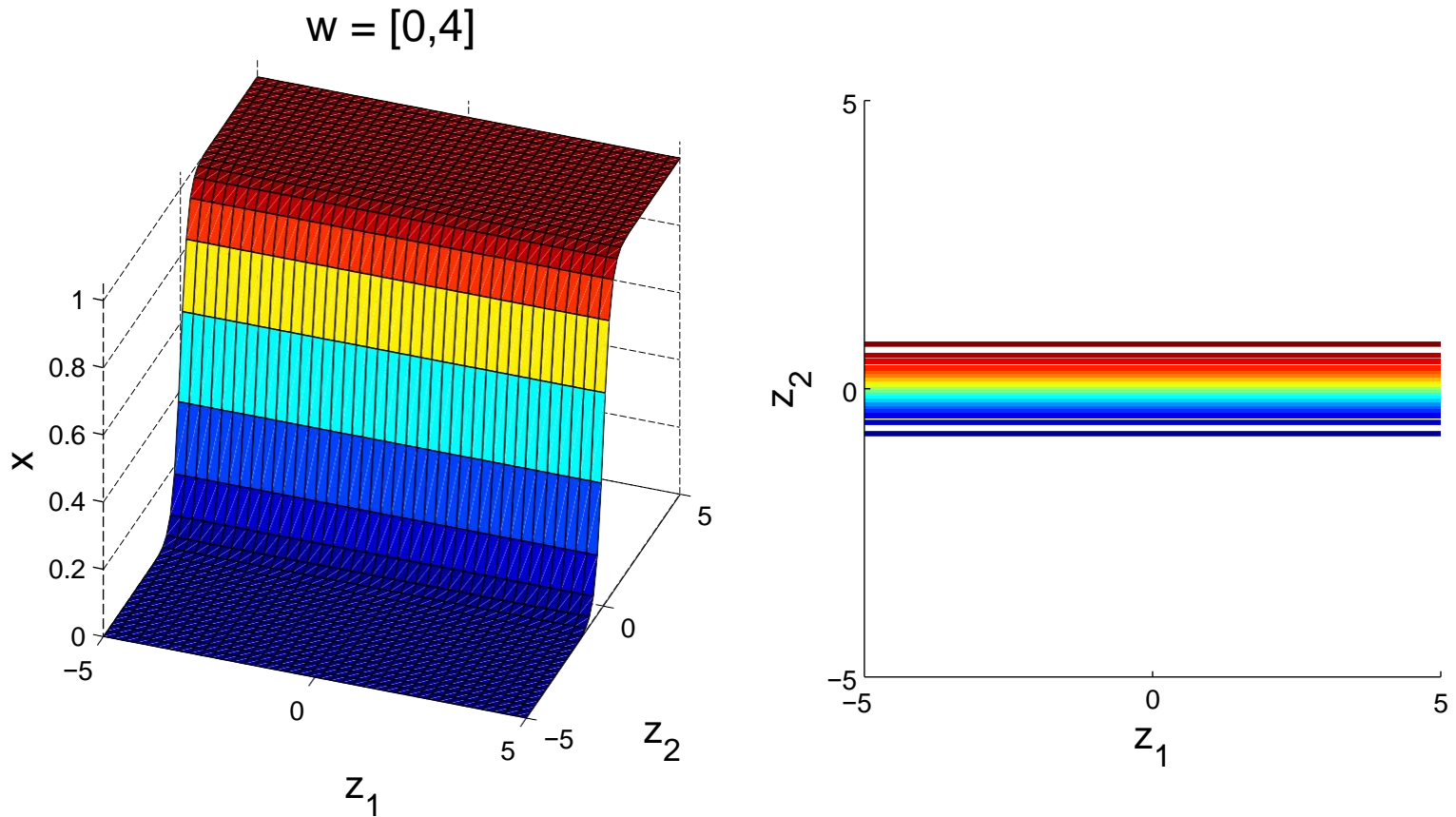
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron (cont'd)



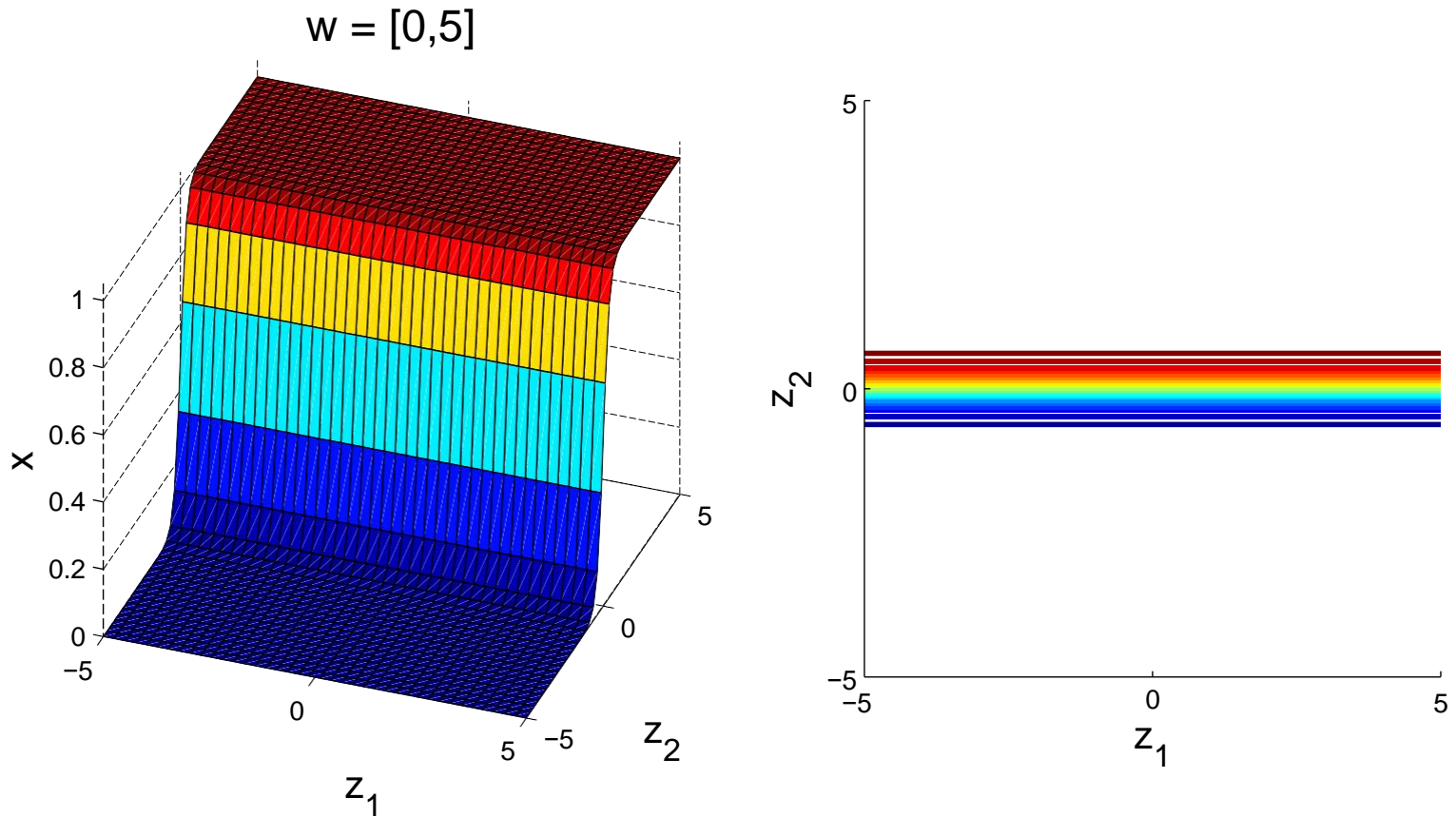
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron (cont'd)



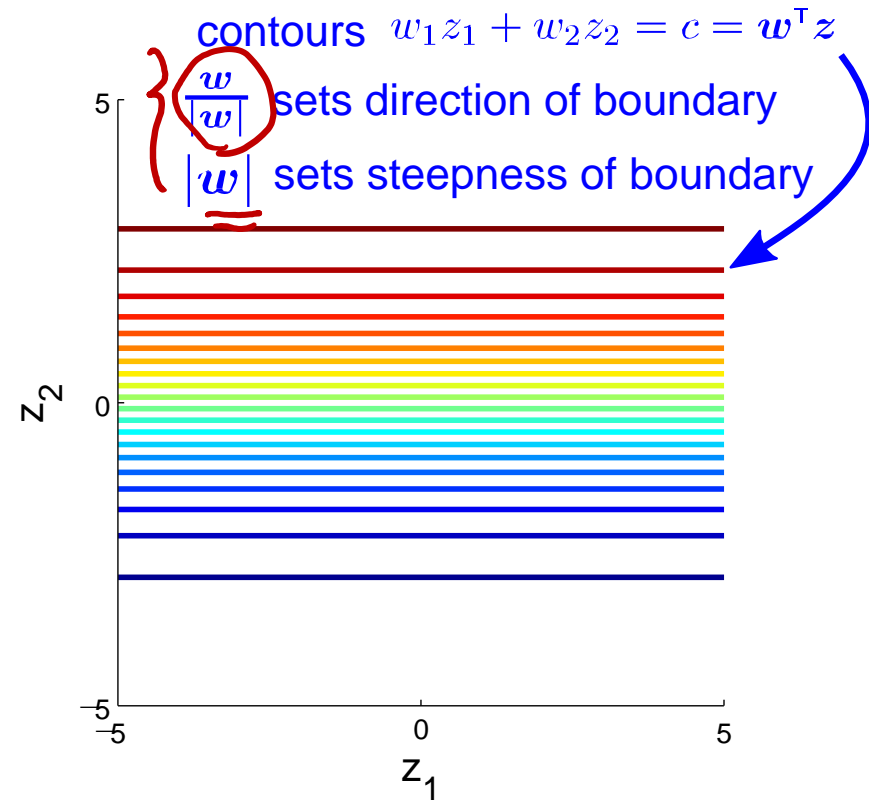
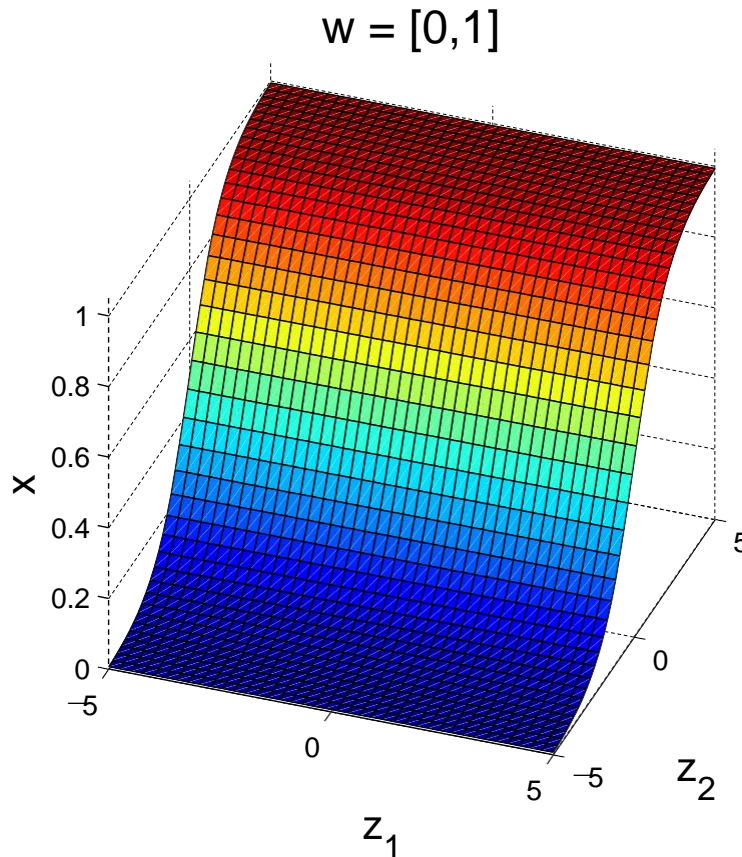
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron (cont'd)



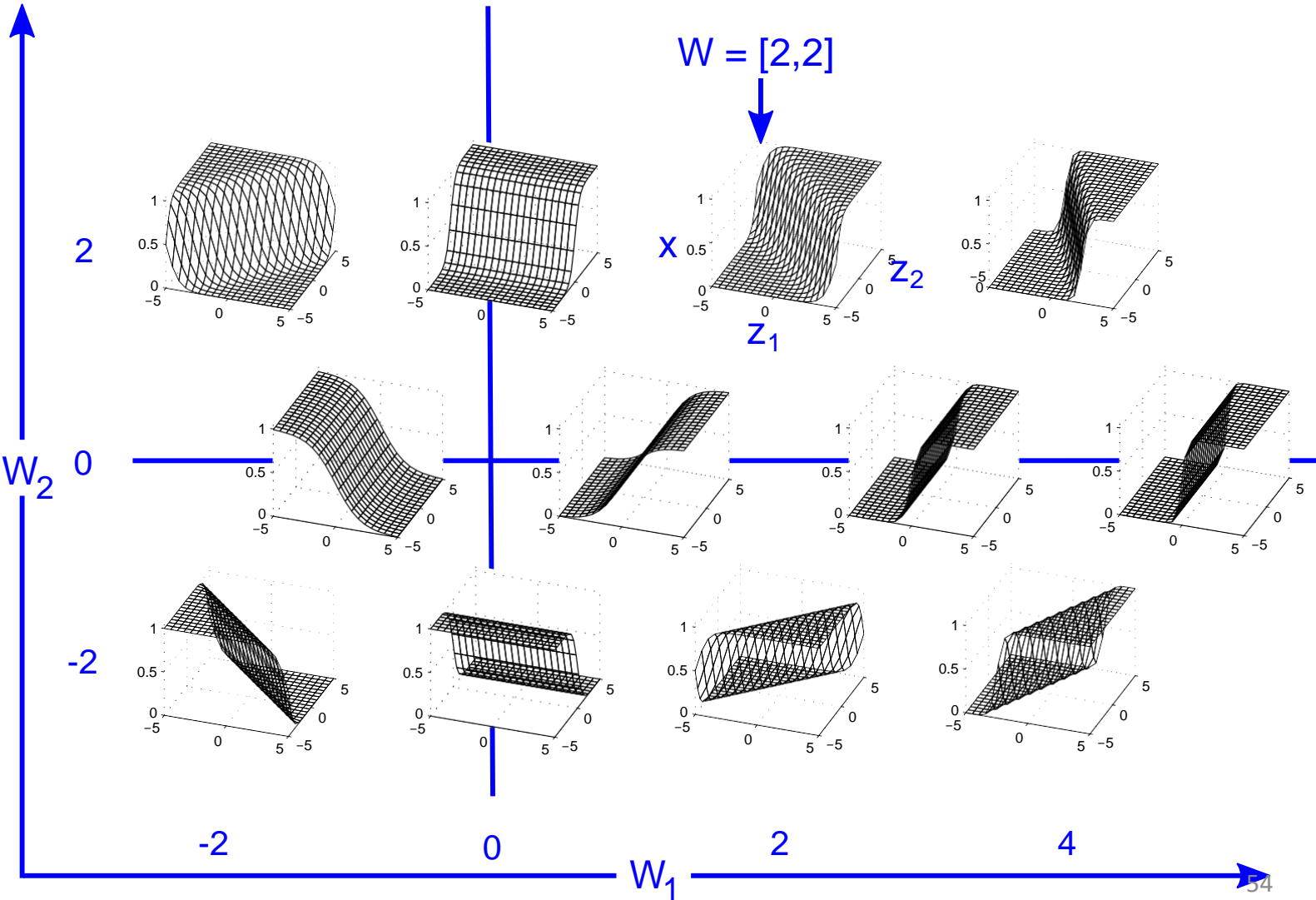
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron (cont'd)

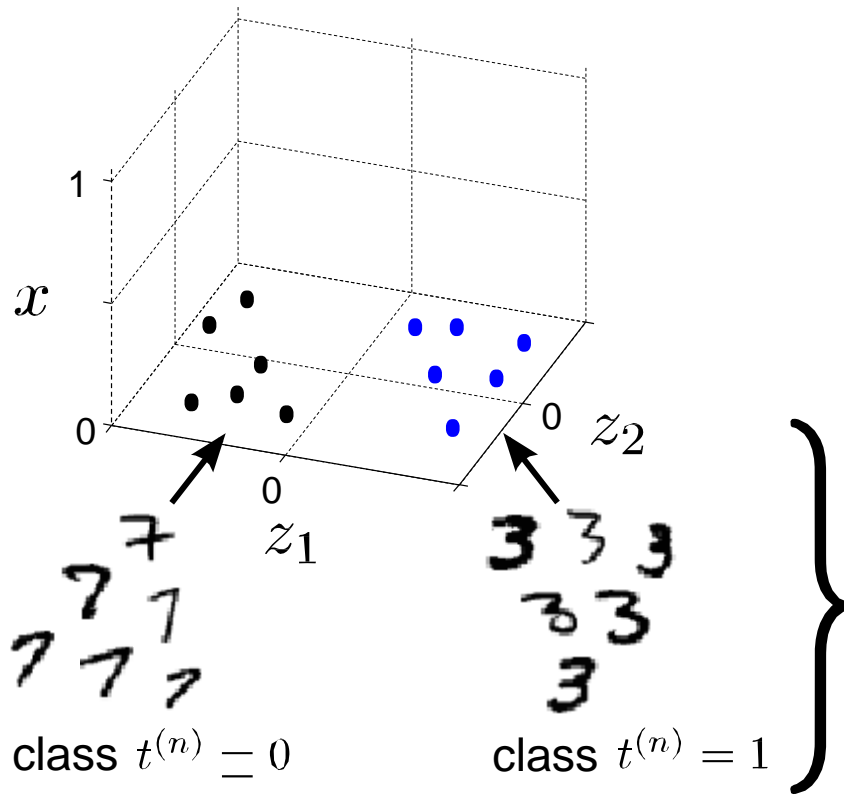


$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Weight Space of a Single Neuron



Training a Single Neuron



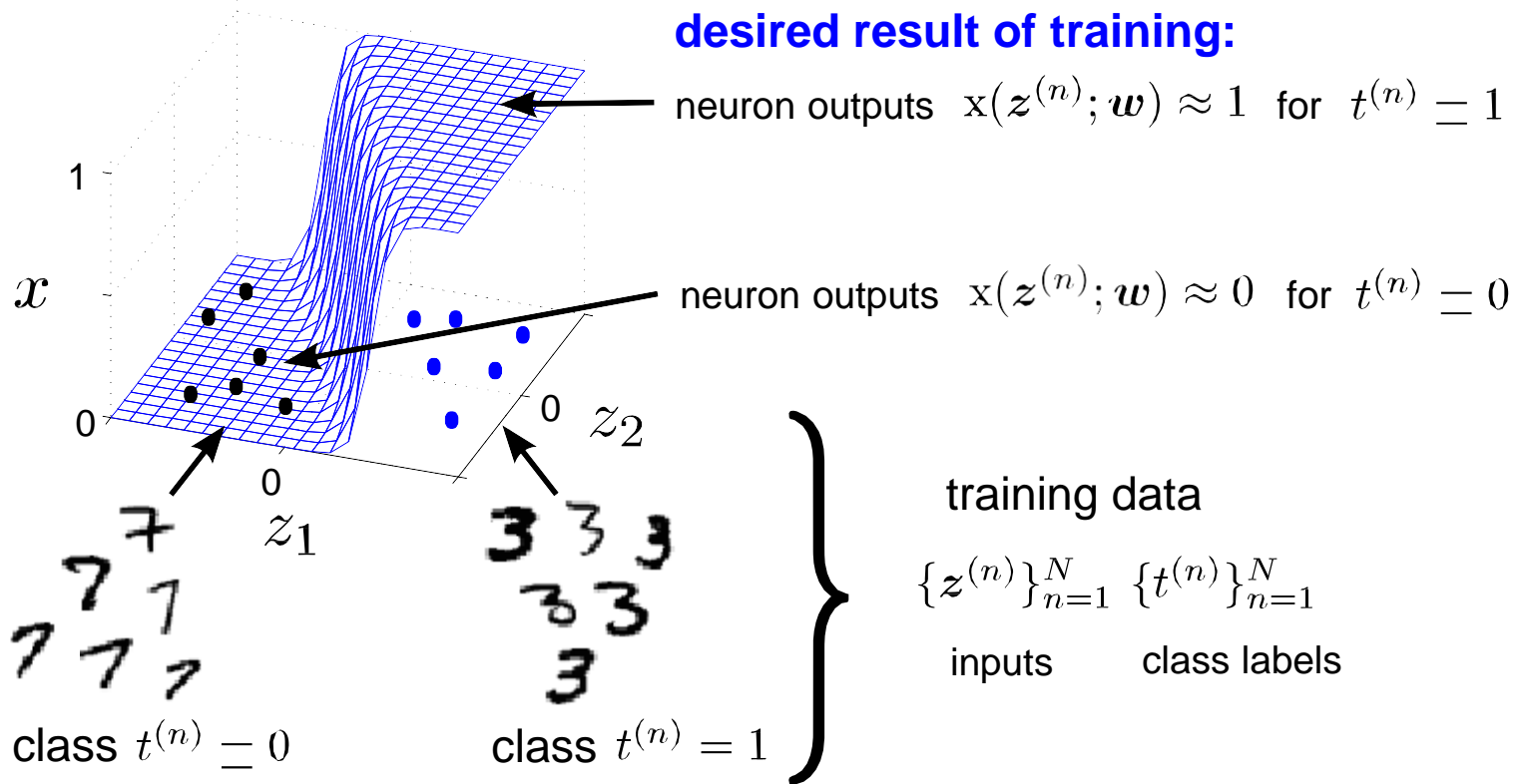
training data

$$\{z^{(n)}\}_{n=1}^N \quad \{t^{(n)}\}_{n=1}^N$$

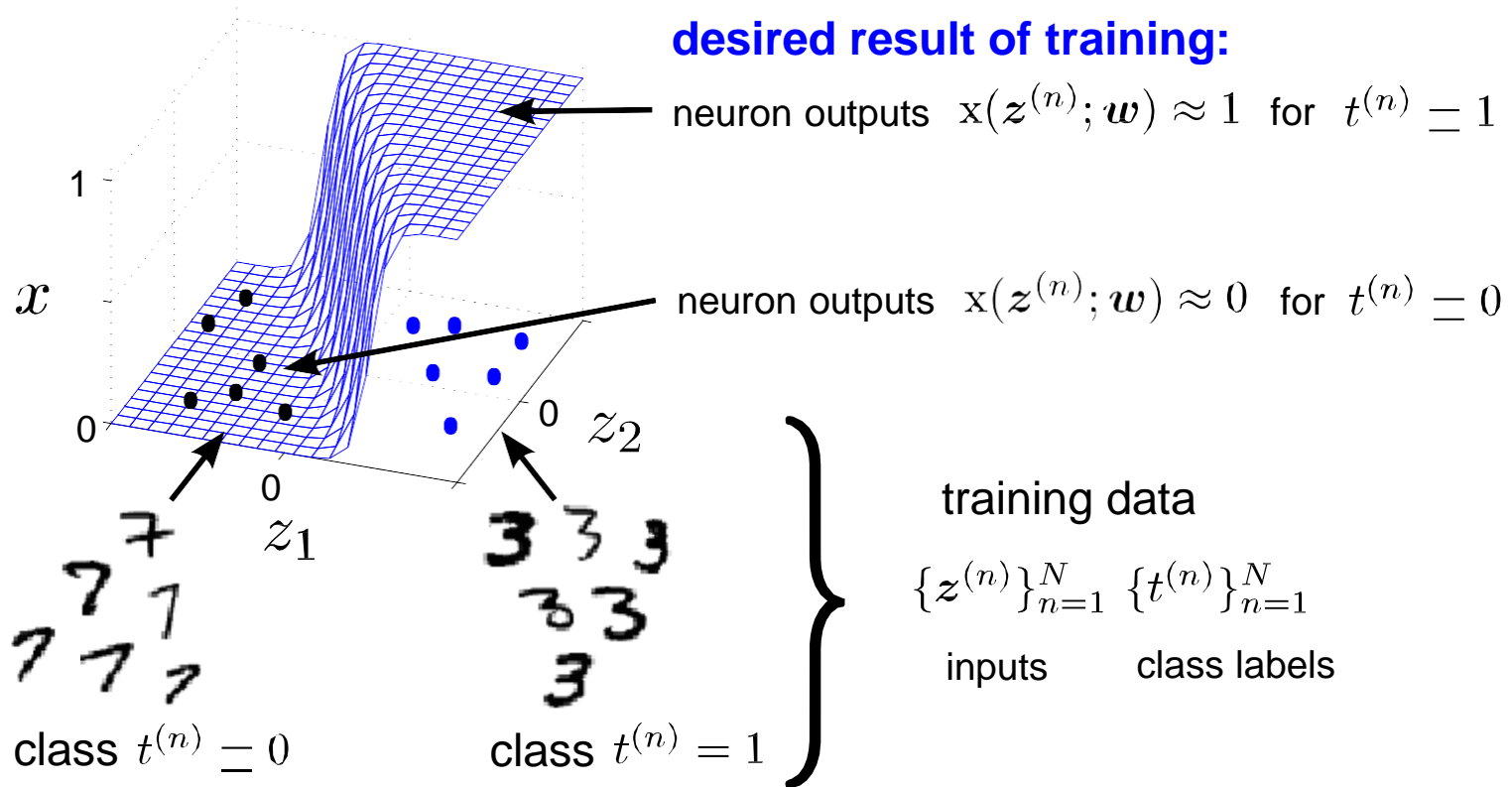
inputs

class labels

Training a Single Neuron



Training a Single Neuron

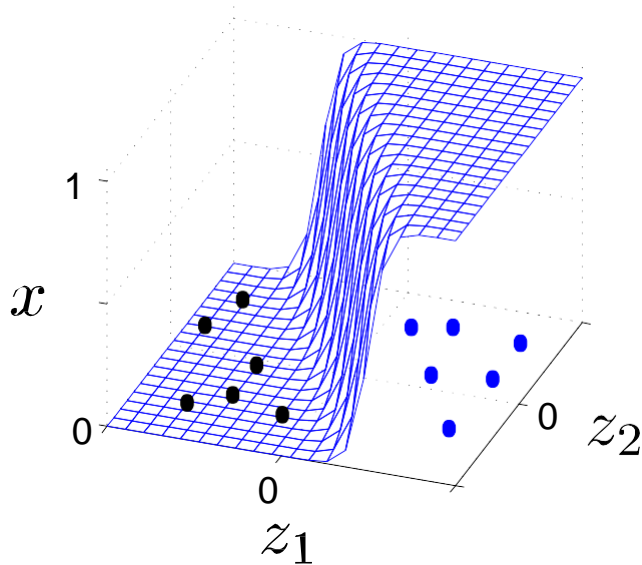


objective function:

$$G(\mathbf{w}) = - \sum_n [t^{(n)} \log x(\mathbf{z}^{(n)}; \mathbf{w}) + (1 - t^{(n)}) \log (1 - x(\mathbf{z}^{(n)}; \mathbf{w}))] \geq 0$$

surprise $-\log p(\text{outcome})$ when observing $t^{(n)}$ } encourages neuron output
 relative entropy between $x(\mathbf{z}^{(n)}; \mathbf{w})$ and $t^{(n)}$ } to match training data

Training a Single Neuron



training data

$$\{\mathbf{z}^{(n)}\}_{n=1}^N \quad \{t^{(n)}\}_{n=1}^N$$

inputs class labels

objective function:

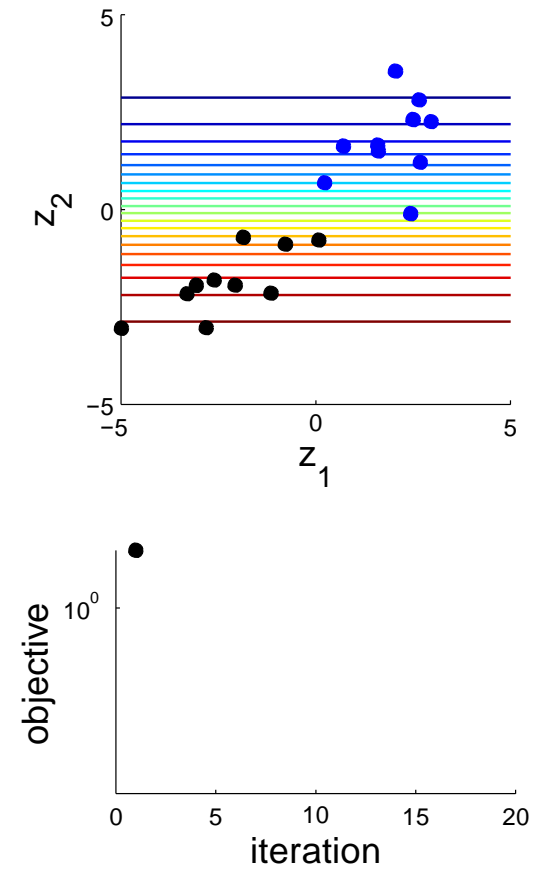
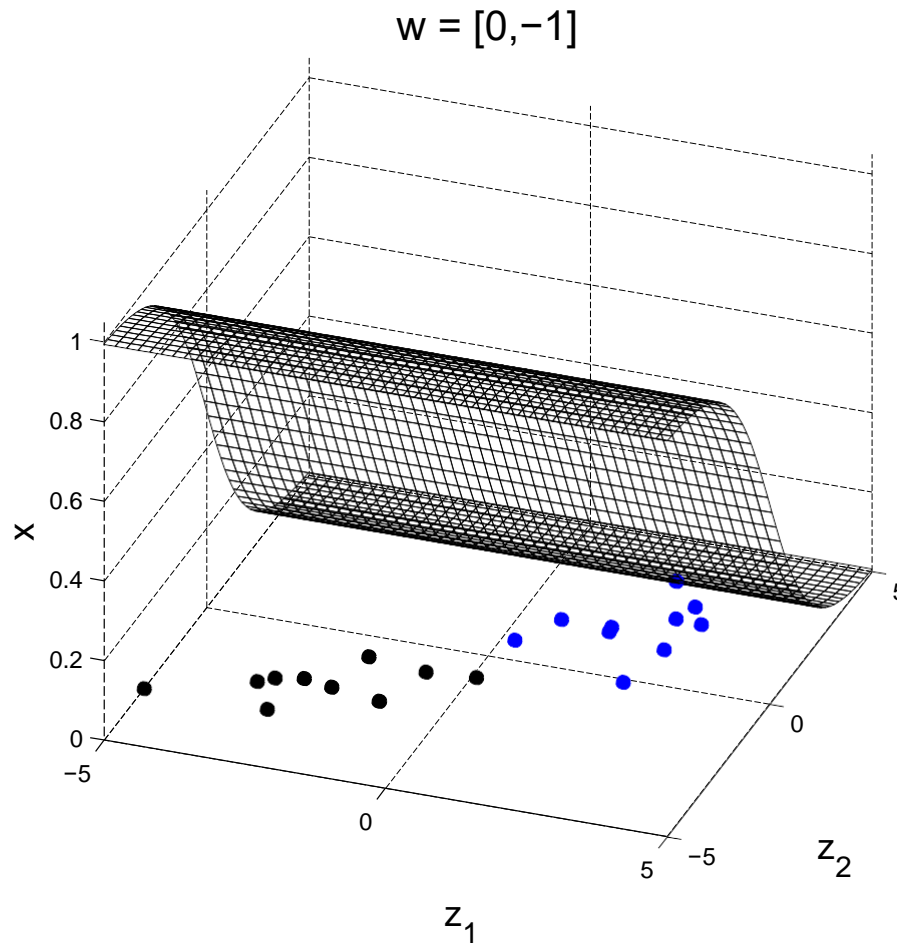
$$G(\mathbf{w}) = - \sum_n [t^{(n)} \log x(\mathbf{z}^{(n)}; \mathbf{w}) + (1 - t^{(n)}) \log (1 - x(\mathbf{z}^{(n)}; \mathbf{w}))] \geq 0$$

$\mathbf{w}^* = \arg \min_{\mathbf{w}} G(\mathbf{w})$ choose the weights that minimise the network's surprise about the training data

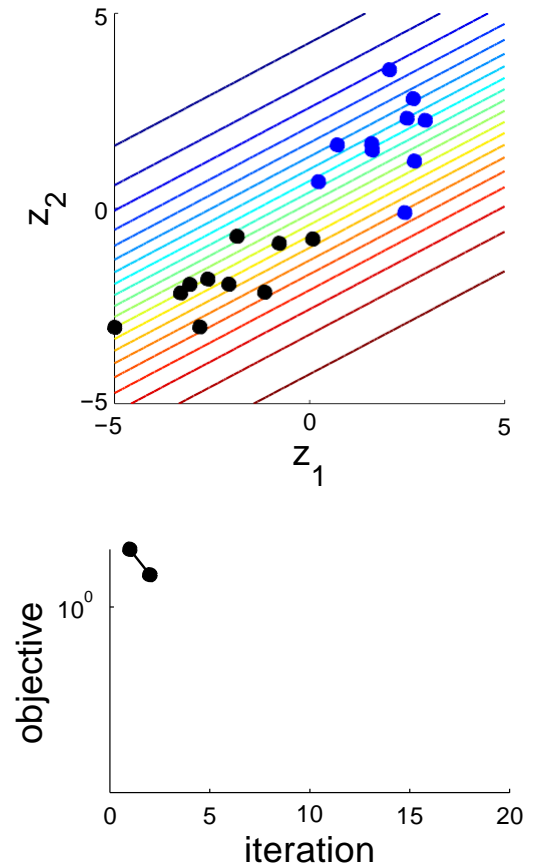
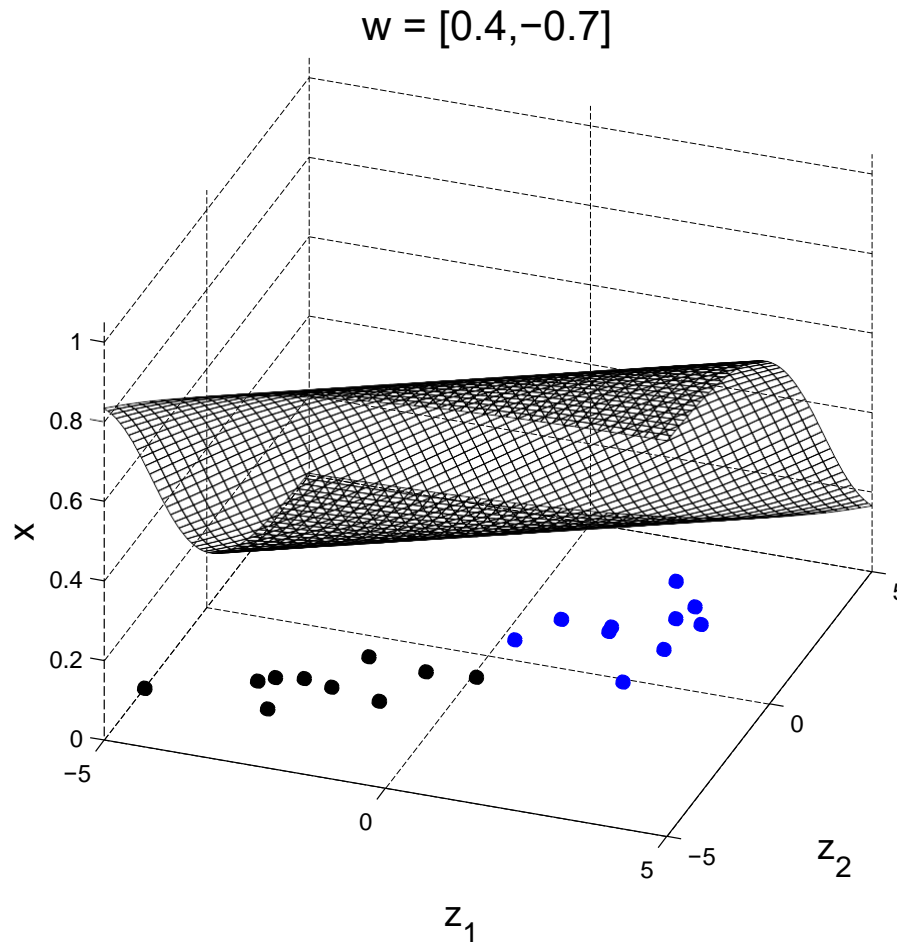
$$\frac{d}{d\mathbf{w}} G(\mathbf{w}) = \sum_n \frac{dG(\mathbf{w})}{dx^{(n)}} \frac{dx^{(n)}}{d\mathbf{w}} = - \sum_n (t^{(n)} - x^{(n)}) \mathbf{z}^{(n)} = \text{prediction error} \times \text{feature}$$

$\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{d}{d\mathbf{w}} G(\mathbf{w})$ iteratively step down the objective (gradient points up hill)

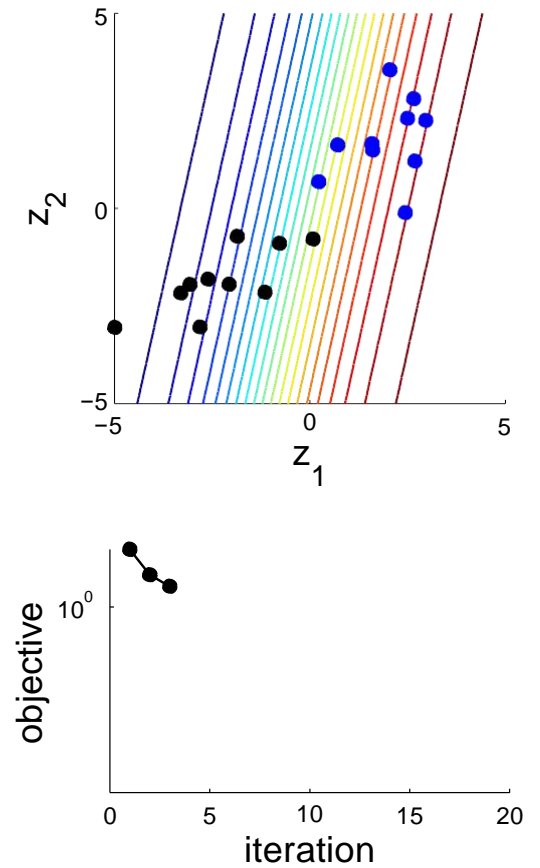
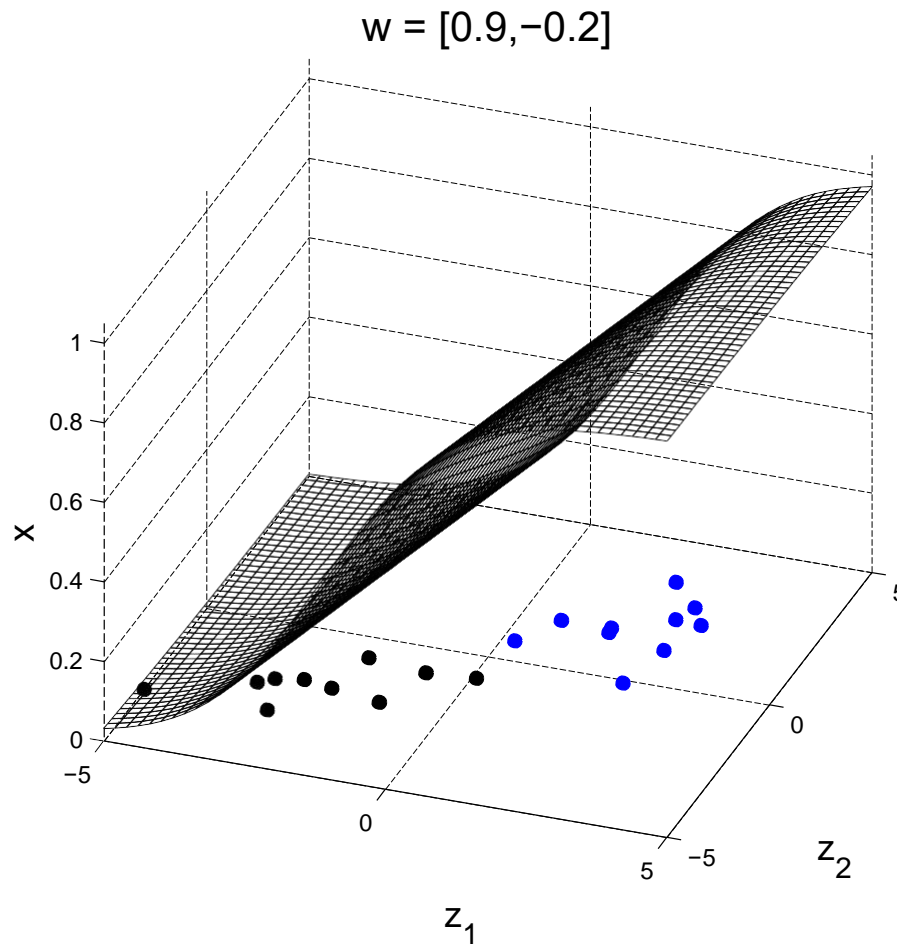
Training a Single Neuron



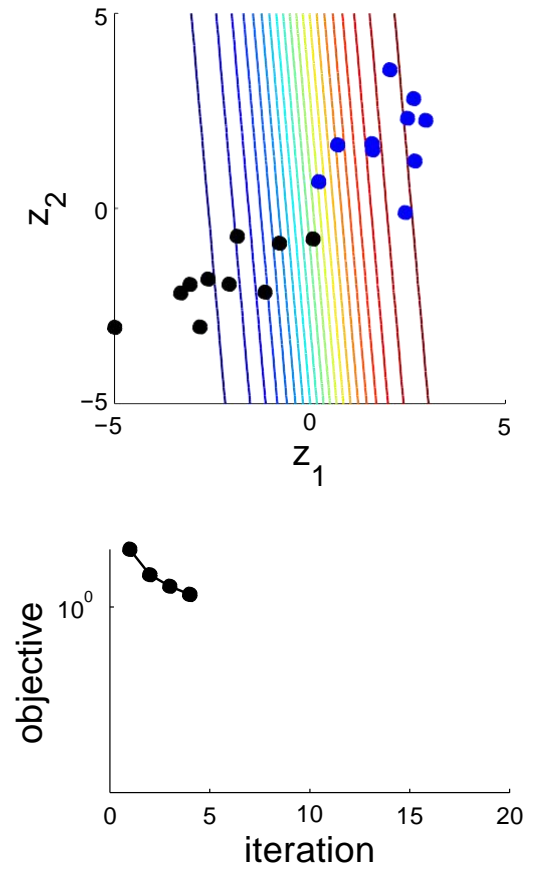
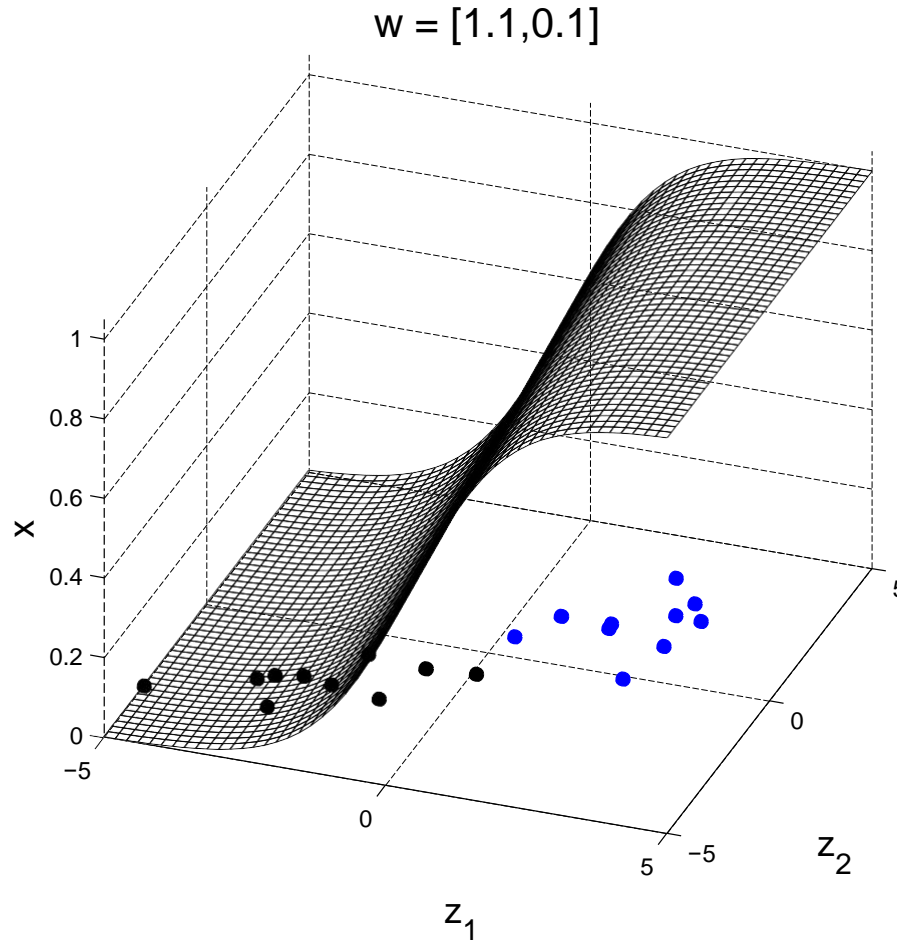
Training a Single Neuron



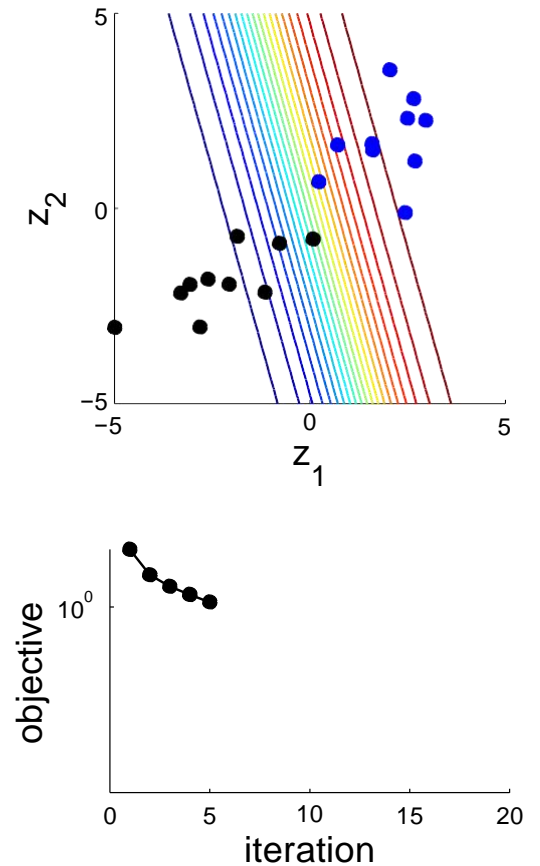
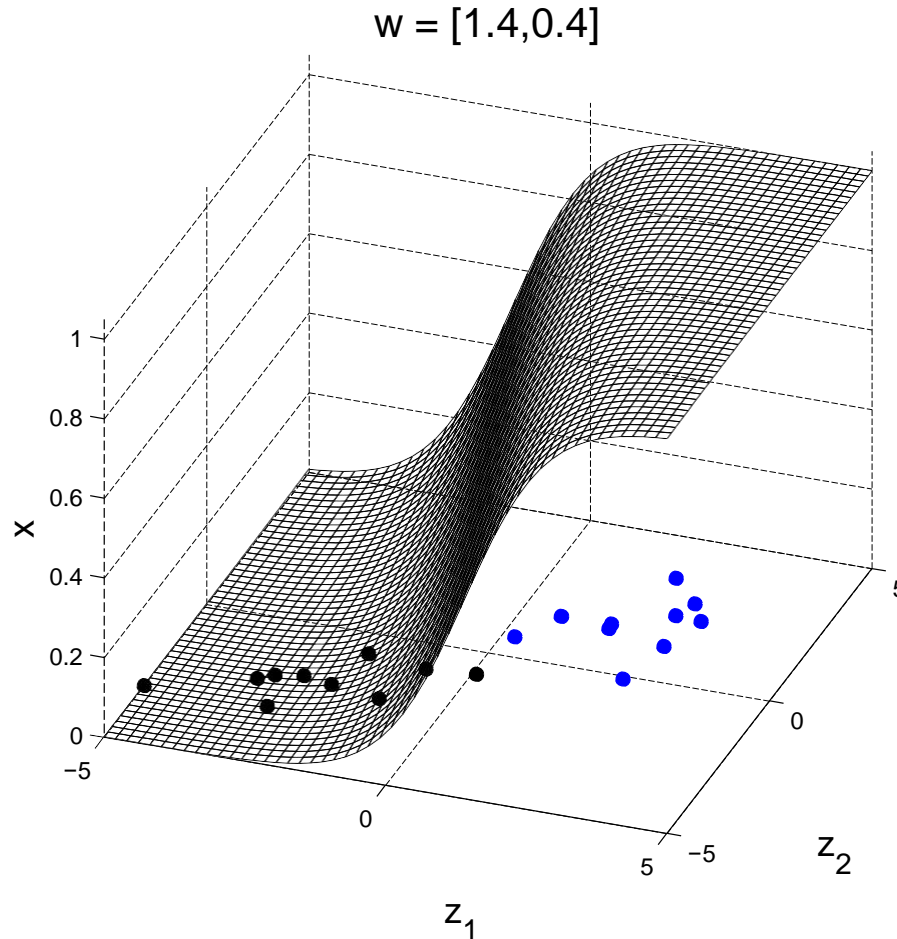
Training a Single Neuron



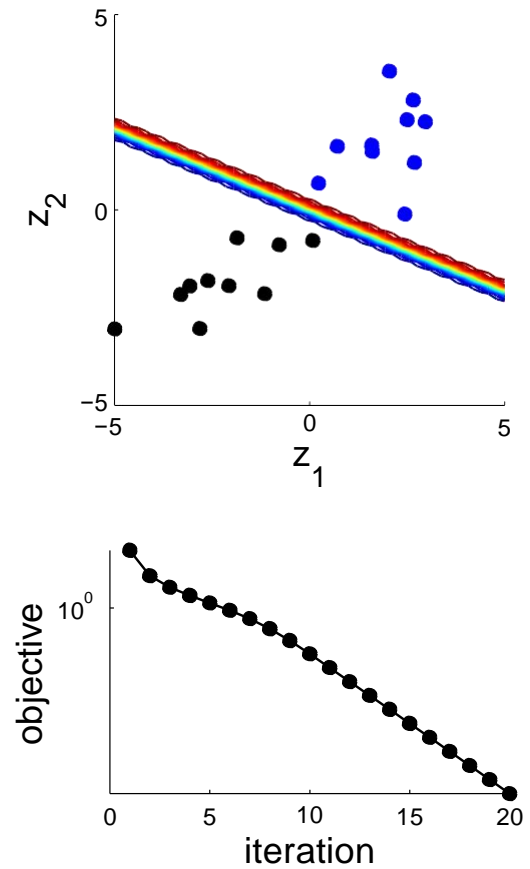
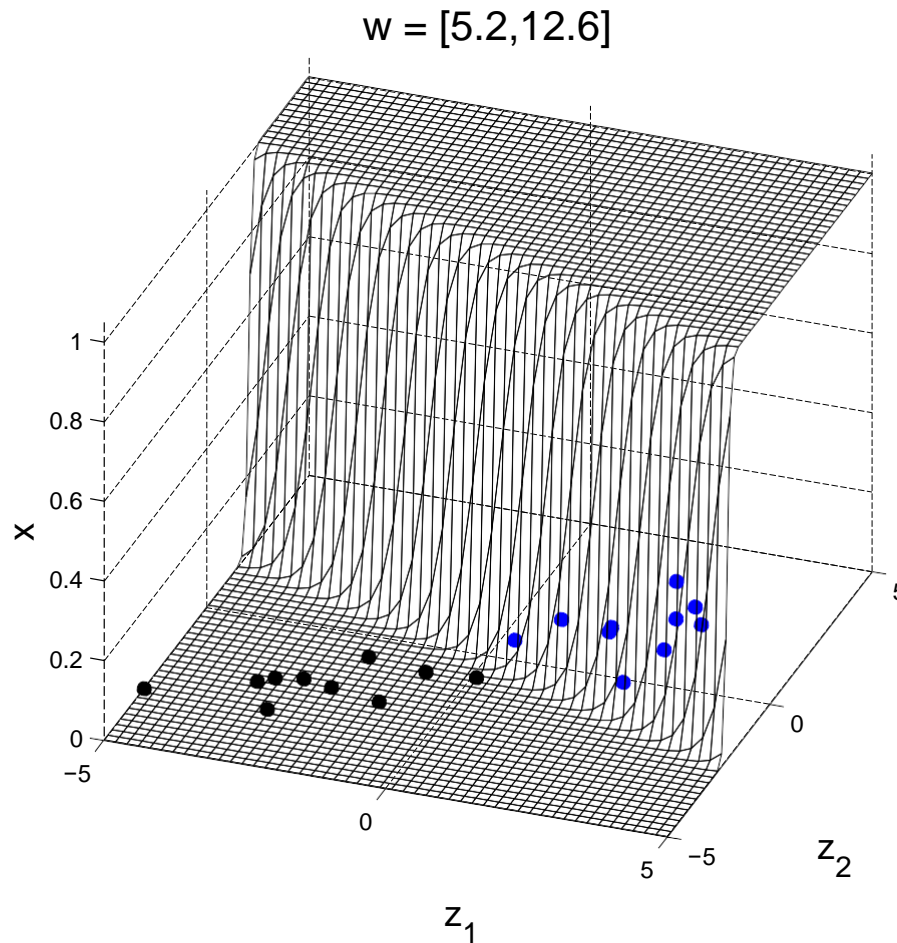
Training a Single Neuron



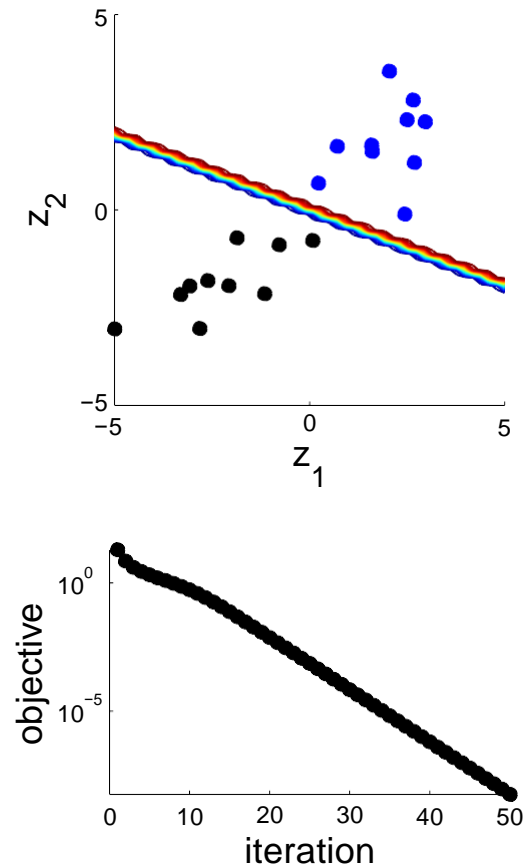
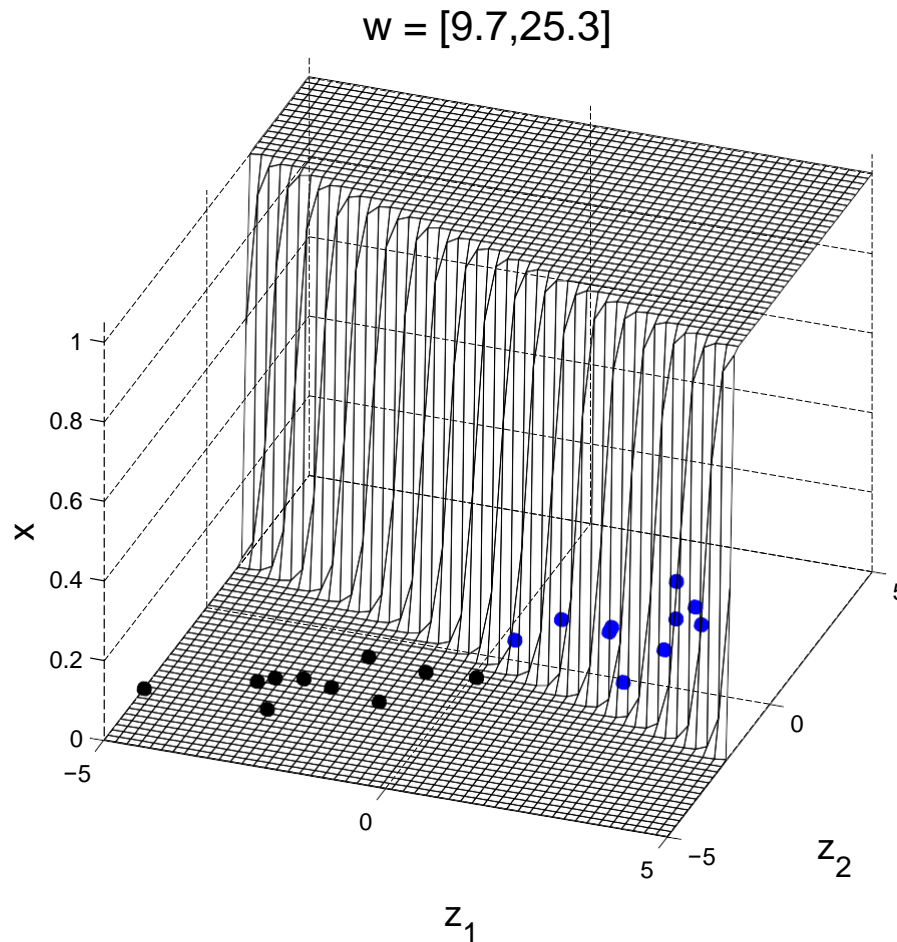
Training a Single Neuron



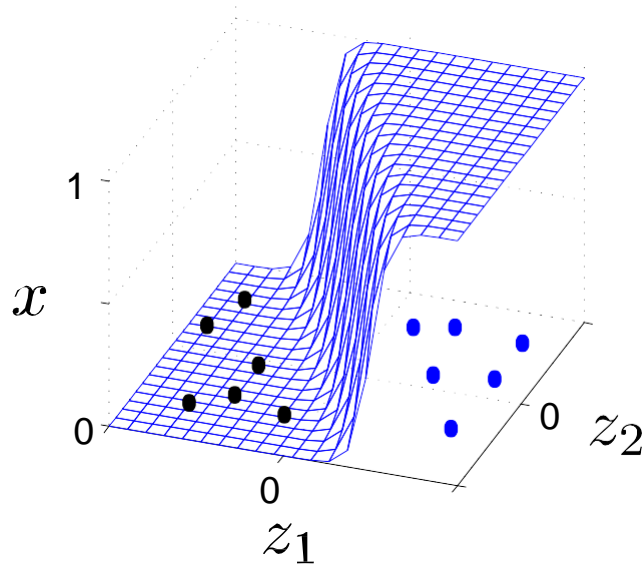
Training a Single Neuron



Training a Single Neuron



Overfitting and Weight Decay



training data

$$\{\mathbf{z}^{(n)}\}_{n=1}^N \quad \{t^{(n)}\}_{n=1}^N$$

inputs

class labels

objective function:

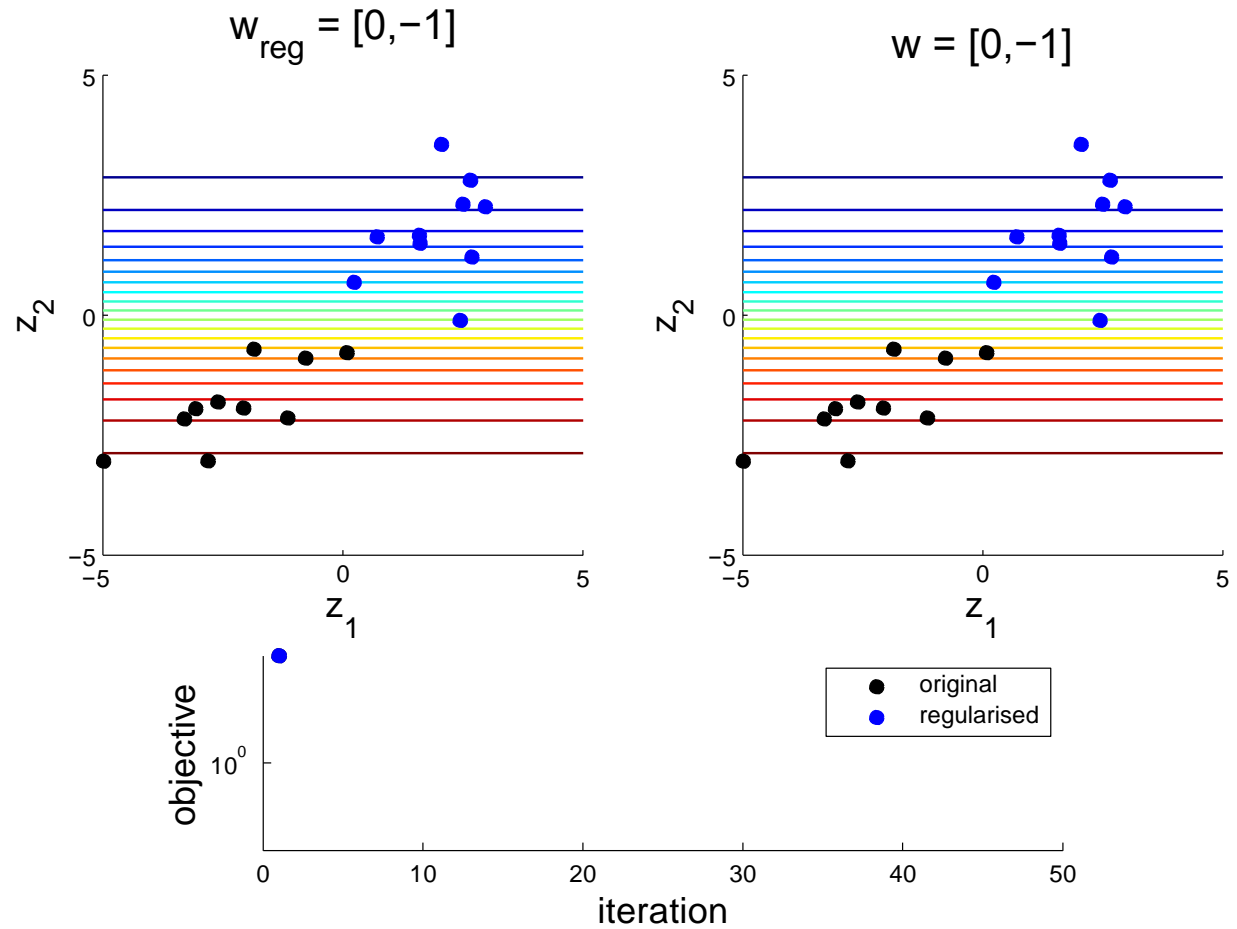
$$G(\mathbf{w}) = - \sum_n [t^{(n)} \log x(\mathbf{z}^{(n)}; \mathbf{w}) + (1 - t^{(n)}) \log (1 - x(\mathbf{z}^{(n)}; \mathbf{w}))]$$

$$E(\mathbf{w}) = \frac{1}{2} \sum_i w_i^2 \quad \text{regulariser discourages the network using extreme weights}$$

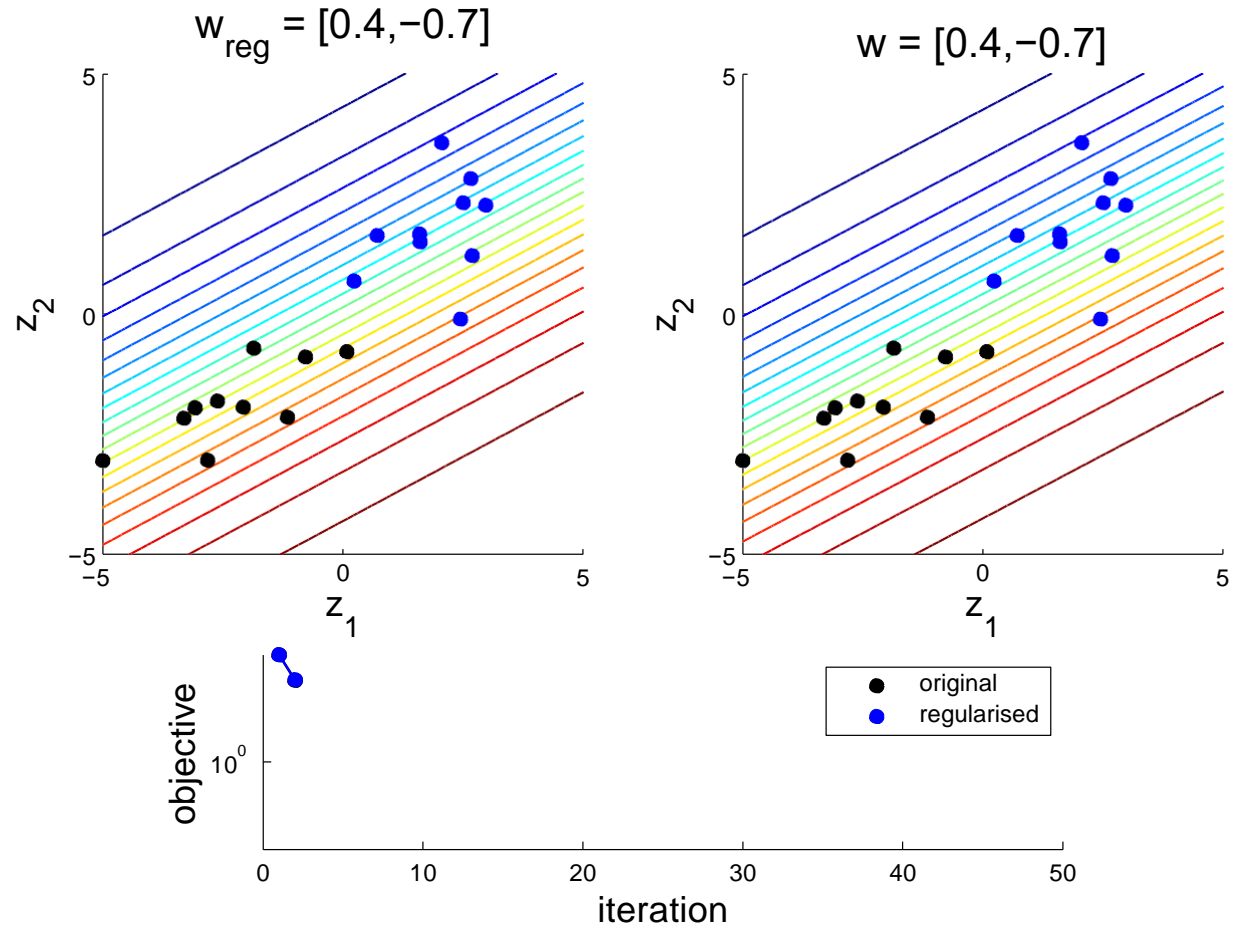
$$\mathbf{w}^* = \arg \min_{\mathbf{w}} M(\mathbf{w}) = \arg \min_{\mathbf{w}} [G(\mathbf{w}) + \alpha E(\mathbf{w})]$$

$$\frac{d}{d\mathbf{w}} M(\mathbf{w}) = - \sum_n (t^{(n)} - x^{(n)}) \mathbf{z}^{(n)} + \alpha \mathbf{w} \quad \text{weight decay - shrinks weights towards zero}$$

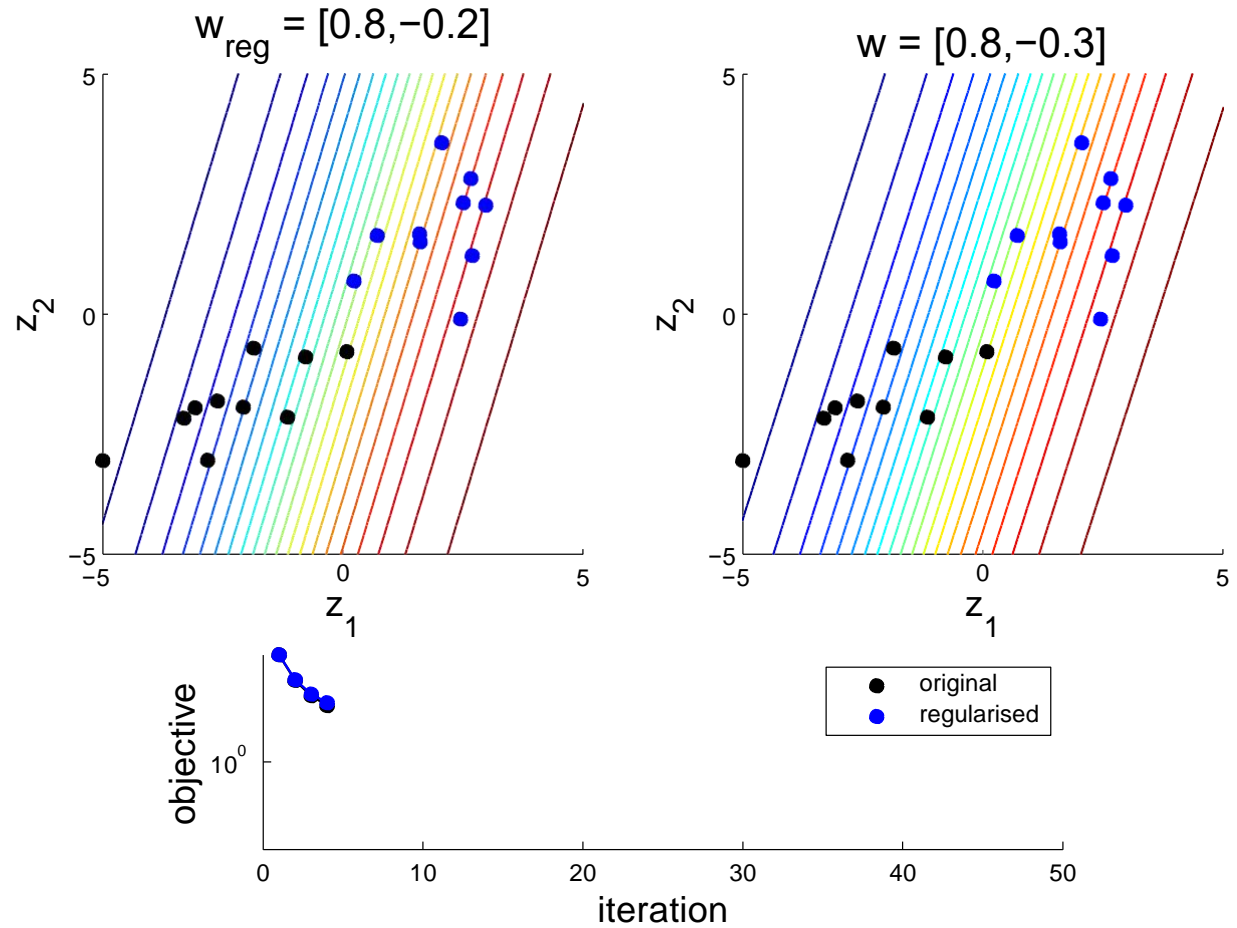
Training a Single Neuron (cont'd)



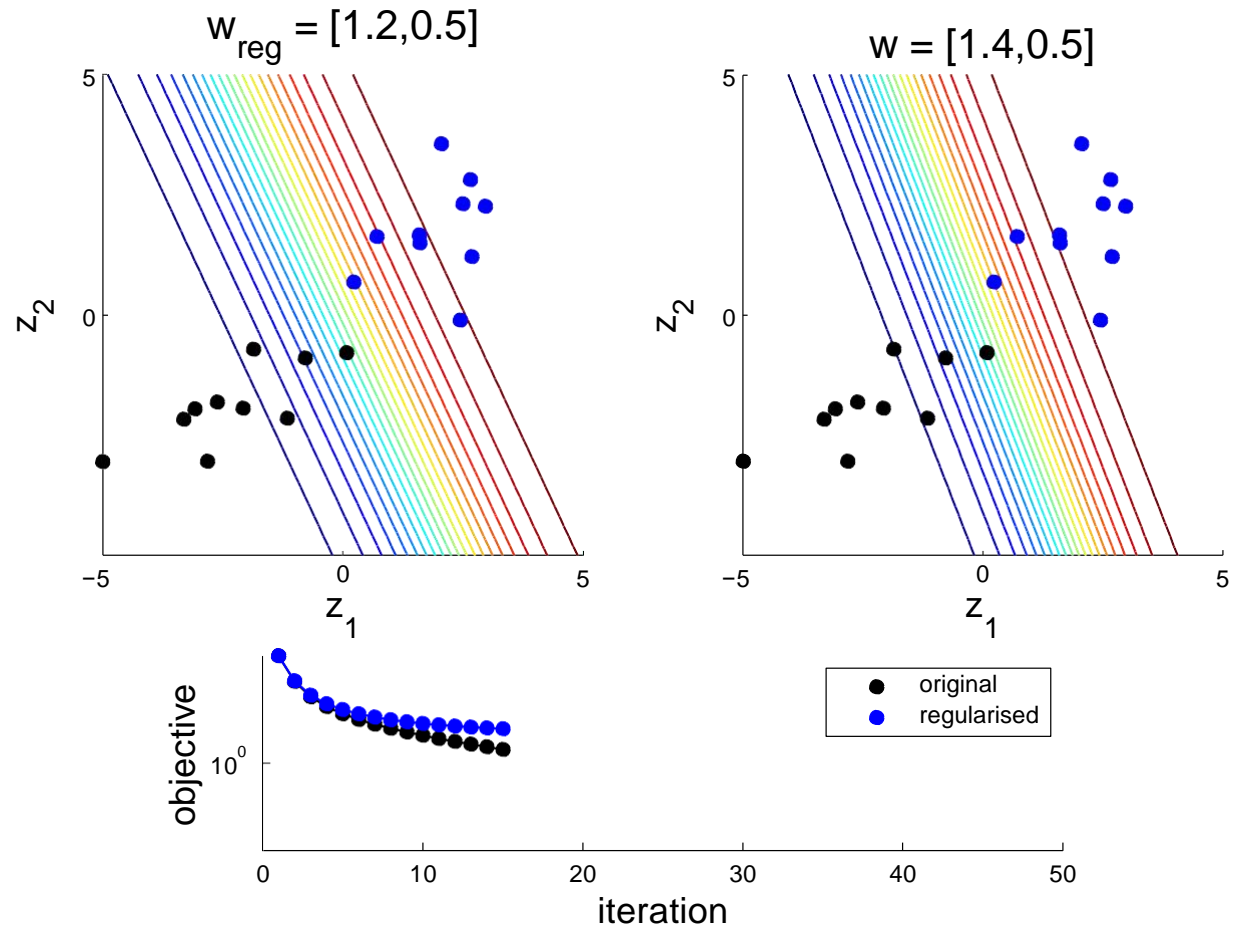
Training a Single Neuron (cont'd)



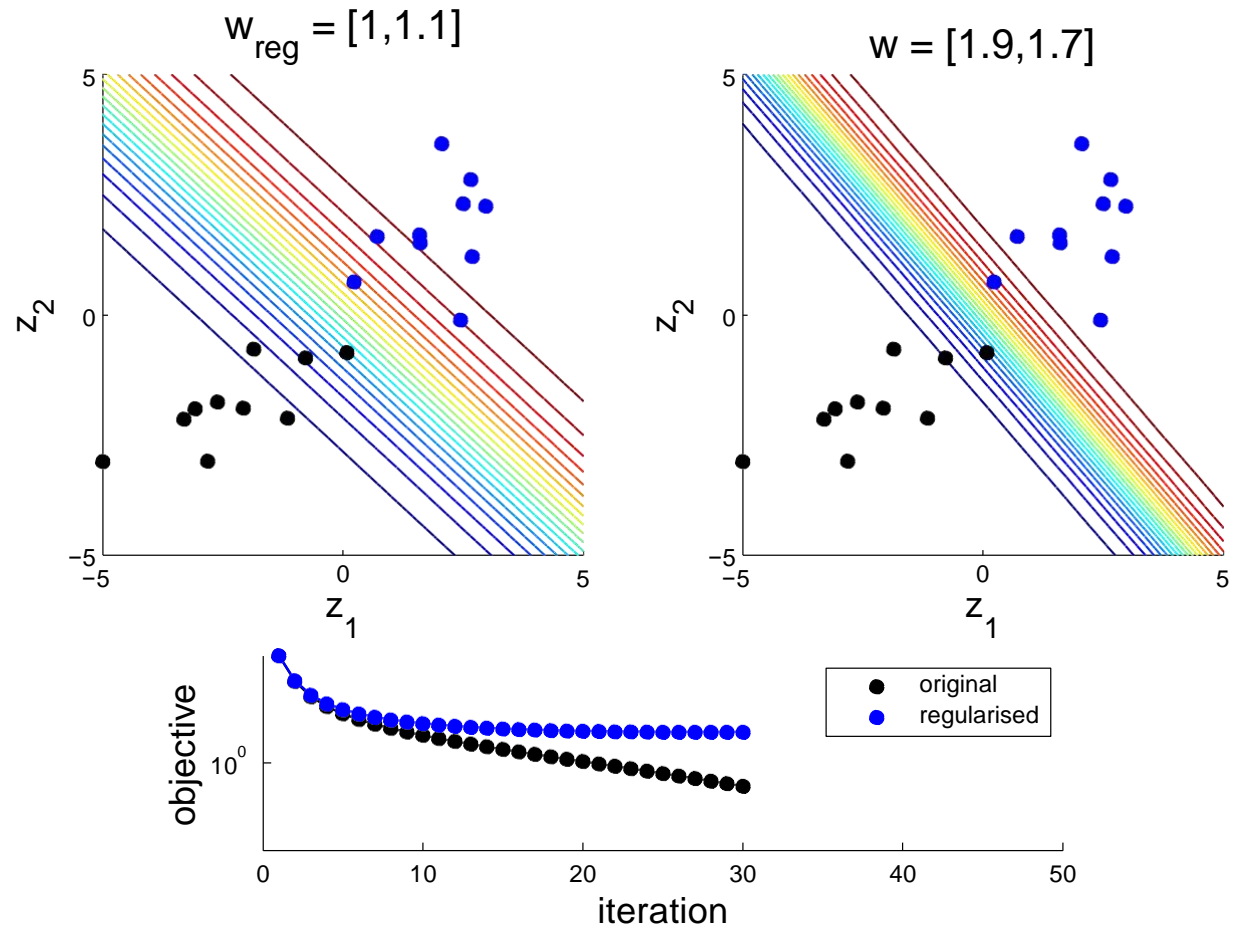
Training a Single Neuron (cont'd)



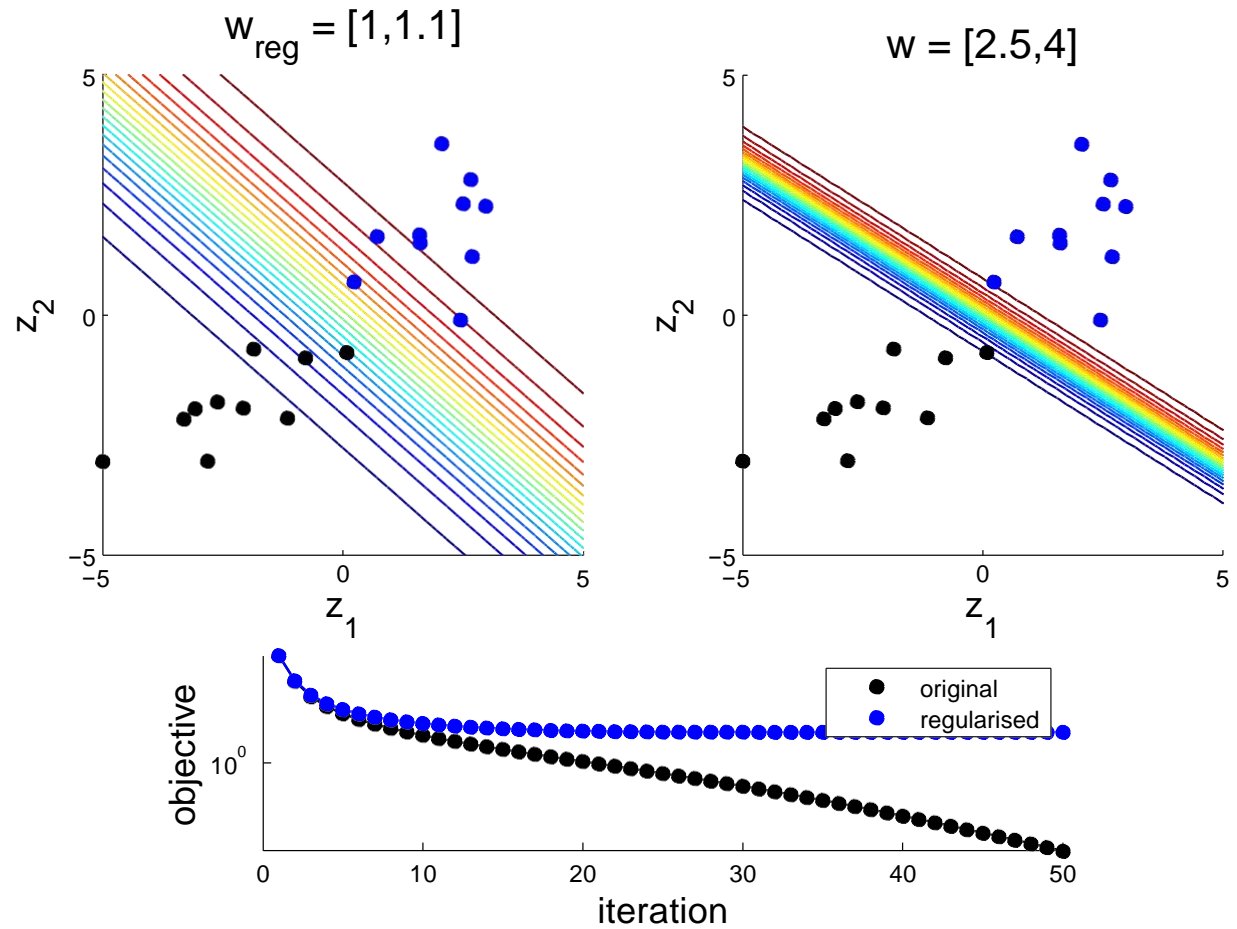
Training a Single Neuron (cont'd)



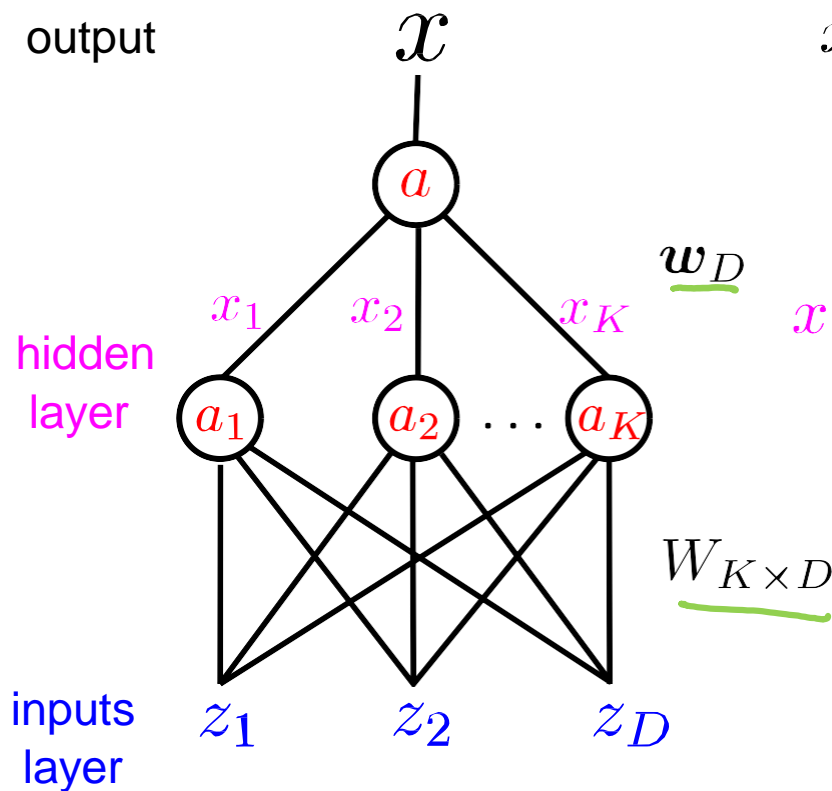
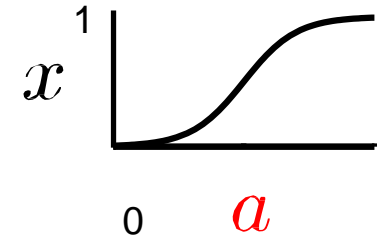
Training a Single Neuron (cont'd)



Training a Single Neuron (cont'd)



Single Hidden Layer Neural Networks



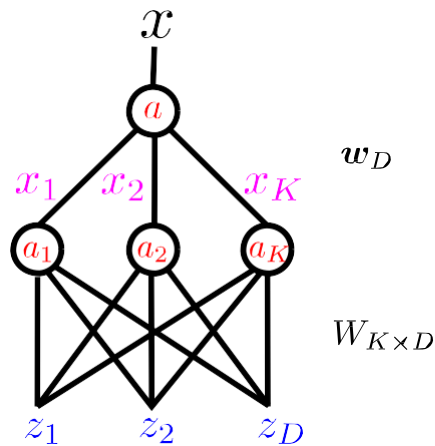
$$x(a) = \frac{1}{1 + \exp(-a)}$$

$$a = \sum_{k=1}^K w_k x_k$$

$$x(a_k) = \frac{1}{1 + \exp(-a_k)}$$

$$a_k = \sum_{d=1}^D W_{k,d} z_d$$

Training a Neural Network with a Single Hidden Layer



$$x(a) = \frac{1}{1 + \exp(-a)}$$

$$a = \sum_{k=1}^K w_k x_k$$

$$x(a_k) = \frac{1}{1 + \exp(-a_k)}$$

$$a_k = \sum_{d=1}^D W_{k,d} z_d$$

objective function:

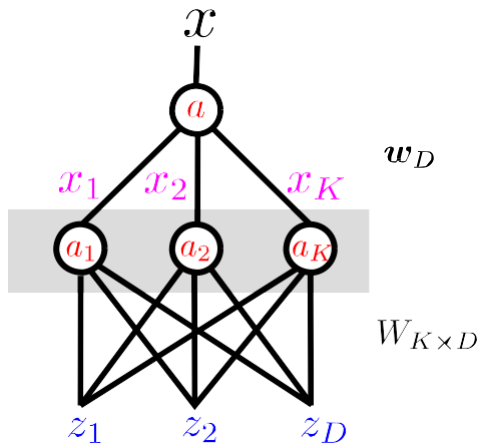
$$G(W, \mathbf{w}) = - \sum_n [t^{(n)} \log x^{(n)} + (1 - t^{(n)}) \log (1 - x^{(n)})] \text{ likelihood same as before}$$

$$E(W, \mathbf{w}) = \frac{1}{2} \sum_i w_i^2 + \frac{1}{2} \sum_{ij} W_{ij}^2 \quad \text{regulariser discourages extreme weights}$$

$$\{W, \mathbf{w}^*\} = \arg \min_{W, \mathbf{w}} M(W, \mathbf{w}) = \arg \min_{W, \mathbf{w}} [G(W, \mathbf{w}) + \alpha E(W, \mathbf{w})]$$

Training a Neural Network with a Single Hidden Layer

Networks with hidden layers can be fit using gradient descent using an algorithm called **back-propagation**.



$$x(a) = \frac{1}{1 + \exp(-a)}$$

$$a = \sum_{k=1}^K w_k x_k$$

$$x(a_k) = \frac{1}{1 + \exp(-a_k)}$$

$$a_k = \sum_{d=1}^D W_{k,d} z_d$$

objective function:

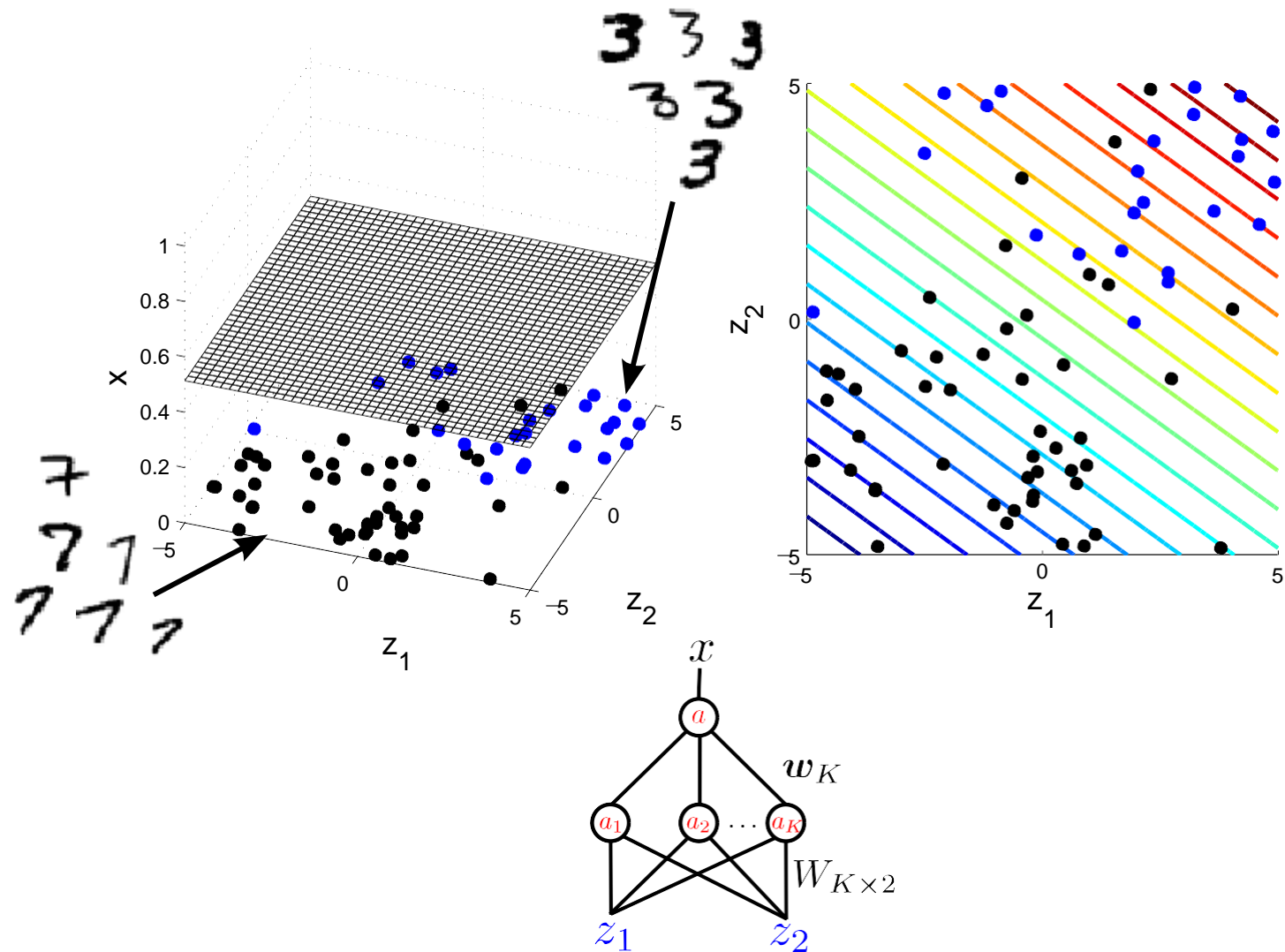
$$G(W, \mathbf{w}) = - \sum_n [t^{(n)} \log x^{(n)} + (1 - t^{(n)}) \log (1 - x^{(n)})] \text{ likelihood same as before}$$

$$E(W, \mathbf{w}) = \frac{1}{2} \sum_i w_i^2 + \frac{1}{2} \sum_{ij} W_{ij}^2 \quad \text{regulariser discourages extreme weights}$$

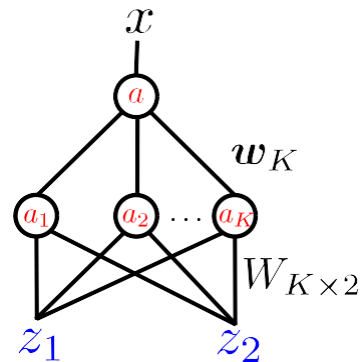
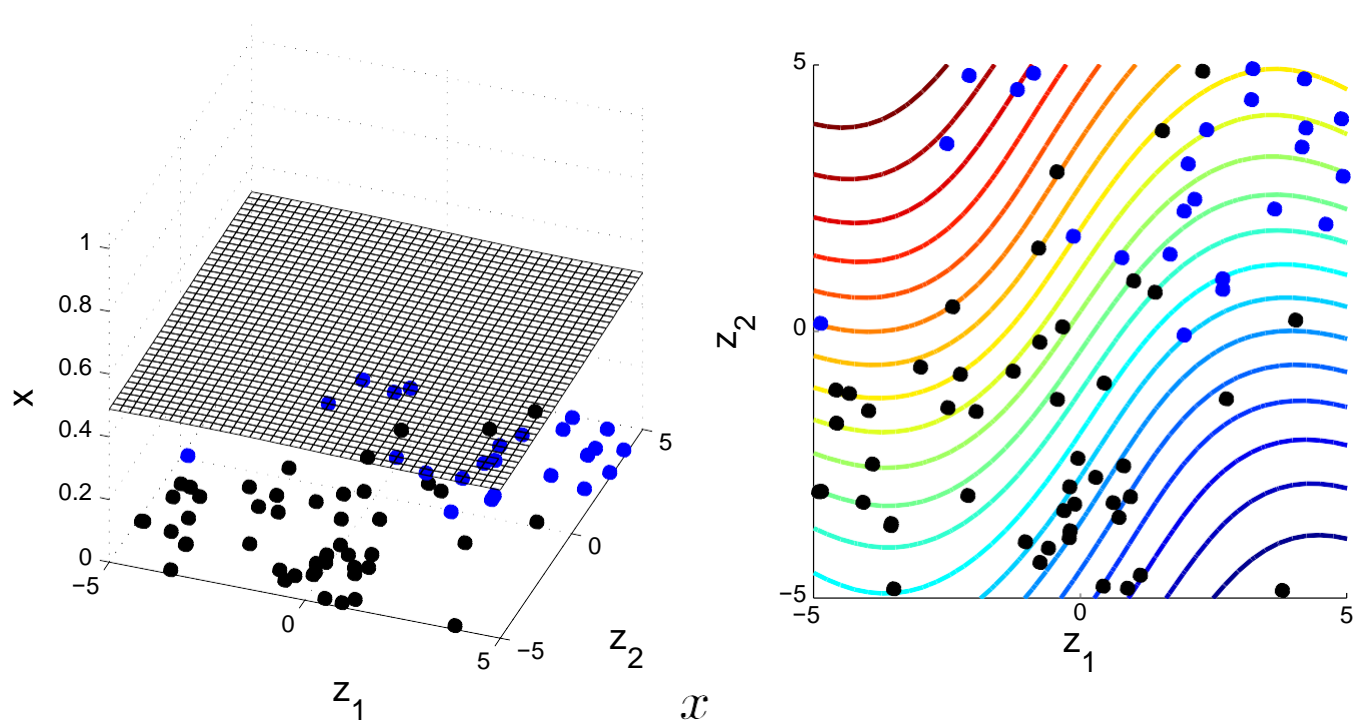
$$\{W, \mathbf{w}^*\} = \arg \min_{W, \mathbf{w}} M(W, \mathbf{w}) = \arg \min_{W, \mathbf{w}} [G(W, \mathbf{w}) + \alpha E(W, \mathbf{w})]$$

$$\begin{aligned} \frac{dG(W, \mathbf{w})}{dW_{ij}} &= \sum_n \frac{dG(W, \mathbf{w})}{dx^{(n)}} \frac{dx^{(n)}}{dW_{ij}} = \sum_n \frac{dG(W, \mathbf{w})}{dx^{(n)}} \frac{da^{(n)}}{dW_{ij}} \\ &= \sum_n \frac{dG(W, \mathbf{w})}{dx^{(n)}} \frac{dx^{(n)}}{da^{(n)}} \frac{da^{(n)}}{dx_i^{(n)}} \frac{dx_i^{(n)}}{dW_{ij}} = \sum_n \frac{dG(W, \mathbf{w})}{dx^{(n)}} \frac{dx^{(n)}}{da^{(n)}} \frac{da^{(n)}}{dx_i^{(n)}} \frac{da_i^{(n)}}{dW_{ij}} \end{aligned}$$

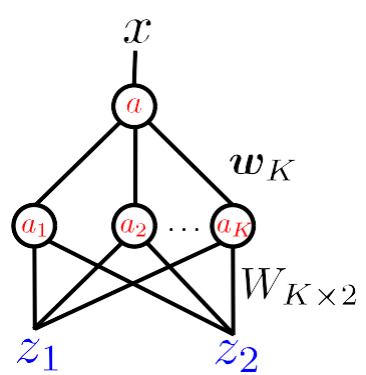
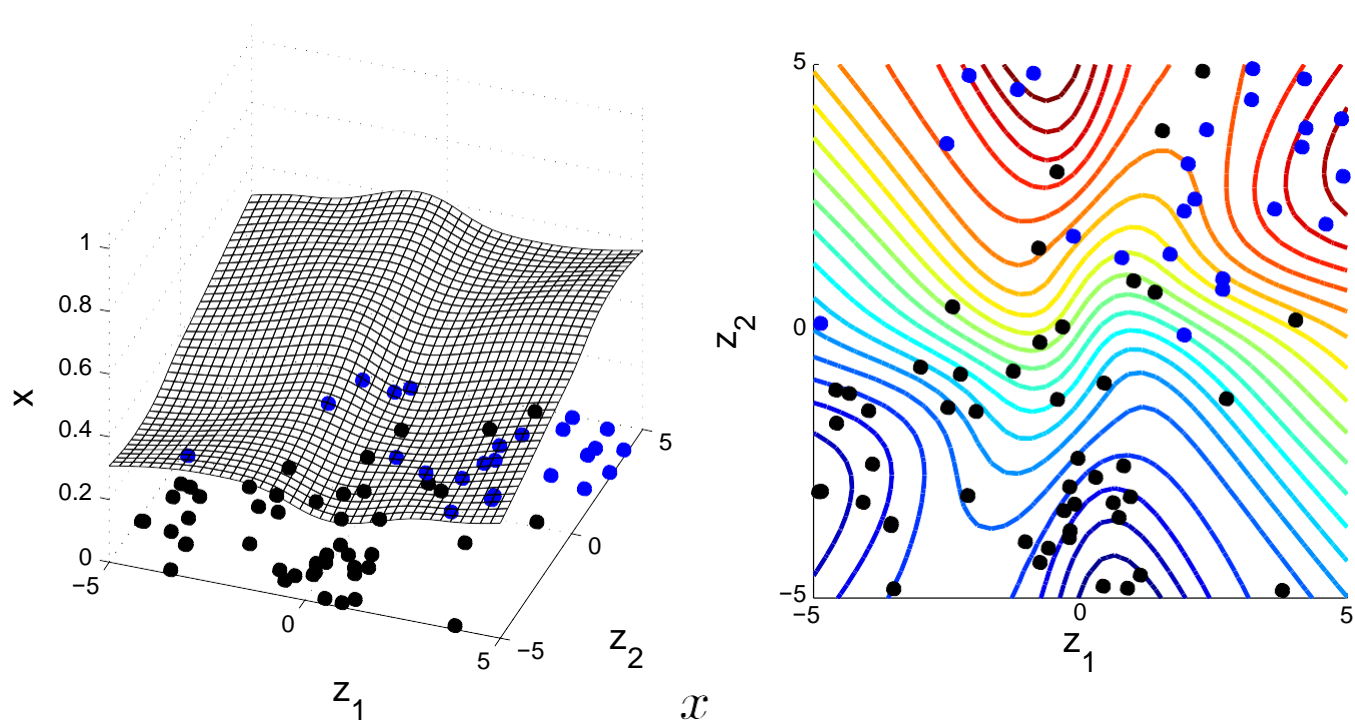
Training a Neural Network with a Single Hidden Layer



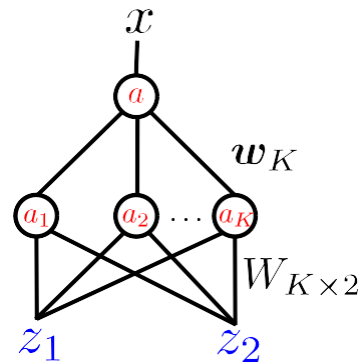
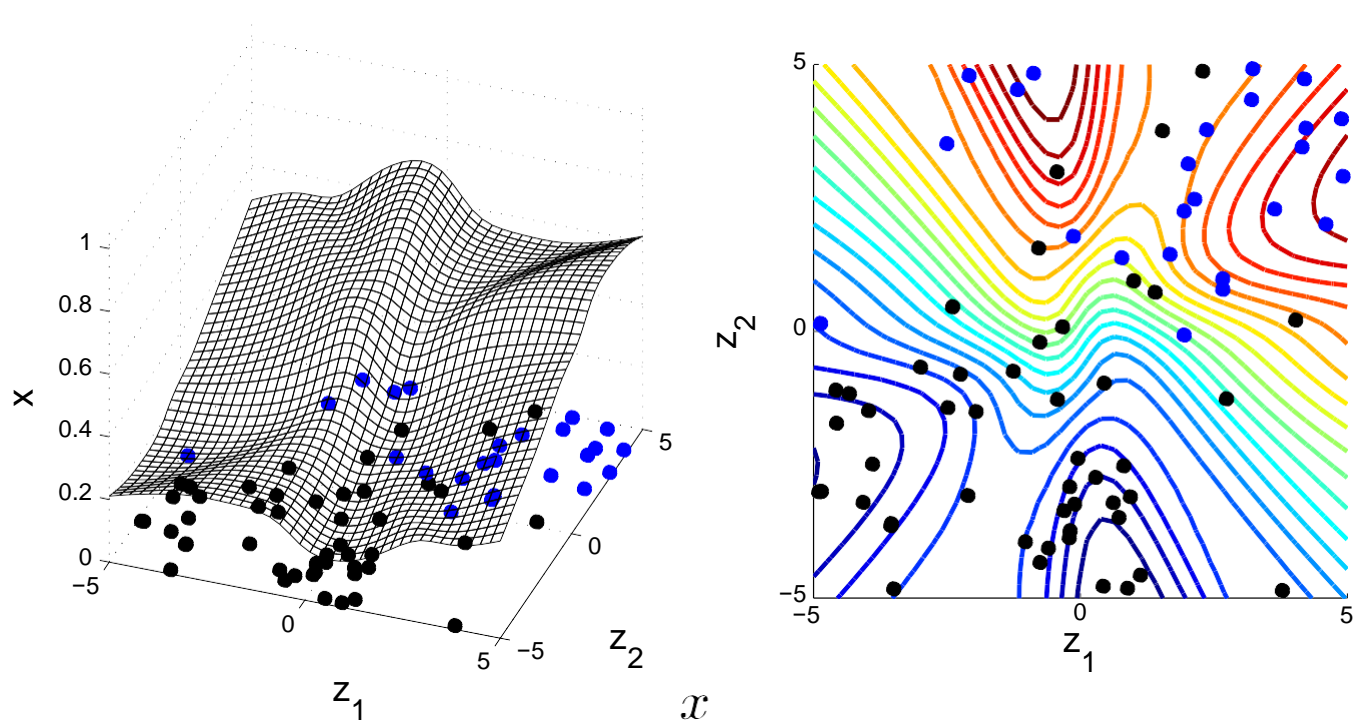
Training a Neural Network with a Single Hidden Layer



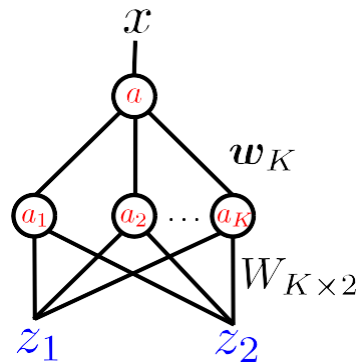
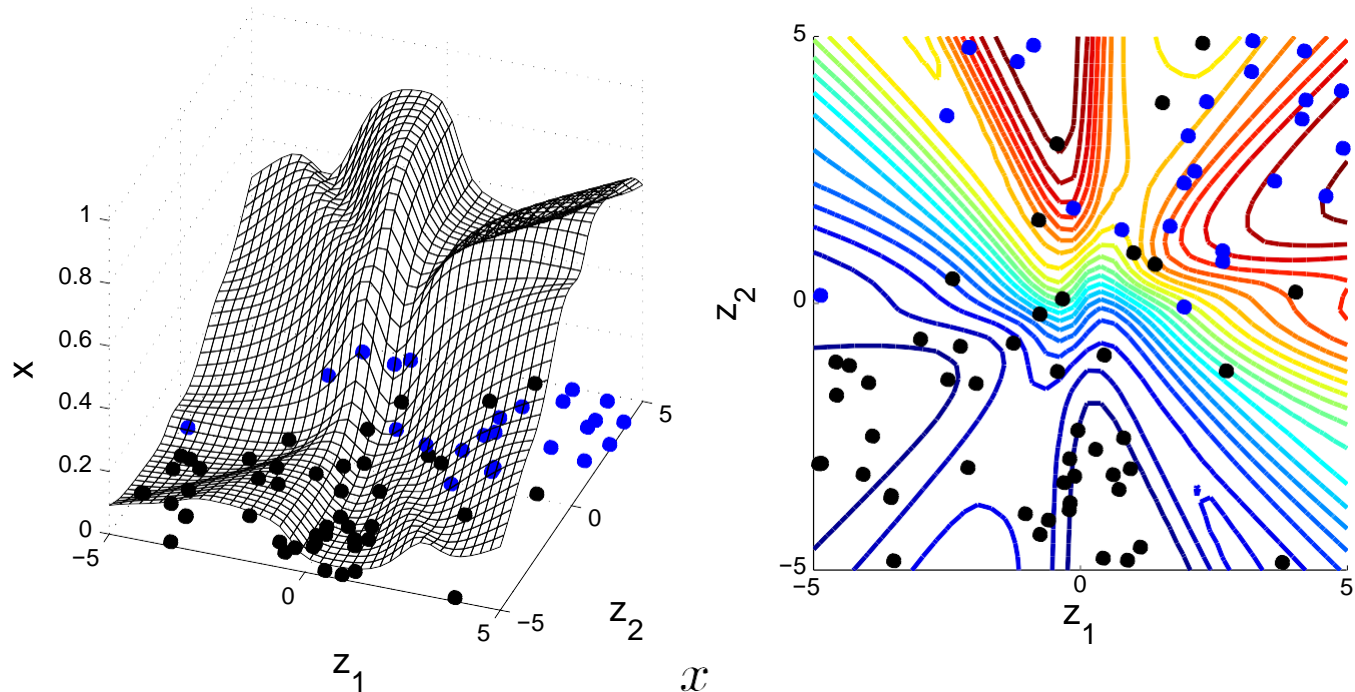
Training a Neural Network with a Single Hidden Layer



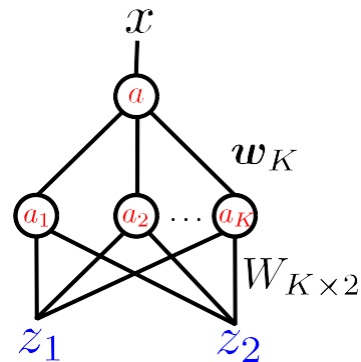
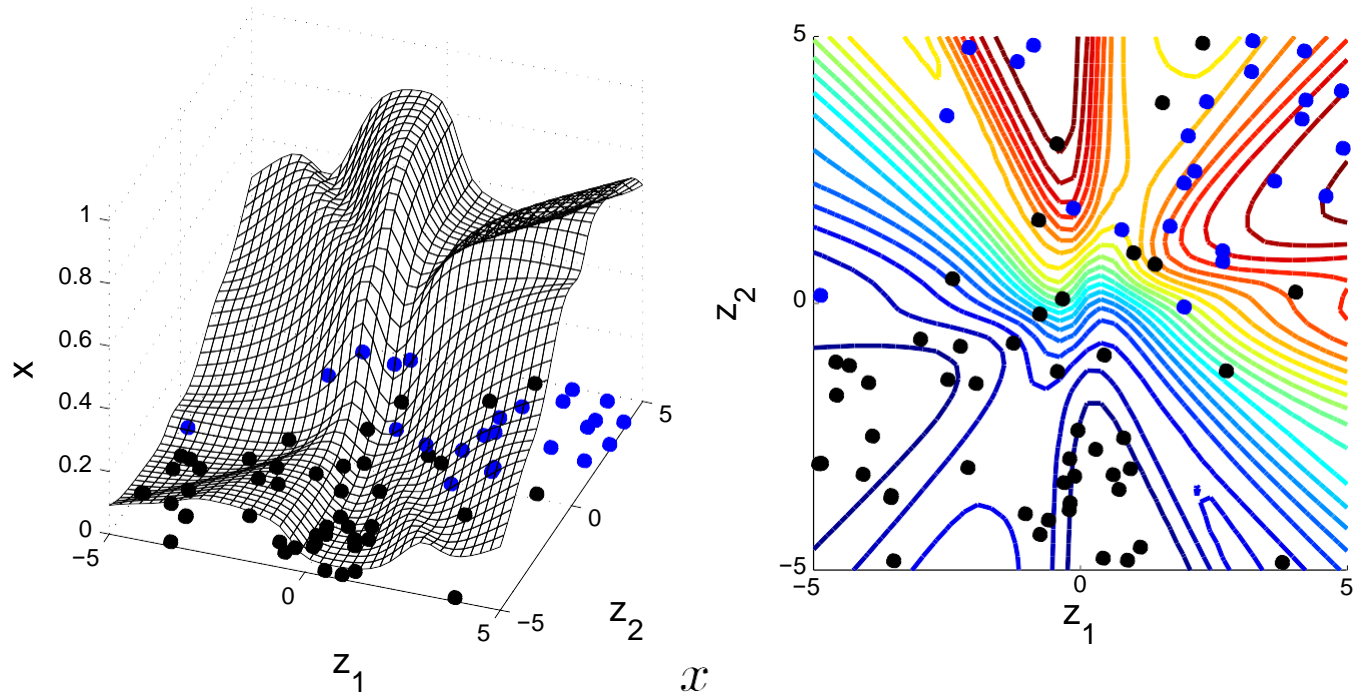
Training a Neural Network with a Single Hidden Layer



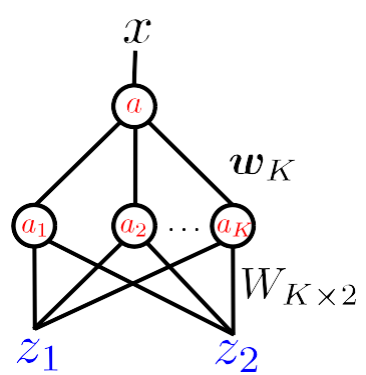
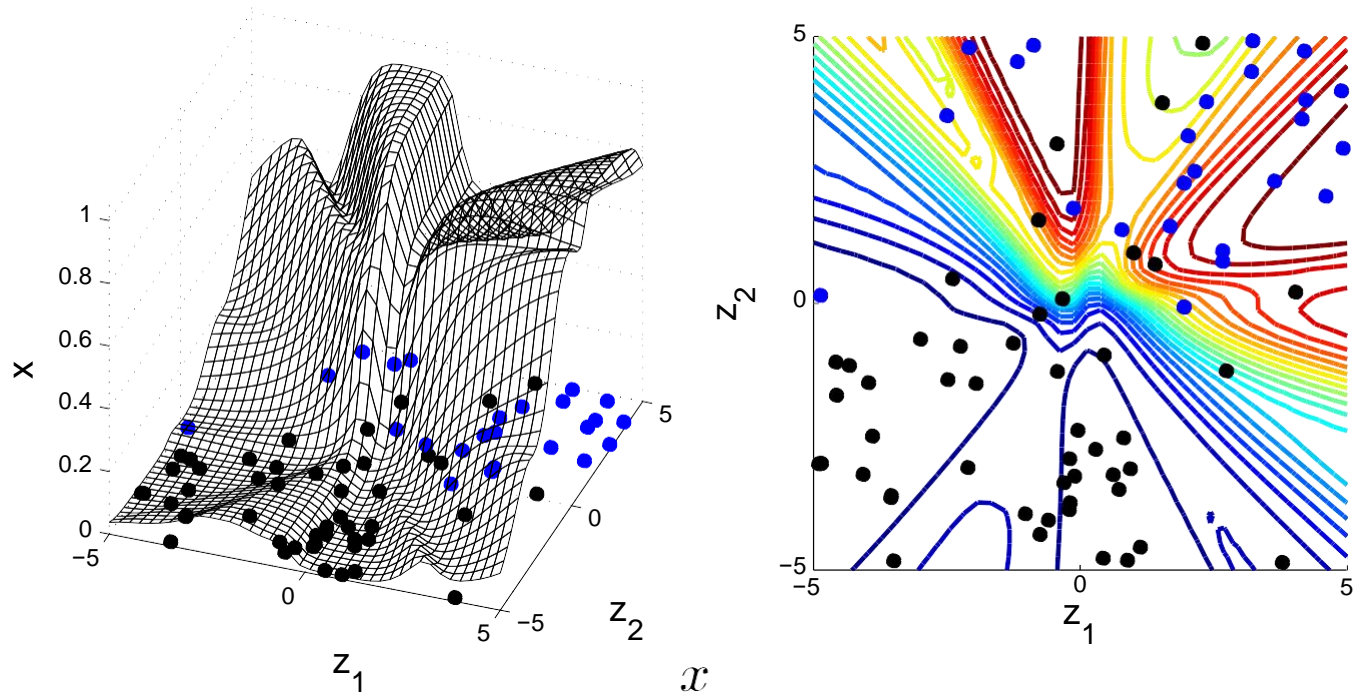
Training a Neural Network with a Single Hidden Layer



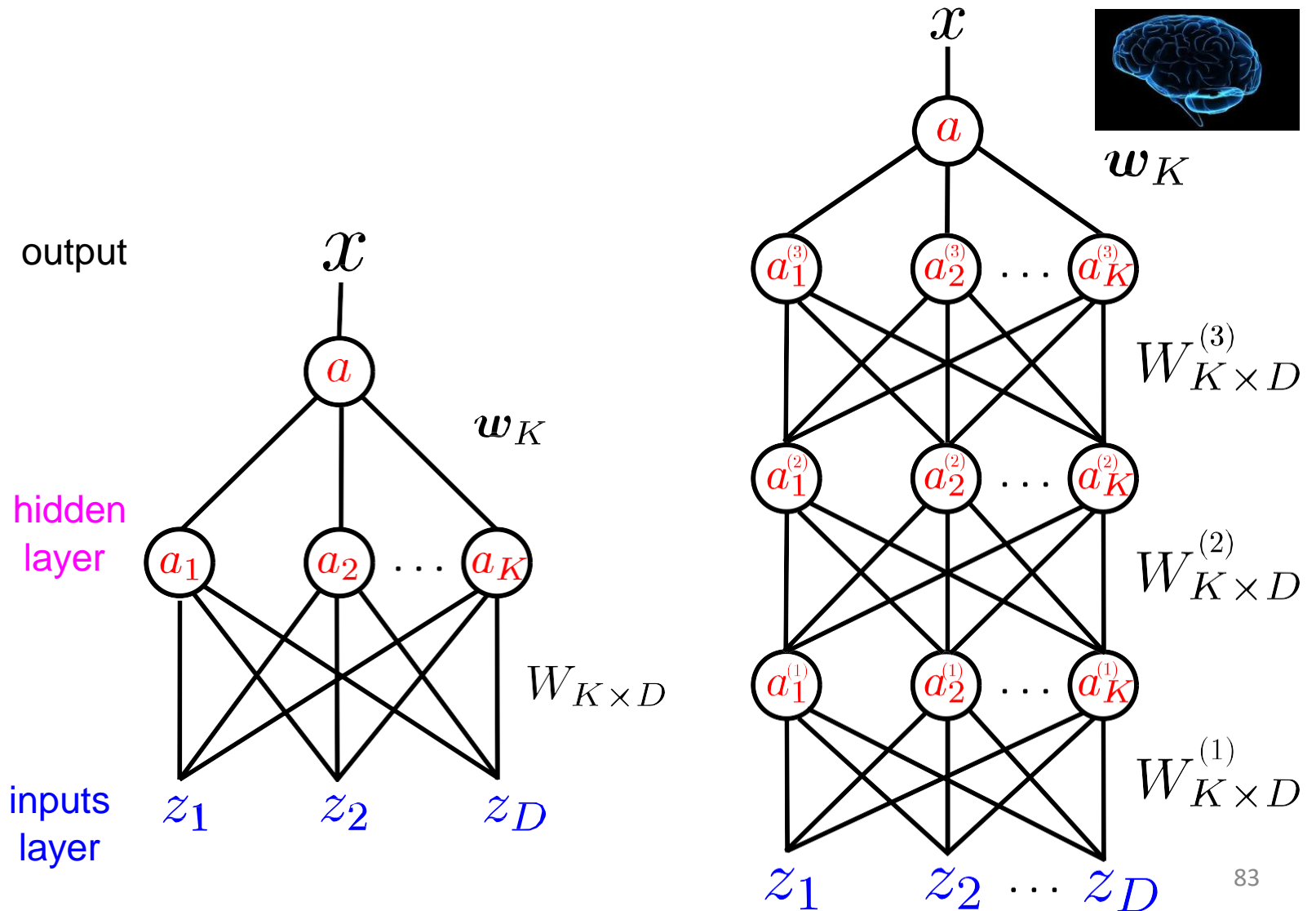
Training a Neural Network with a Single Hidden Layer



Training a Neural Network with a Single Hidden Layer



Hierarchical Models with Many Layers



What We've Covered Today...

- Unsupervised vs. Supervised Learning
 - Clustering & Dimension Reduction
 - Training, testing, & validation
 - Linear Classification
 - From Linear Classifier to Neural Nets

