

Deep Learning for Computer Vision

Spring 2019

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

<https://ceiba.ntu.edu.tw/1072CommE5052> (grade, etc.)

FB: [DLCV Spring 2019](#)

Yu-Chiang Frank Wang 王鈺強, Associate Professor

Dept. Electrical Engineering, National Taiwan University

Distance Metric of p-Norm

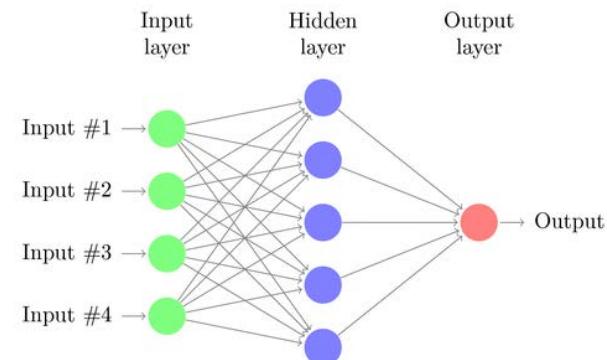
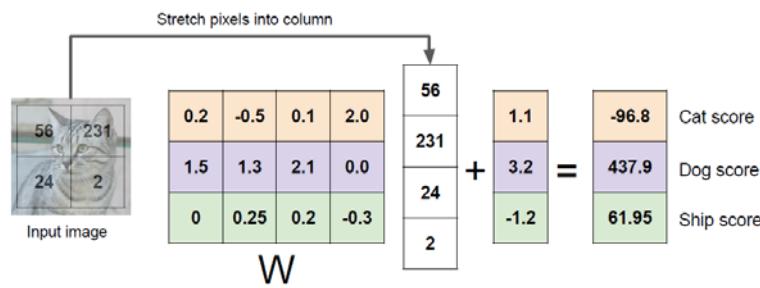
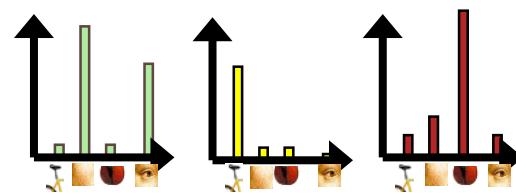
- Remarks
 - We have $\|x\|_p = \left(\sum_{i=1}^d |x_i|^p\right)^{1/p}$ and L_0 norm: $\|x\|_0$
 - Strictly speaking, L_0 norm is not a norm.
 - It is a cardinality function which has its definition in the form of L_p norm.
 - In practice, for L_0 norm, most engineers use the following definition instead:

$$\|x\|_0 = \lim_{p \rightarrow 0} \sum_{i=1}^d |x_i|^p \text{ or } \|x\|_0 = \#\{i | x_i \neq 0\}$$

In other words, it returns the total number of non-zero entries in x .

What's to Be Covered Today...

- General Framework for Visual Classification
 - Bag-of-Words Models as Image Features
- Intro to Neural Networks & CNN
 - Linear Classification
 - Neural Network for Machine Vision
 - Multi-Layer Perceptron
 - Convolutional Neural Networks



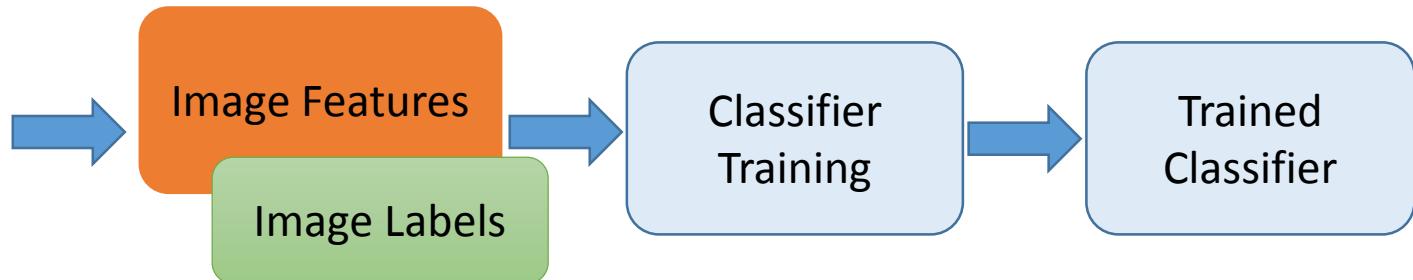
Supervised Learning for Visual Classification

- Training vs. Testing Phases

Training Images



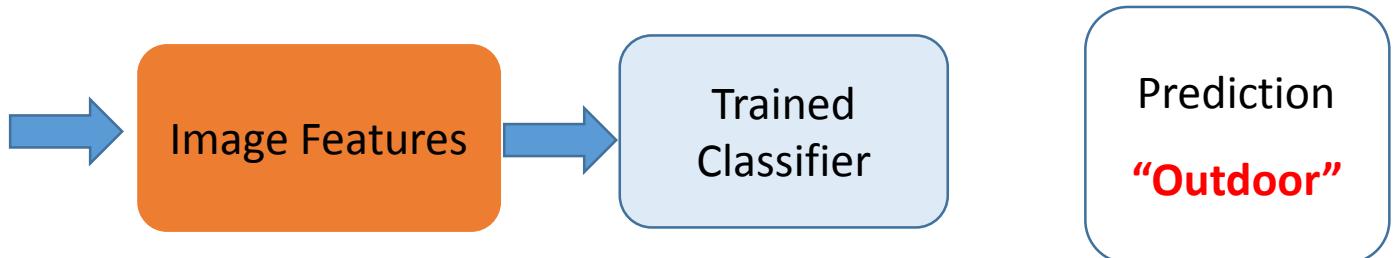
Training



Testing



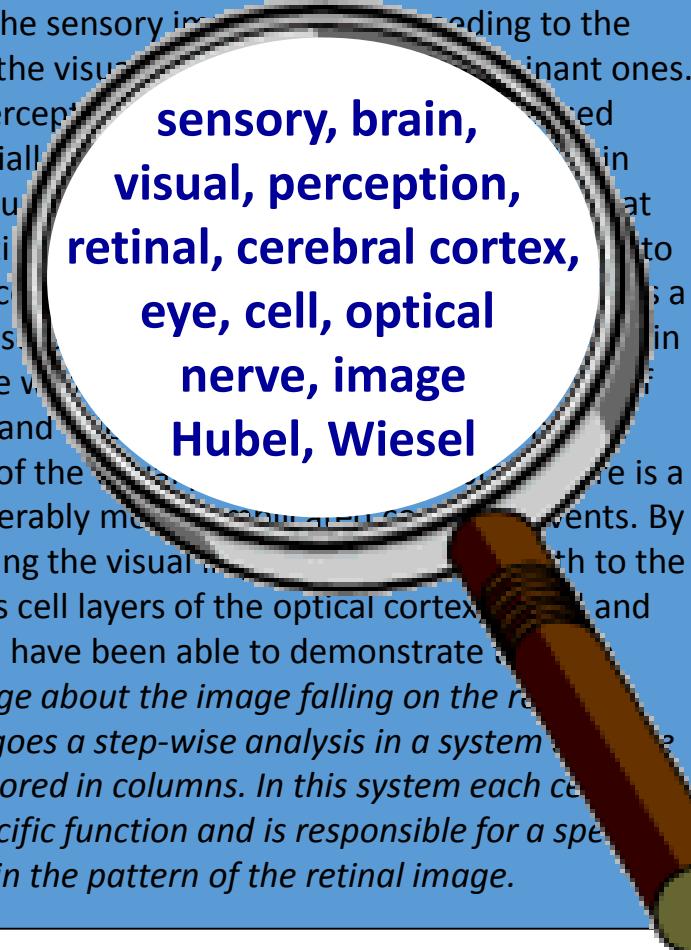
Test Image



Bag-of-Words Models for Image Classification

- Analogy to document categorization

Of all the sensory information entering the brain, the visual system is the dominant one. Our perception of the world is based essentially on what we see. Light from our environment enters the eye via the retina and is converted into visual cues such as movie screens, faces, the eye wall, Hubel and Wiesel have shown that the origin of the visual system is considerably more complex than previously thought. By following the visual pathway from the eye to the various cell layers of the optical cortex, Hubel and Wiesel have been able to demonstrate a *step-wise analysis*. A message about the image falling on the retina undergoes a step-wise analysis in a system of cells stored in columns. In this system each cell has its specific function and is responsible for a specific detail in the pattern of the retinal image.



China is forecasting a trade surplus of \$90bn (£51bn) to \$100bn this year, up from \$75bn in 2004 and \$32bn. The Chinese government is concerned that the surplus would be causing inflation. Exports to the US are imports from China. This has annoyed the US, which exports more than it imports. Undervaluation of the yuan is high, but the Chinese government also needed to increase the value of the yuan against the dollar by 2.1% in July and permitted it to trade within a narrow band, but the US wants the yuan to be allowed to trade freely. However, Beijing has made it clear that it will take its time and tread carefully before allowing the yuan to rise further in value.



Bag of Words (or Visual Words)

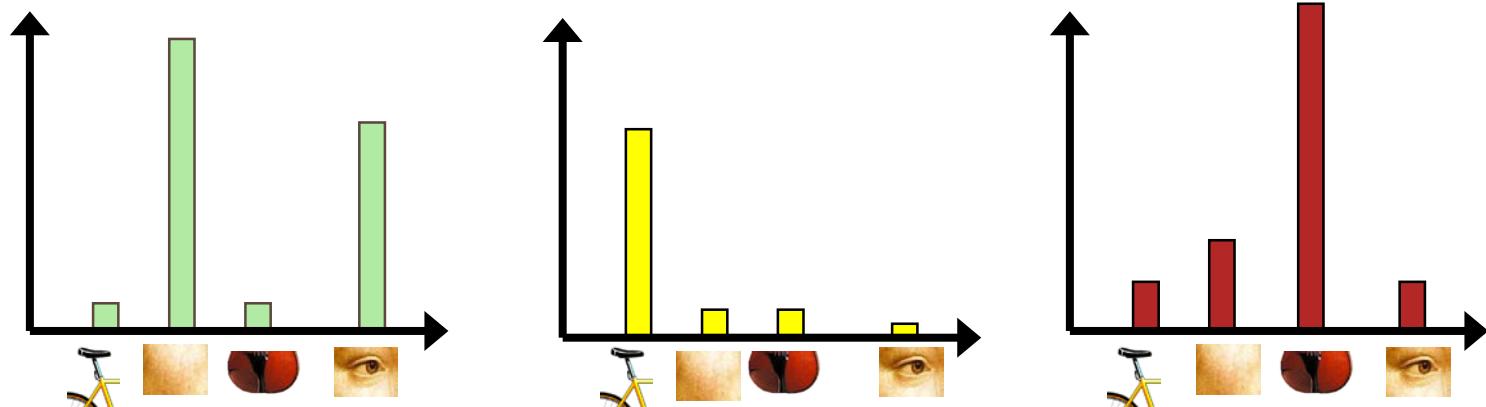


Image Representation: Histograms

- Take images with 2D features/descriptors as an example

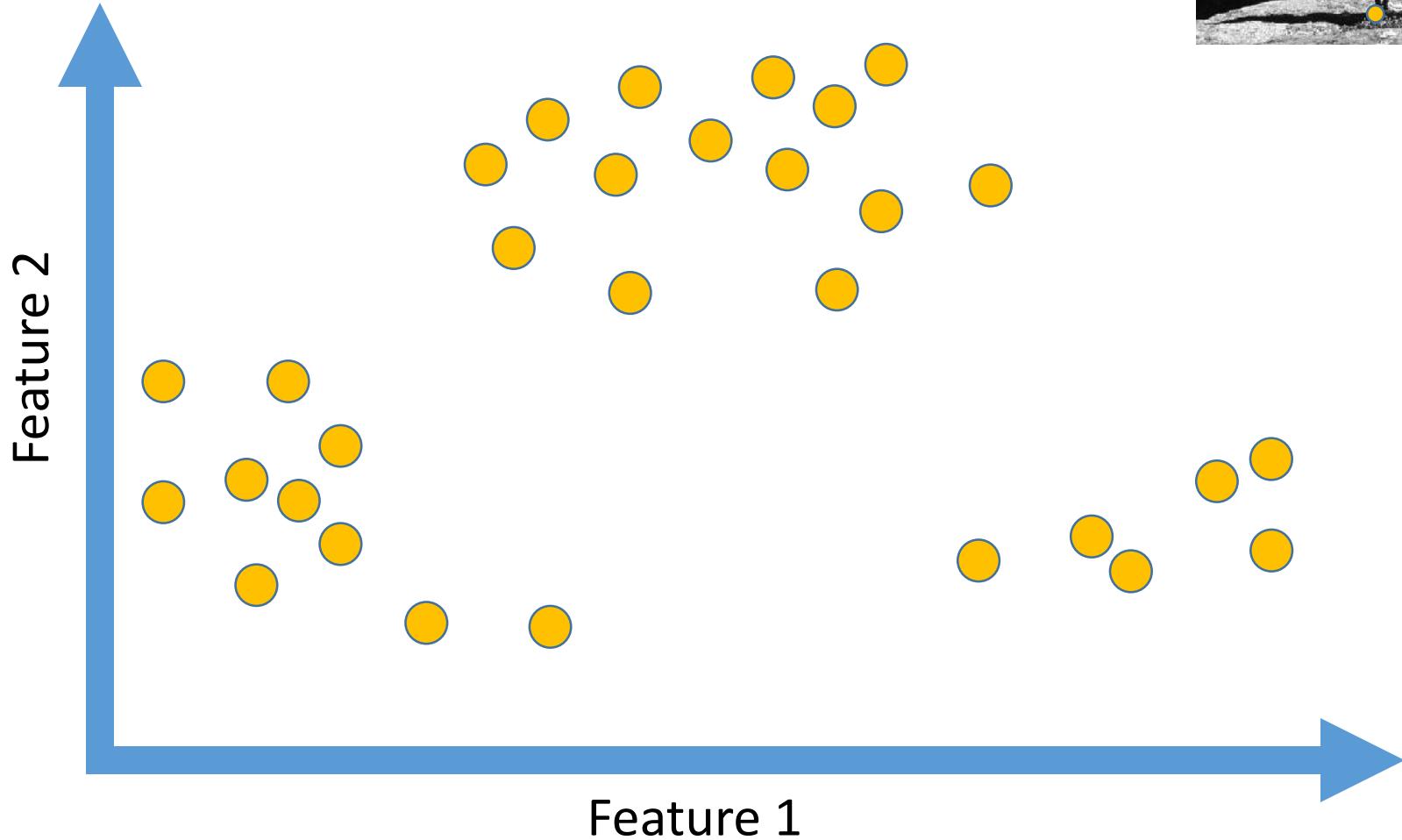
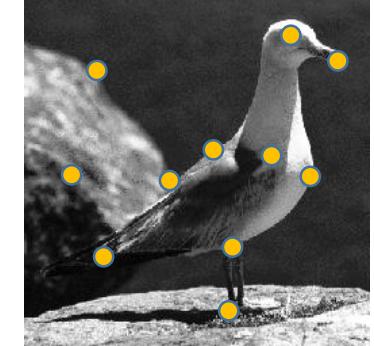


Image Representation: Histograms

- # of occurrence of data in each bin
- Marginal histogram of feature 1

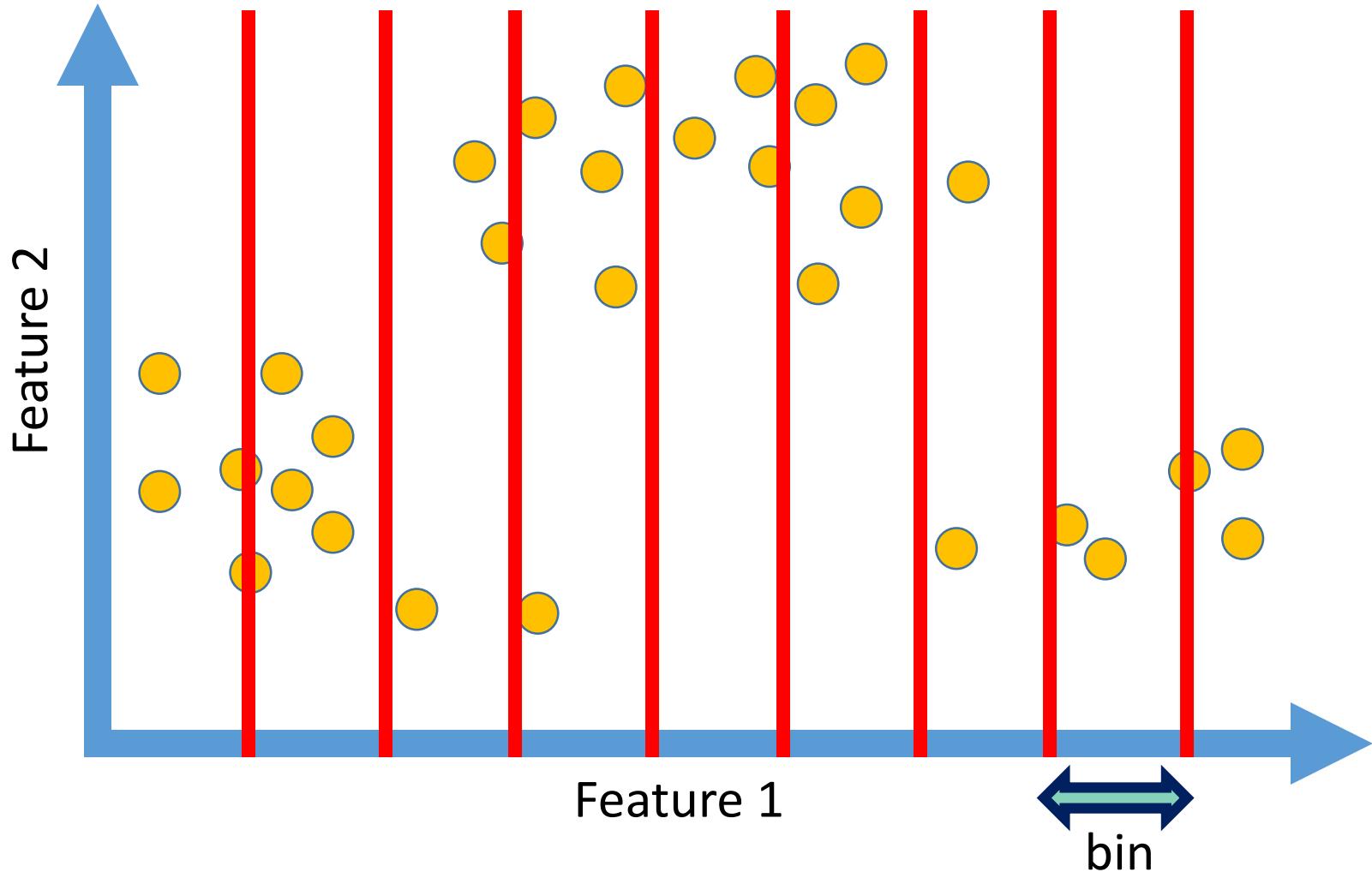


Image Representation: Histograms

- # of occurrence of data in each bin
- Marginal histogram of feature 2

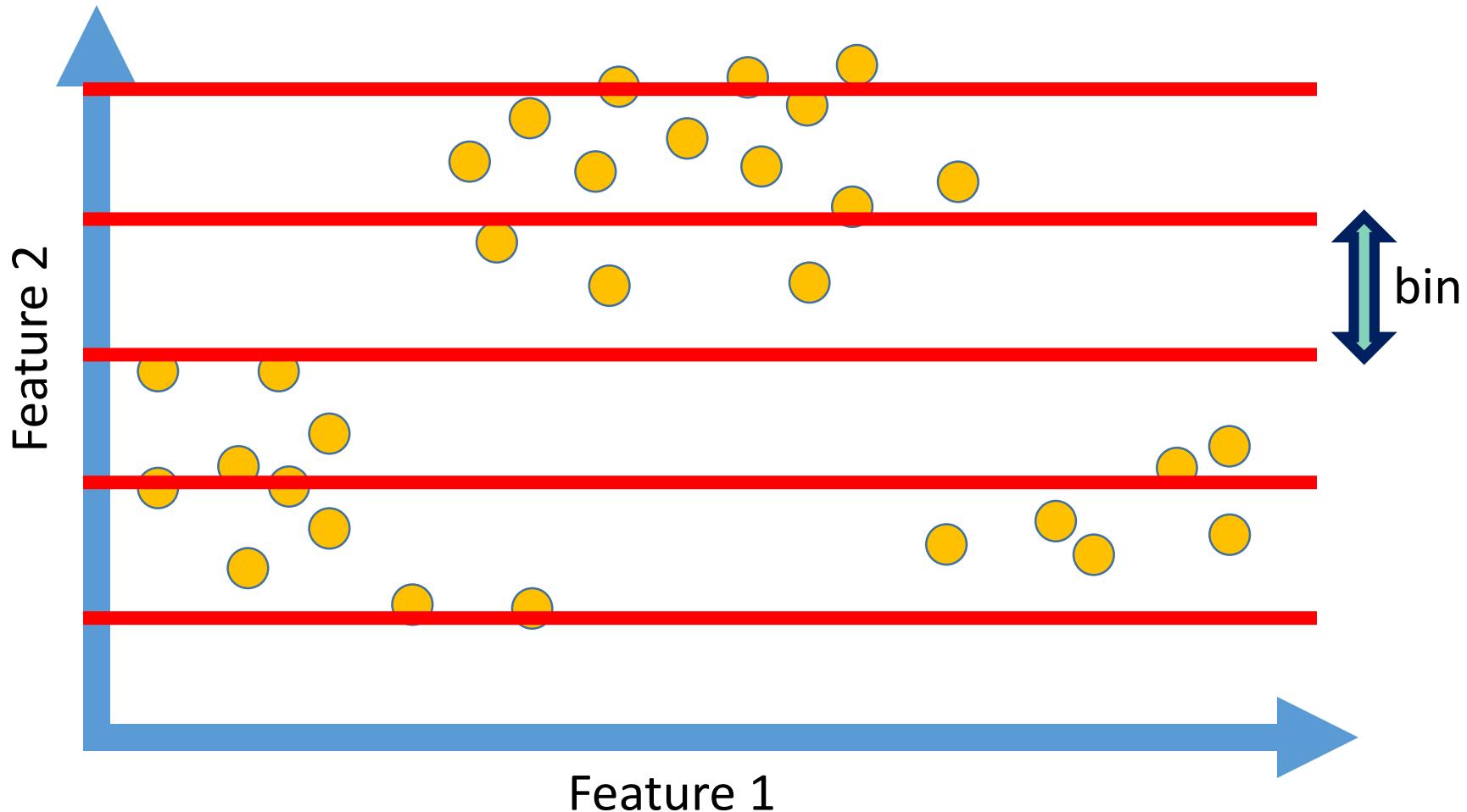


Image Representation: Histograms

- Better modeling (quantization) of multi-dimensional data
- Clustering

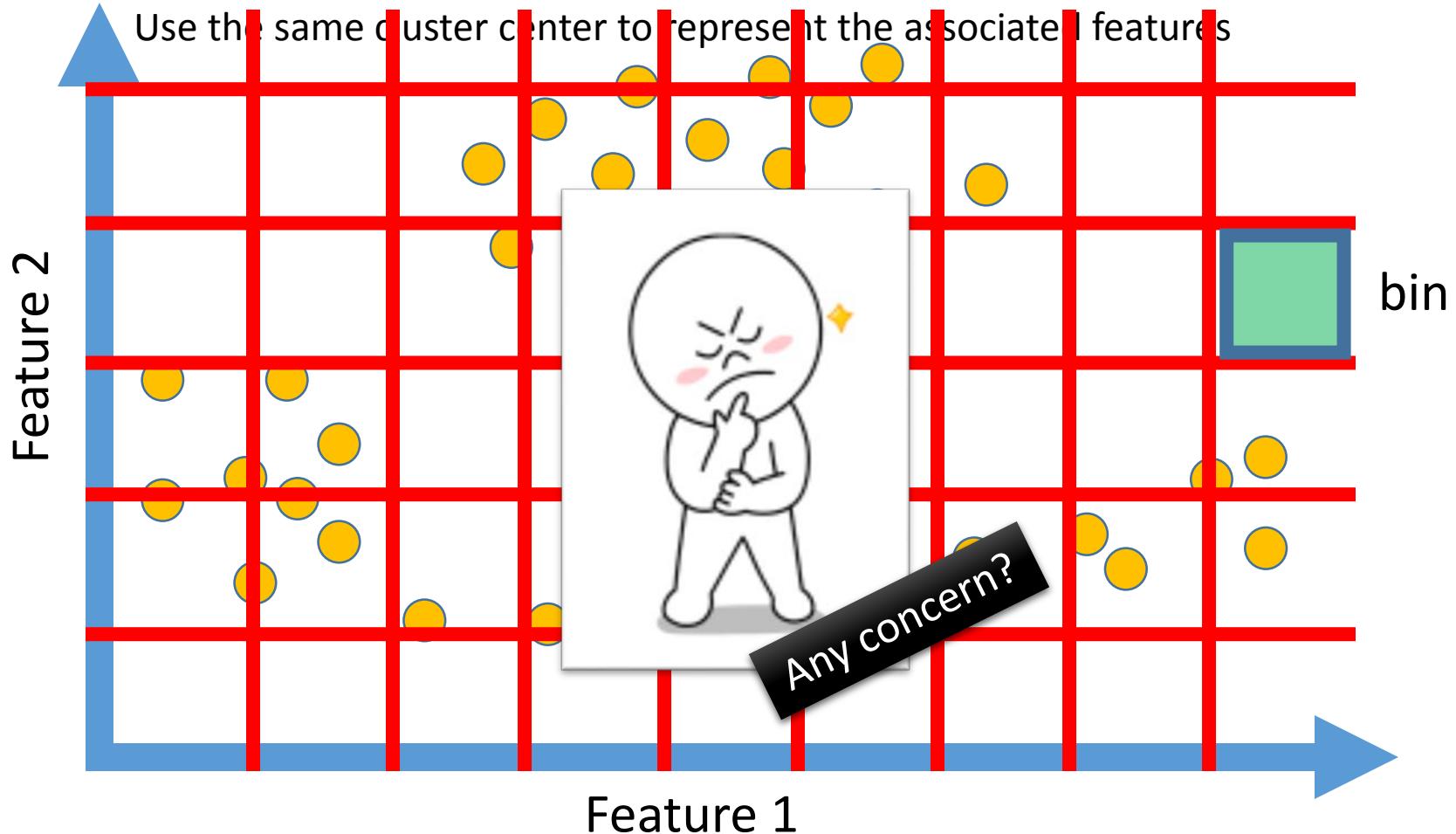
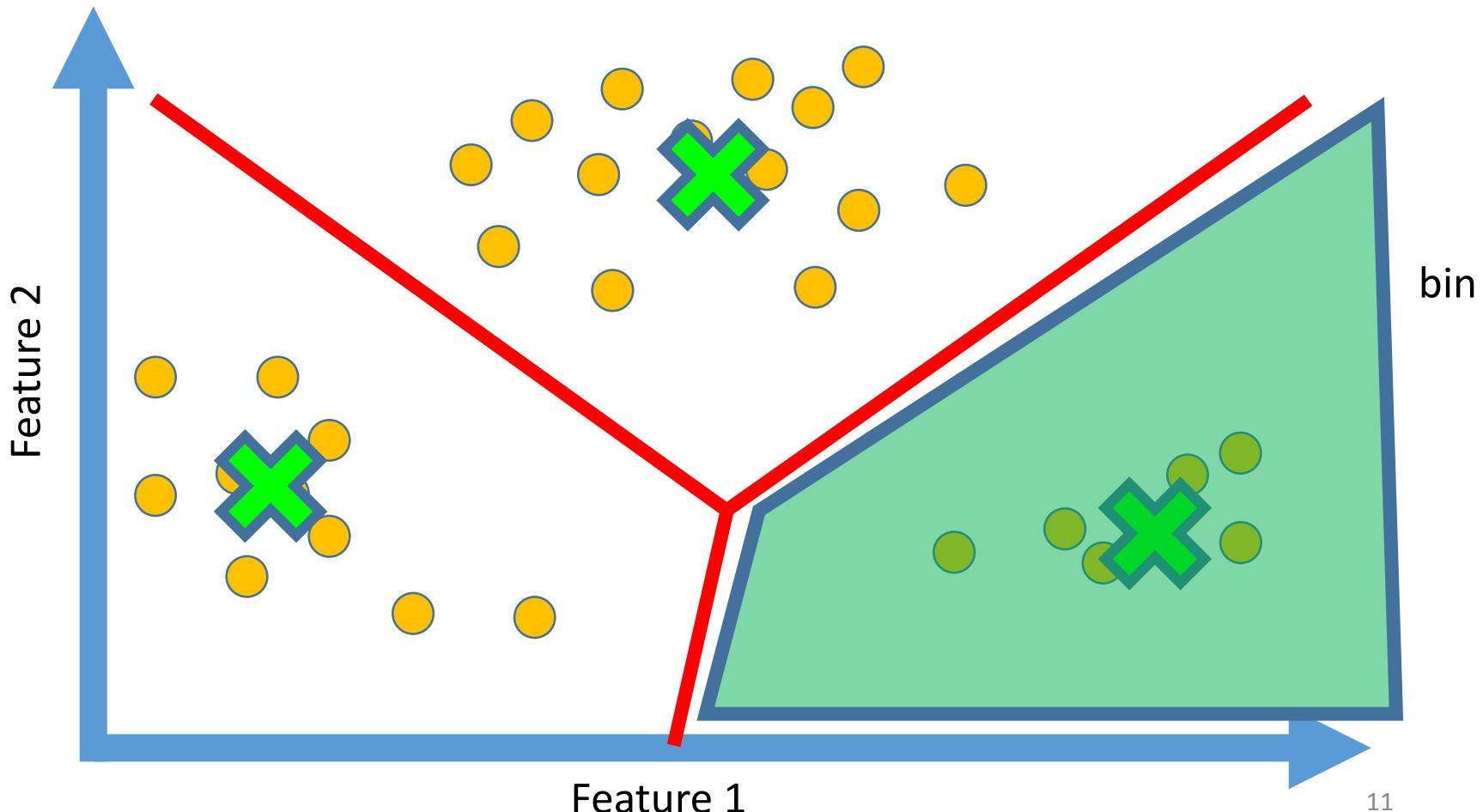


Image Representation: Histograms

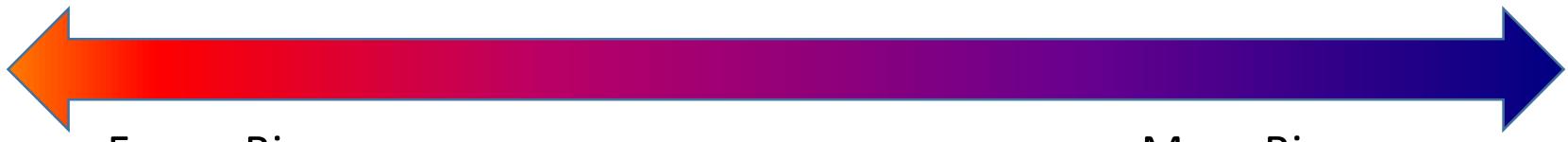
- Better modeling (quantization) of multi-dimensional data
- Clustering
 - Use the same cluster center to represent the associated features



Remarks on Histogram-Based Image Representation

- Quantization

- Grids vs. clusters



Fewer Bins

Need less data

Coarser representation

More Bins

Need more data

Finer representation

- Possible distance metrics

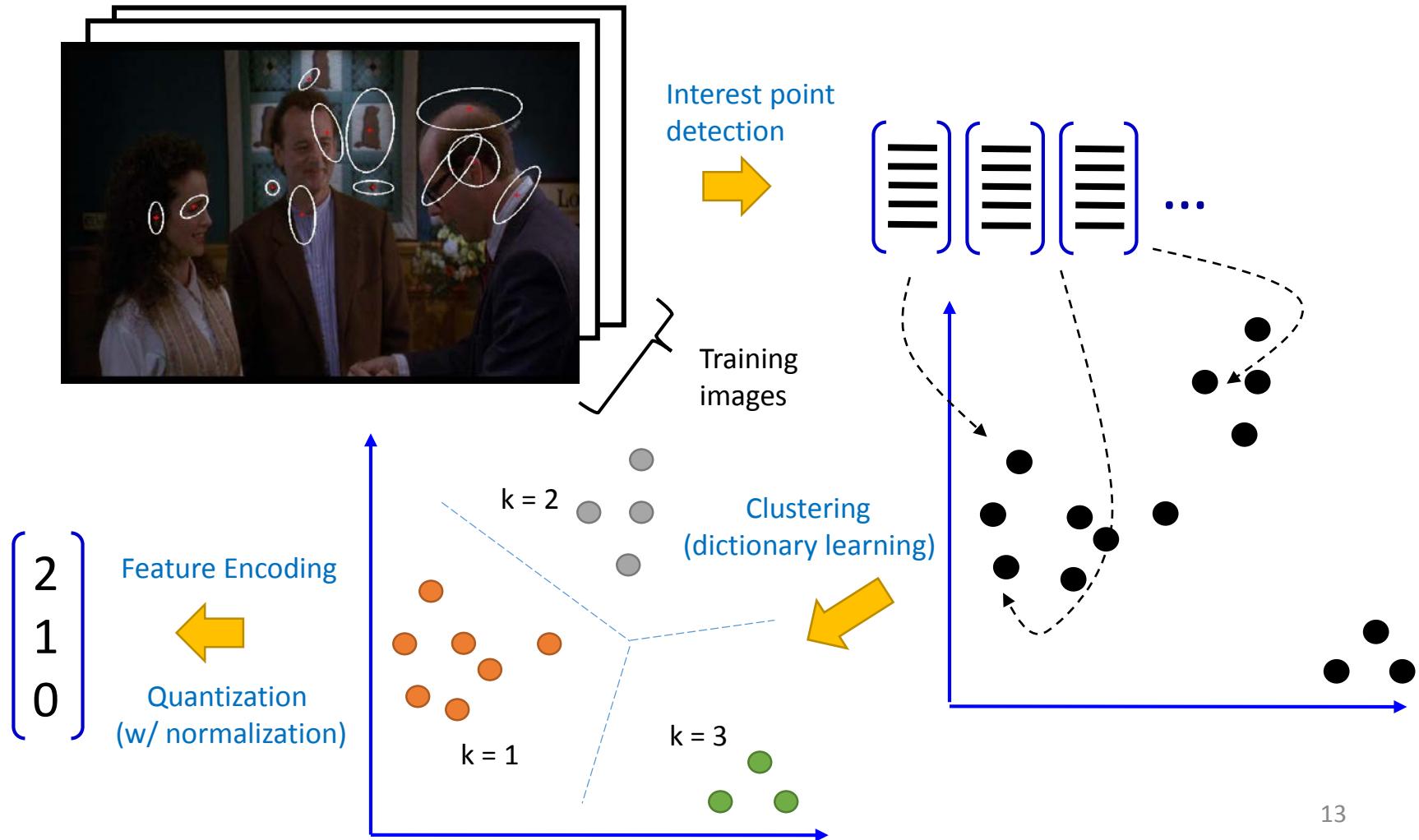
- Euclidean distance
- Histogram intersection kernel
- Chi-squared distance
- Earth mover's distance
(min cost to transform one distribution to another)

$$\text{histint}(h_i, h_j) = 1 - \sum_{m=1}^K \min(h_i(m), h_j(m))$$

$$\chi^2(h_i, h_j) = \frac{1}{2} \sum_{m=1}^K \frac{[h_i(m) - h_j(m)]^2}{h_i(m) + h_j(m)}$$

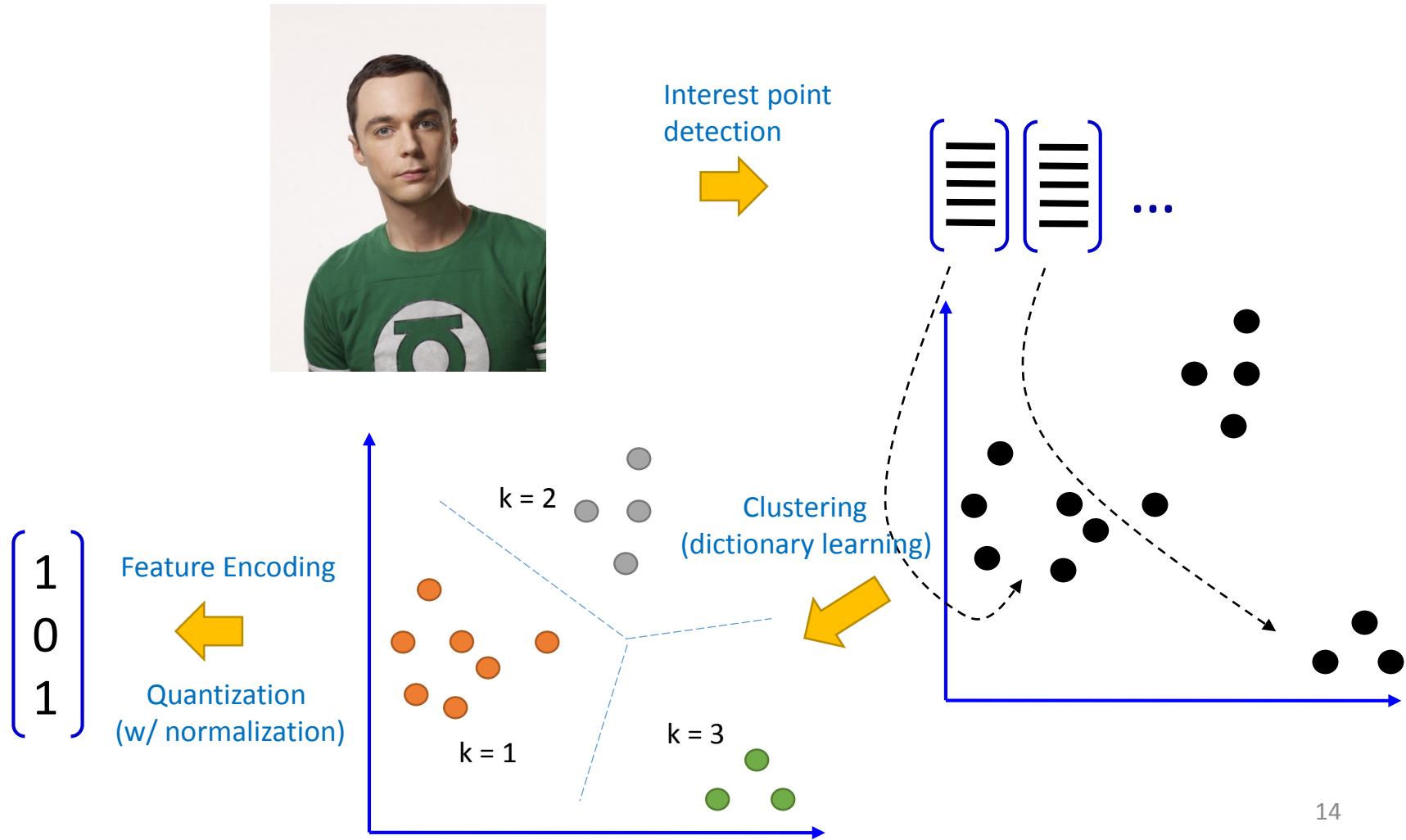
Bag-of-Words for Image Classification

- Training



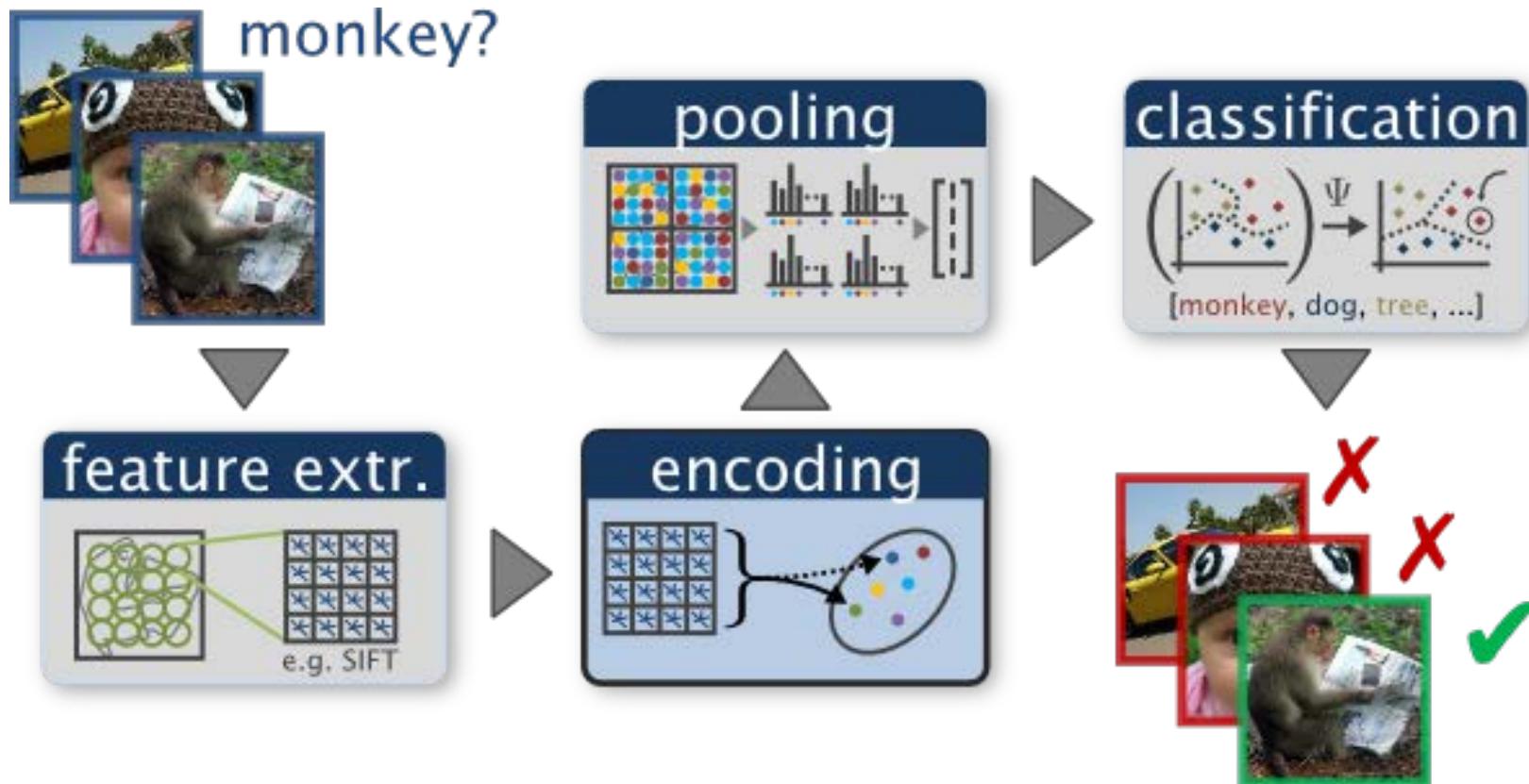
Bag-of-Words for Image Classification

- Testing

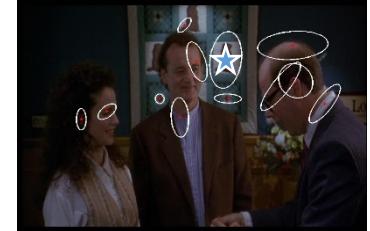


Bag-of-Words for Image Classification

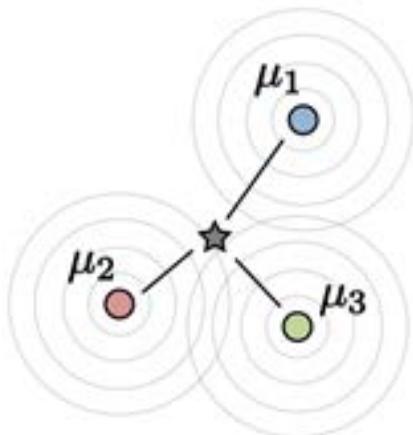
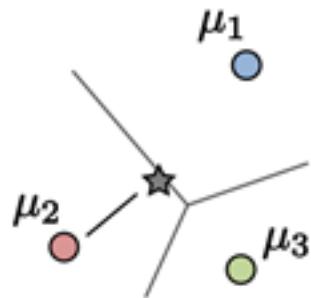
- Overview



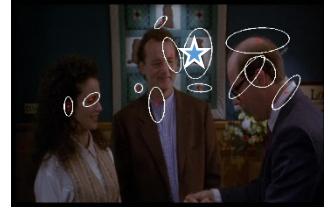
About Feature Encoding for Bag-of-Words



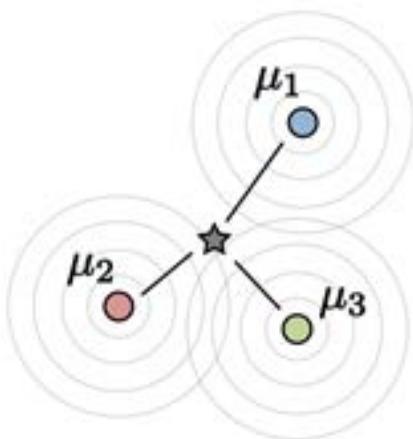
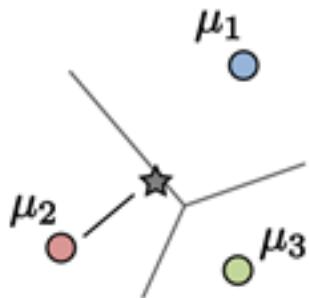
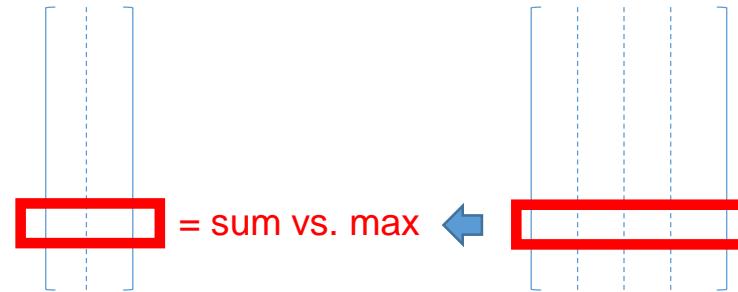
- Hard vs. soft assignments to clusters



About Feature Encoding for Bag-of-Words

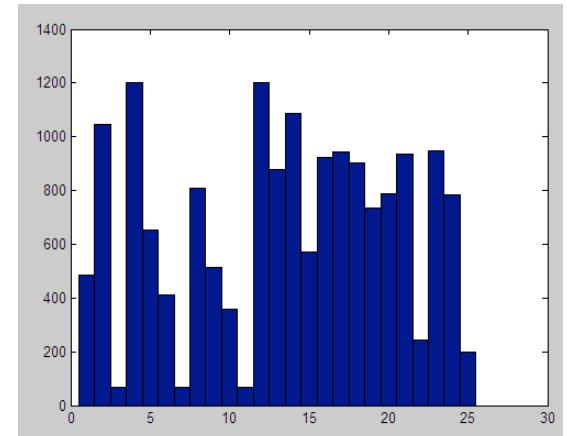
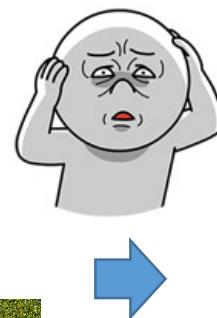


- Sum vs. Max Pooling



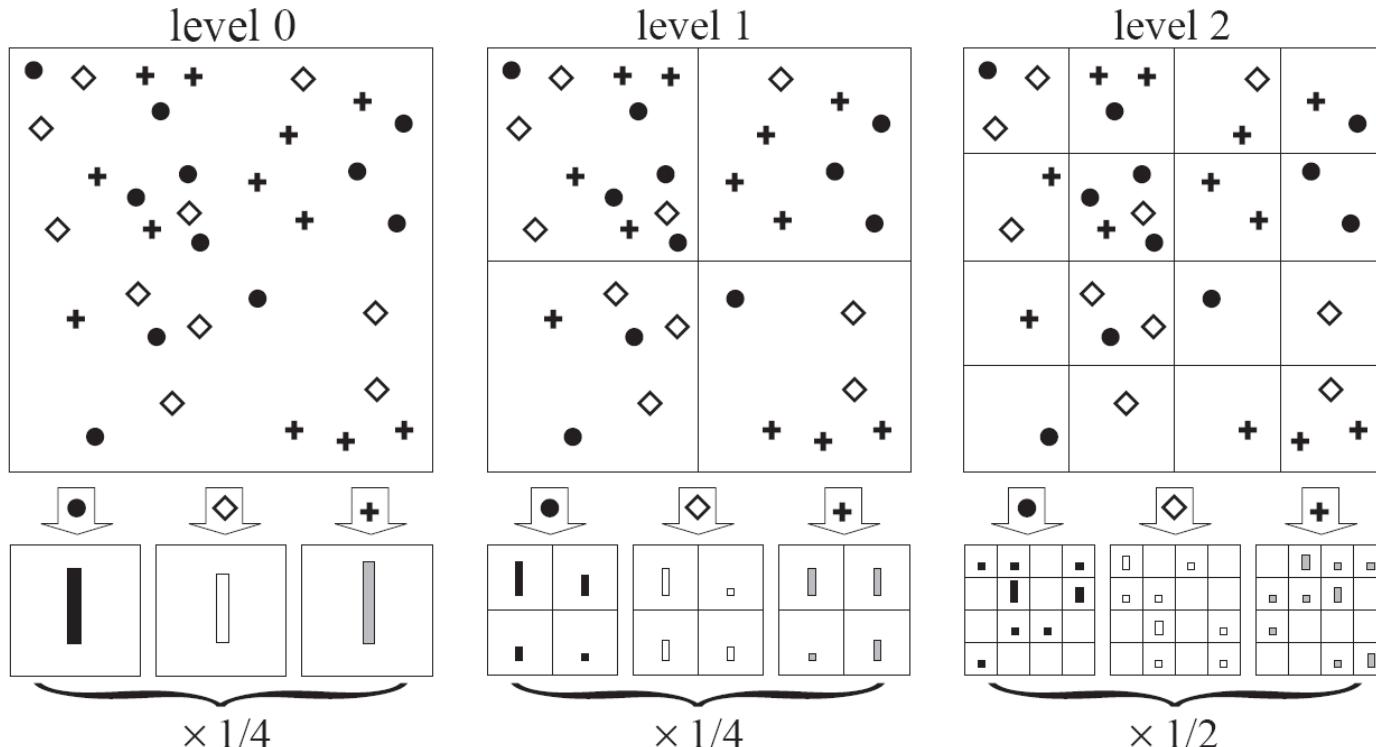
Final Remarks on BoW

- What's the limitation?
 - Loss of...
 - What's the possible solution?



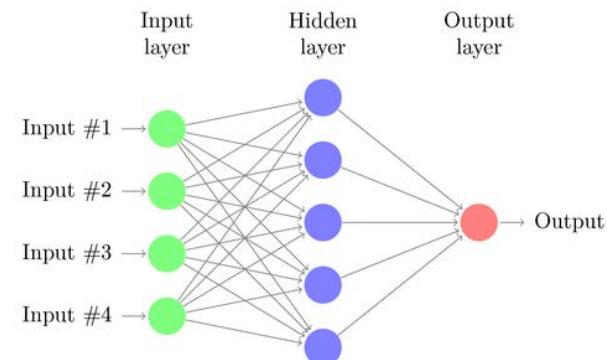
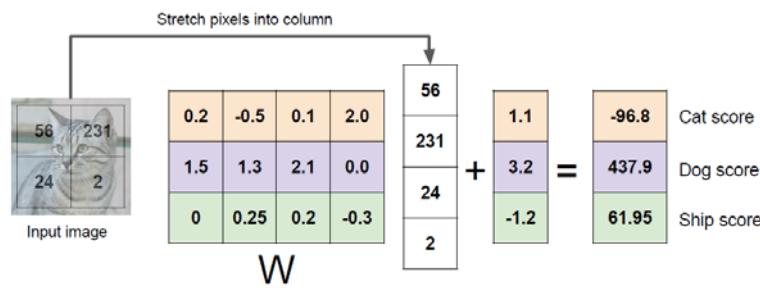
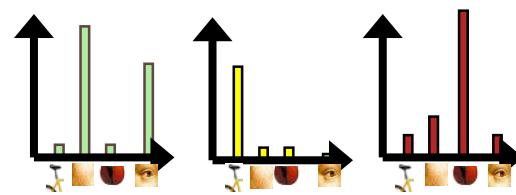
Final Remarks on BoW

- Spatial pyramid
 - Compute BoW in each spatial grid + concatenation



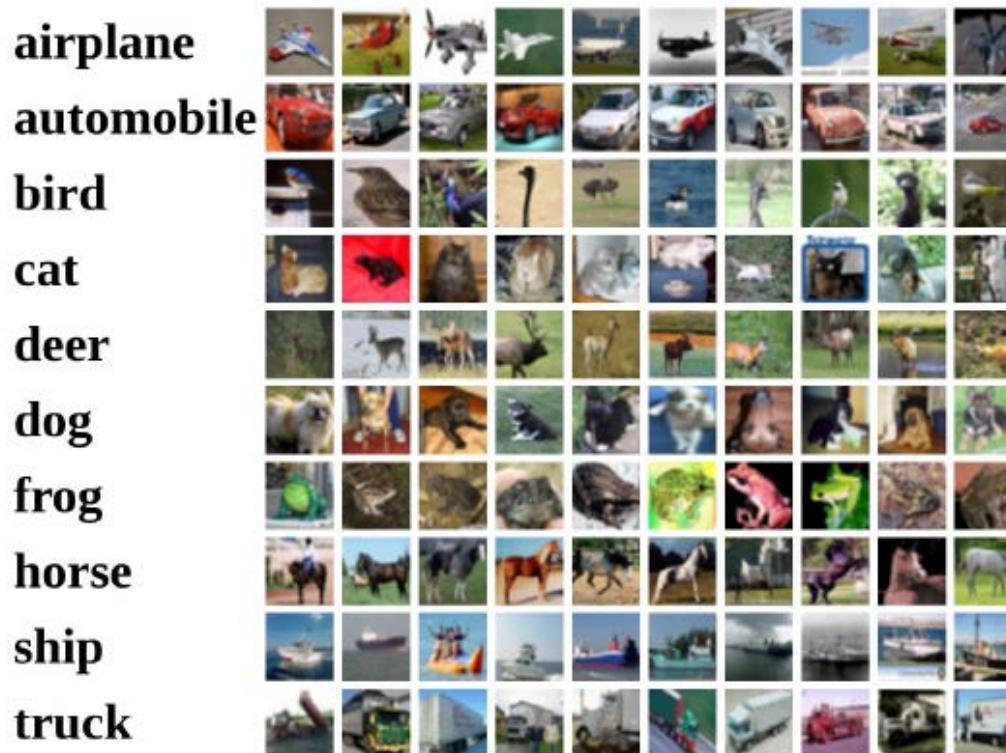
What's to Be Covered Today...

- General Framework for Visual Classification
 - Bag-of-Words Models as Image Features
- Intro to Neural Networks & CNN
 - Linear Classification
 - Neural Network for Machine Vision
 - Multi-Layer Perceptron
 - Convolutional Neural Networks



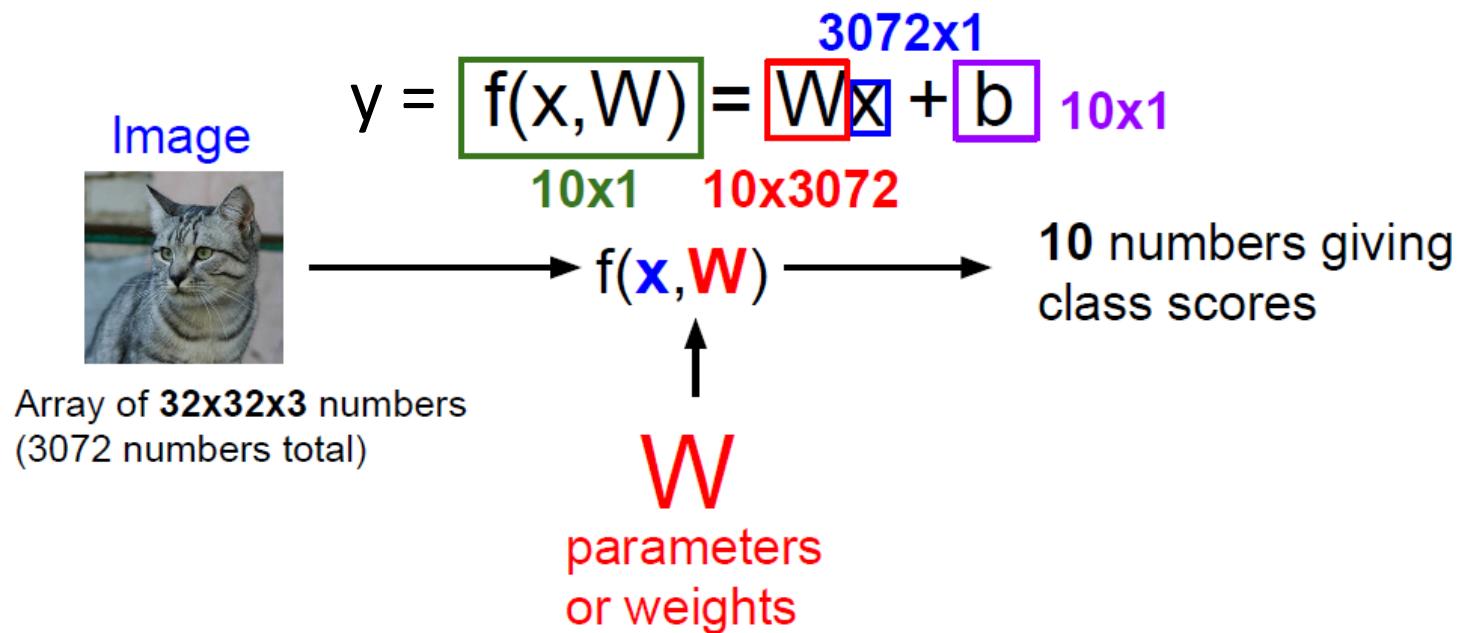
Linear Classification

- Linear Classifier
 - Can be viewed as a **parametric approach**. Why?
 - Assuming that we need to recognize 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.



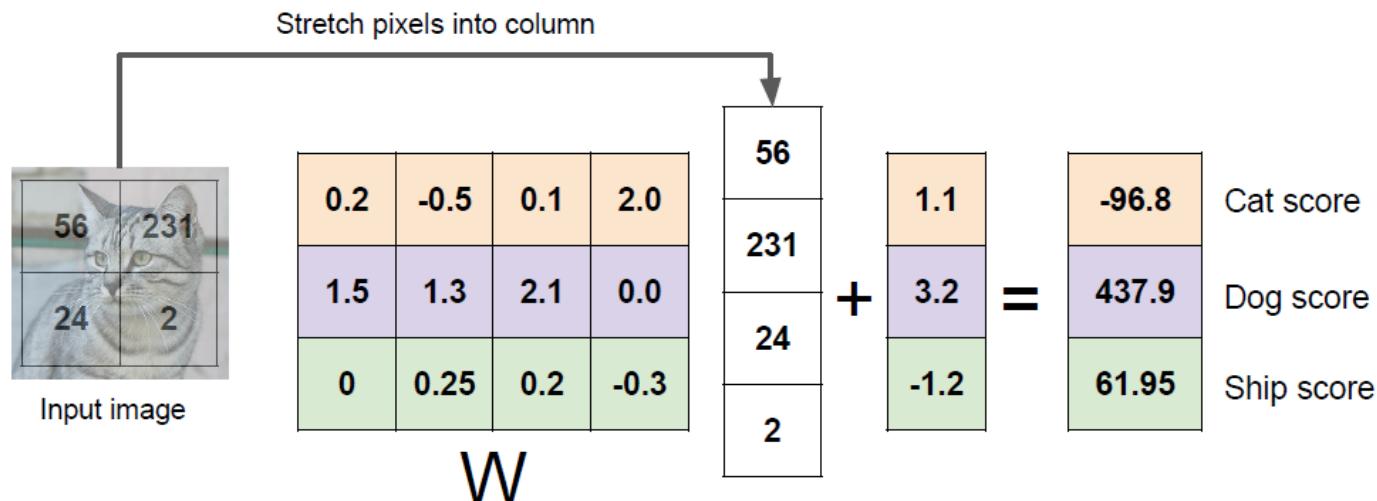
Linear Classification (cont'd)

- Linear Classifier
 - Can be viewed as a **parametric approach**. Why?
 - Assuming that we need to recognize 10 object categories of interest (e.g., CIFAR10).
 - Let's take the input image as \mathbf{x} , and the linear classifier as \mathbf{W} . We hope to see that $\mathbf{y} = \mathbf{Wx} + \mathbf{b}$ as a 10-dimensional output indicating the score for each class.



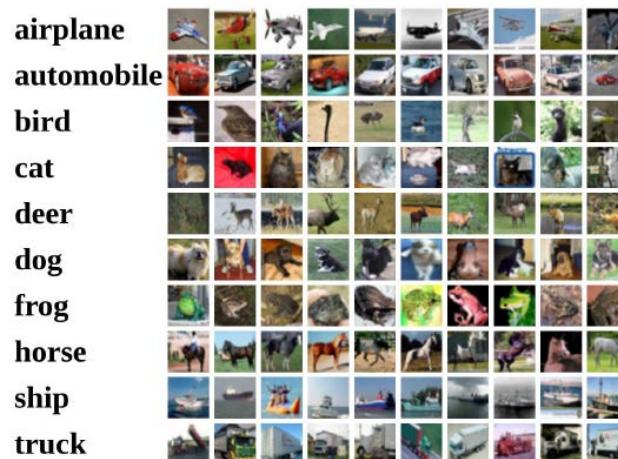
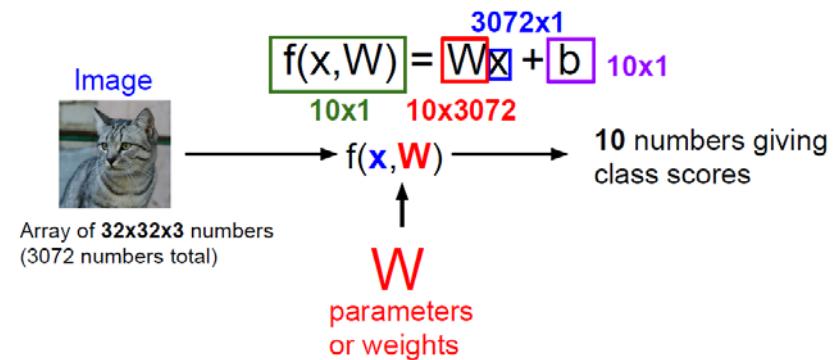
Linear Classification (cont'd)

- Linear Classifier
 - Can be viewed as a **parametric approach**. Why?
 - Assuming that we need to recognize 10 object categories of interest (e.g., CIFAR10).
 - Let's take the input image as \mathbf{x} , and the linear classifier as \mathbf{W} . We hope to see that $\mathbf{y} = \mathbf{Wx} + \mathbf{b}$ as a 10-dimensional output indicating the score for each class.
 - Take an image with 2×2 pixels & 3 classes of interest as example: we need to learn linear transformation/classifier \mathbf{W} and bias \mathbf{b} , so that desirable outputs $\mathbf{y} = \mathbf{Wx} + \mathbf{b}$ can be expected.



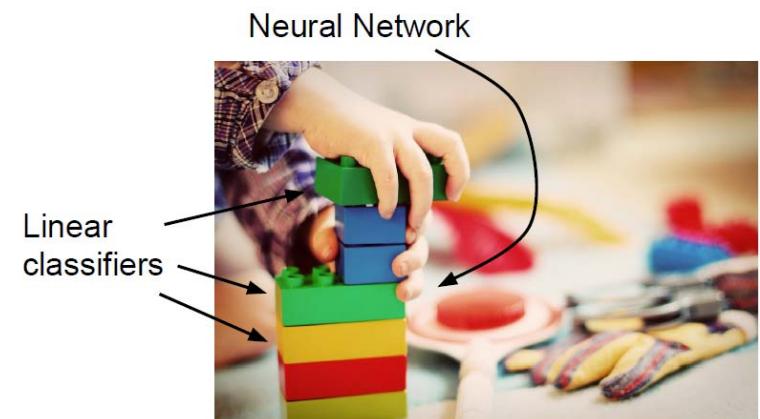
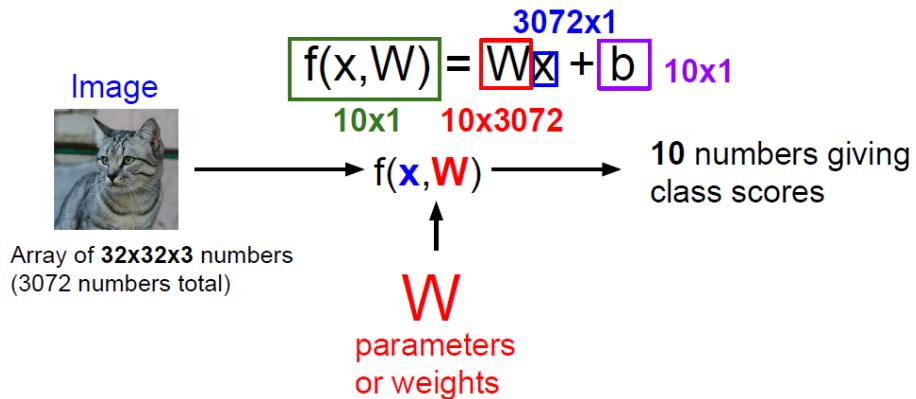
Some Remarks

- Interpreting $y = Wx + b$
 - What can we say about the learned W ?
 - The weights in W are trained by observing training data X and their ground truth Y .
 - Each column in W can be viewed as an exemplar of the corresponding class.
 - Thus, Wx basically performs **inner product** (or **correlation**) between the input x and the exemplar of each class. (Signal & Systems!)



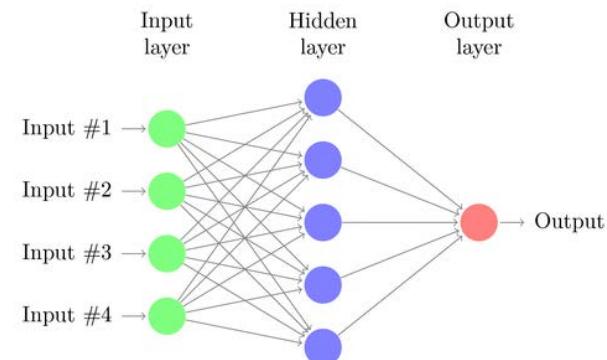
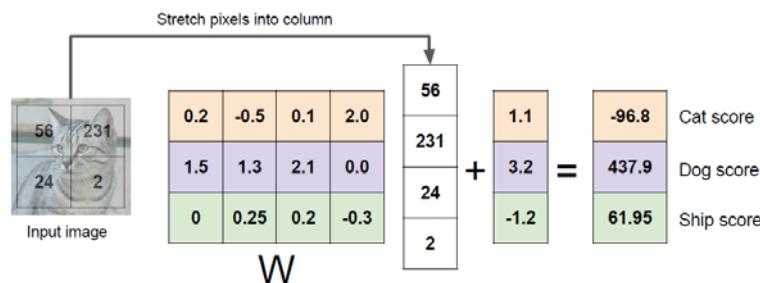
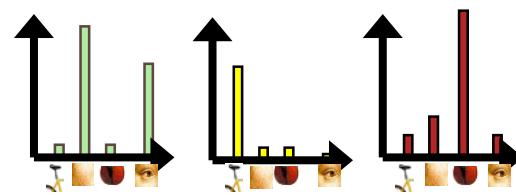
Linear Classification

- Remarks
 - Starting points for many multi-class or complex/nonlinear classifier
 - How to determine a proper loss function for matching \mathbf{y} and $\mathbf{Wx} + \mathbf{b}$, and thus how to learn the model \mathbf{W} (including the bias \mathbf{b}), are the keys to the learning of an effective classification model.

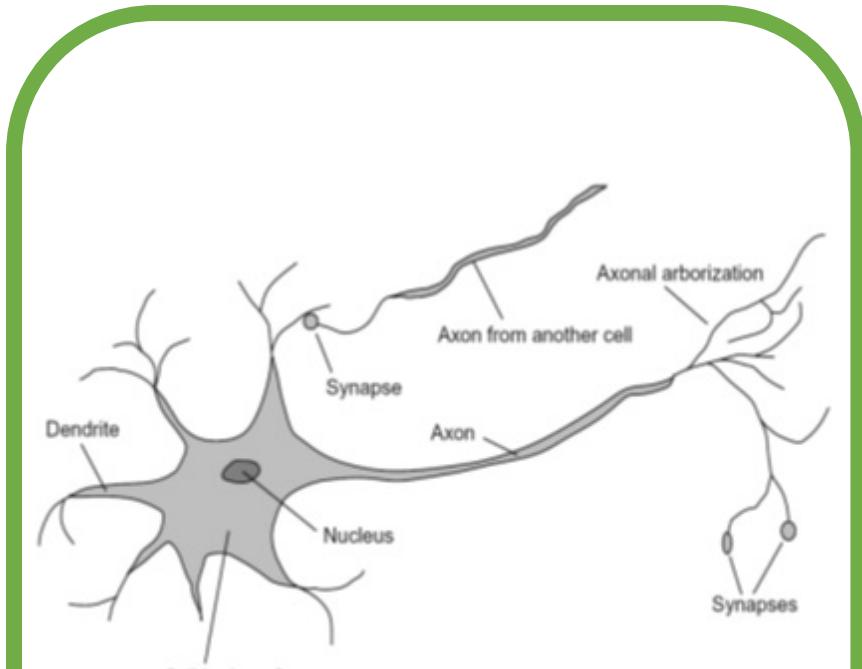


What's to Be Covered Today...

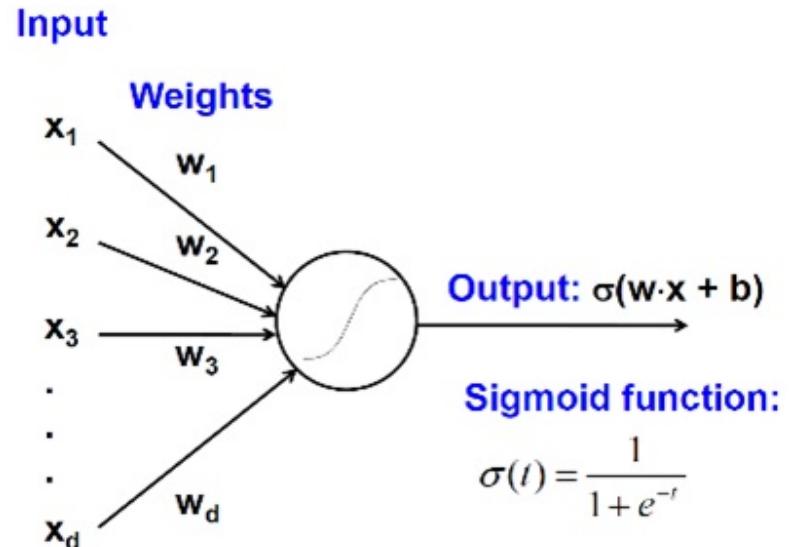
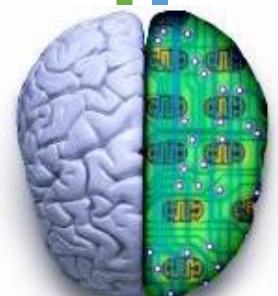
- General Framework for Visual Classification
 - Bag-of-Words Models as Image Features
- Intro to Neural Networks & CNN
 - Linear Classification
 - Neural Network for Machine Vision
 - Multi-Layer Perceptron
 - Convolutional Neural Networks



Biological neuron and Perceptrons



A biological neuron

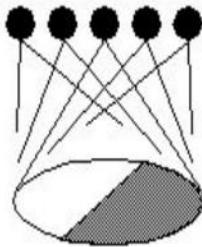


An artificial neuron (Perceptron) - a linear classifier

Hubel/Wiesel Architecture and Multi-layer Neural Network

Hubel & Weisel

topographical mapping

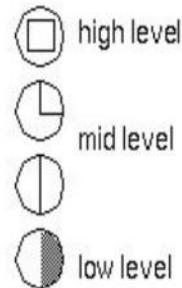
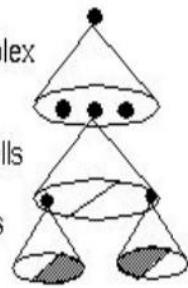


featural hierarchy

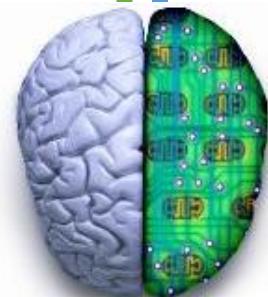
hyper-complex cells

complex cells

simple cells



Hubel and Weisel's architecture

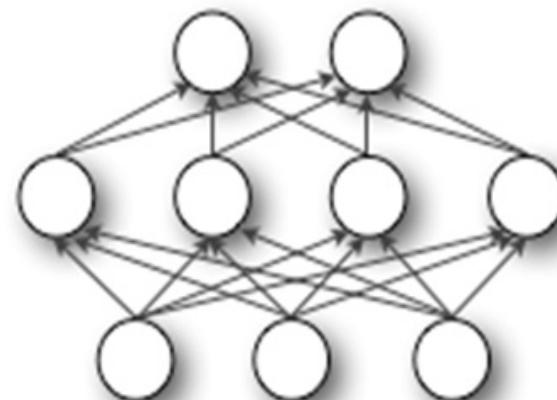


Multi-layer Neural Network
- A *non-linear* classifier

output layer

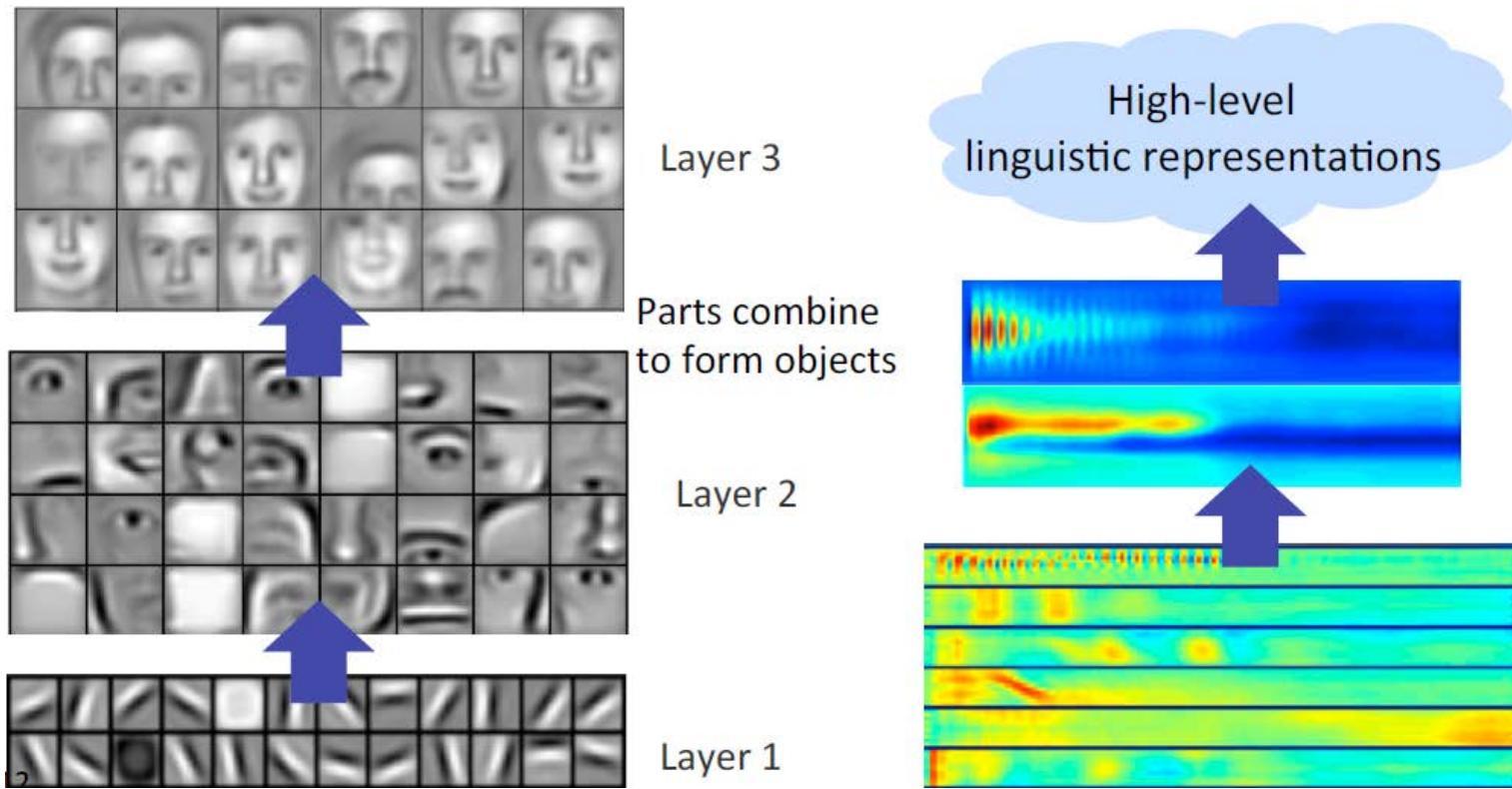
hidden layer

input layer



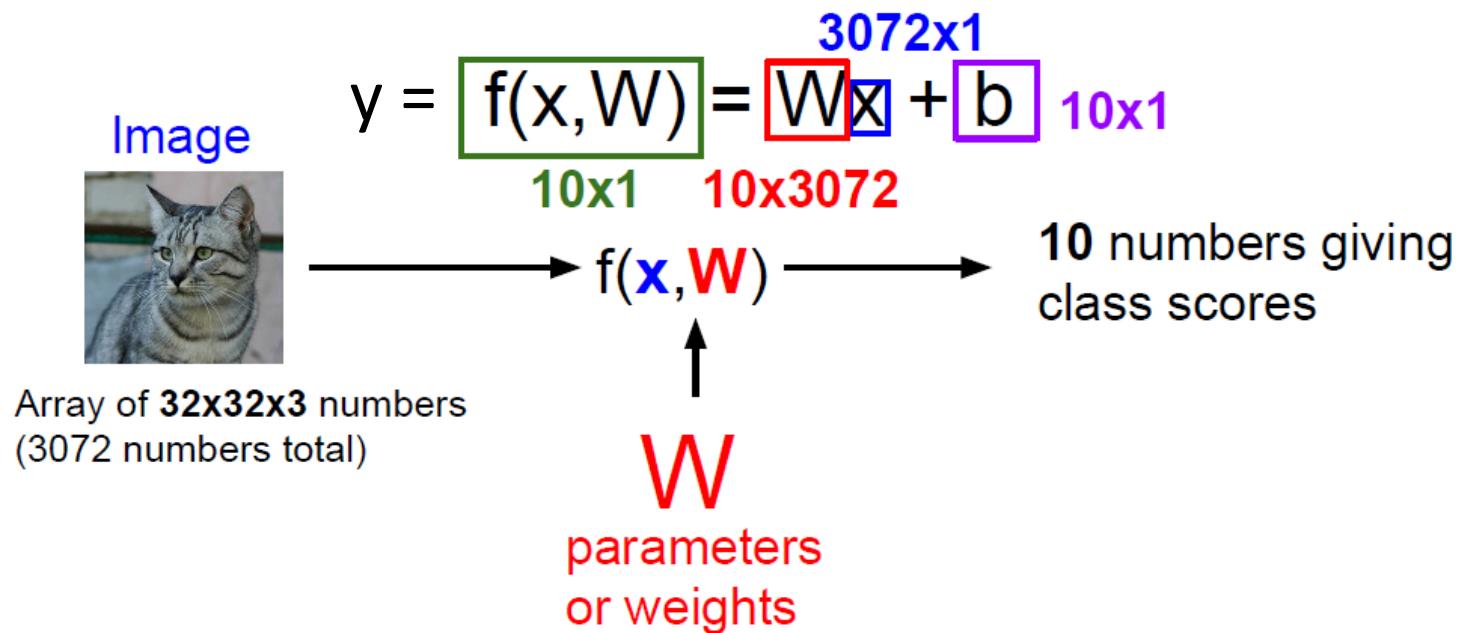
Hierarchical Learning

- Successive model layers learn deeper intermediate representations.

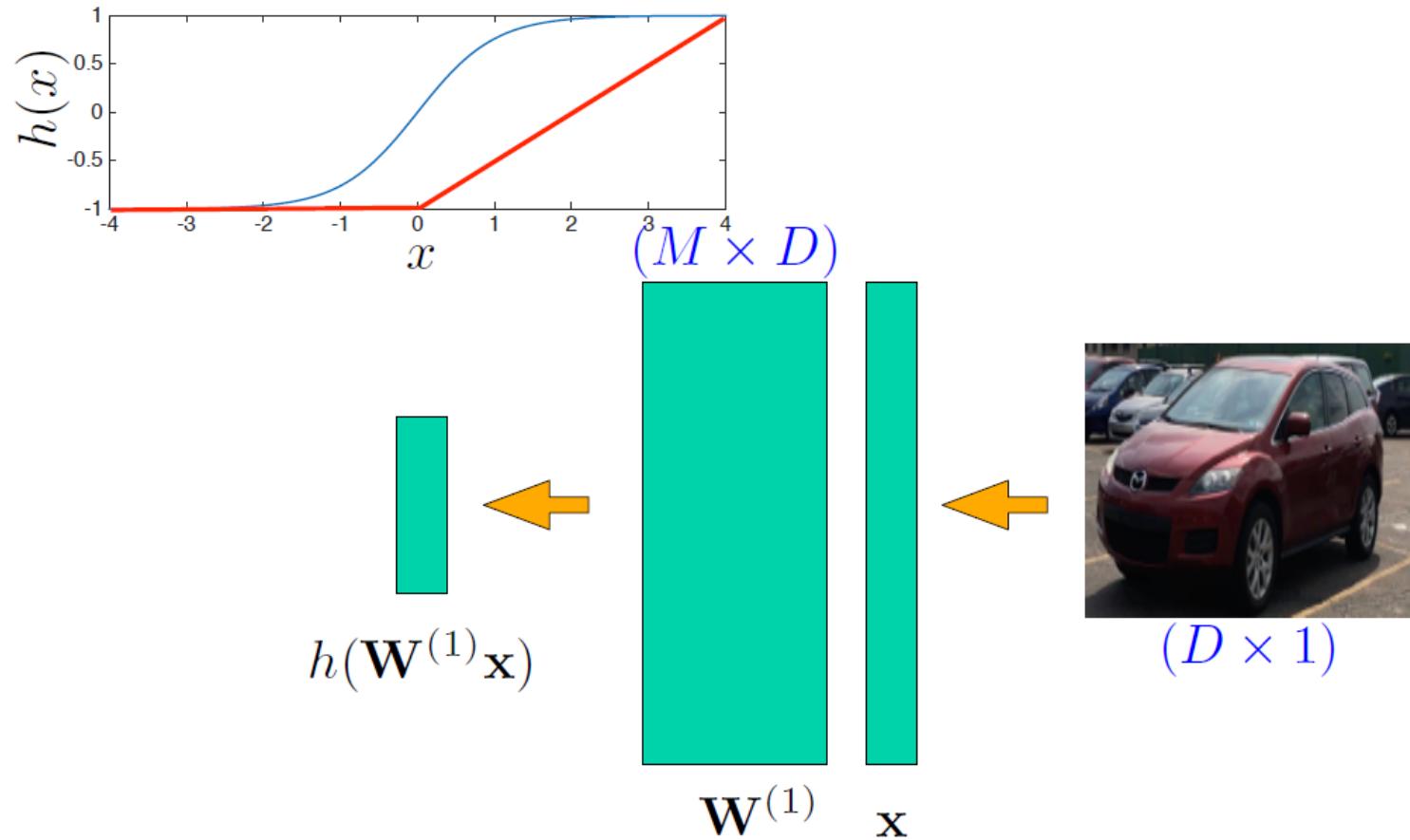


Revisit of Linear Classification

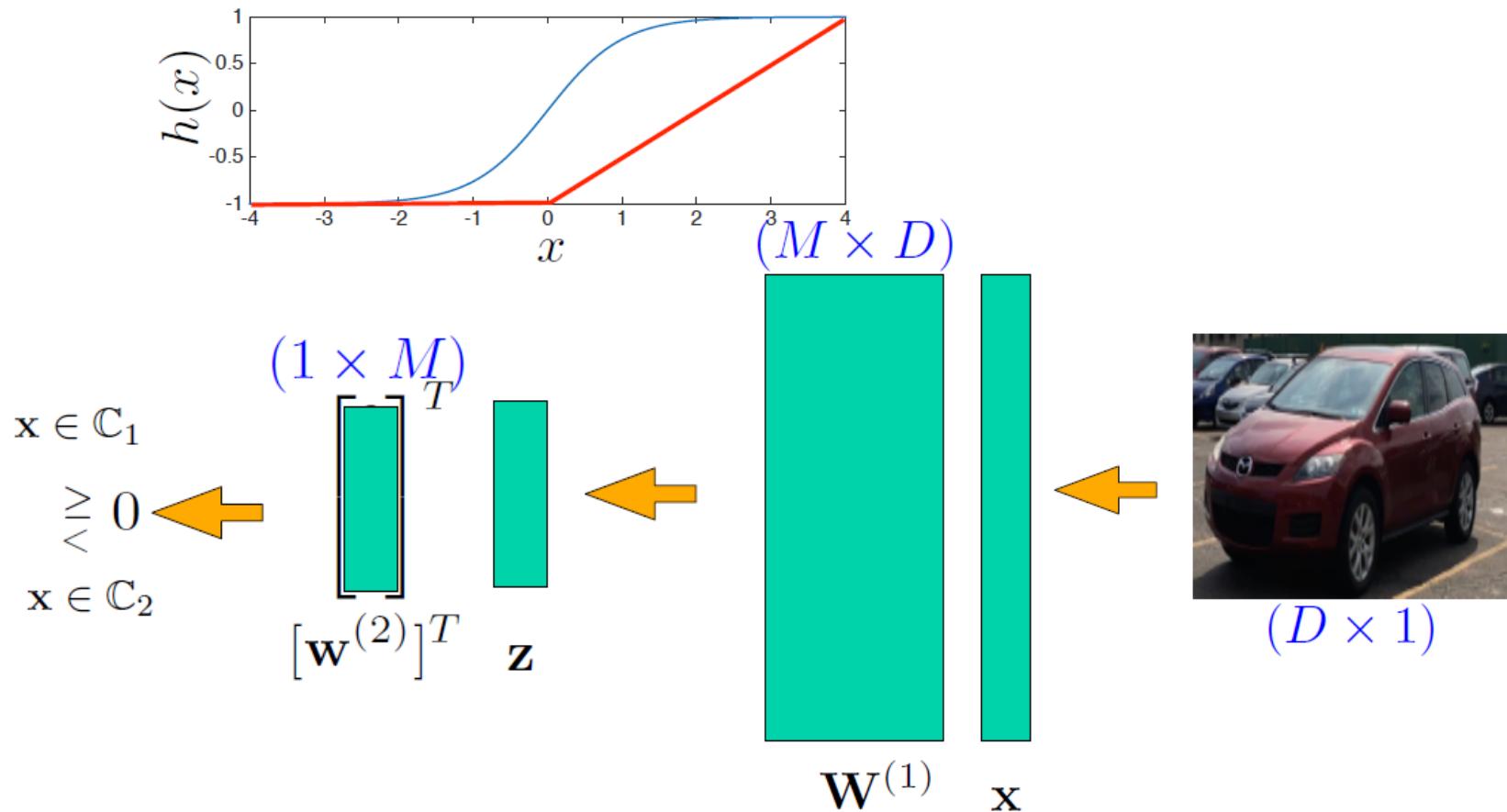
- Linear Classifier
 - Can be viewed as a **parametric approach**. Why?
 - Assuming that we need to recognize 10 object categories of interest (e.g., CIFAR10).
 - Let's take the input image as \mathbf{x} , and the linear classifier as \mathbf{W} . We hope to see that $\mathbf{y} = \mathbf{Wx} + \mathbf{b}$ as a 10-dimensional output indicating the score for each class.



Multi-Layer Perceptron: A Nonlinear Classifier

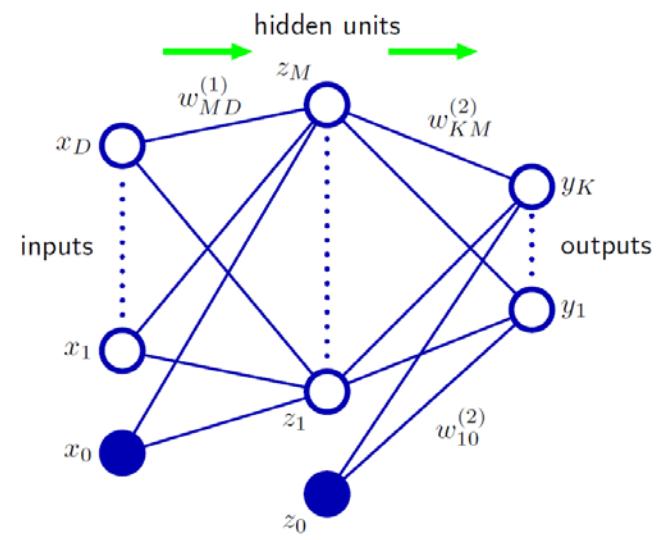


Multi-Layer Perceptron: A Nonlinear Classifier (cont'd)



Layer 1 in MLP

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

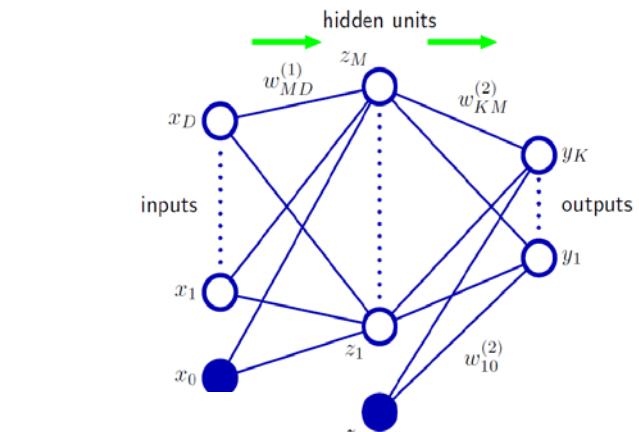


$h()$ = non-linear function

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

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

Layer 2 in MLP



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

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

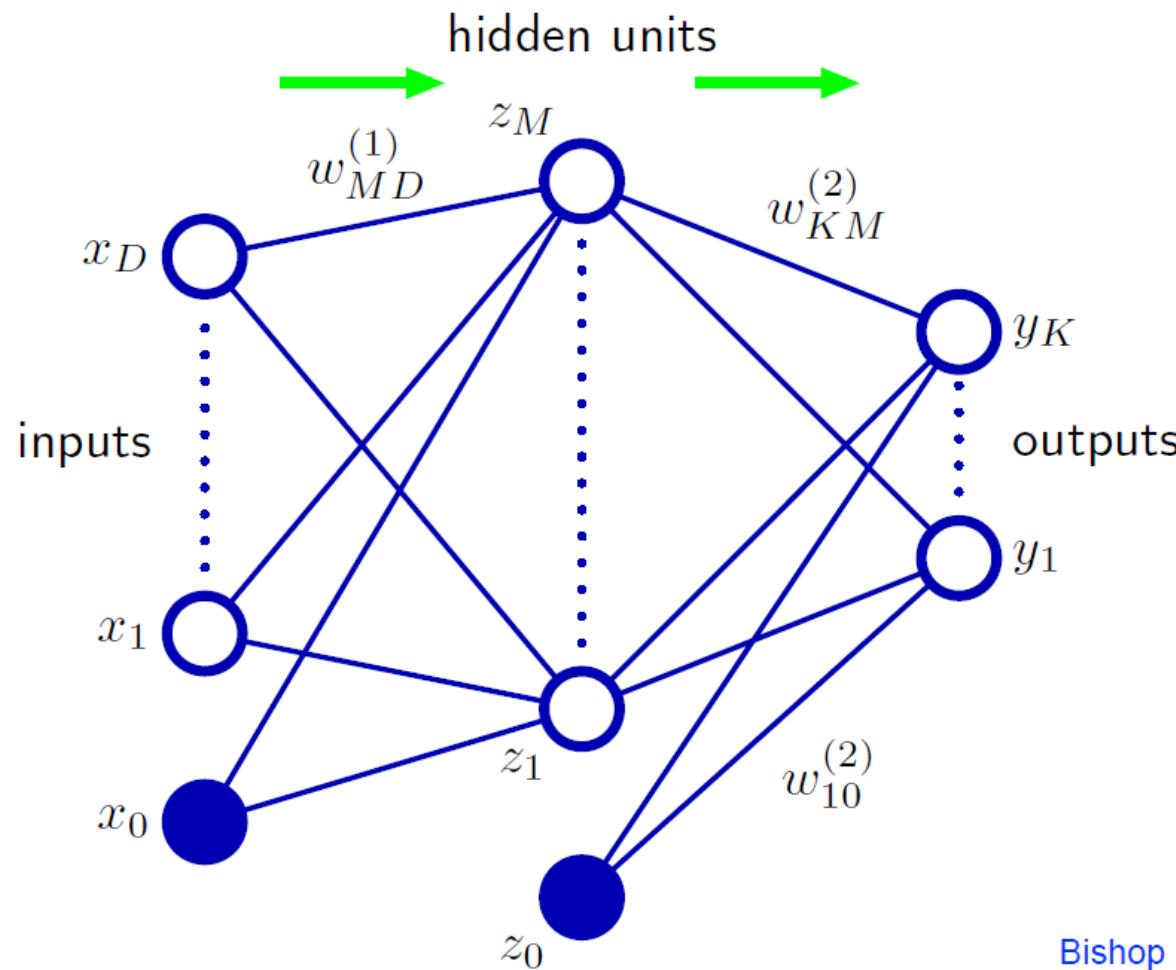
$$\mathbf{z}^T \mathbf{w}^{(2)} \geq 0$$

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

$\mathbf{z} = M \times 1$ output of layer 1

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

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

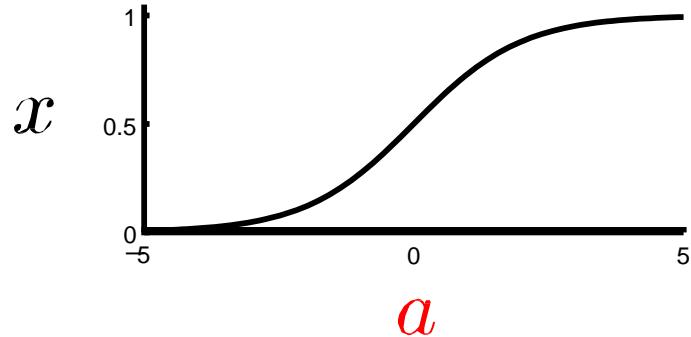
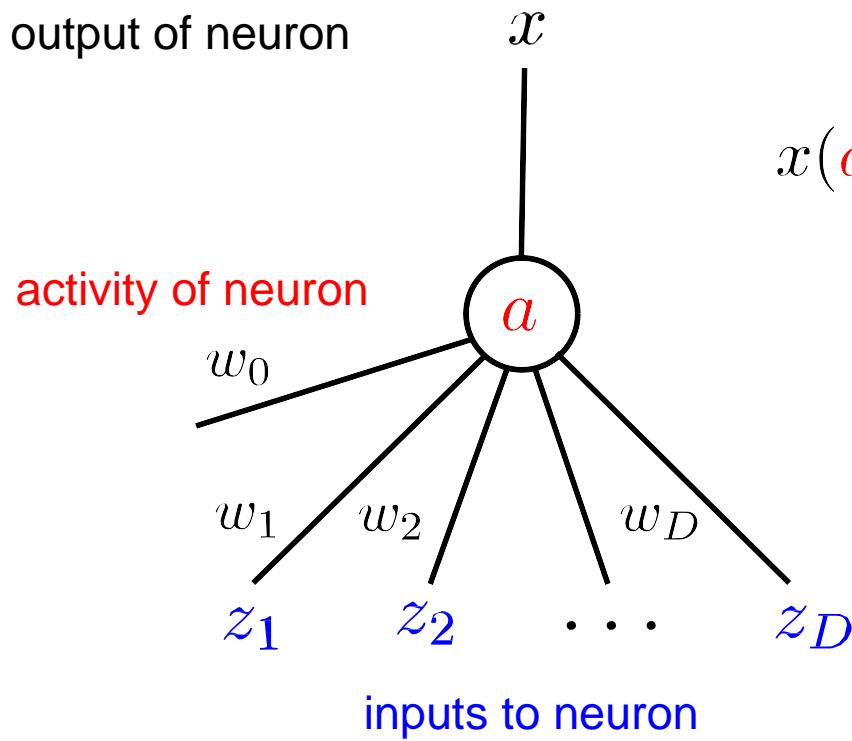


Bishop 2006

Let's Get a Closer Look...

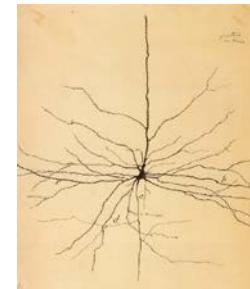
- A single neuron

output of neuron

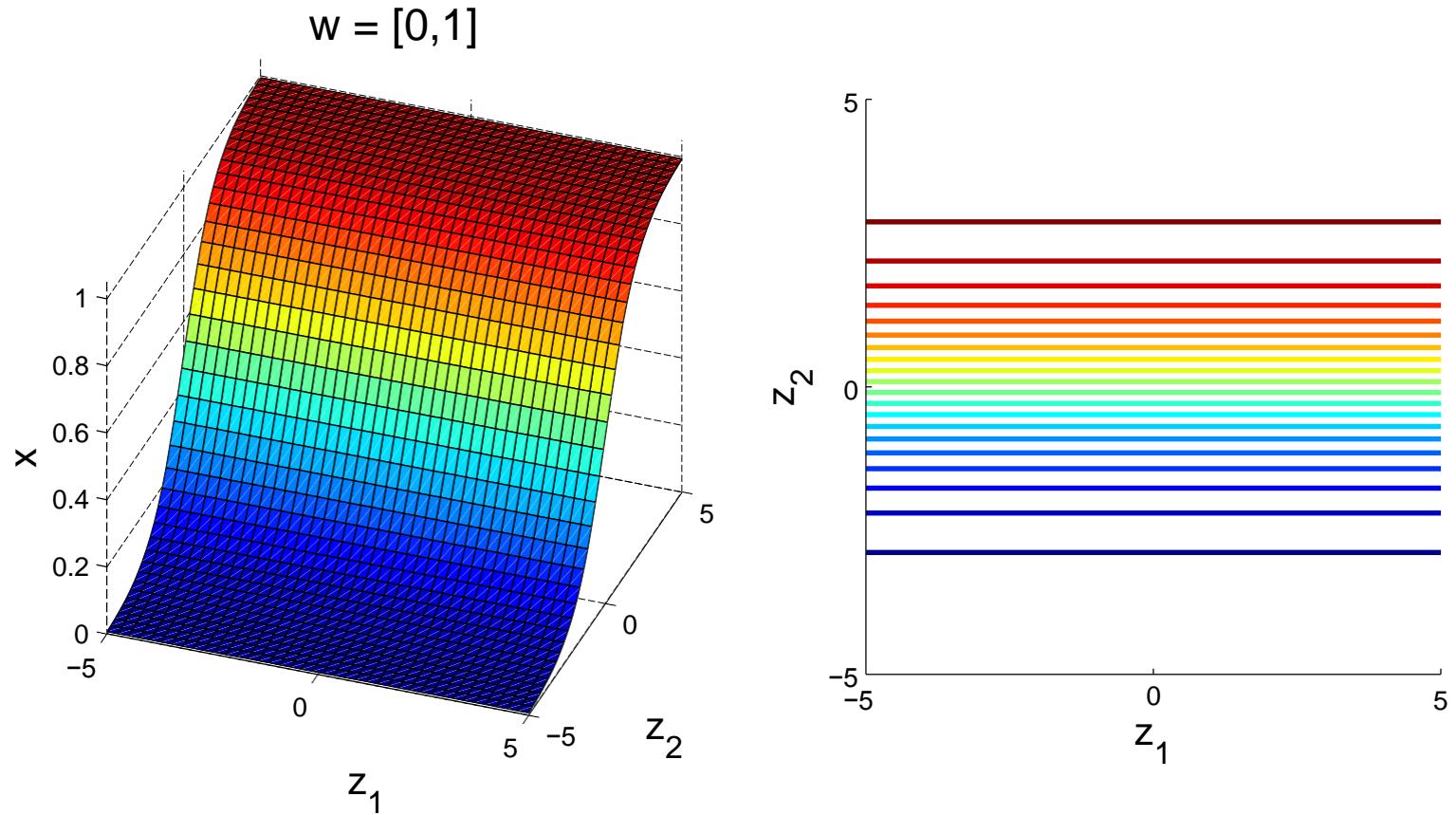


$$x(\textcolor{red}{a}) = \frac{1}{1+\exp(-\textcolor{red}{a})} \quad x \in (0, 1)$$

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

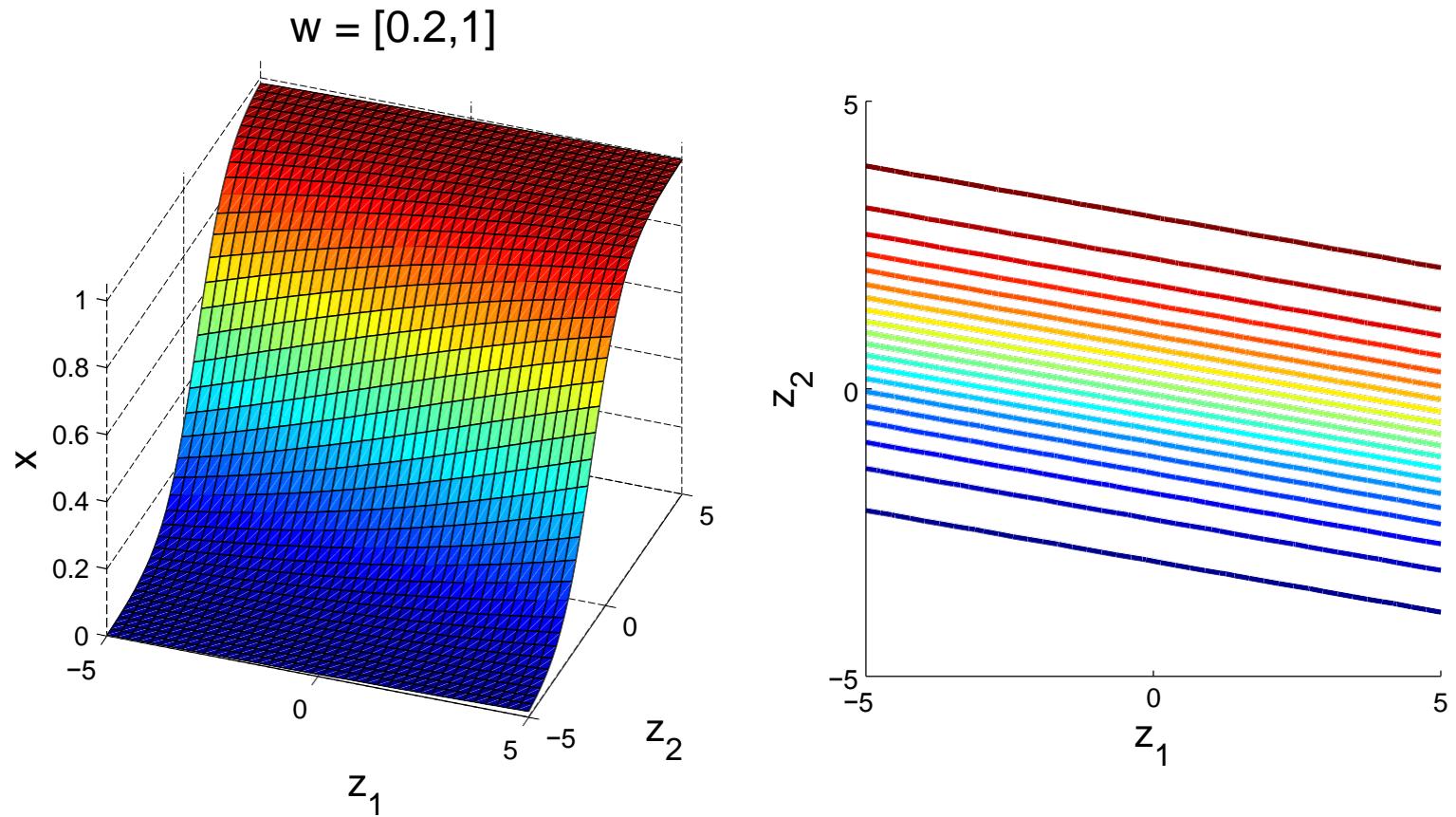


Input-Output Function of a Single Neuron



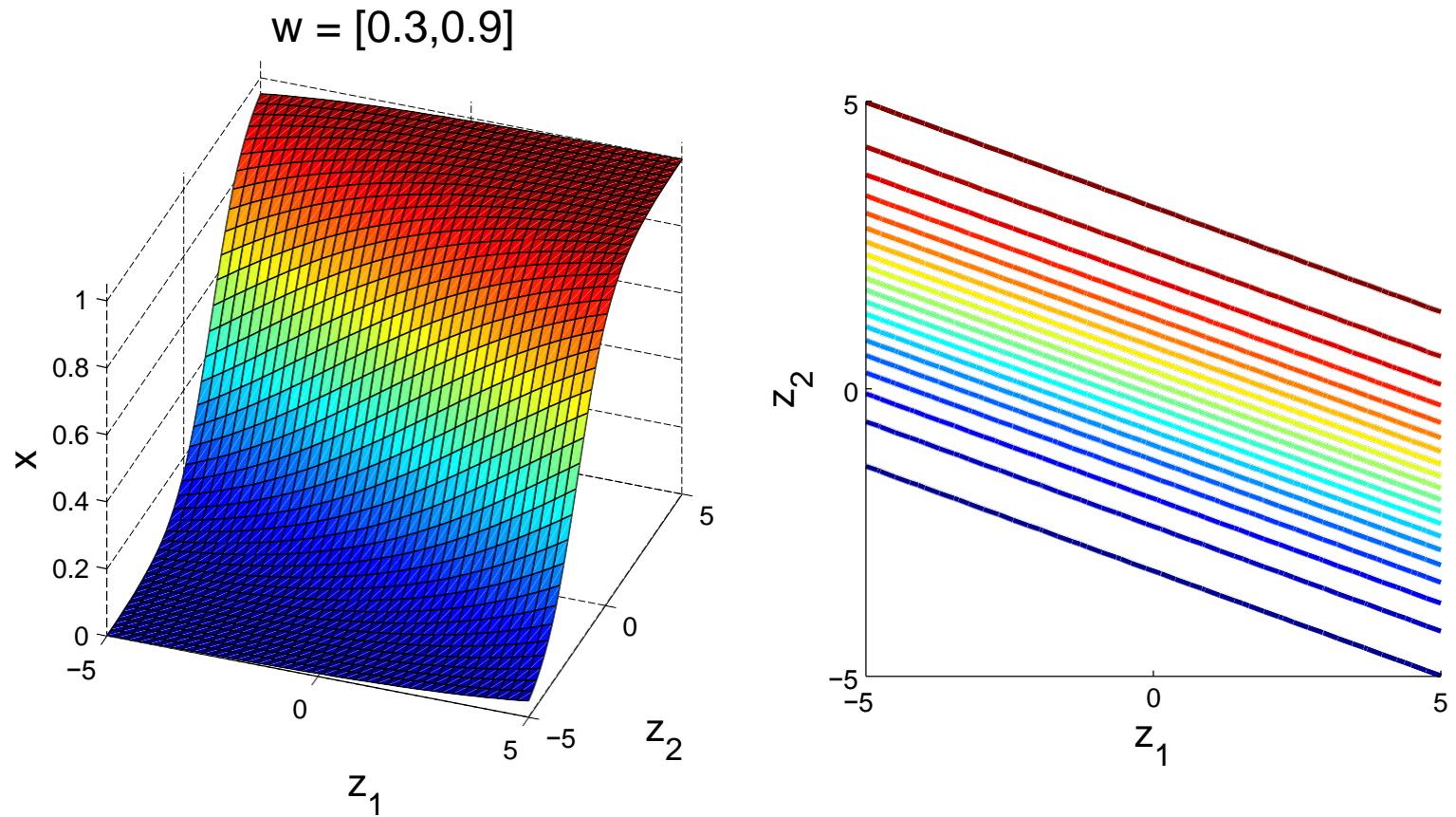
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron



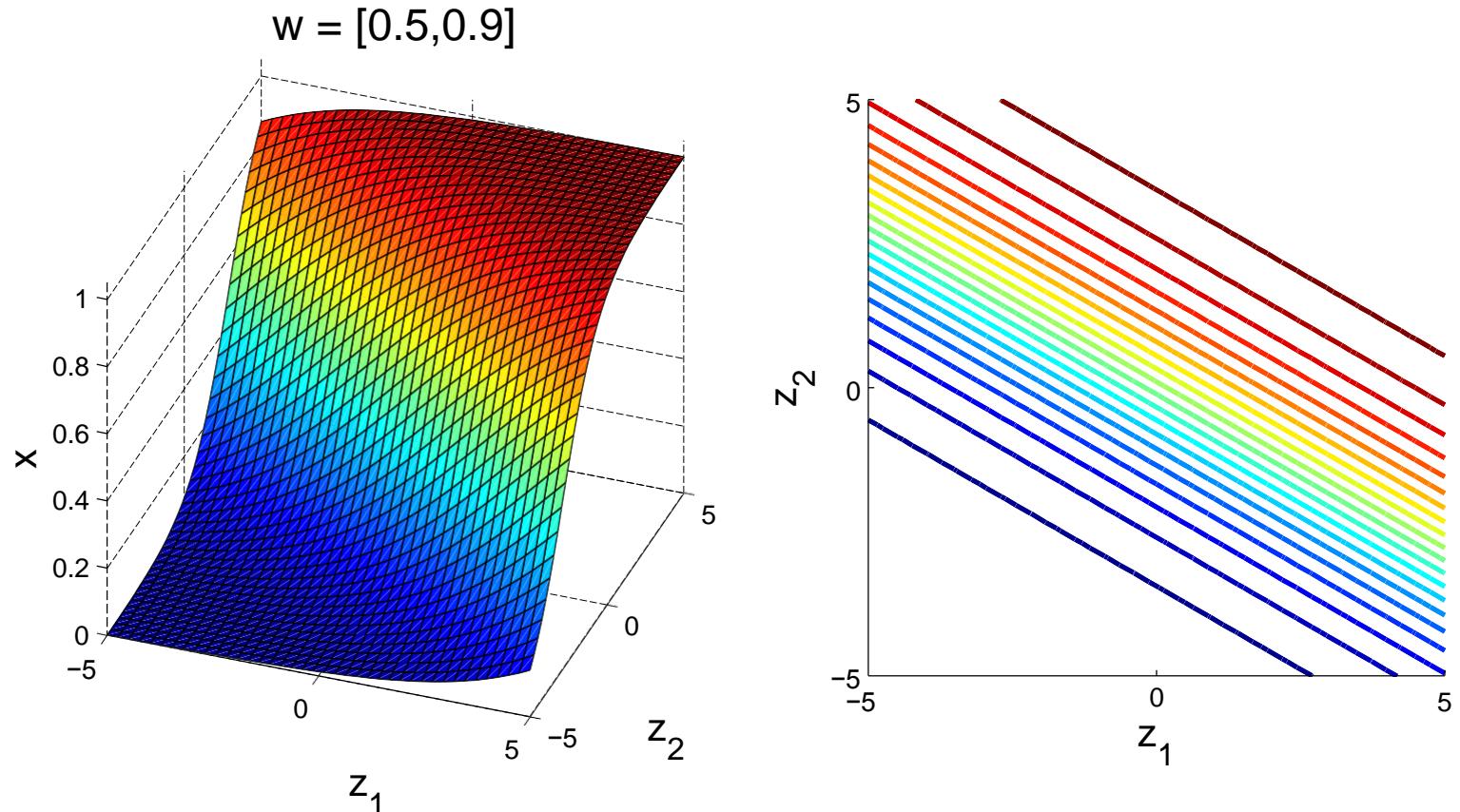
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron



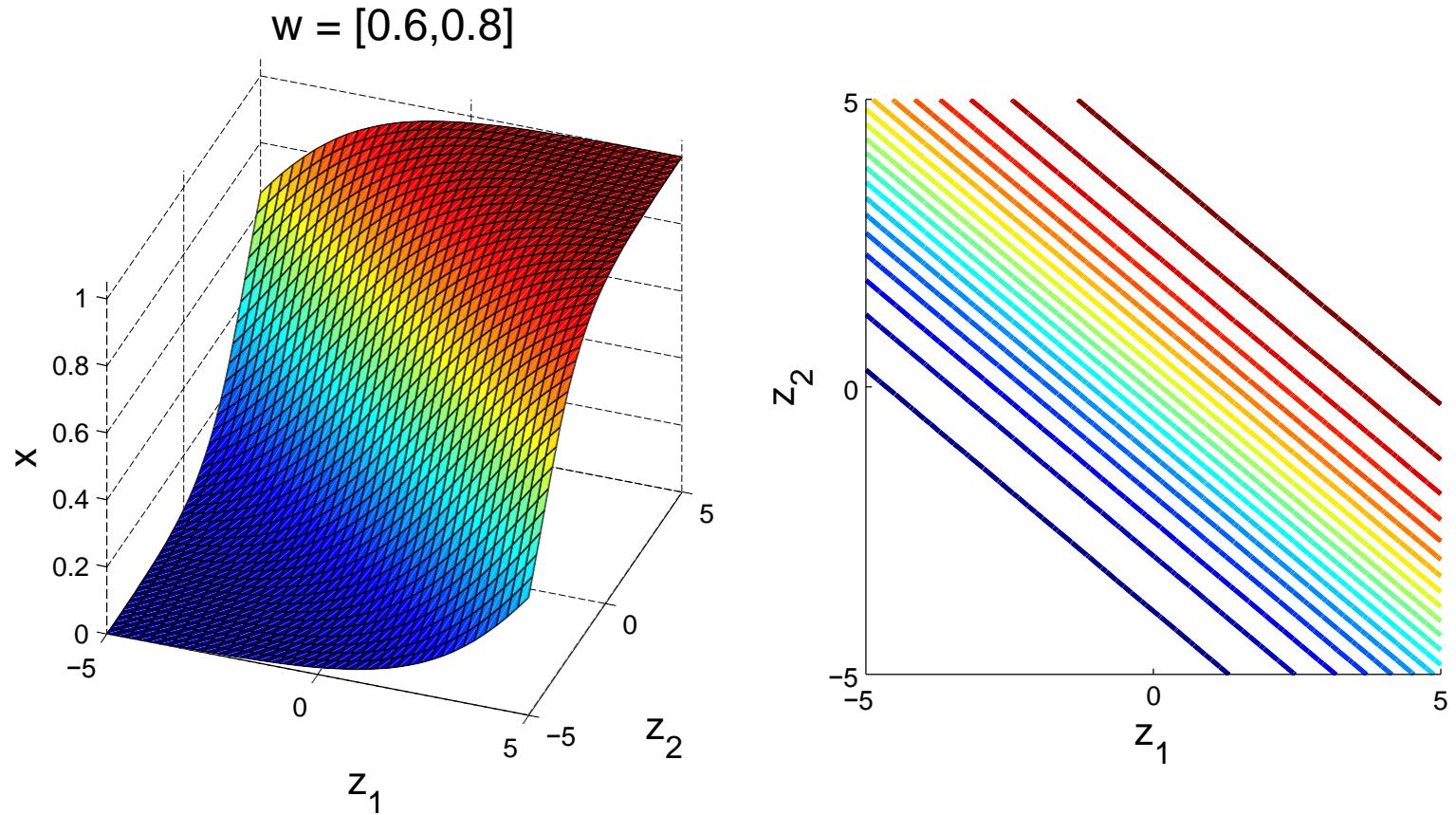
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron



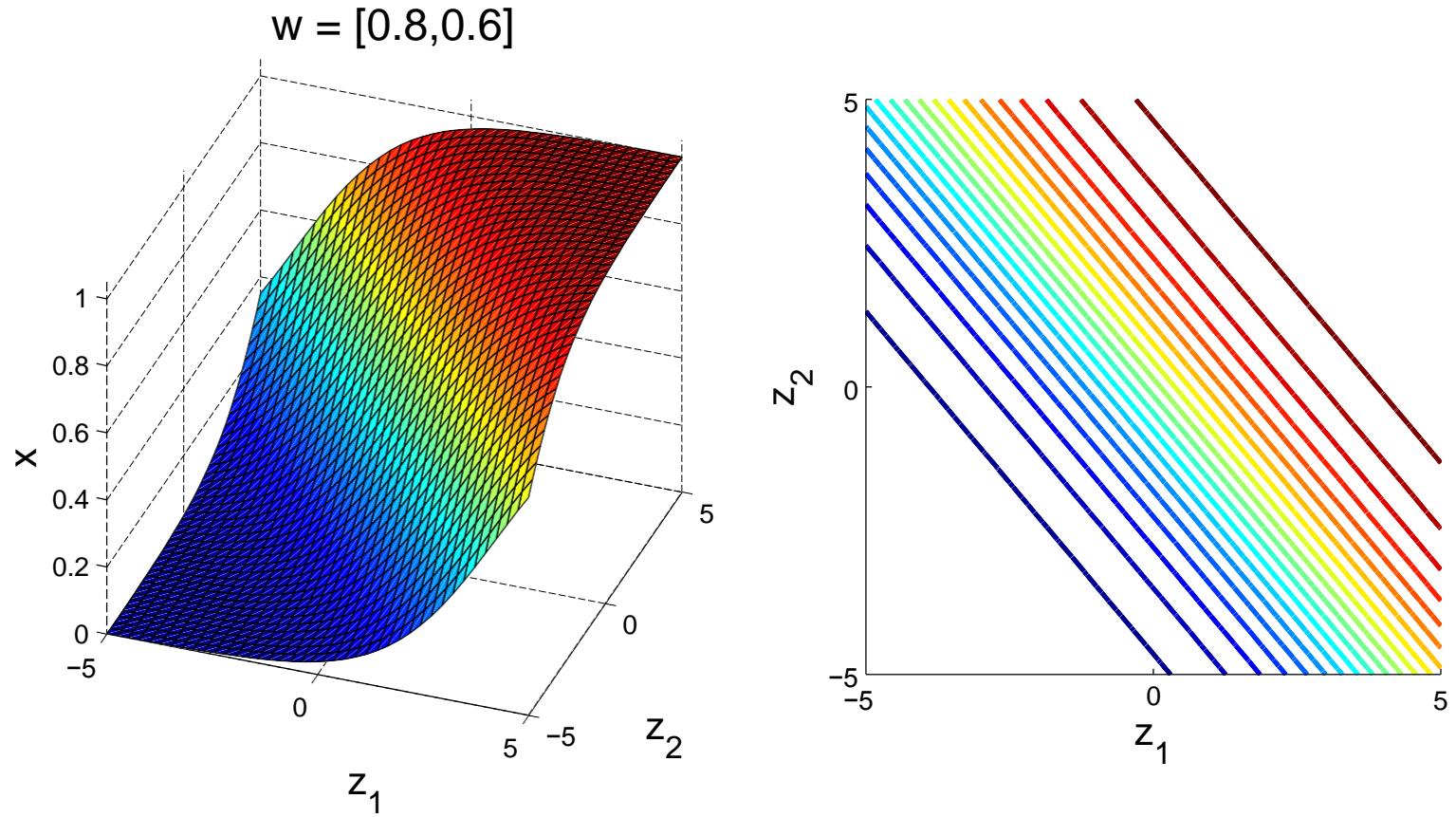
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron



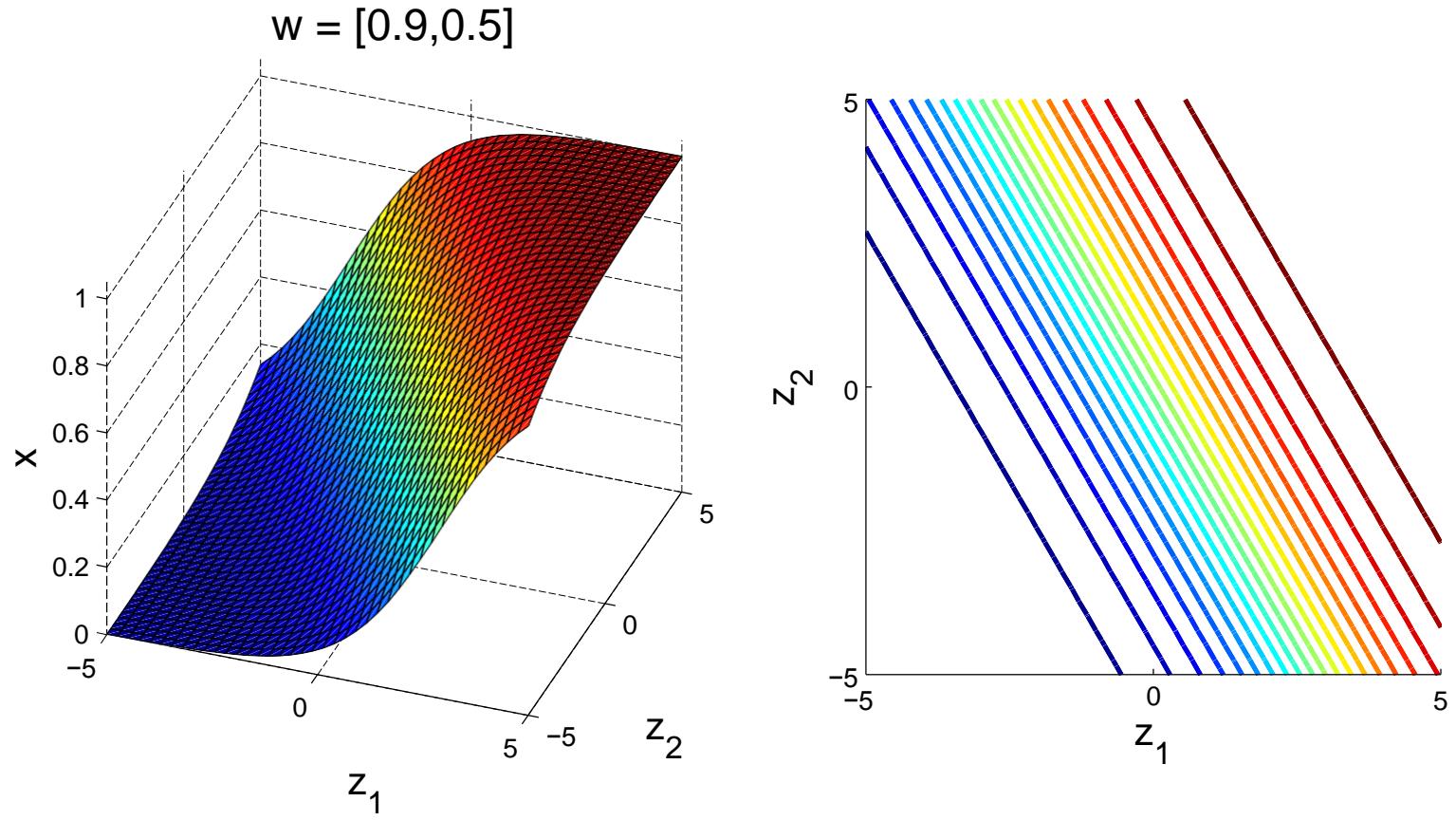
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron



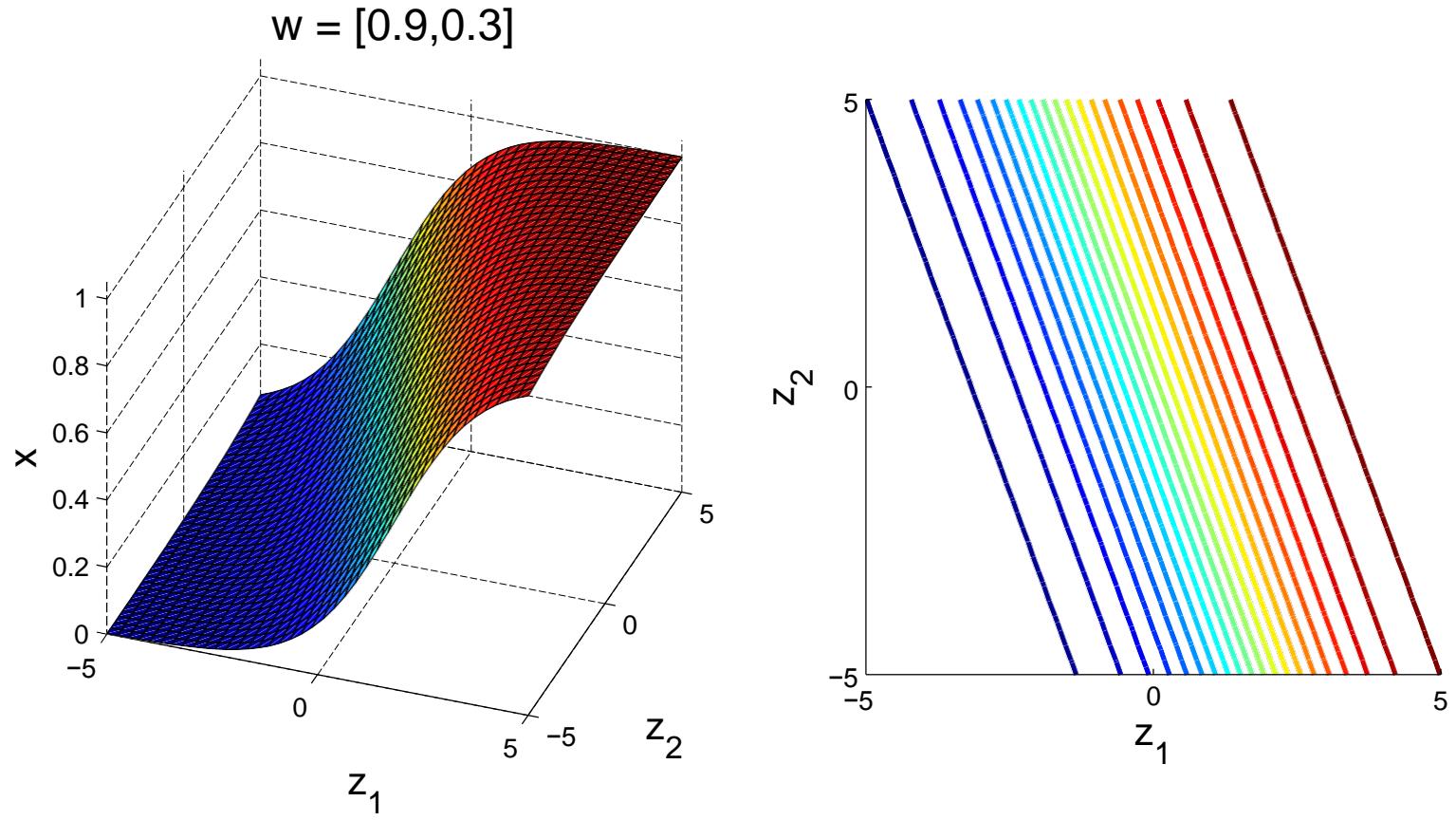
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron



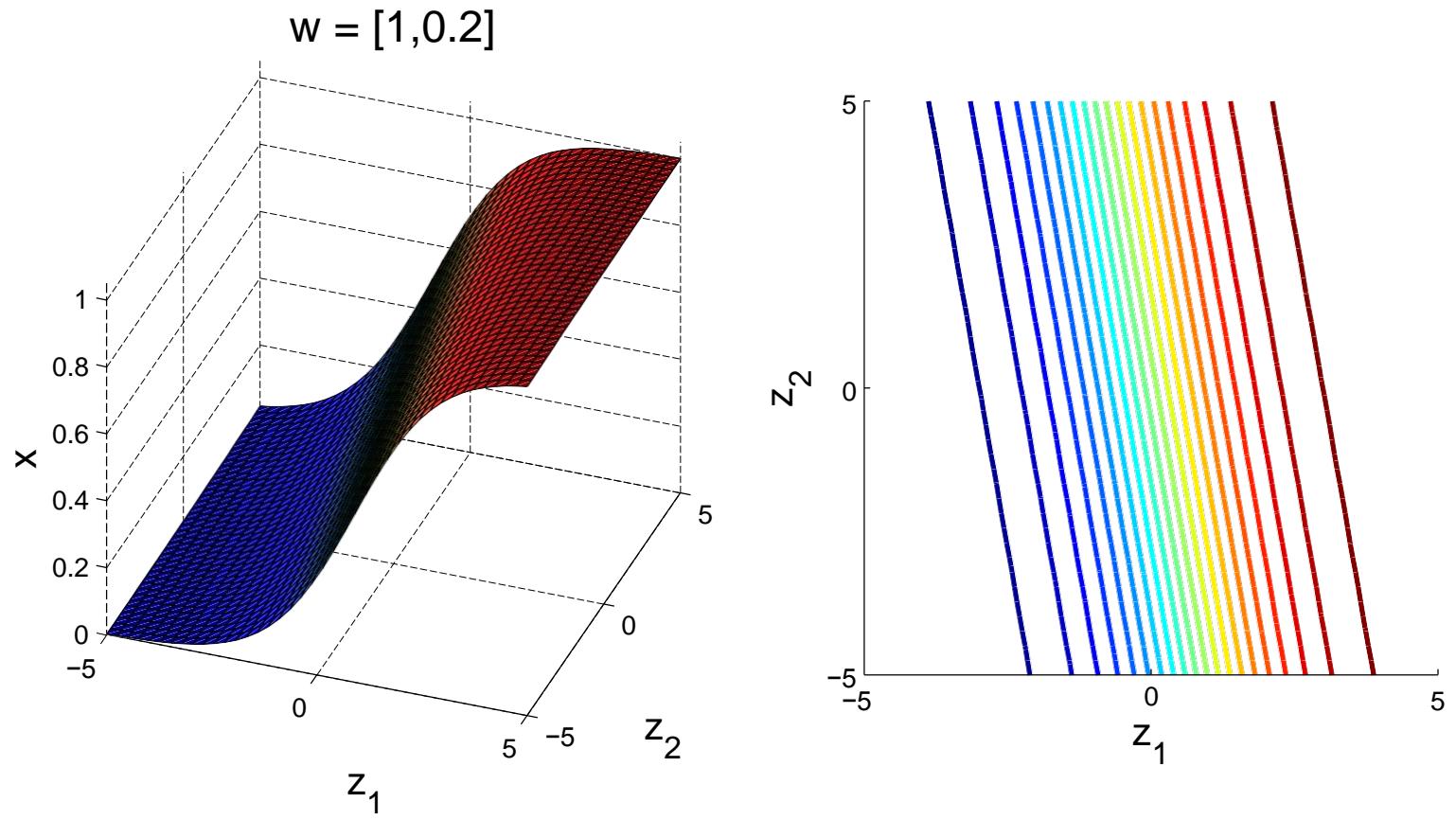
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron



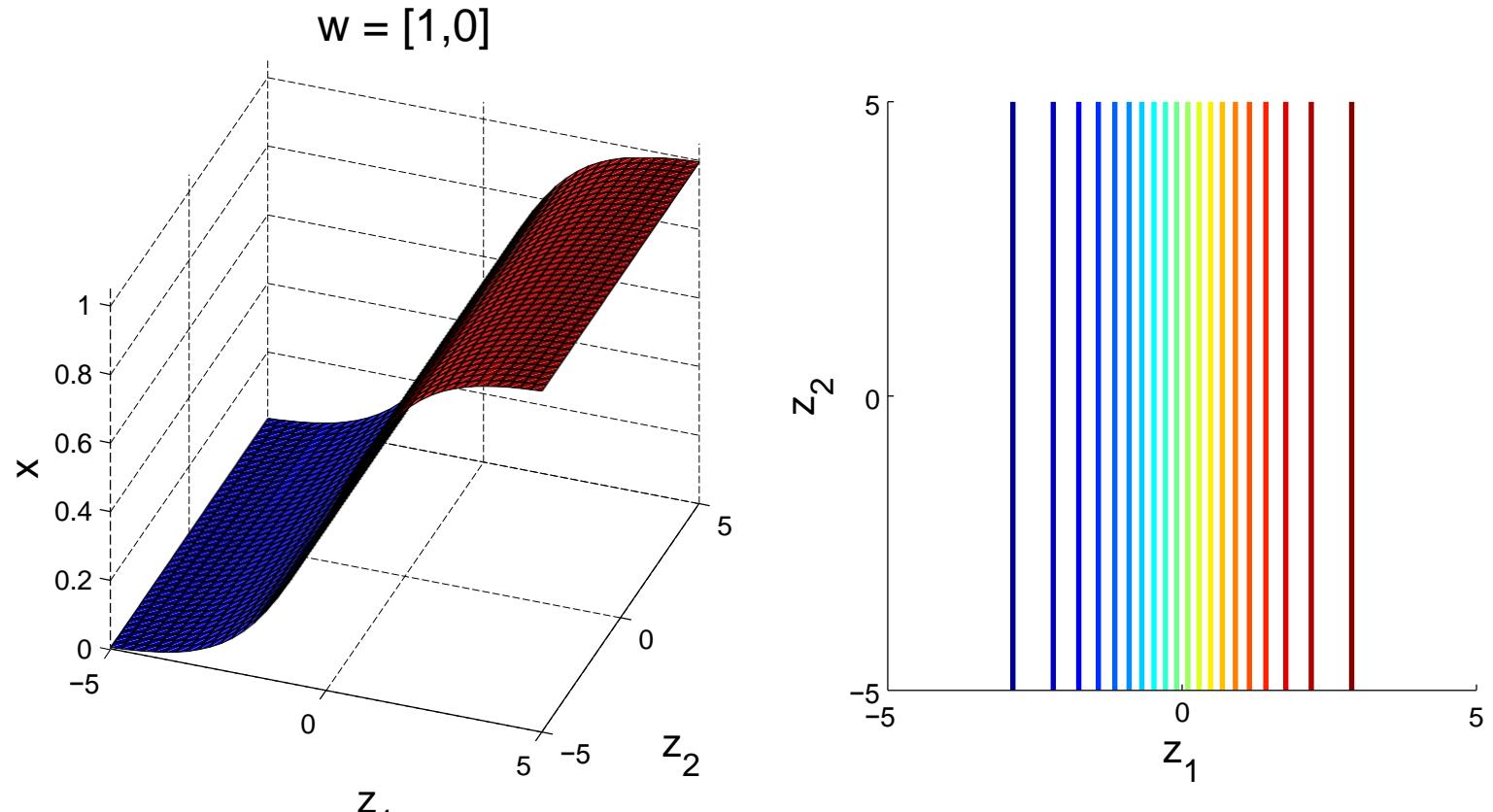
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron



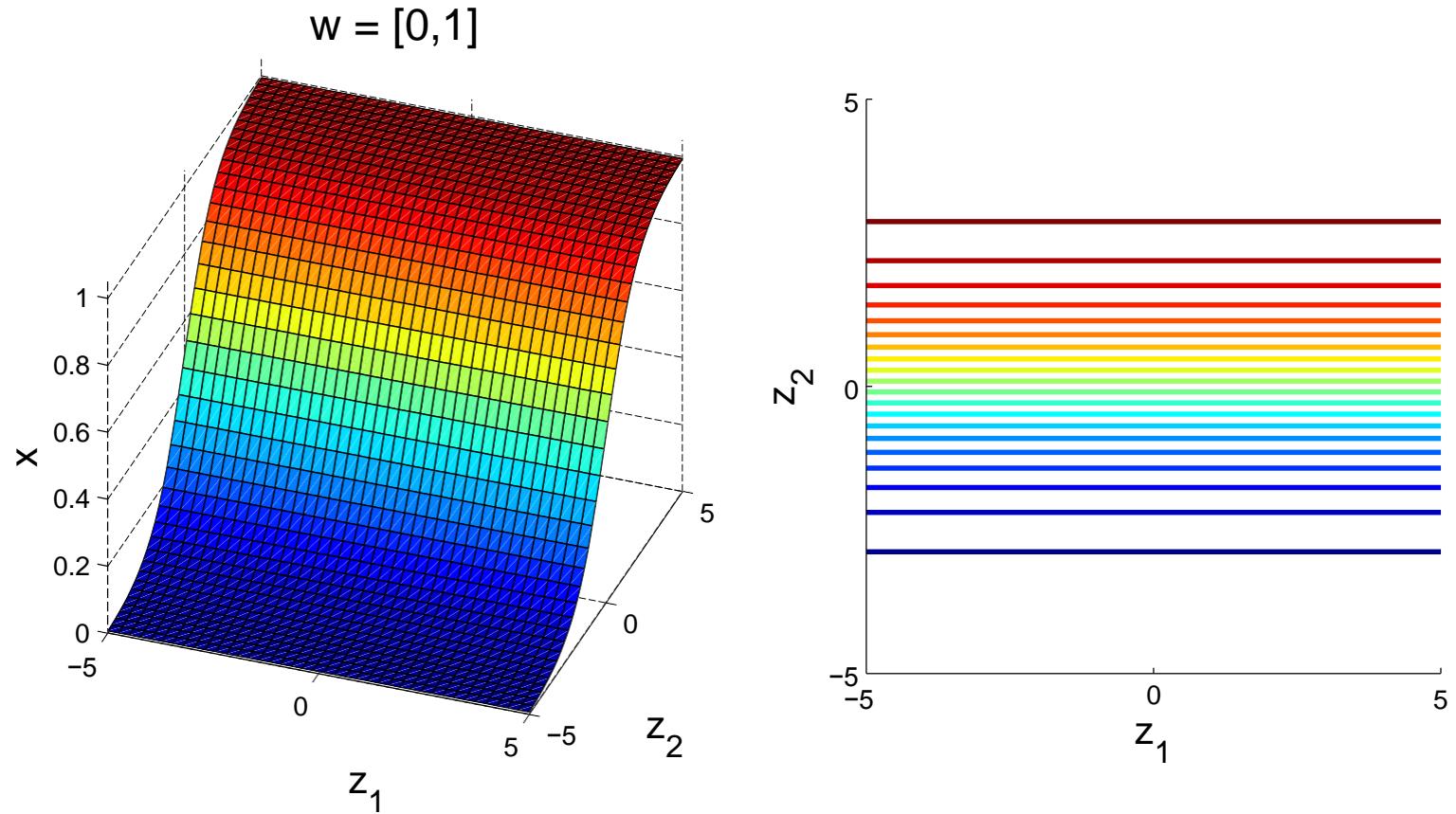
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron



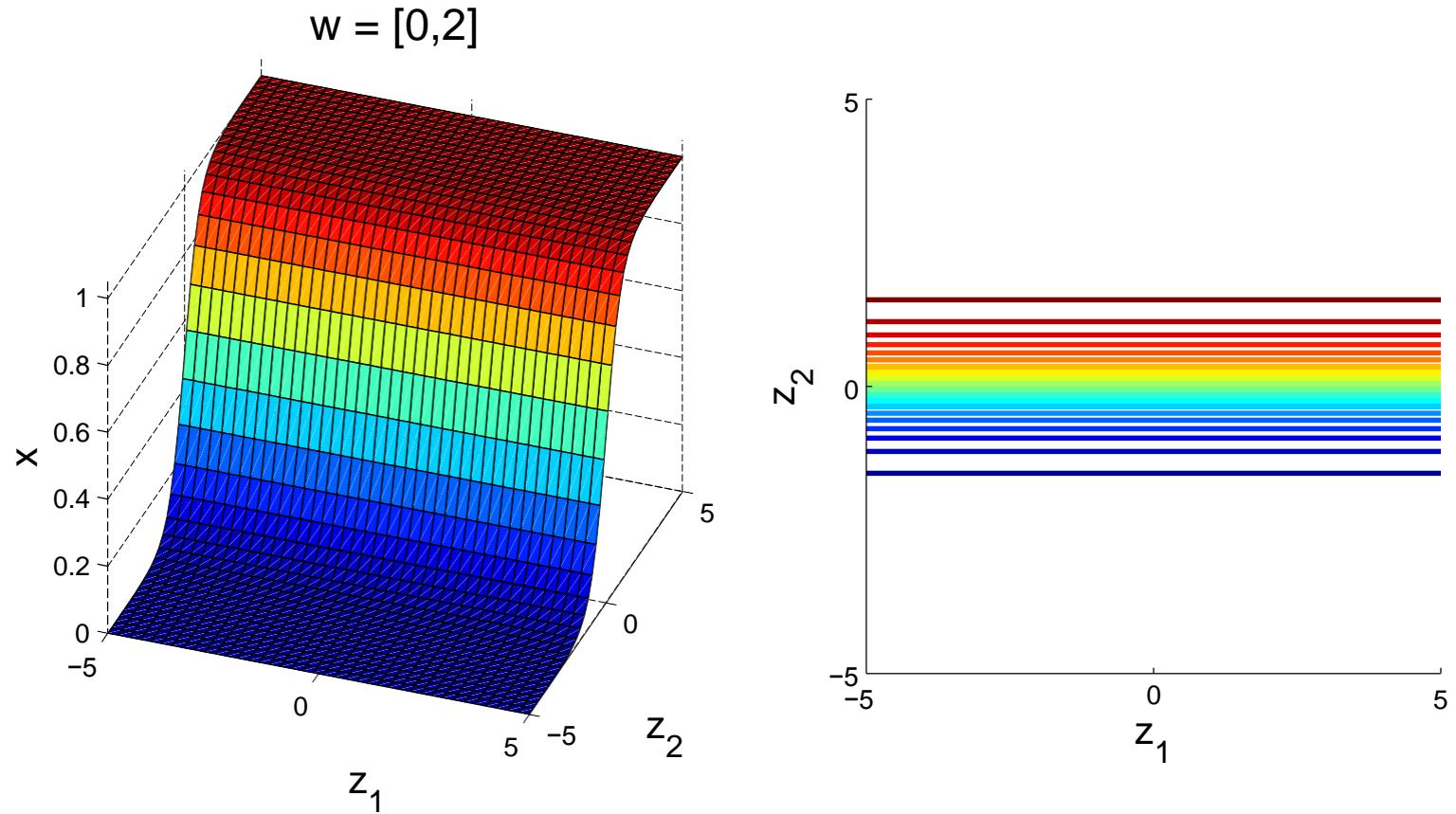
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

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



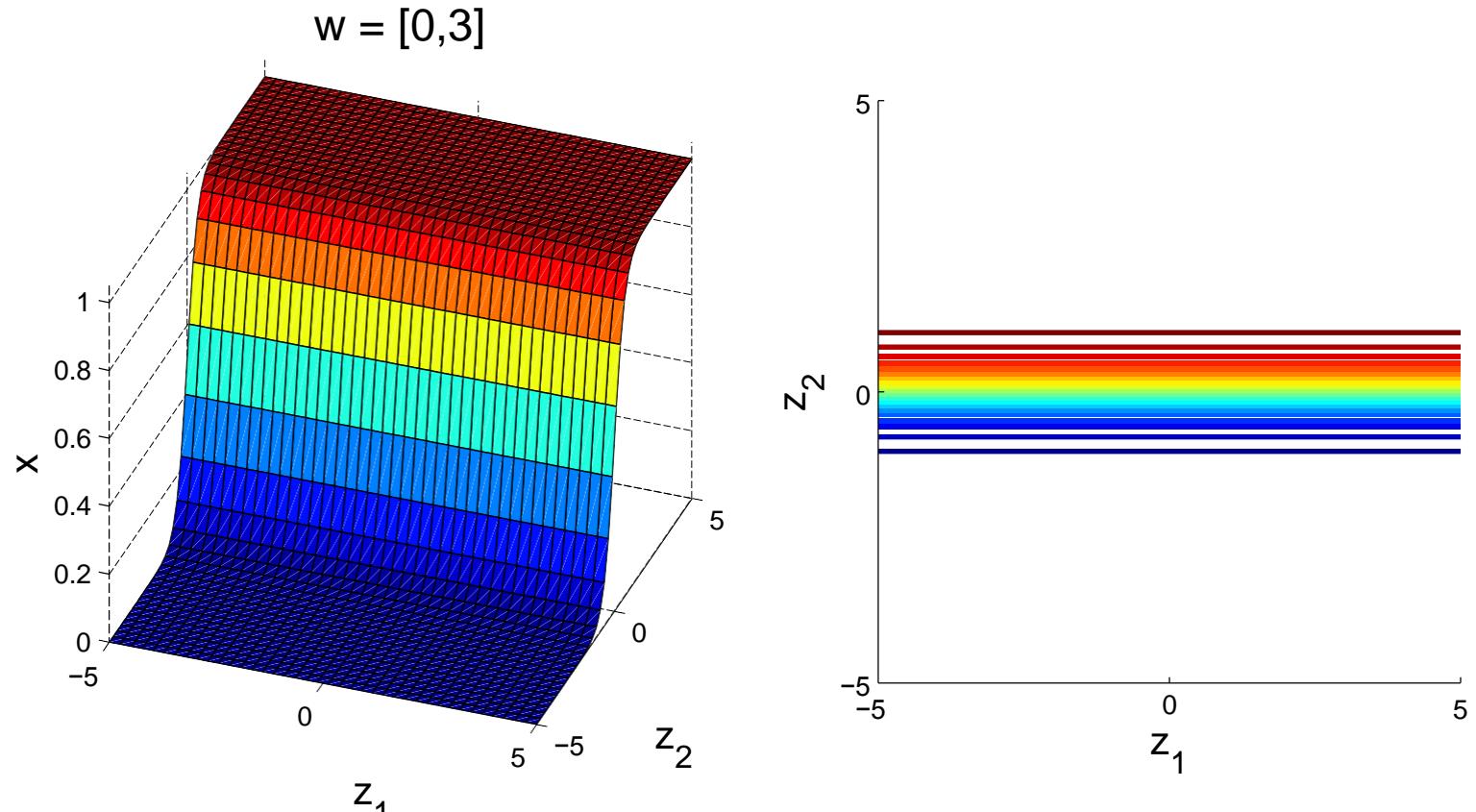
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

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



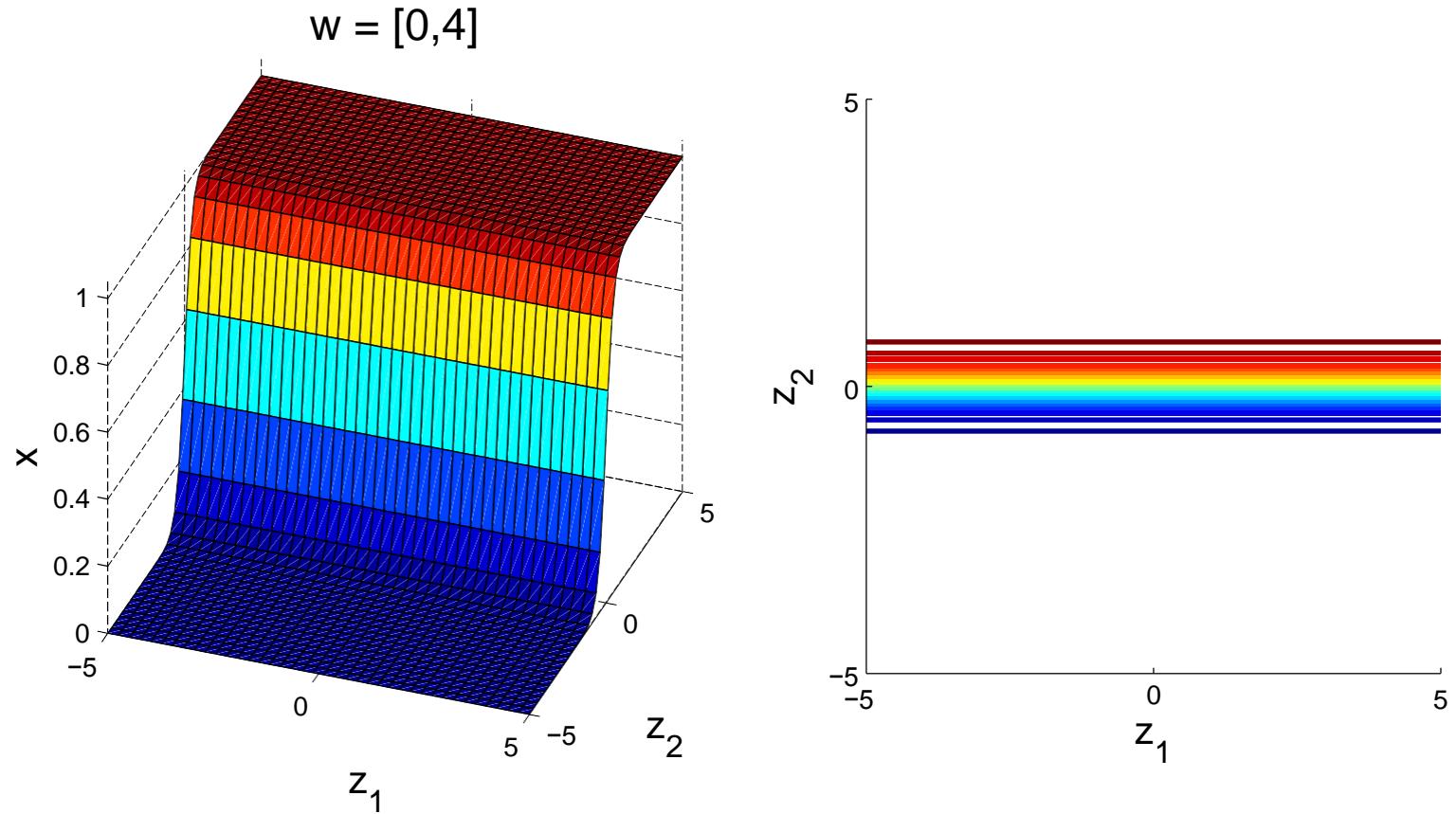
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

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



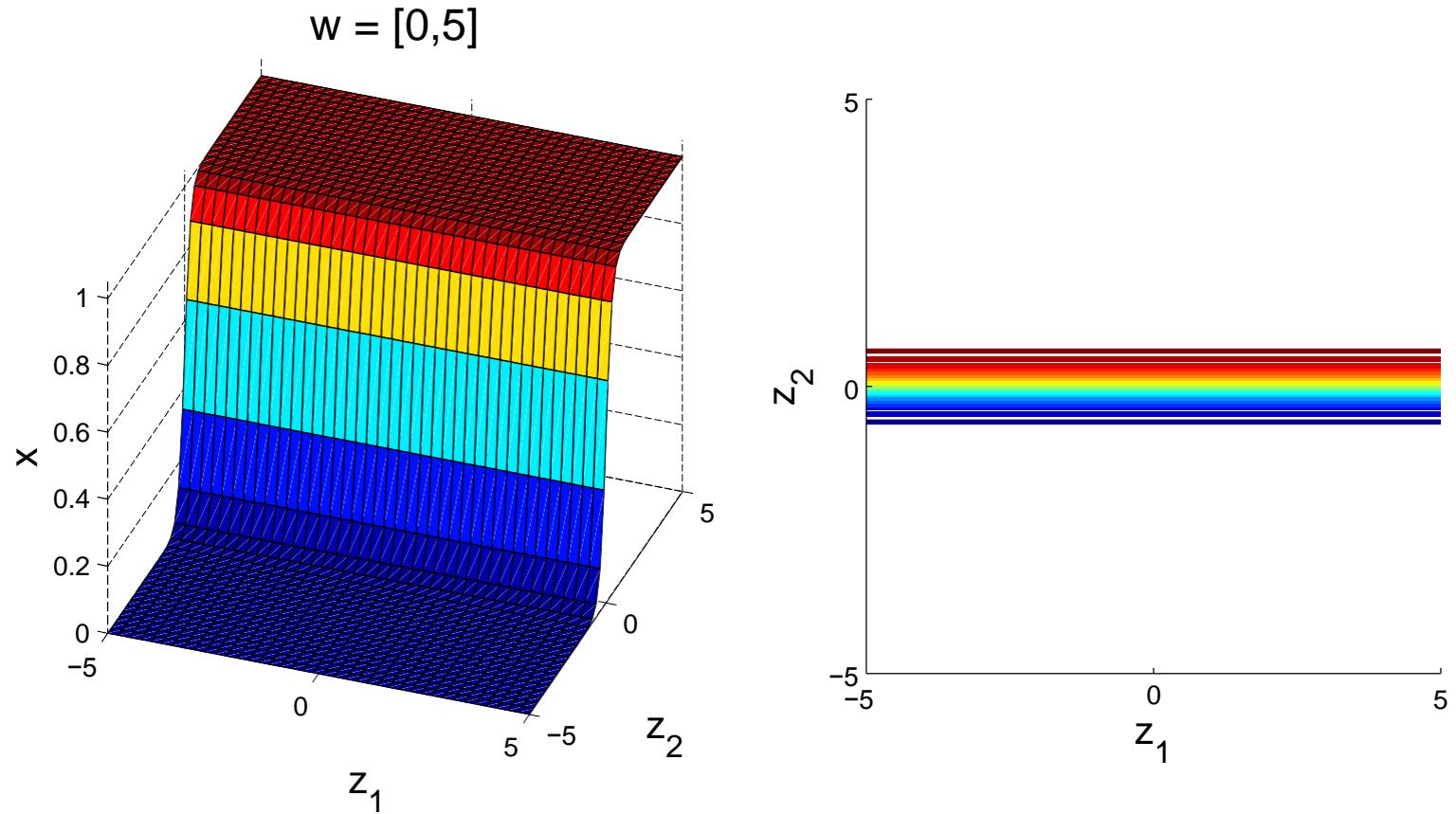
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

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



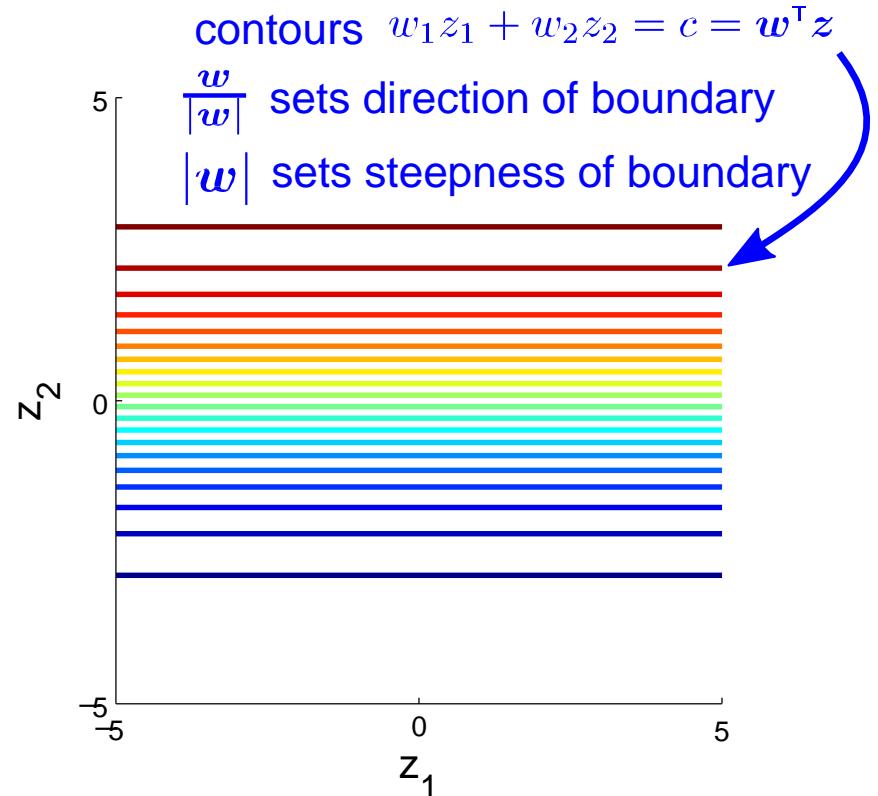
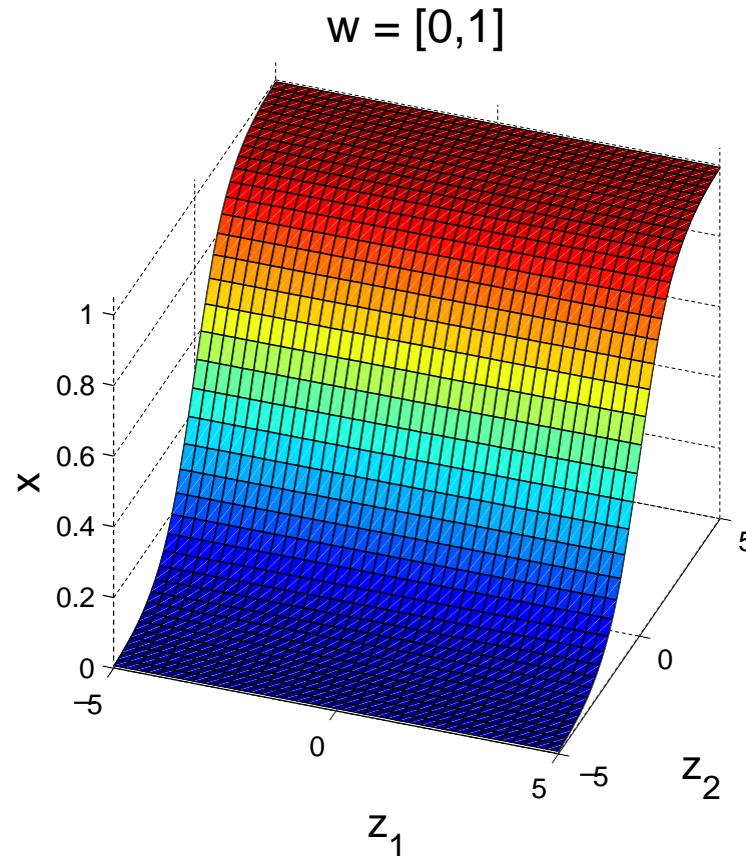
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

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



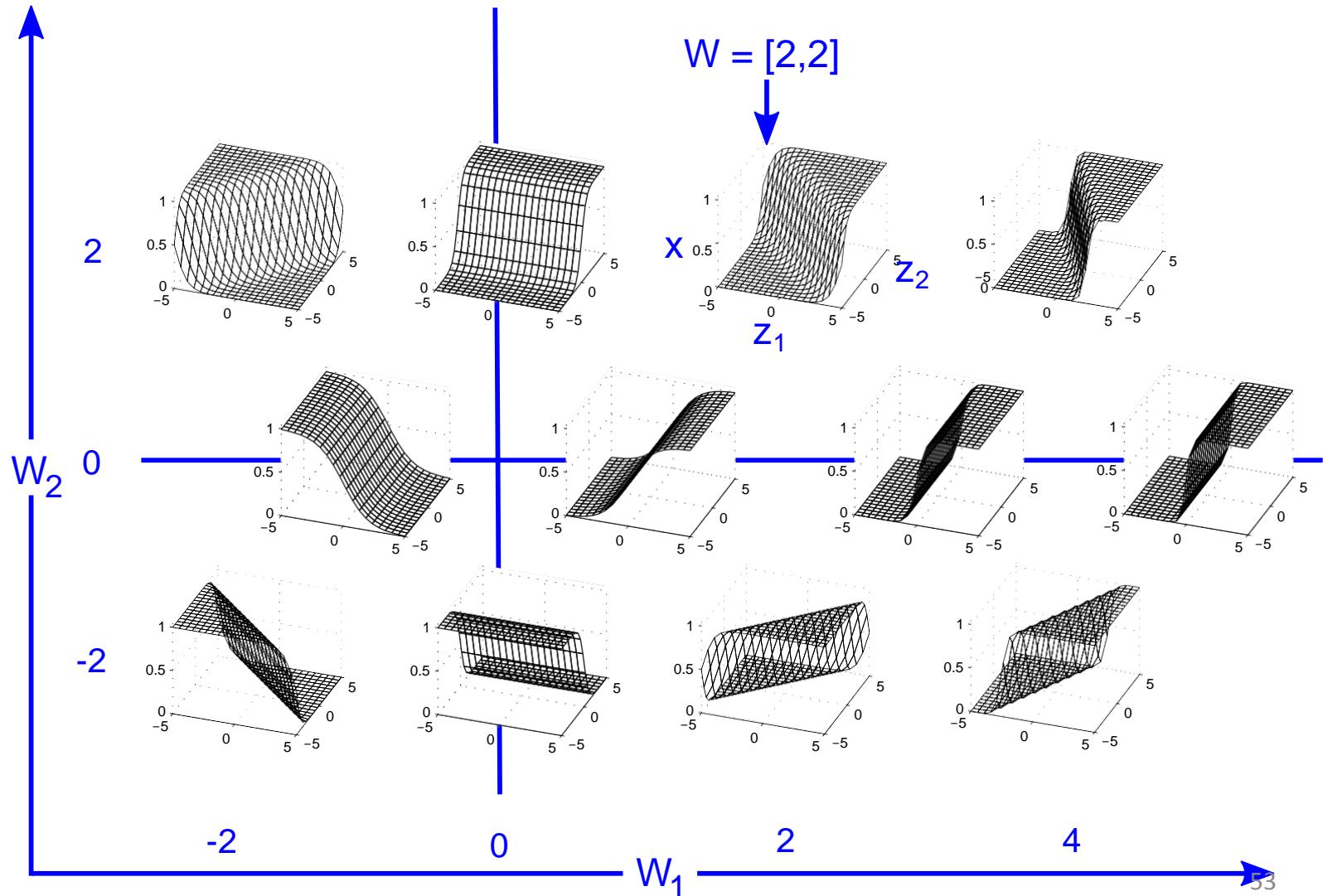
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

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

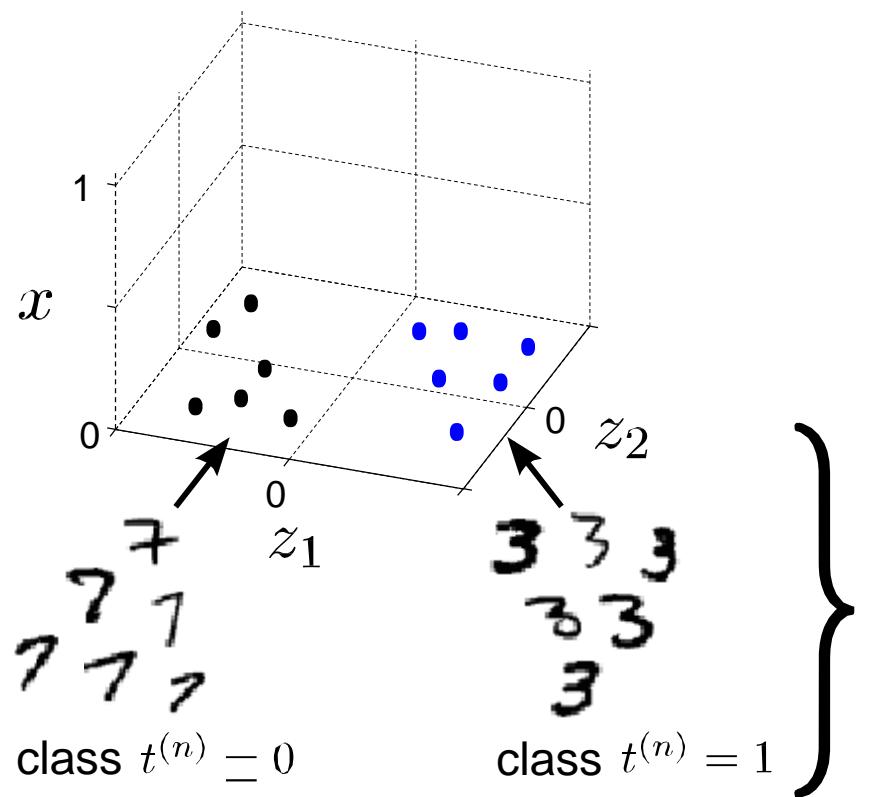


$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Weight Space of a Single Neuron



Training a Single Neuron

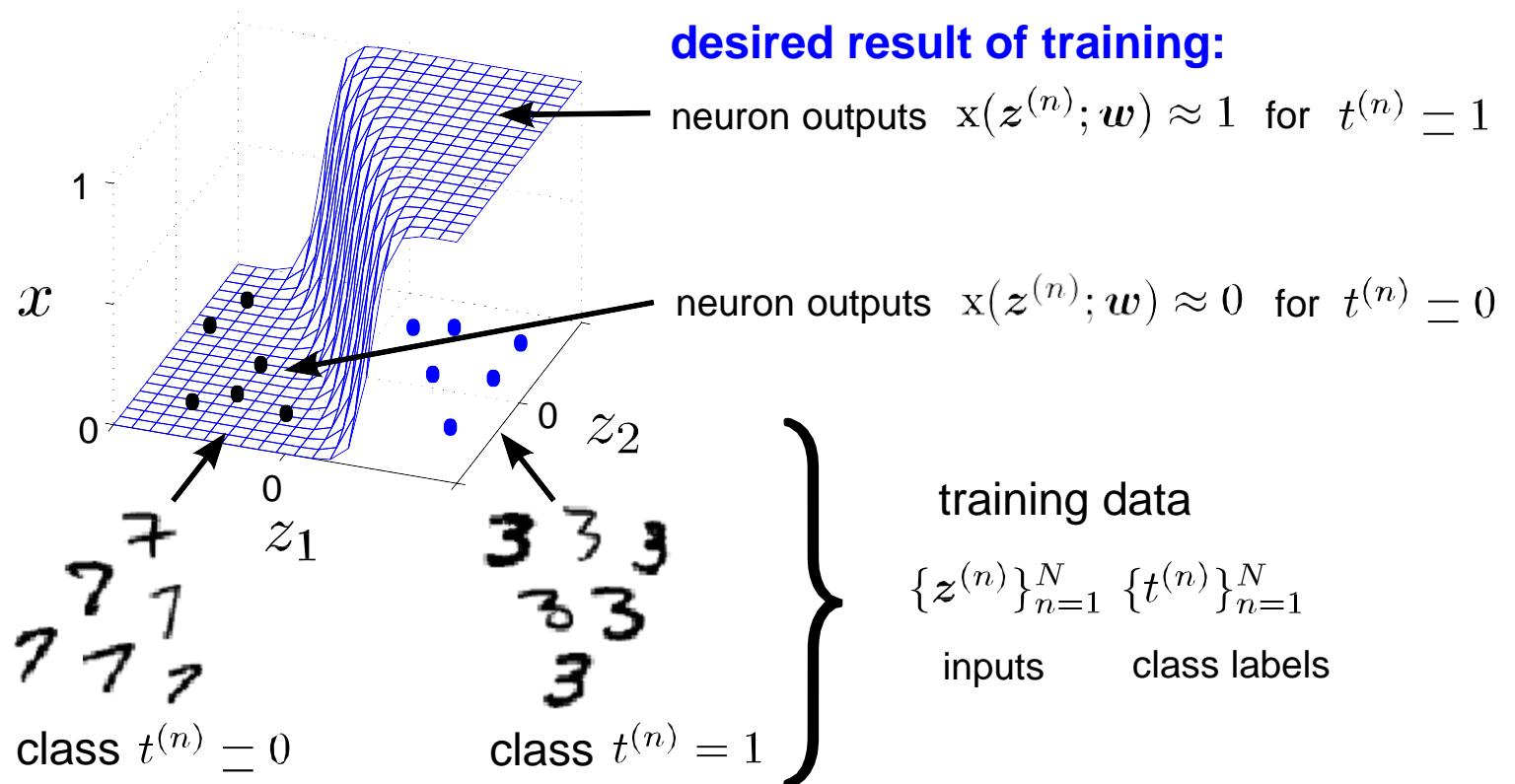


training data

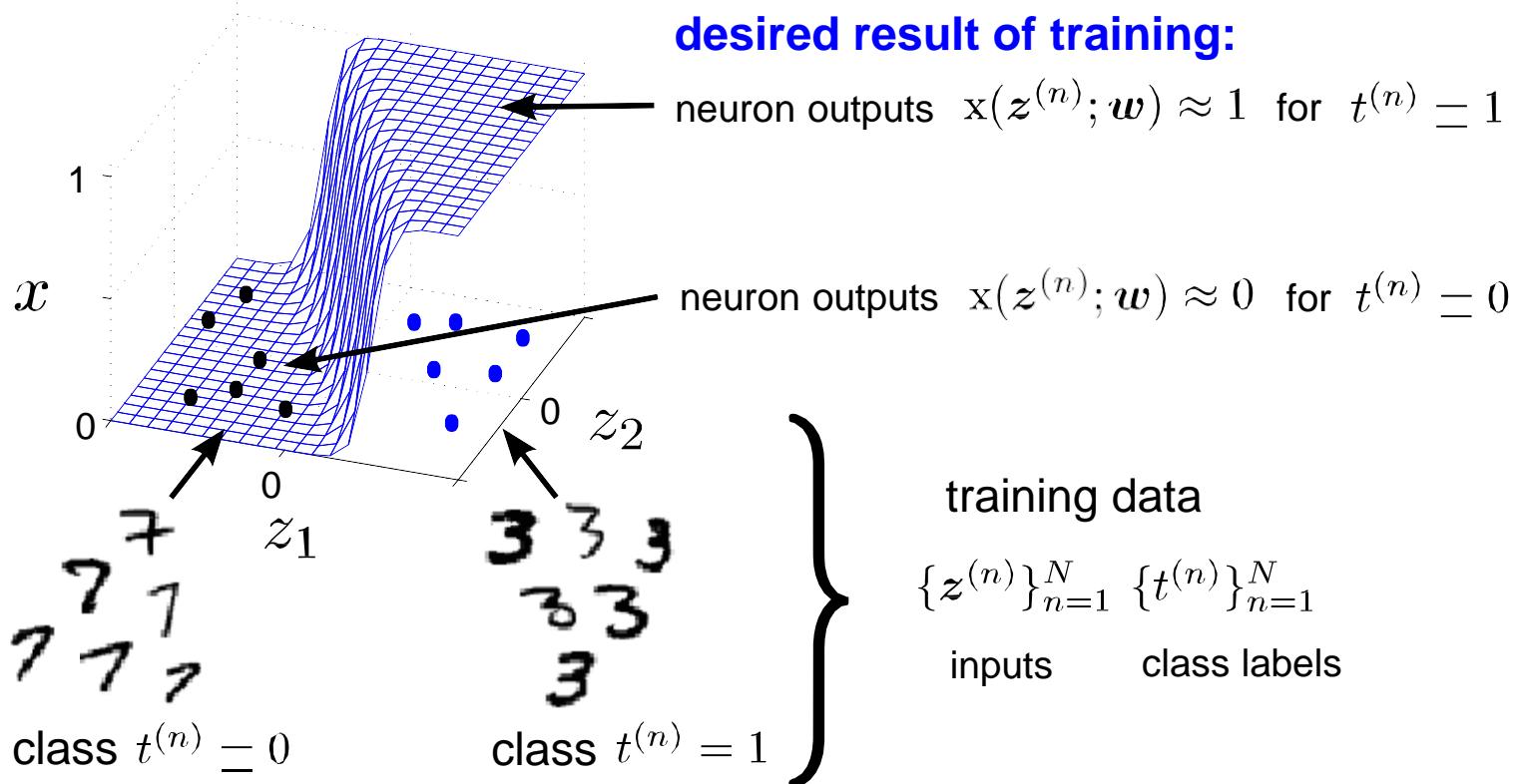
$$\{z^{(n)}\}_{n=1}^N \quad \{t^{(n)}\}_{n=1}^N$$

inputs class labels

Training a Single Neuron



Training a Single Neuron

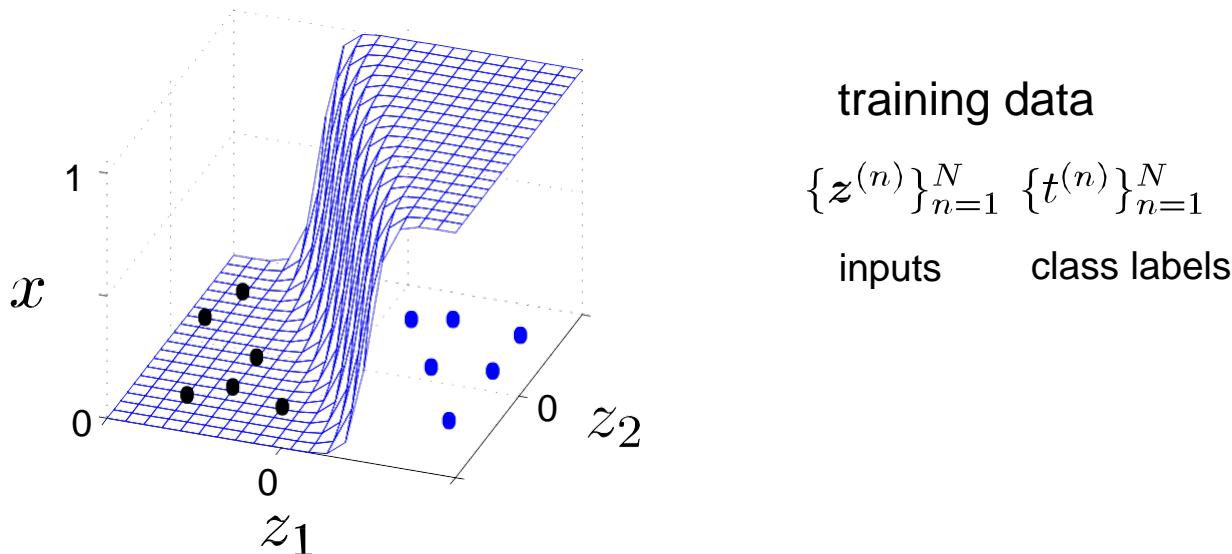


objective function:

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

surprise $-\log p(\text{outcome})$ when observing $t^{(n)}$ } encourages neuron output
relative entropy between $x(z^{(n)}; \mathbf{w})$ and $t^{(n)}$ } to match training data 56

Training a Single Neuron



training data

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

inputs

class labels

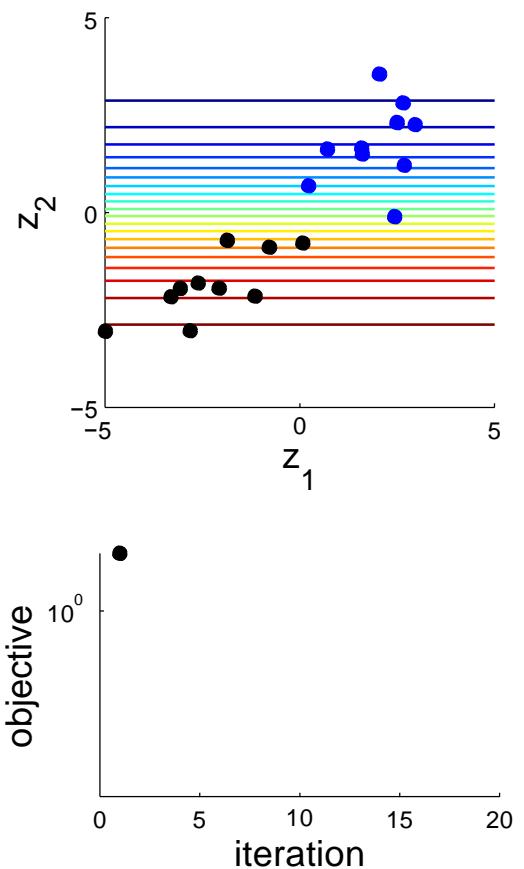
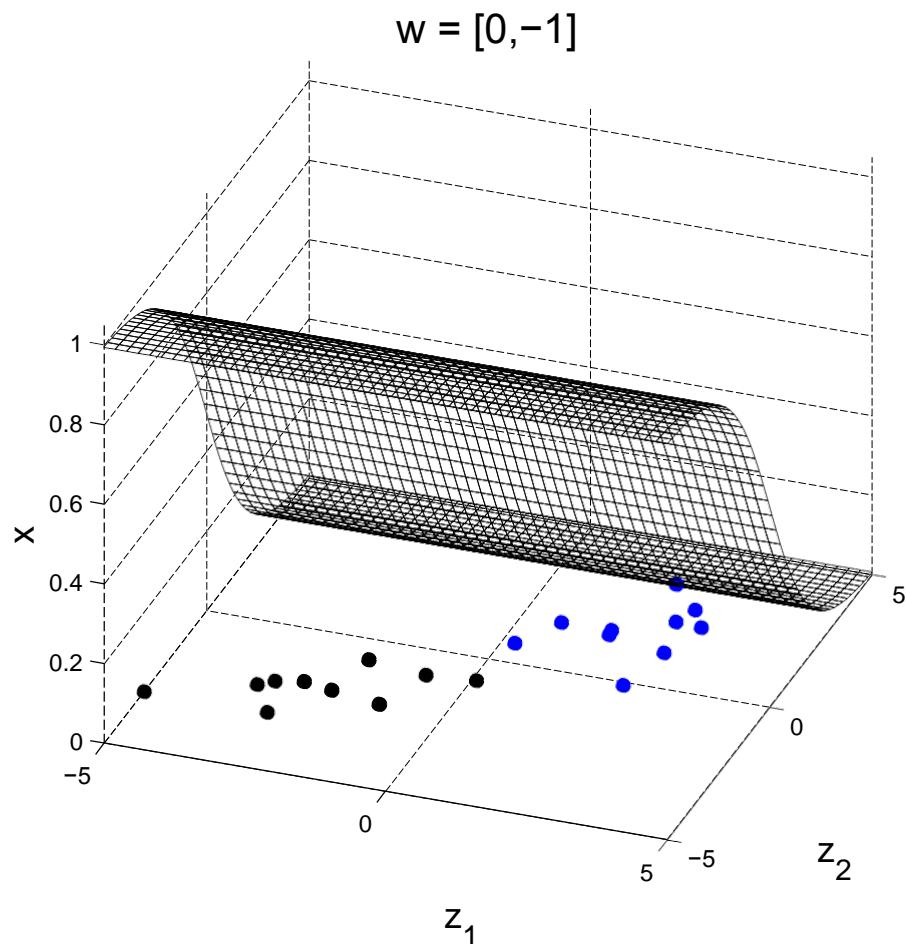
objective function:

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

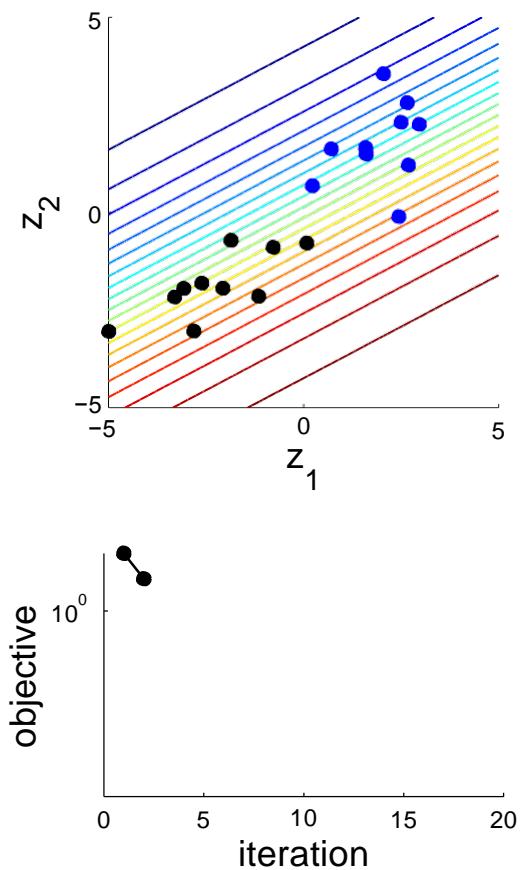
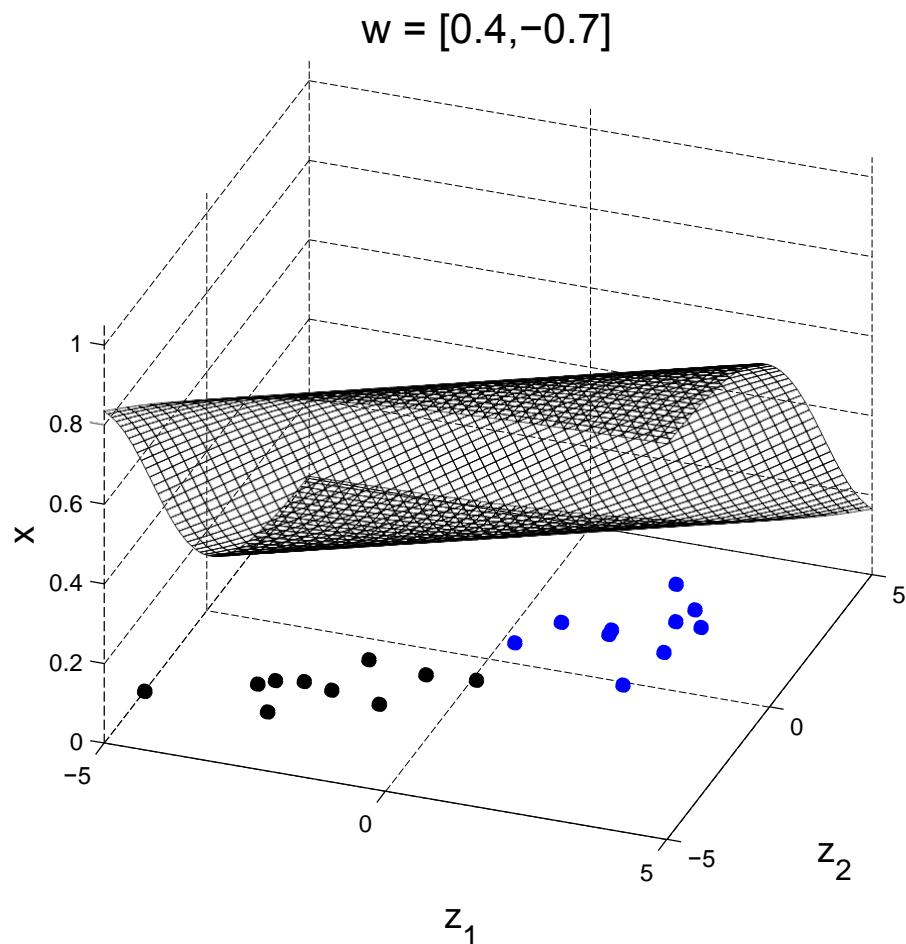
$$\mathbf{w}^* = \arg \min_{\mathbf{w}} G(\mathbf{w}) \quad \text{choose the weights that minimise the network's surprise about the training data}$$

$$\frac{d}{d\mathbf{w}} G(\mathbf{w}) = \sum_n \frac{dG(\mathbf{w})}{dx^{(n)}} \frac{dx^{(n)}}{d\mathbf{w}} = - \sum_n (t^{(n)} - x^{(n)}) \mathbf{z}^{(n)} = \text{prediction error} \times \text{feature}$$
$$\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{d}{d\mathbf{w}} G(\mathbf{w}) \quad \text{iteratively step down the objective (gradient points up hill)}$$

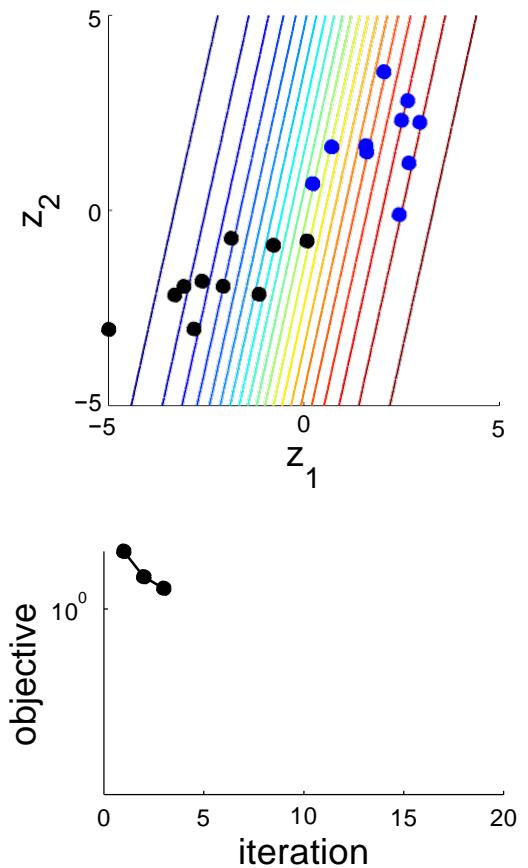
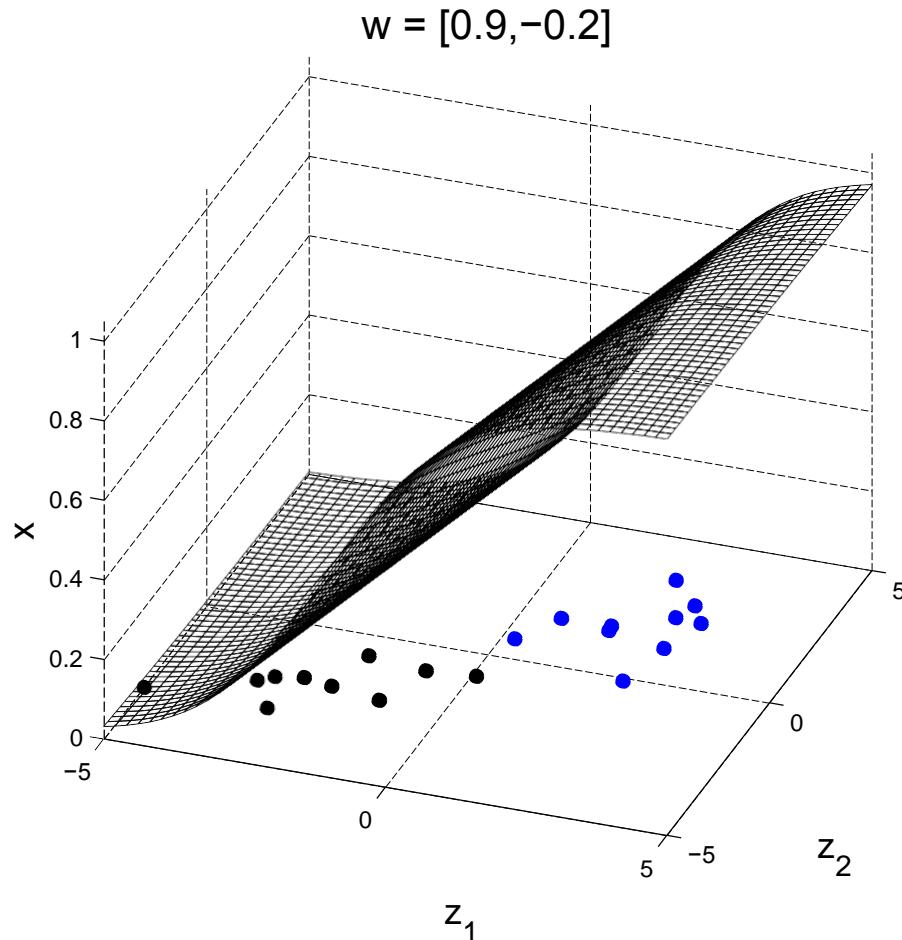
Training a Single Neuron



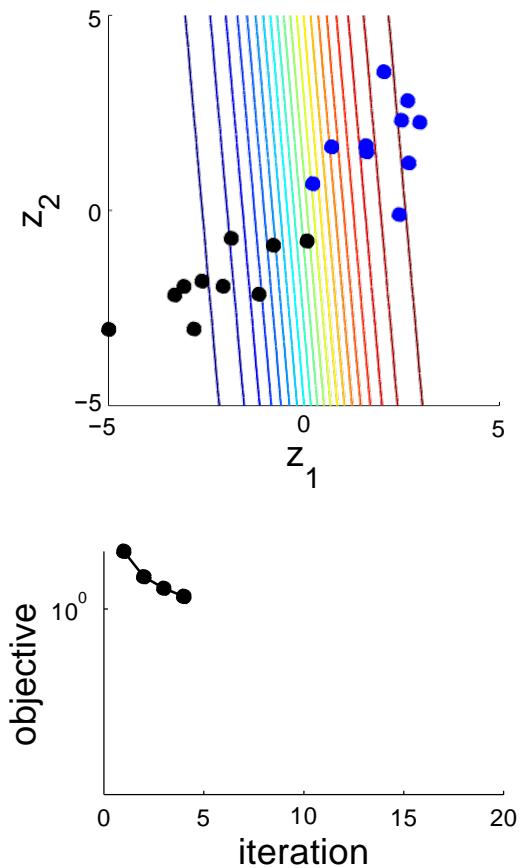
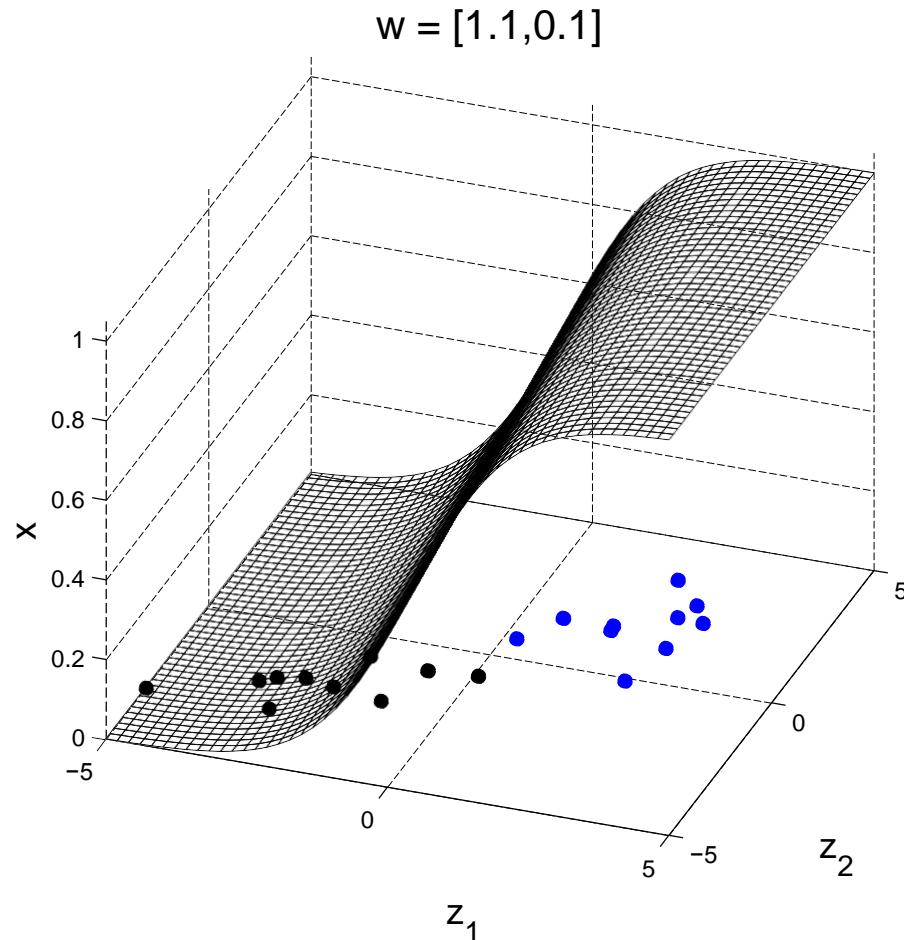
Training a Single Neuron



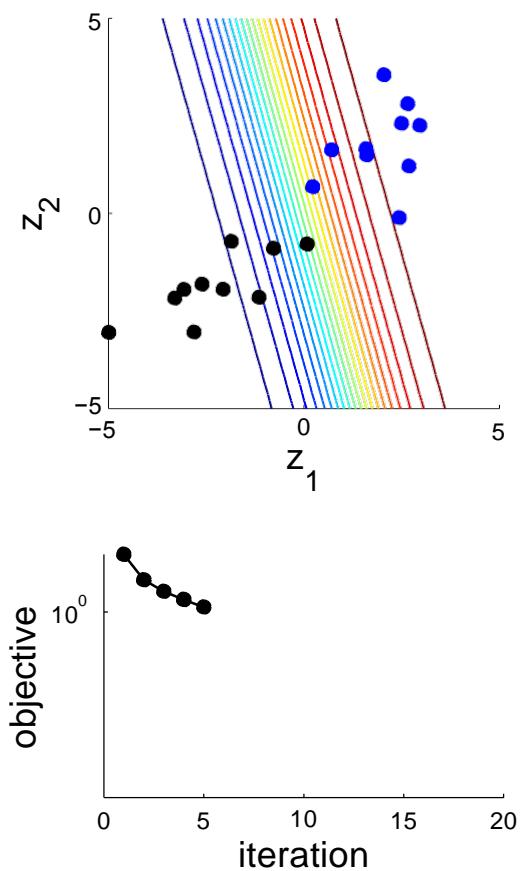
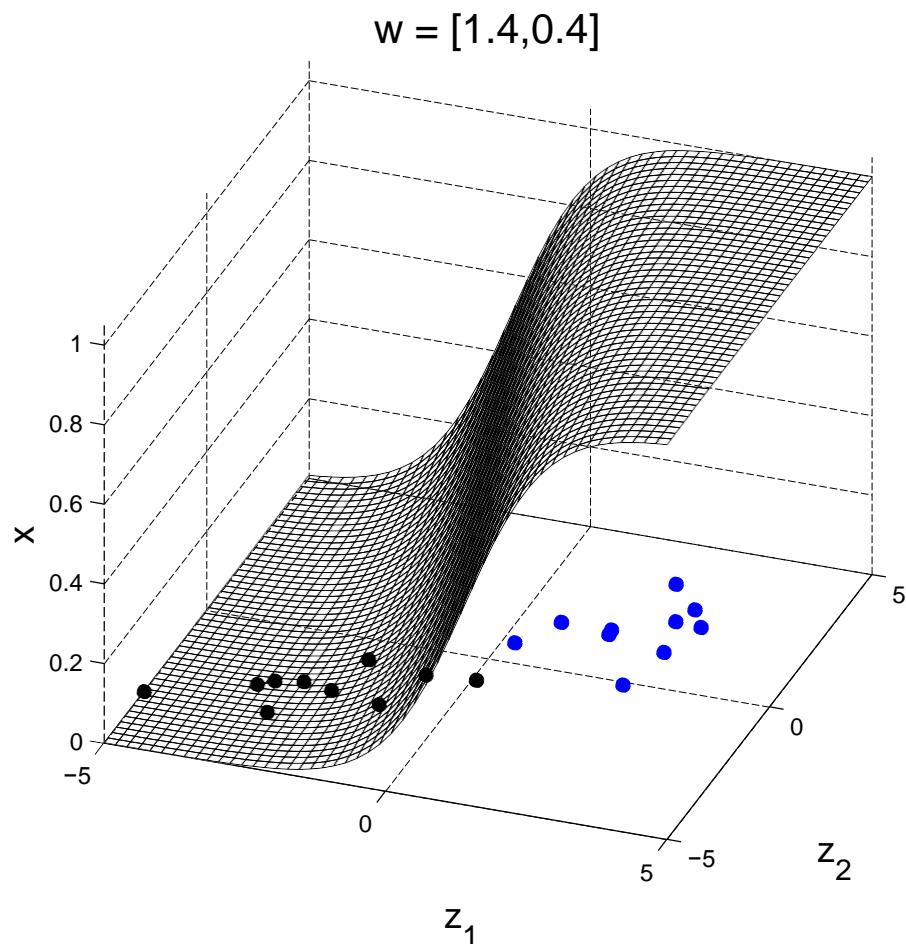
Training a Single Neuron



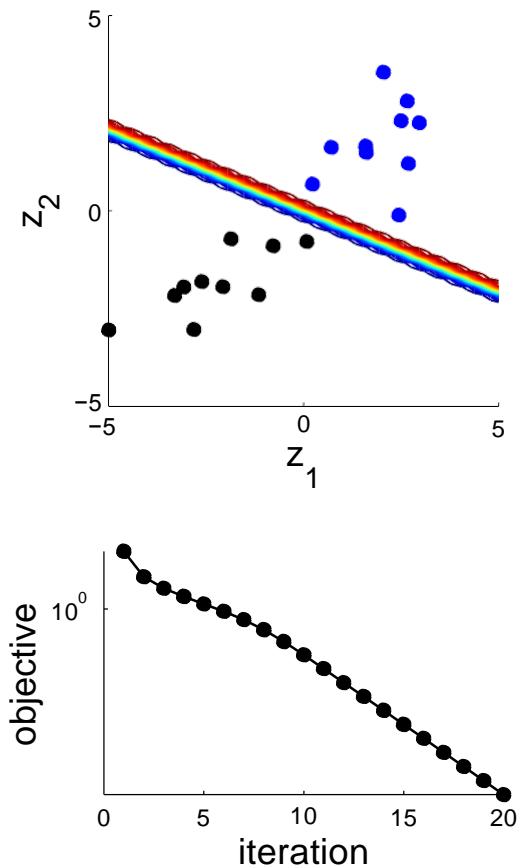
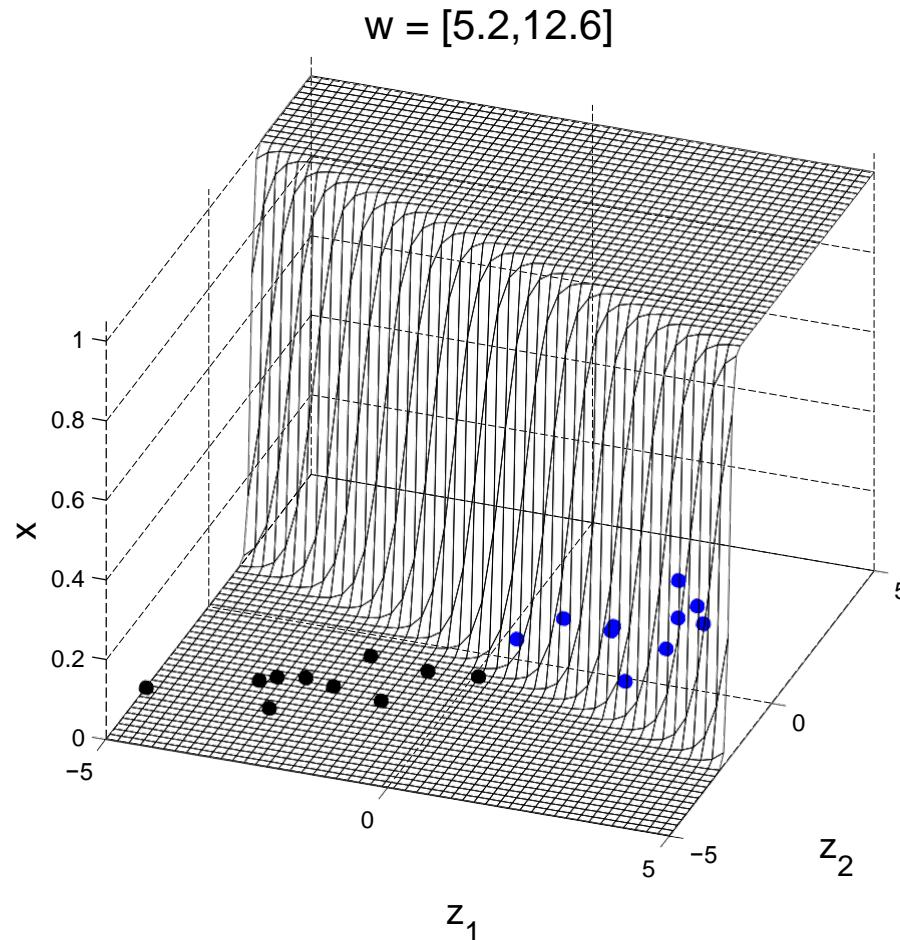
Training a Single Neuron



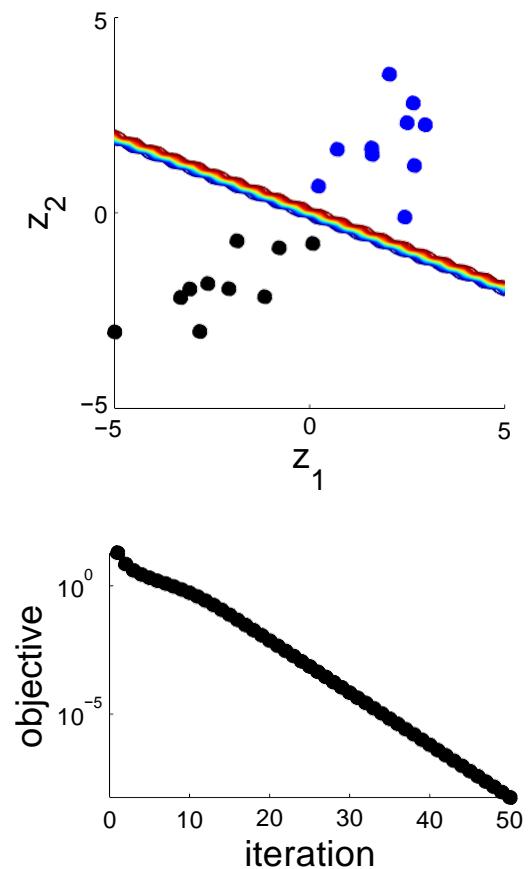
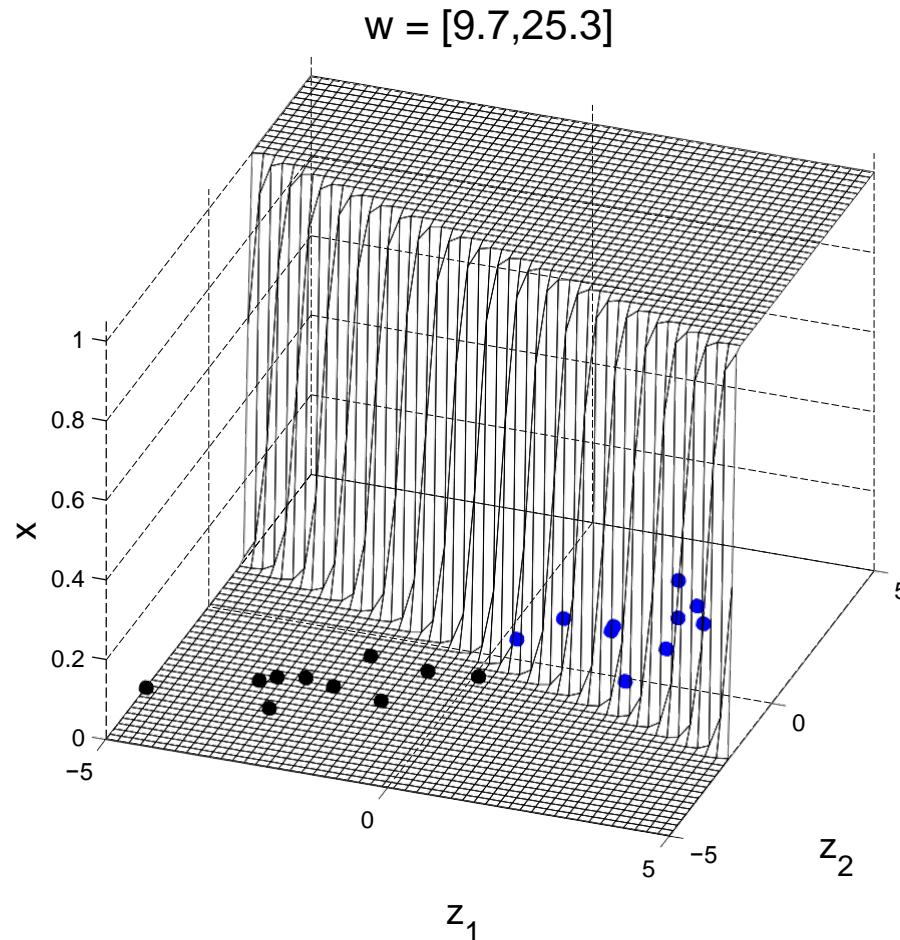
Training a Single Neuron



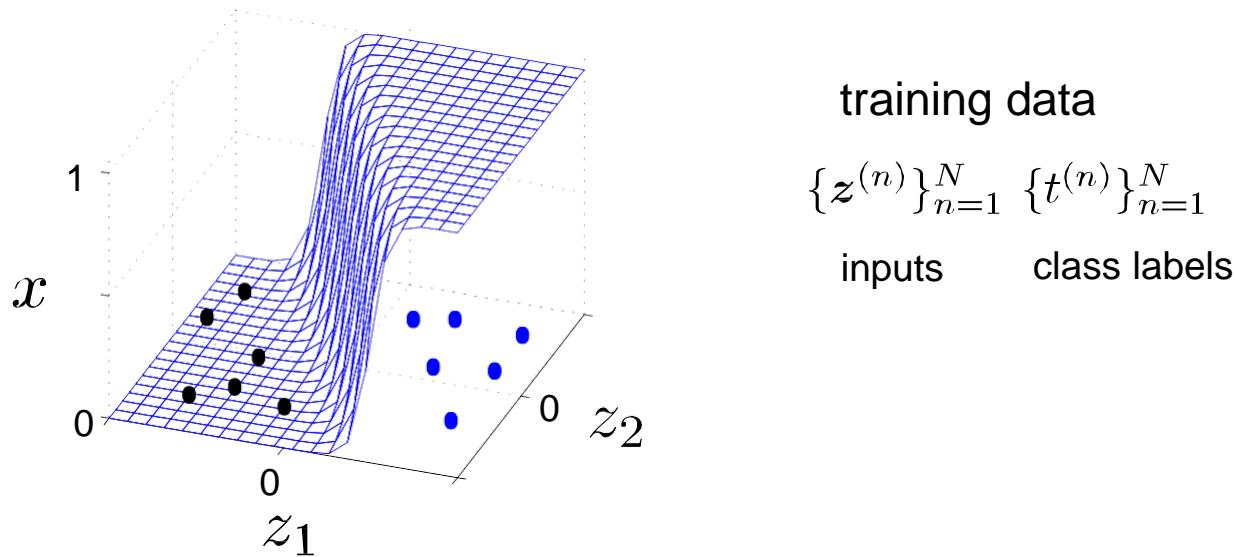
Training a Single Neuron



Training a Single Neuron



Overfitting and Weight Decay



objective function:

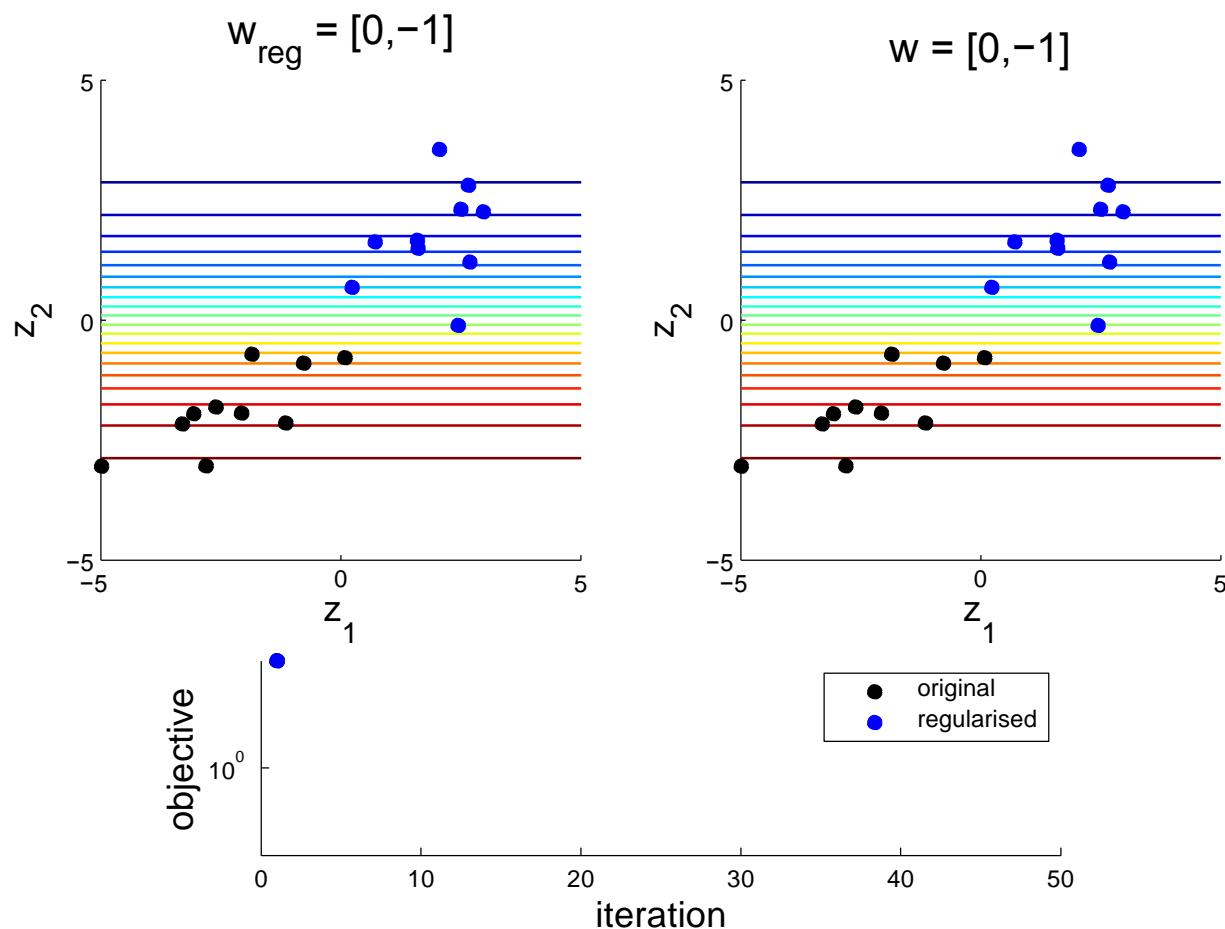
$$G(\mathbf{w}) = - \sum_n [t^{(n)} \log x(\mathbf{z}^{(n)}; \mathbf{w}) + (1 - t^{(n)}) \log (1 - x(\mathbf{z}^{(n)}; \mathbf{w}))]$$

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

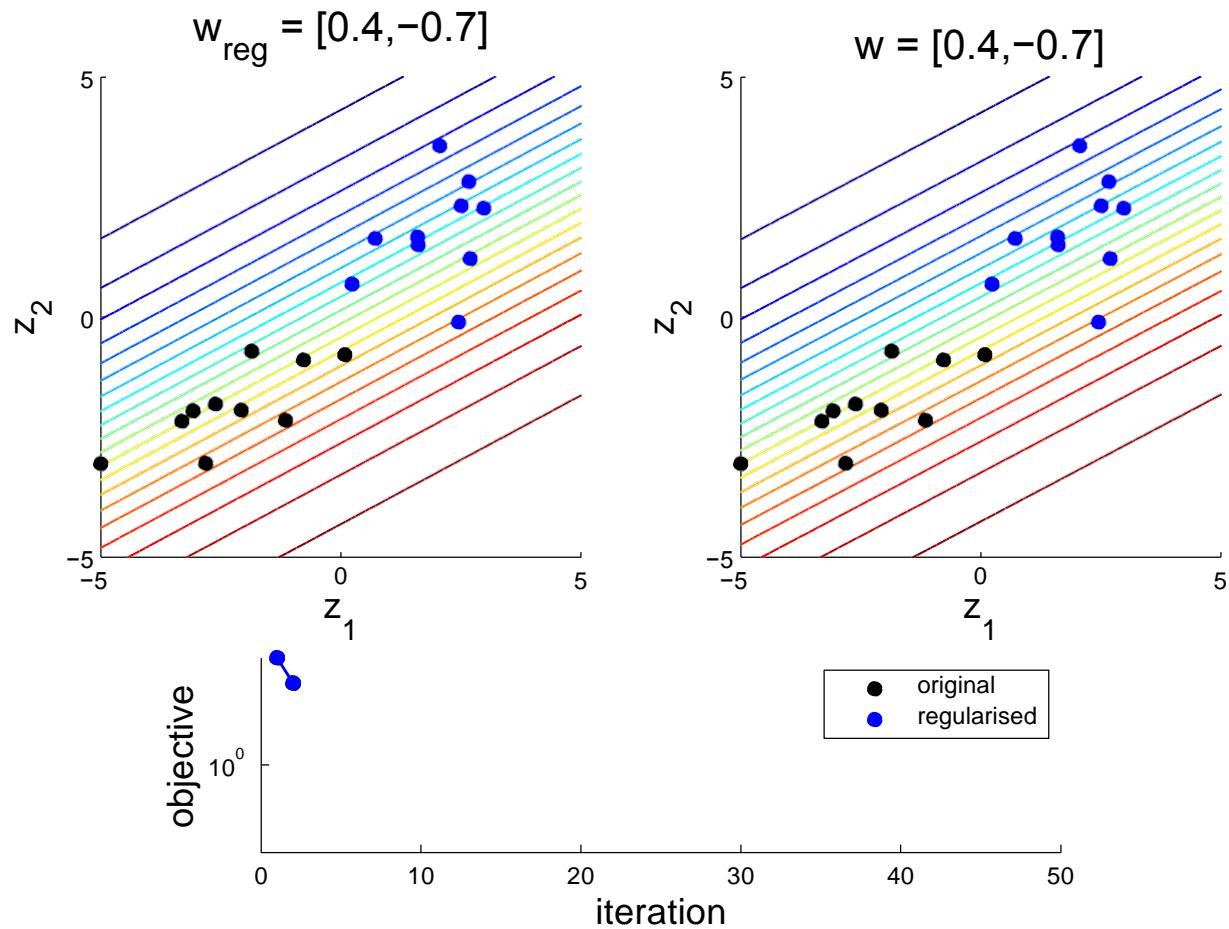
$$\mathbf{w}^* = \arg \min_{\mathbf{w}} M(\mathbf{w}) = \arg \min_{\mathbf{w}} [G(\mathbf{w}) + \alpha E(\mathbf{w})]$$

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

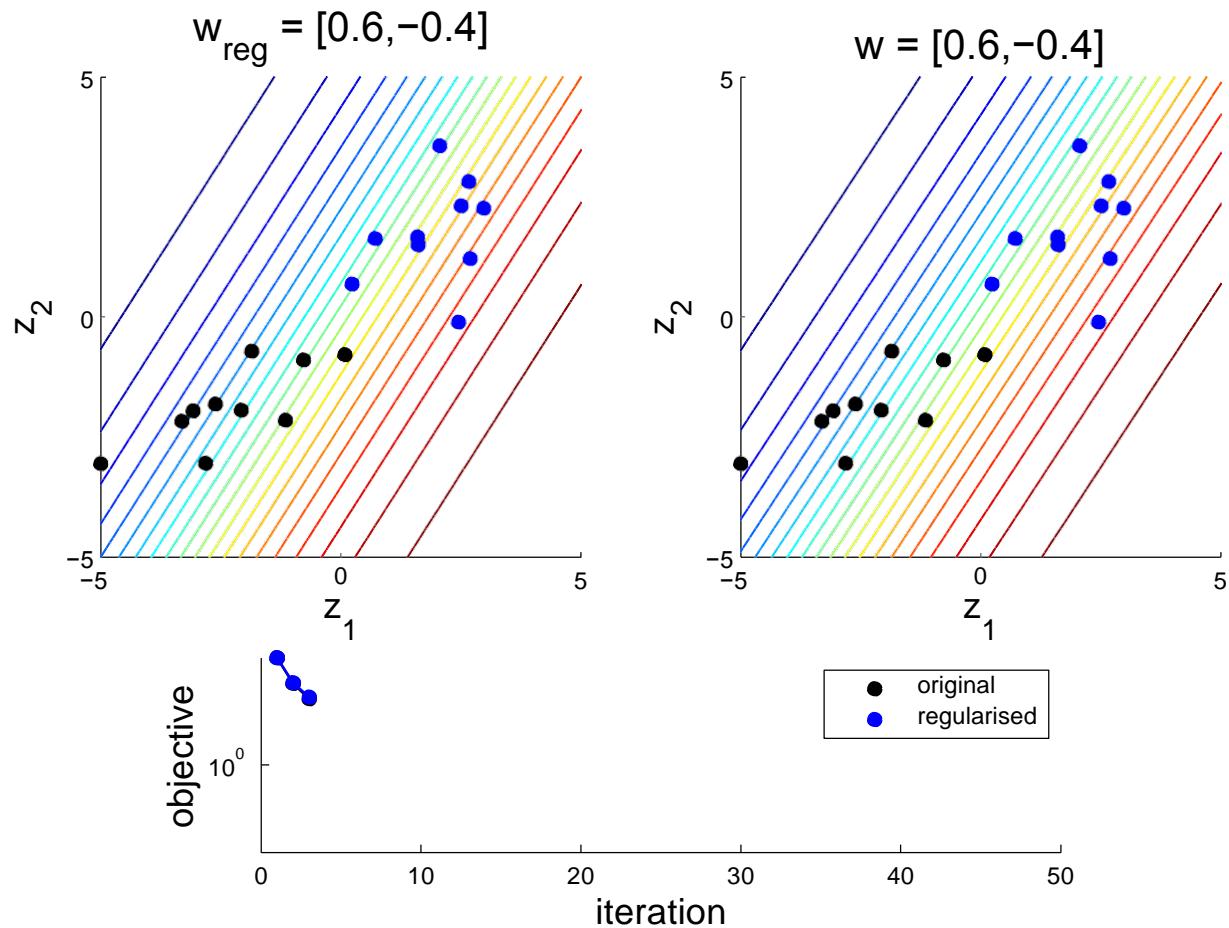
Training a Single Neuron (cont'd)



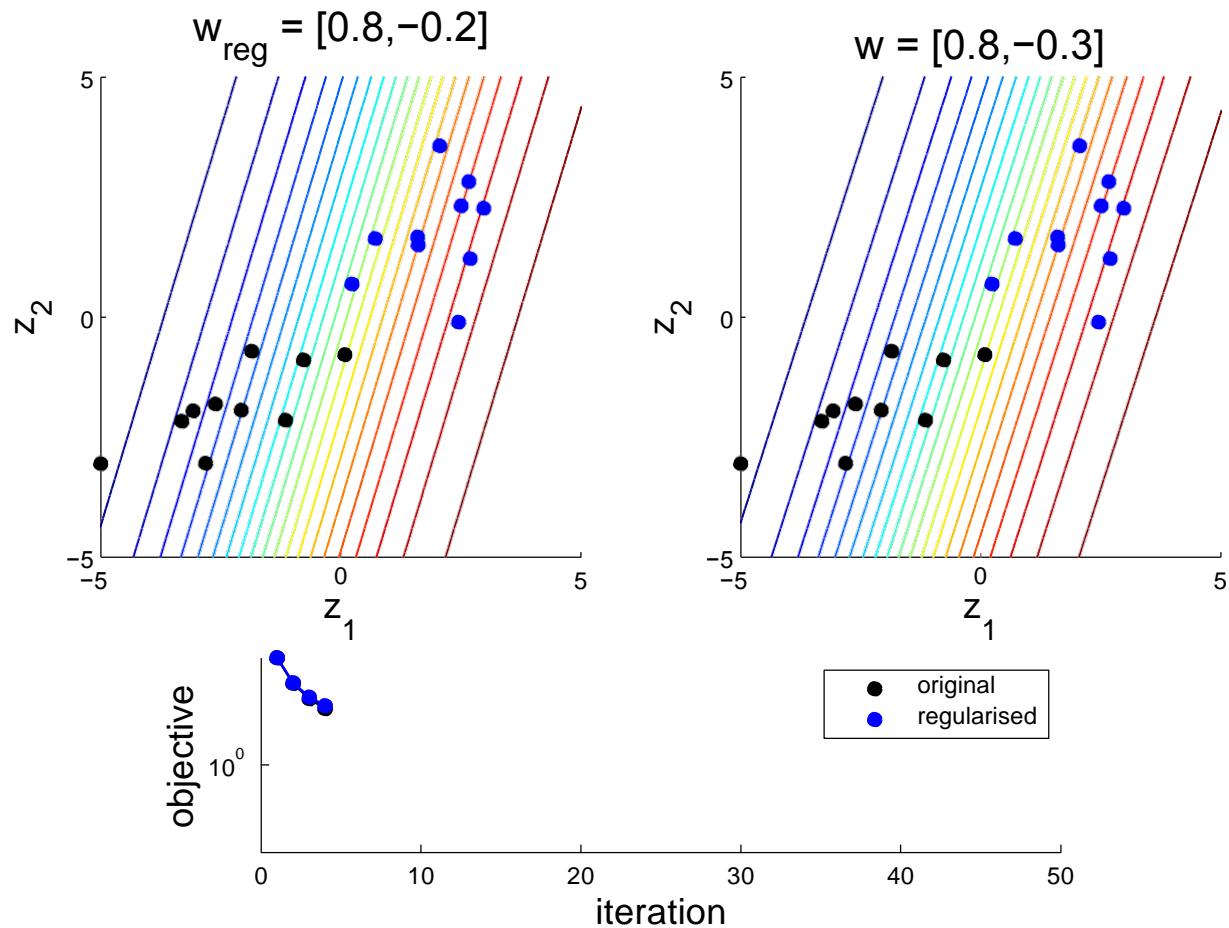
Training a Single Neuron (cont'd)



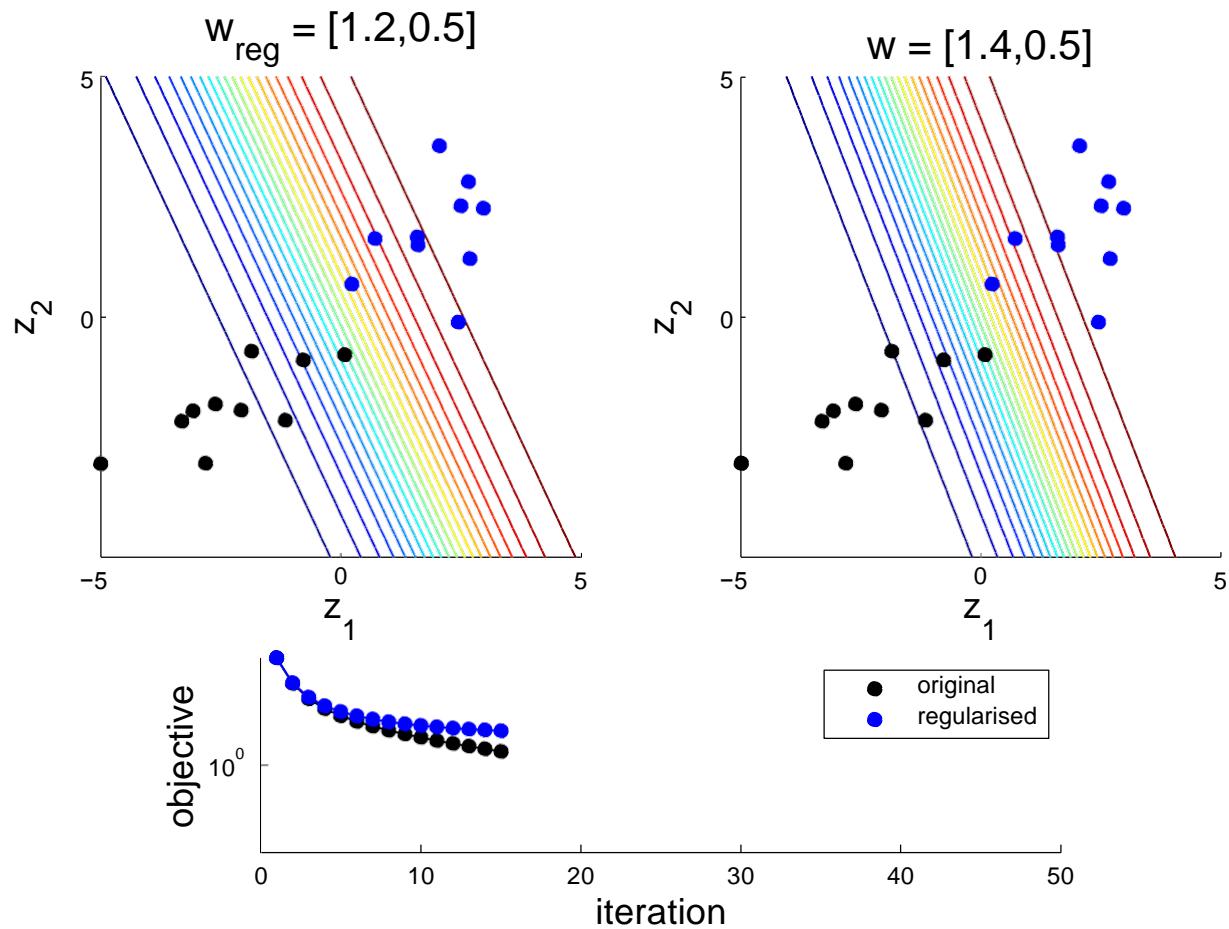
Training a Single Neuron (cont'd)



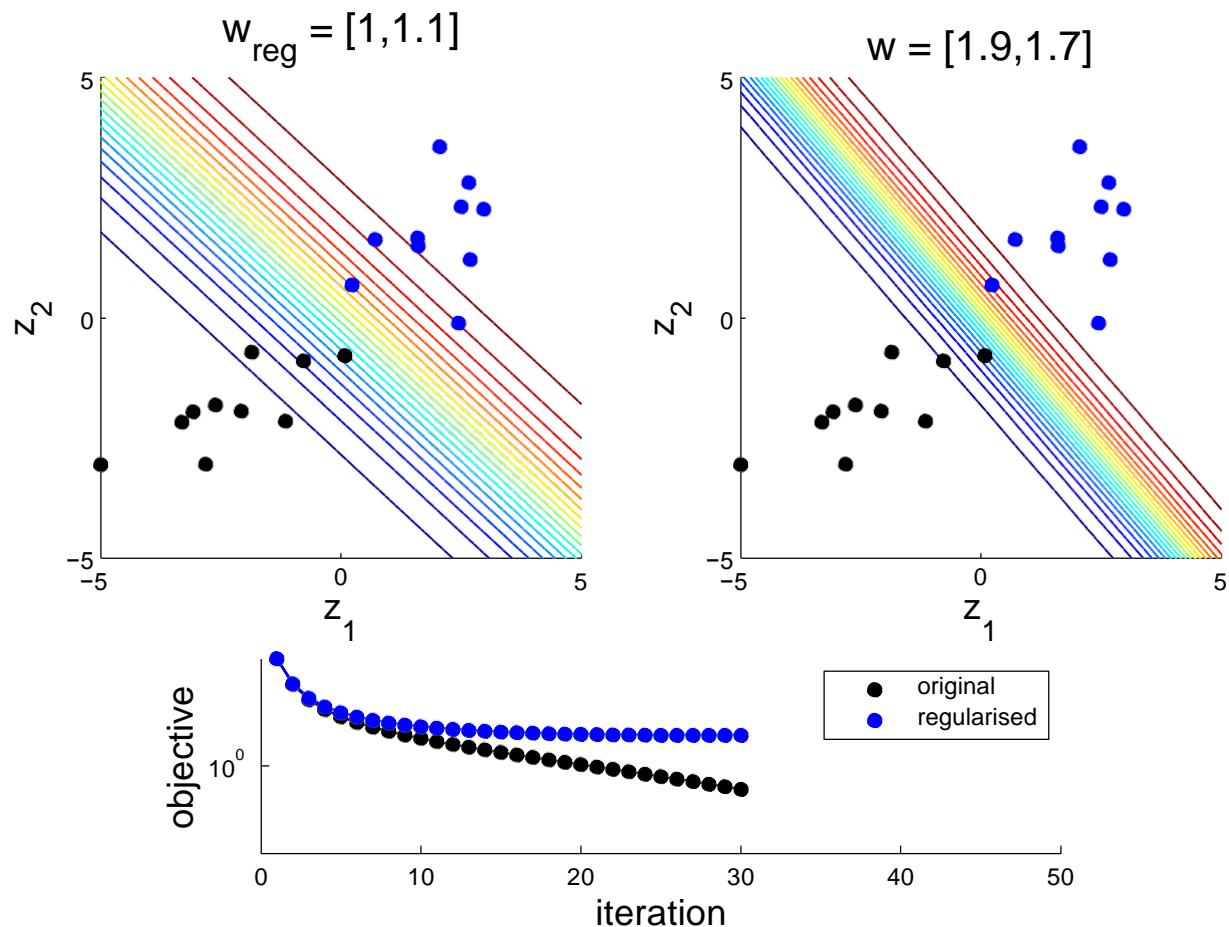
Training a Single Neuron (cont'd)



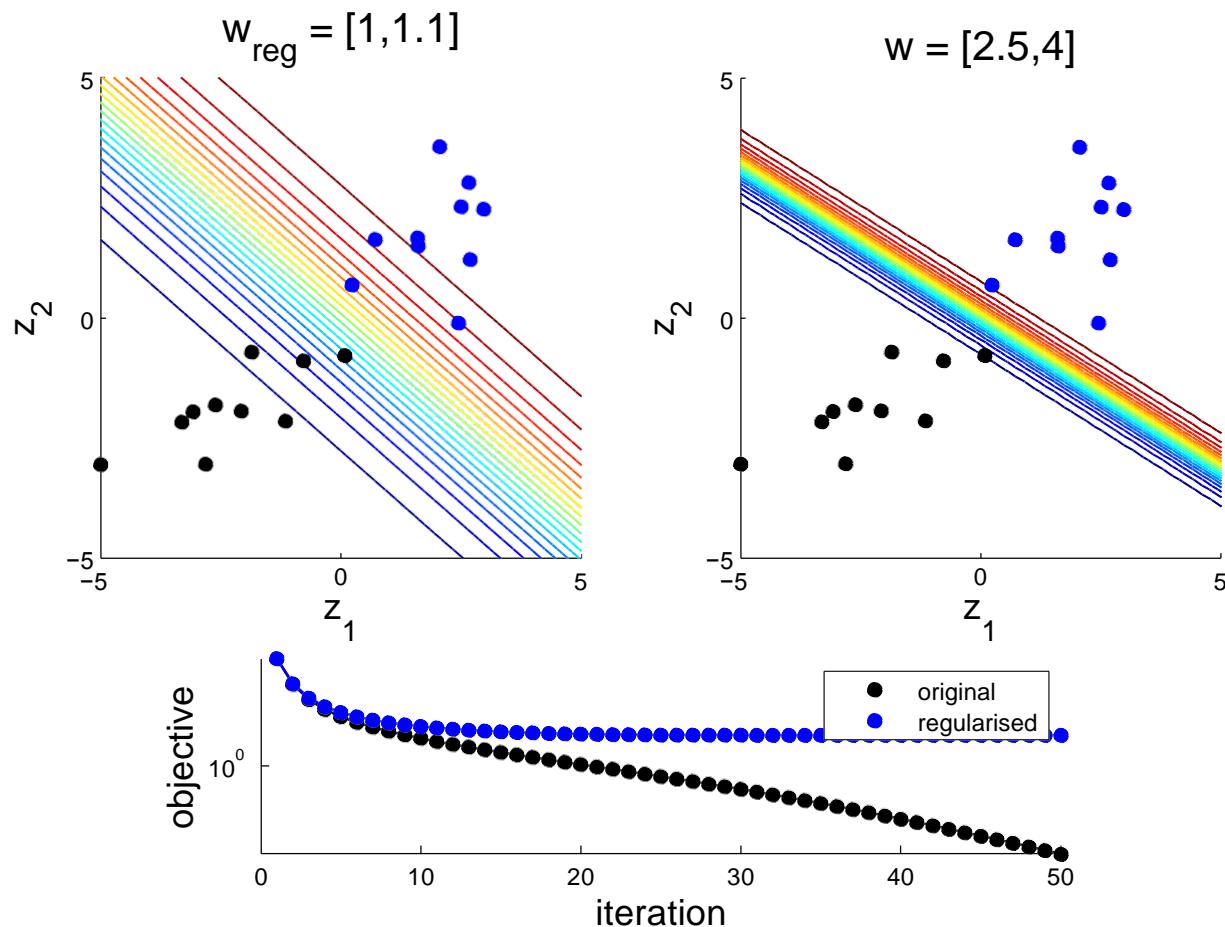
Training a Single Neuron (cont'd)



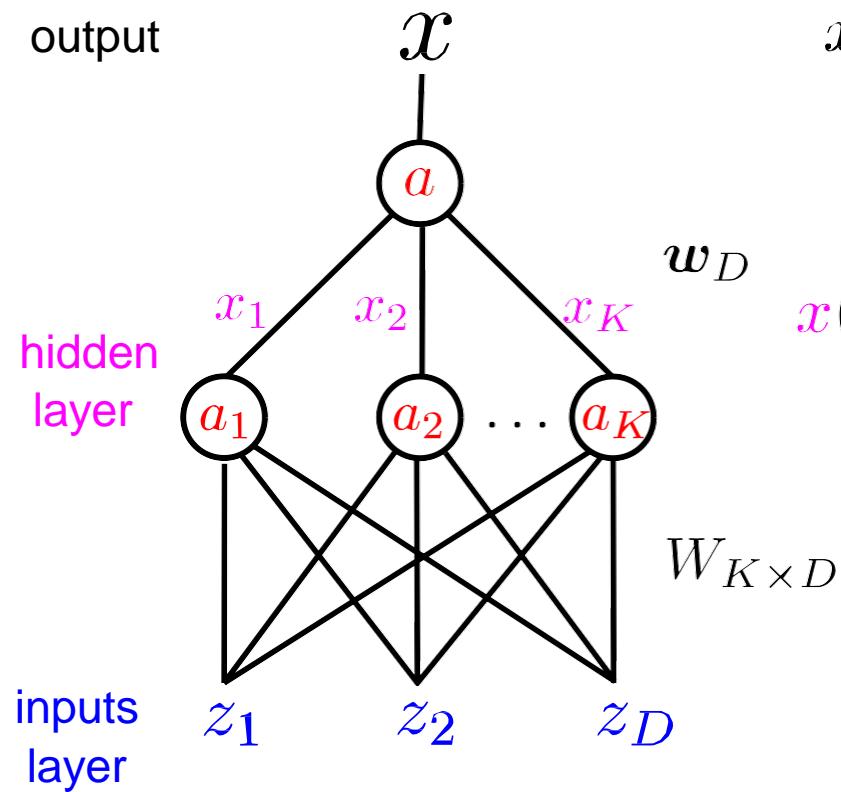
Training a Single Neuron (cont'd)



Training a Single Neuron (cont'd)



Single Hidden Layer Neural Networks

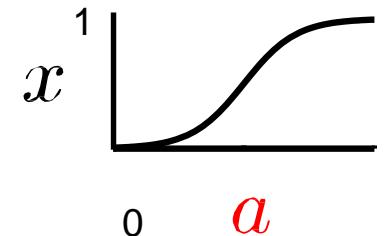


$$x(\textcolor{red}{a}) = \frac{1}{1 + \exp(-\textcolor{red}{a})}$$

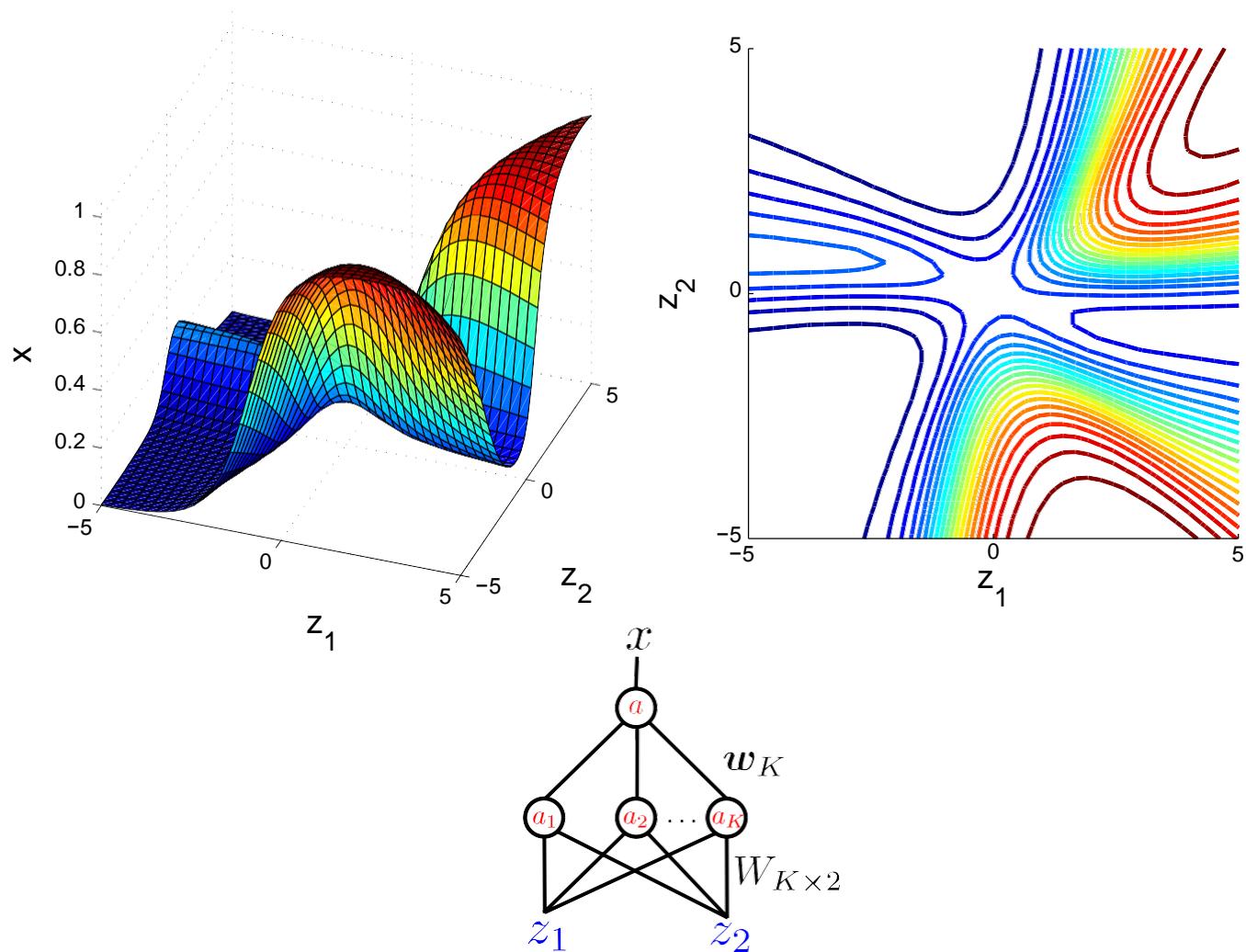
$$\textcolor{red}{a} = \sum_{k=1}^K w_k \textcolor{violet}{x}_k$$

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

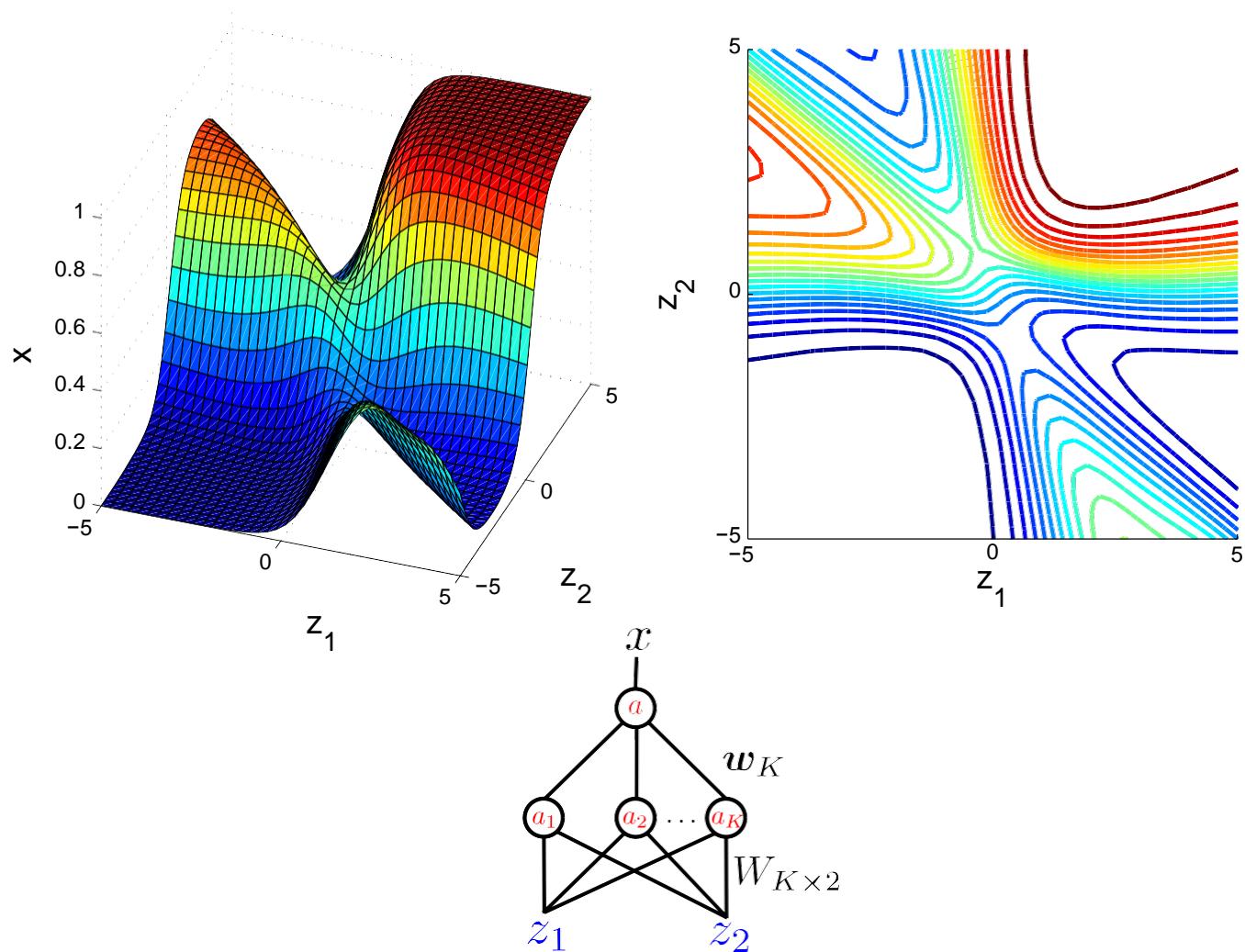
$$a_k = \sum_{d=1}^D W_{k,d} \textcolor{blue}{z}_d$$



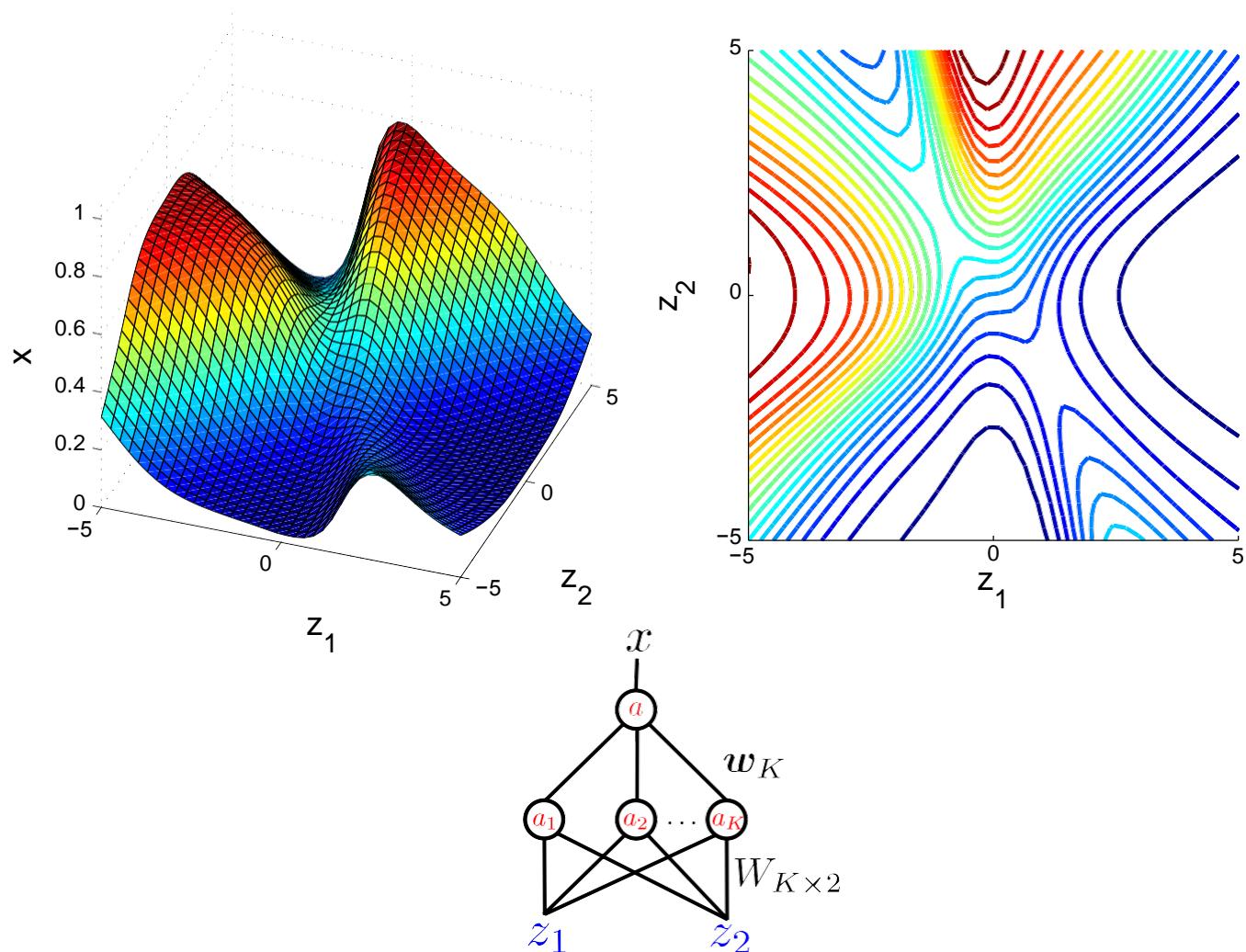
Sampling Random Neural Network Classifiers



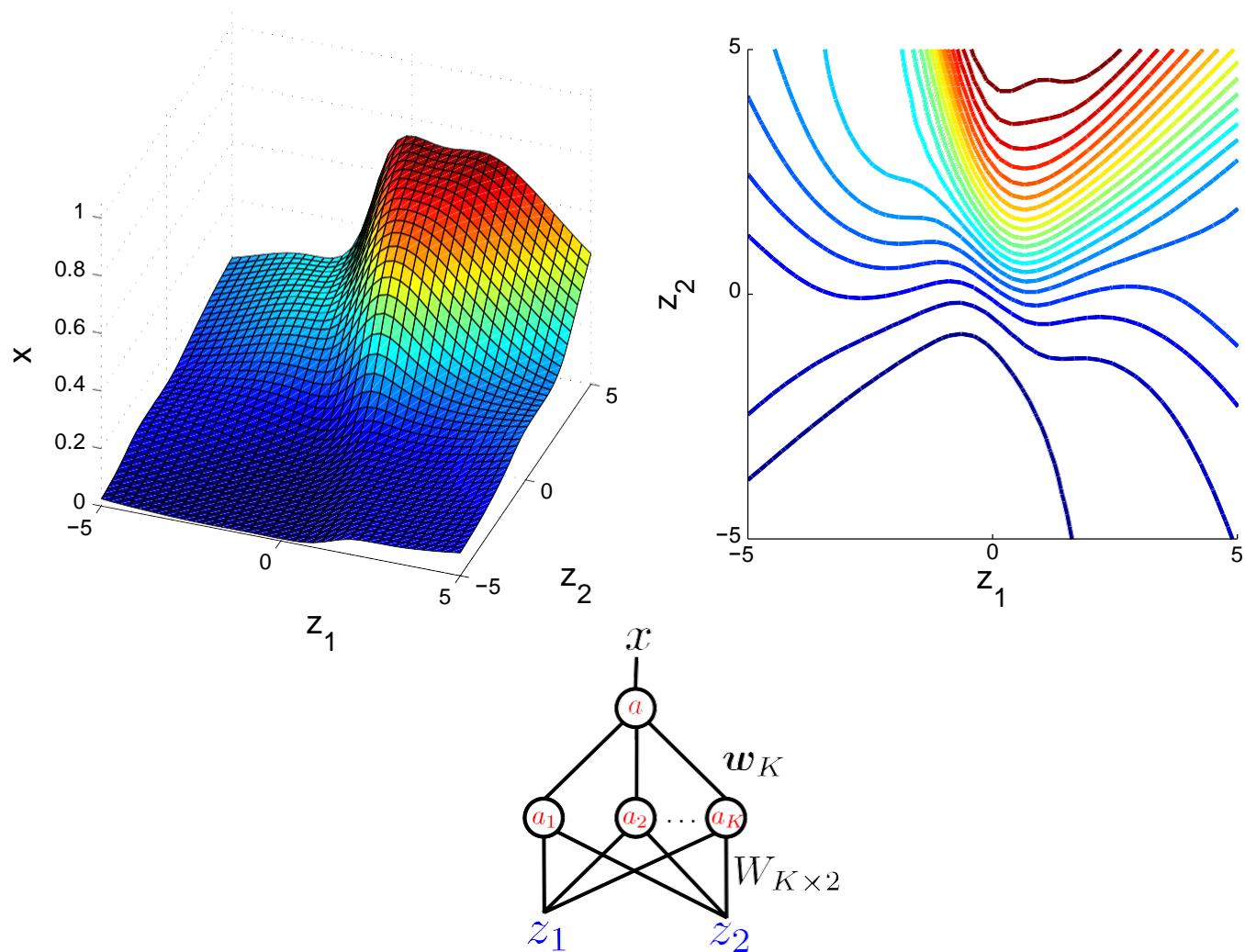
Sampling Random Neural Network Classifiers



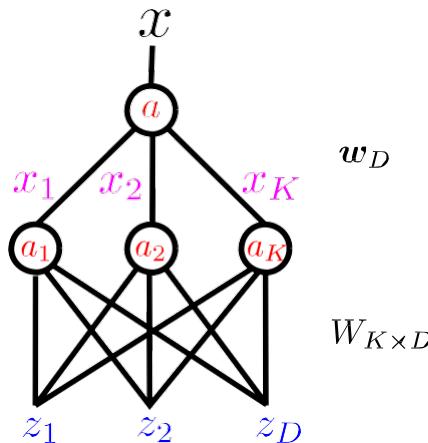
Sampling Random Neural Network Classifiers



Sampling Random Neural Network Classifiers



Training a Neural Network with a Single Hidden Layer



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

$$\textcolor{red}{a} = \sum_{k=1}^K w_k \textcolor{violet}{x}_k$$

$$\textcolor{violet}{x}(\textcolor{red}{a}_k) = \frac{1}{1 + \exp(-\textcolor{red}{a}_k)}$$

$$a_k = \sum_{d=1}^D W_{k,d} \textcolor{blue}{z}_d$$

objective function:

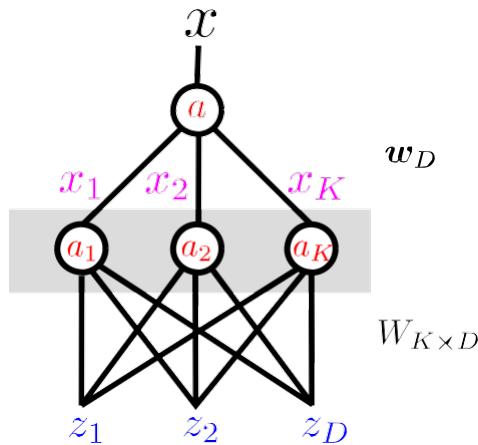
$$G(W, \mathbf{w}) = - \sum_n [t^{(n)} \log x^{(n)} + (1 - t^{(n)}) \log (1 - x^{(n)})] \quad \text{likelihood same as before}$$

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

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

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**.



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

$$\mathbf{a} = \sum_{k=1}^K w_k \mathbf{x}_k$$

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

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

objective function:

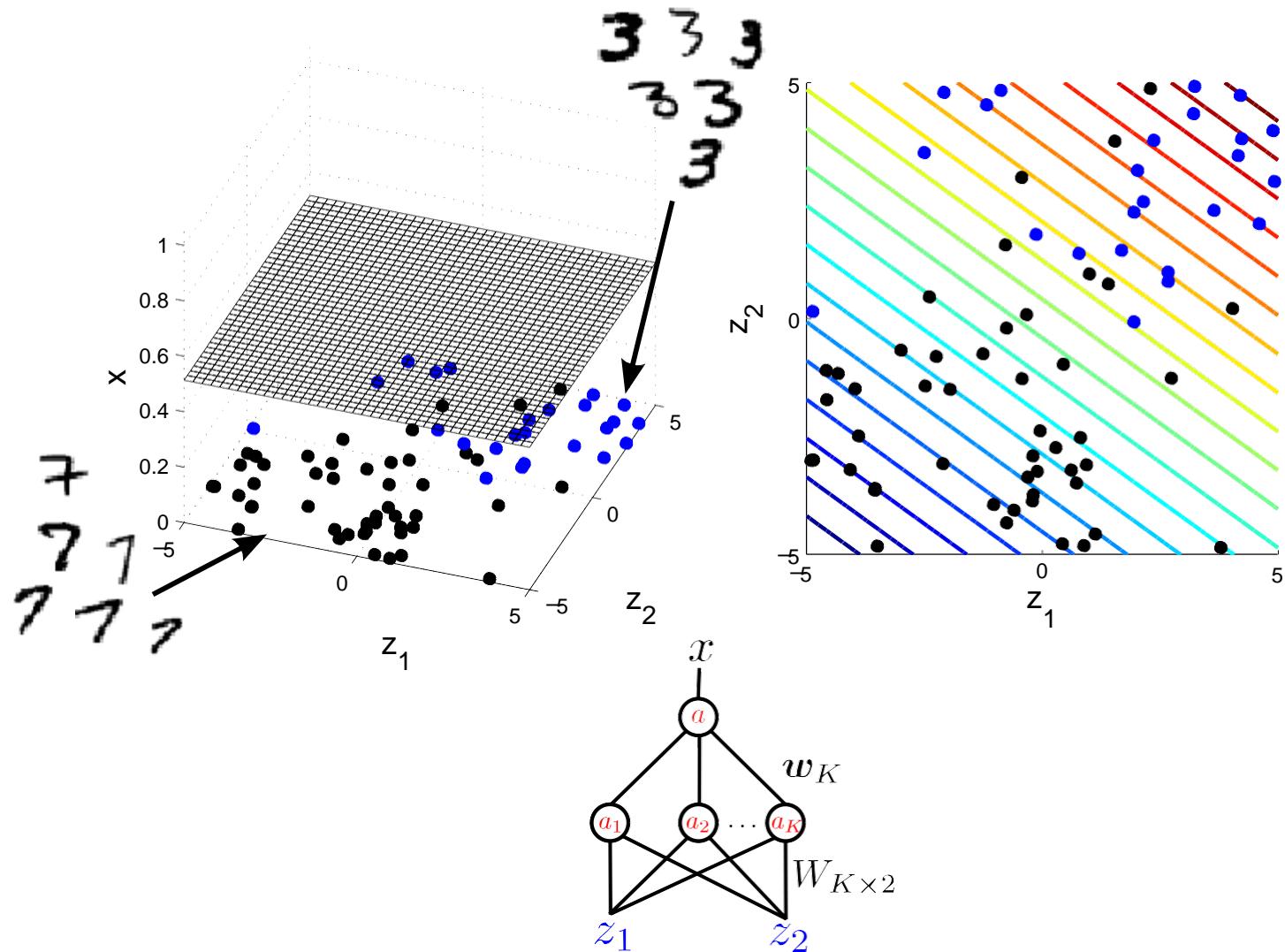
$$G(W, \mathbf{w}) = - \sum_n [t^{(n)} \log x^{(n)} + (1 - t^{(n)}) \log (1 - x^{(n)})] \quad \text{likelihood same as before}$$

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

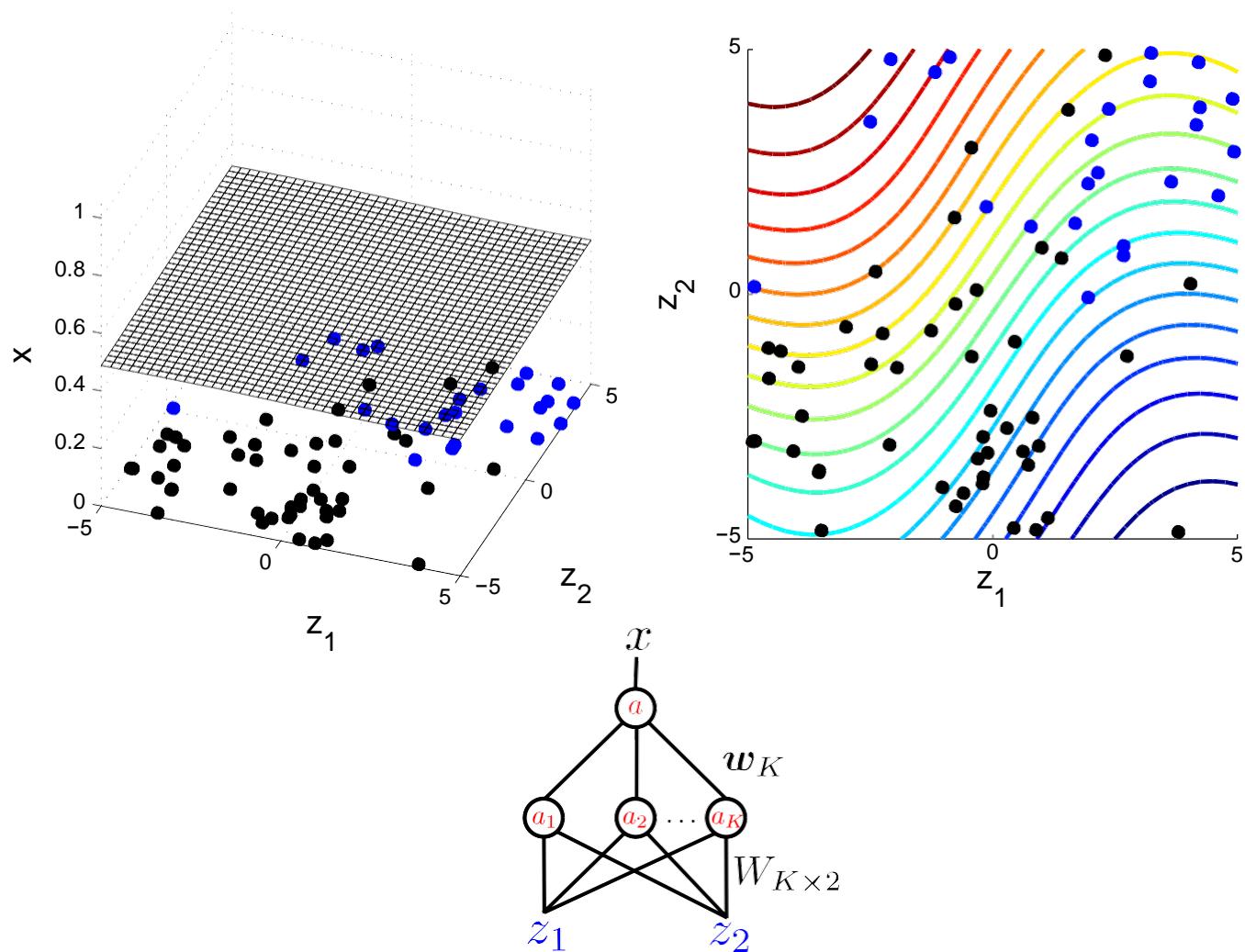
$$\{W, \mathbf{w}^*\} = \arg \min_{W, \mathbf{w}} M(W, \mathbf{w}) = \arg \min_{W, \mathbf{w}} [G(W, \mathbf{w}) + \alpha E(W, \mathbf{w})]$$

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

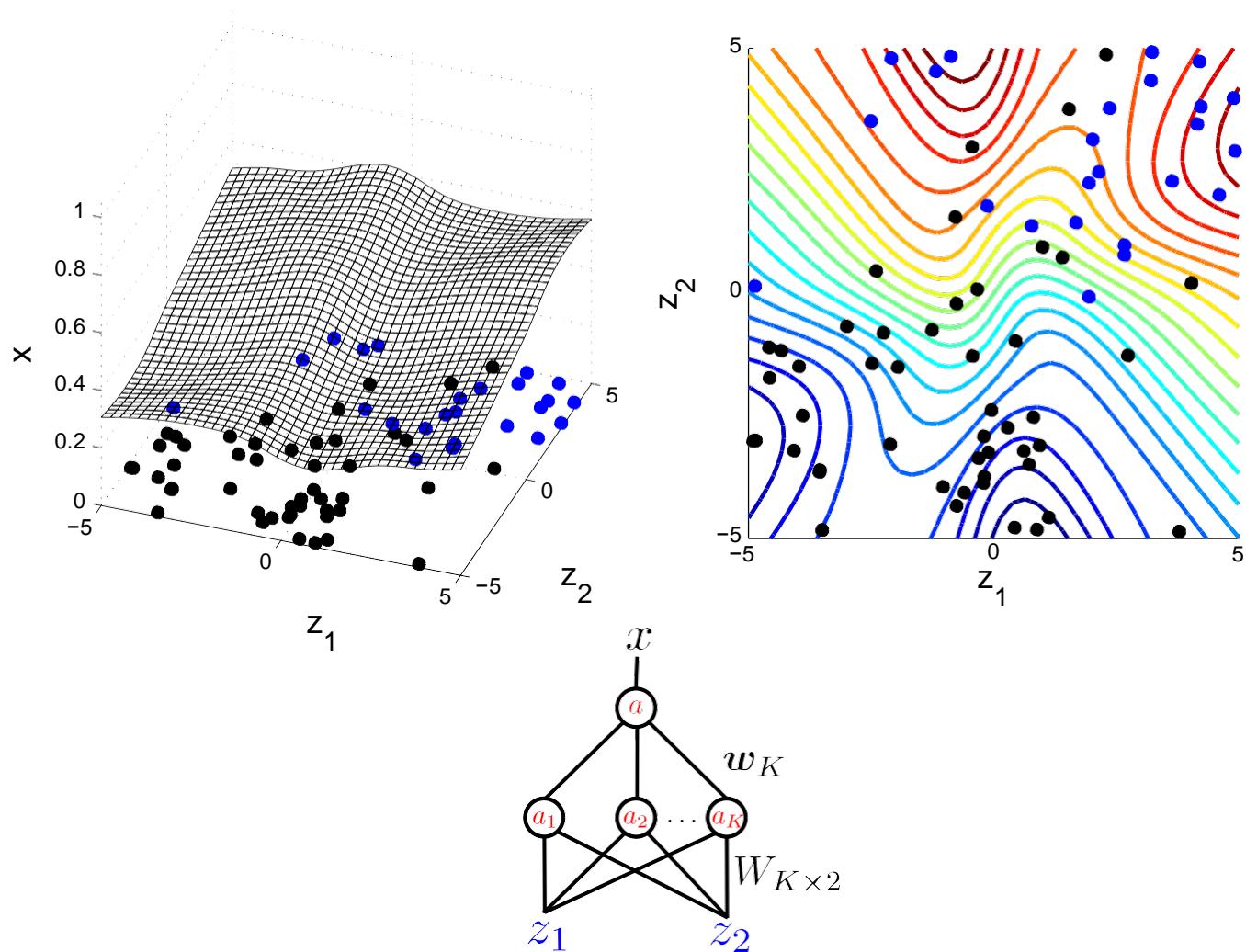
Training a Neural Network with a Single Hidden Layer



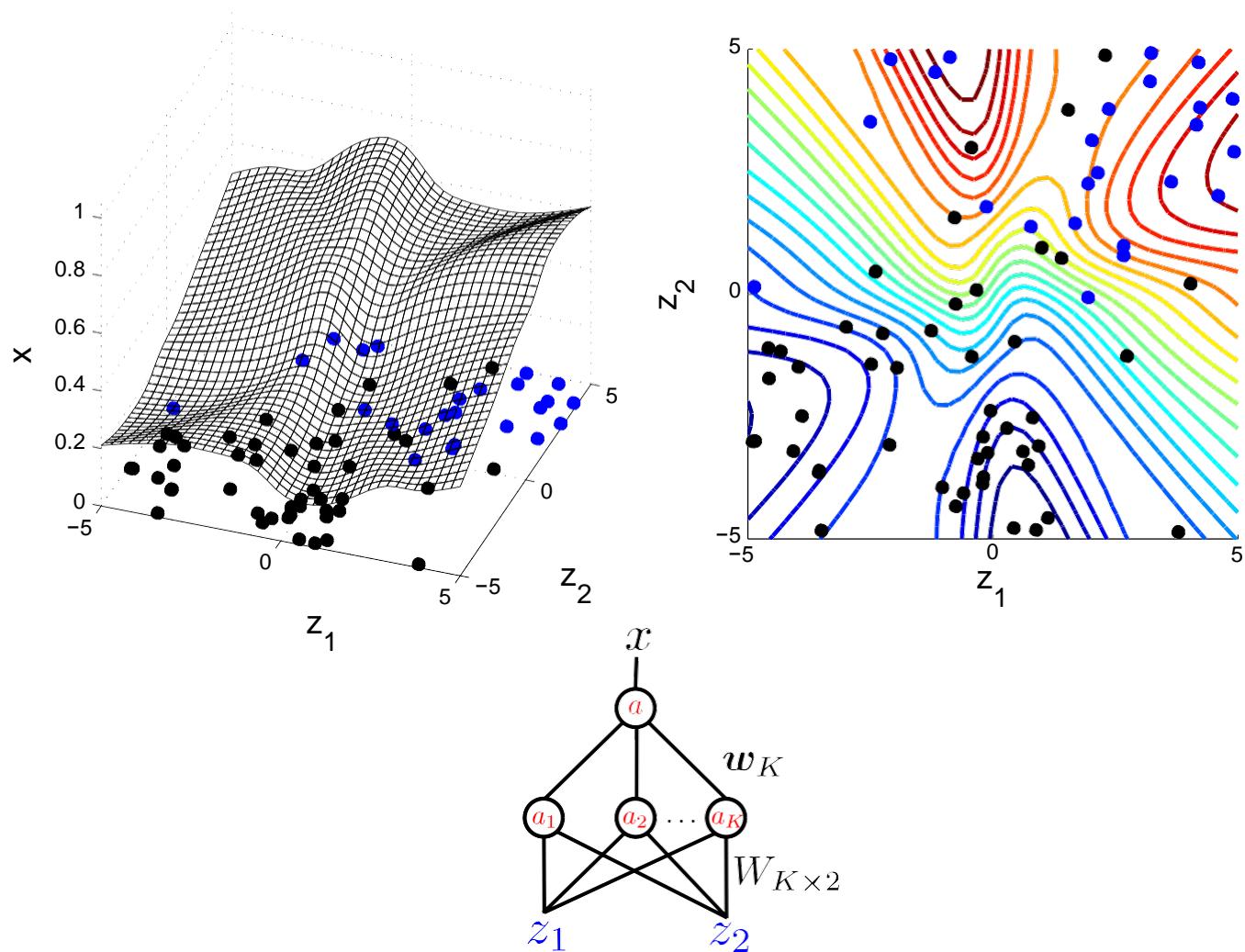
Training a Neural Network with a Single Hidden Layer



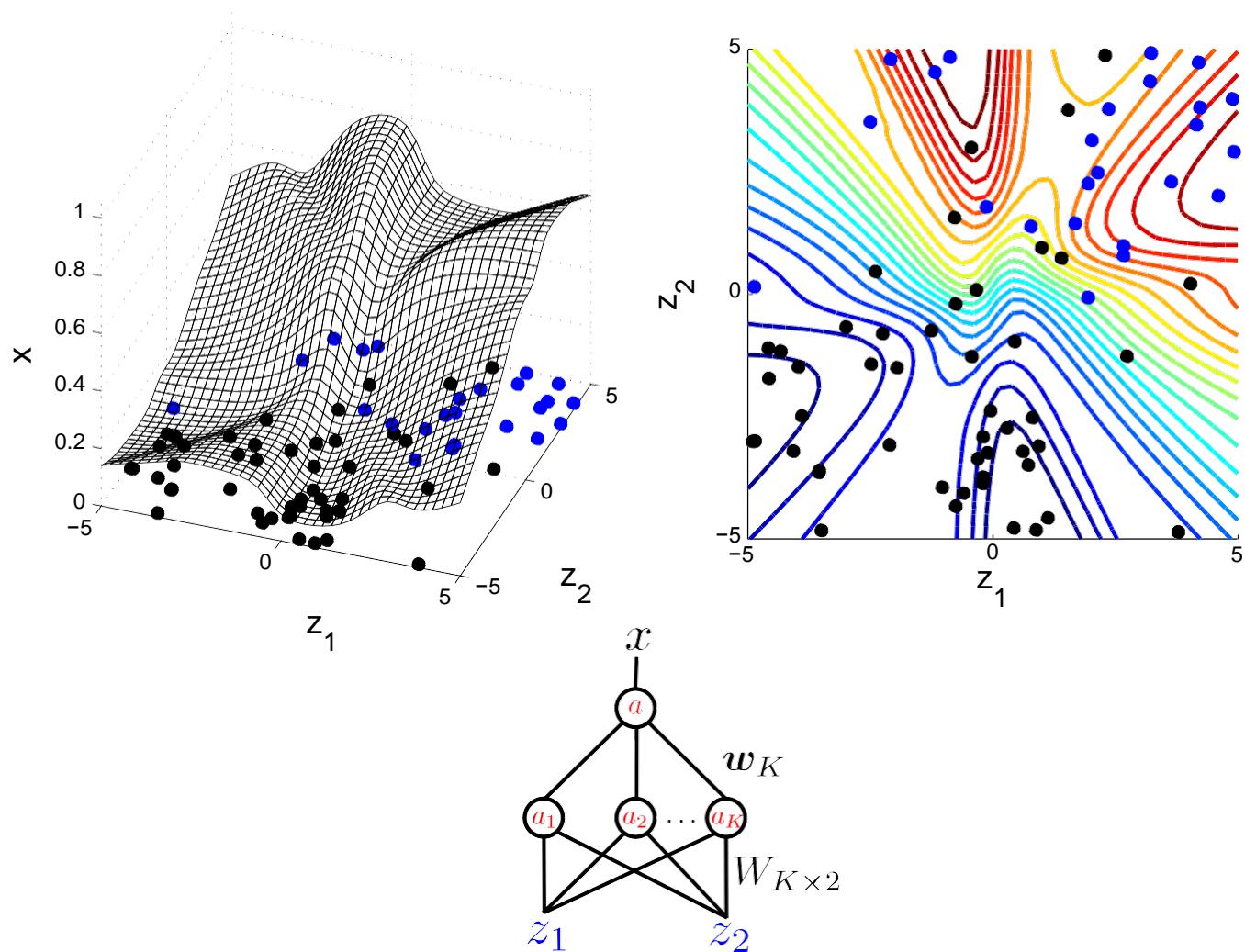
Training a Neural Network with a Single Hidden Layer



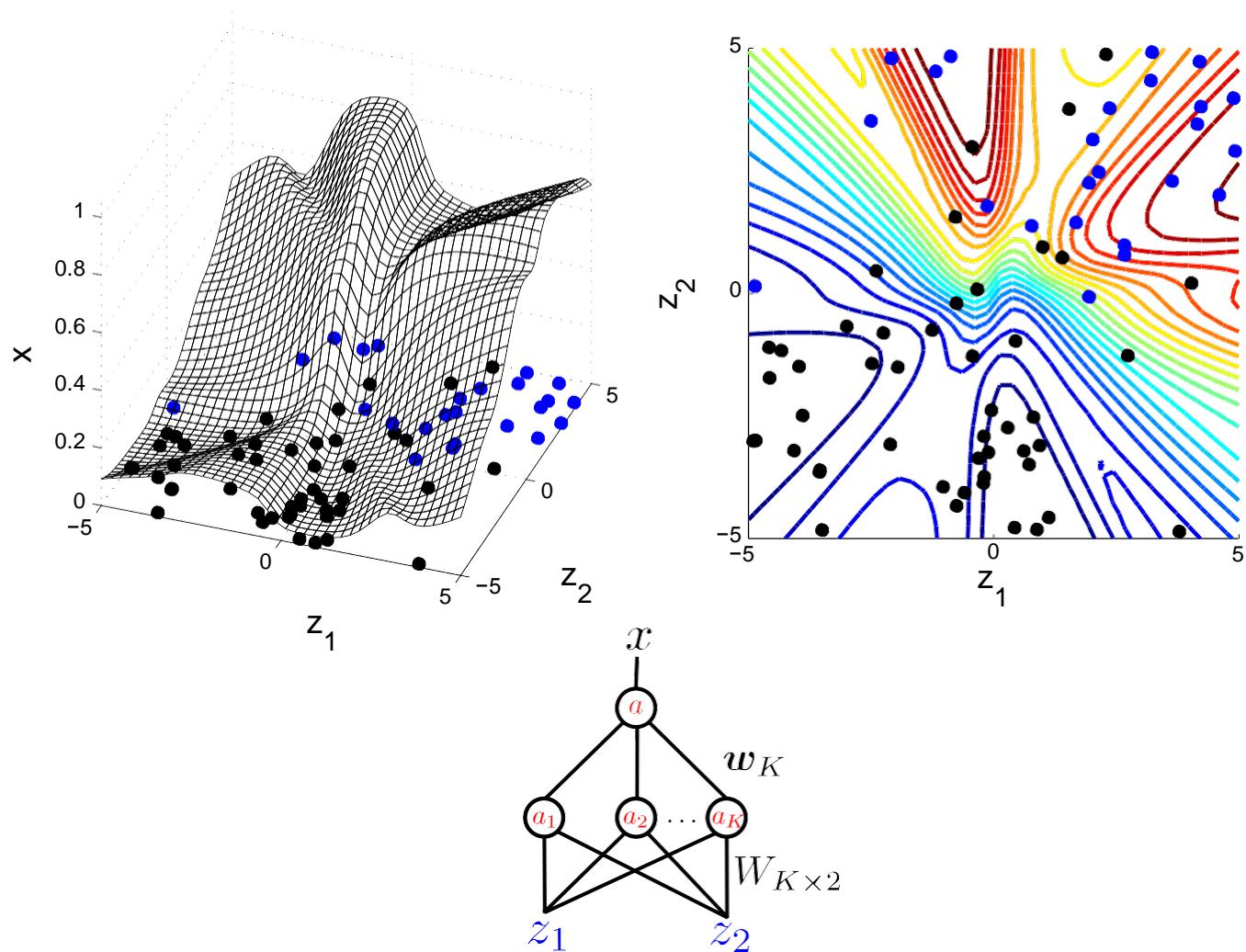
Training a Neural Network with a Single Hidden Layer



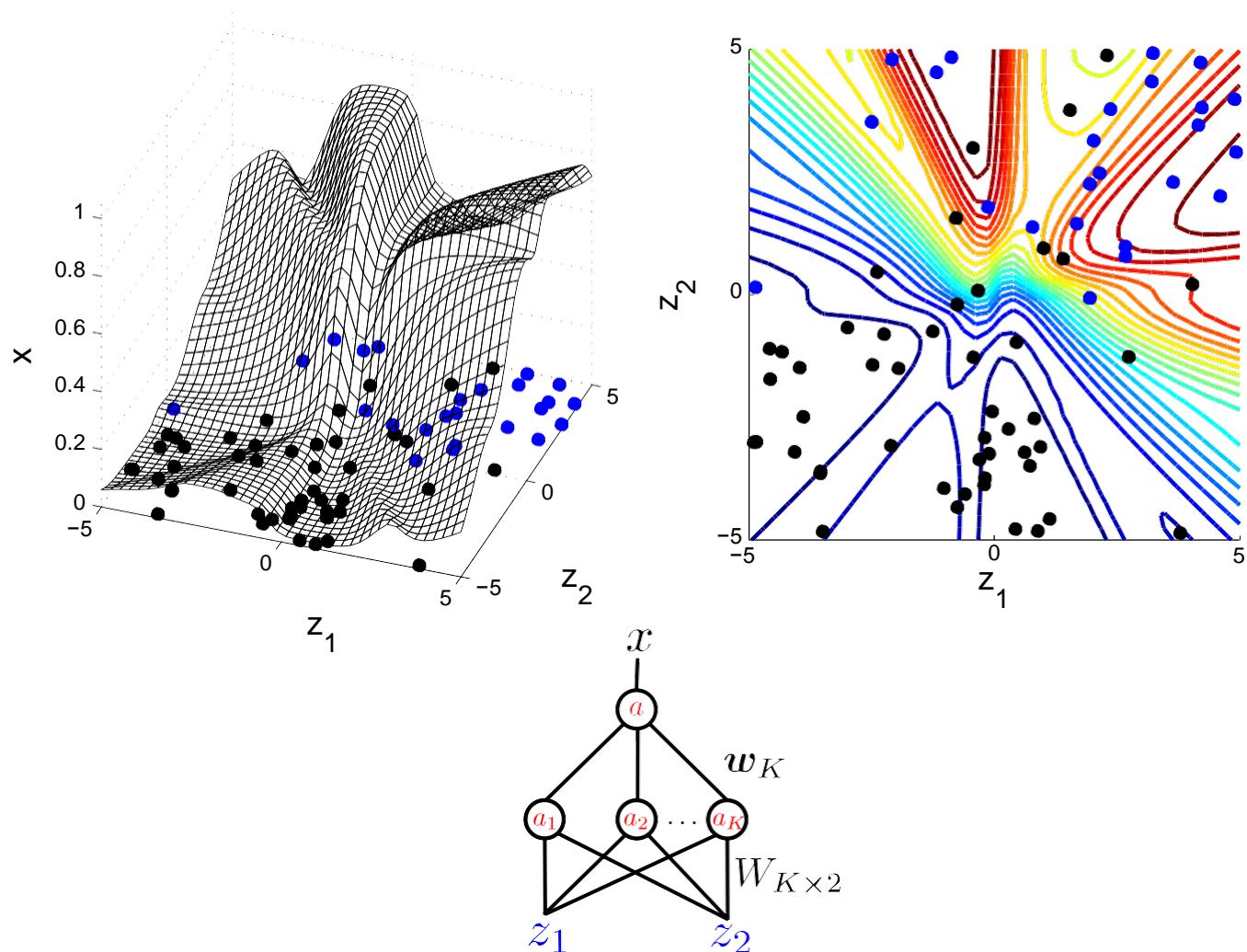
Training a Neural Network with a Single Hidden Layer



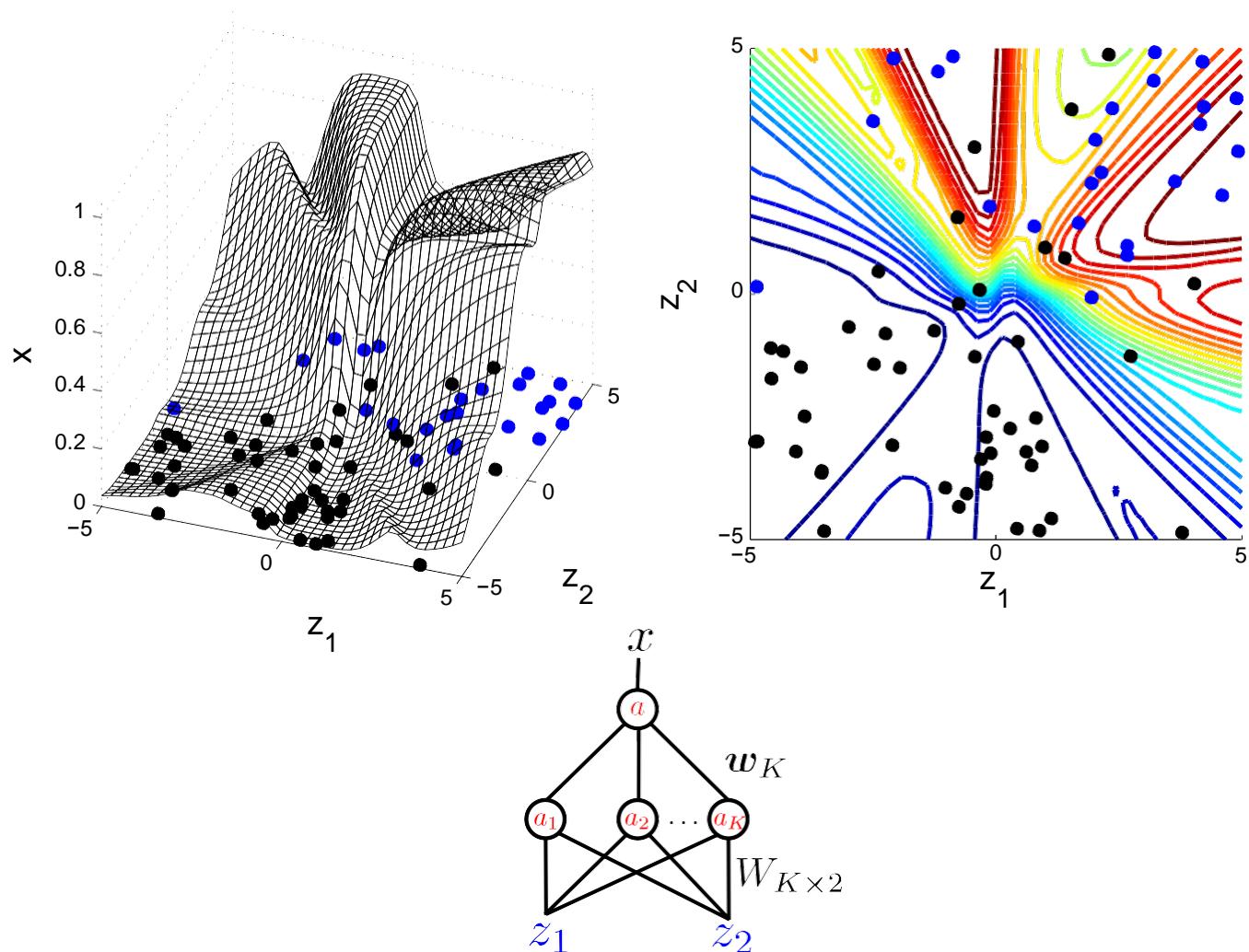
Training a Neural Network with a Single Hidden Layer



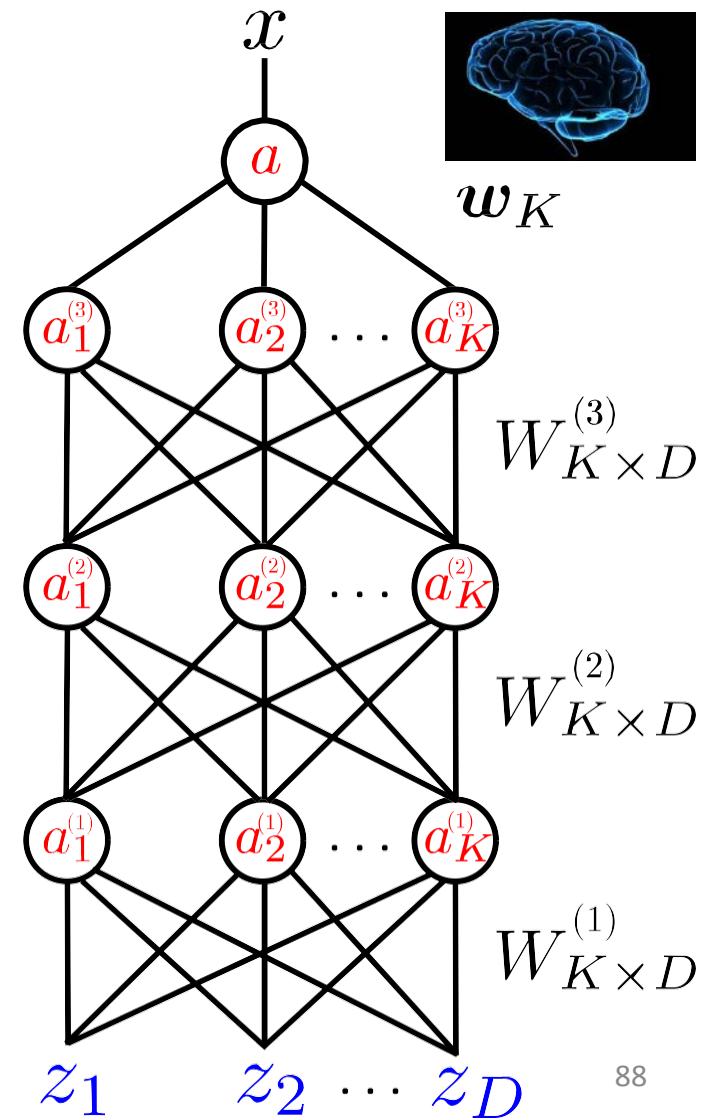
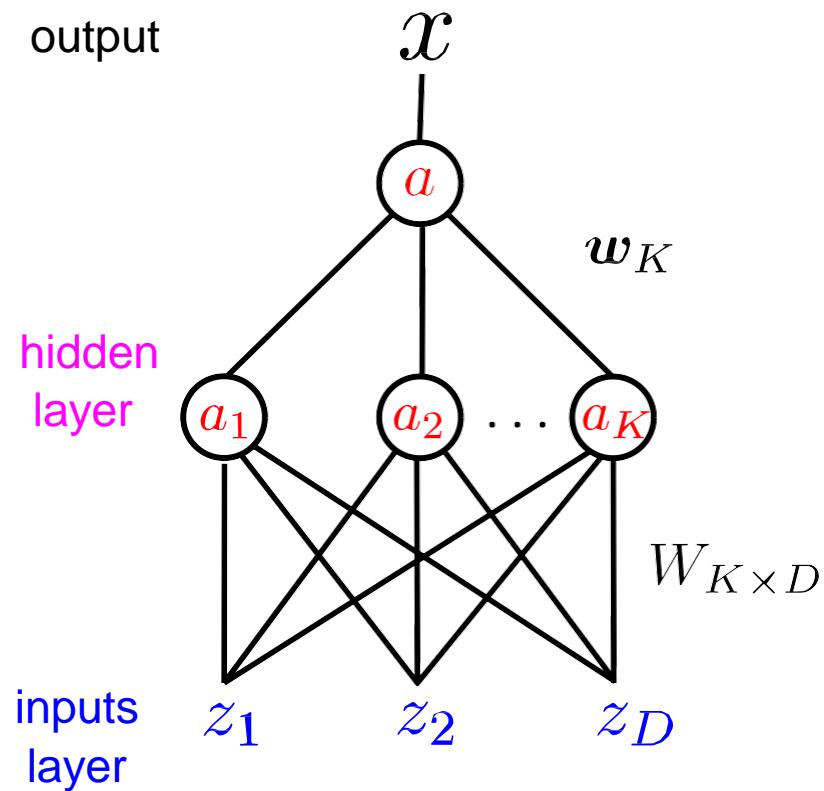
Training a Neural Network with a Single Hidden Layer



Training a Neural Network with a Single Hidden Layer

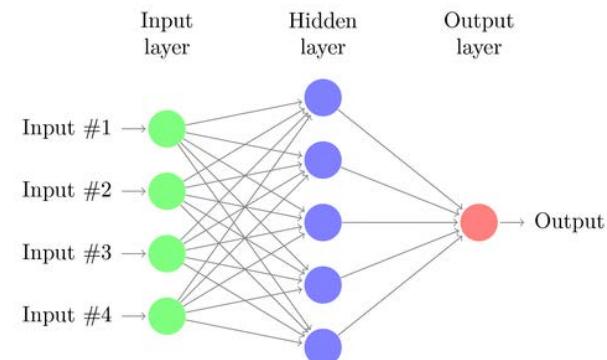
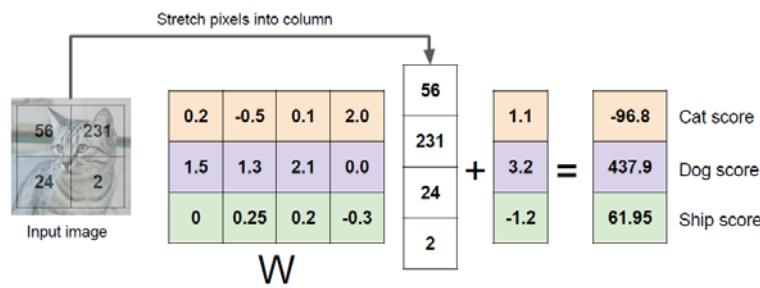
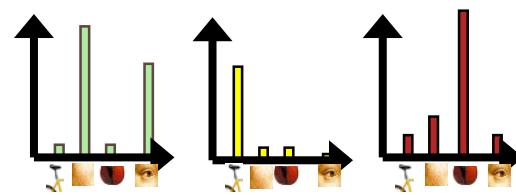


Hierarchical Models with Many Layers



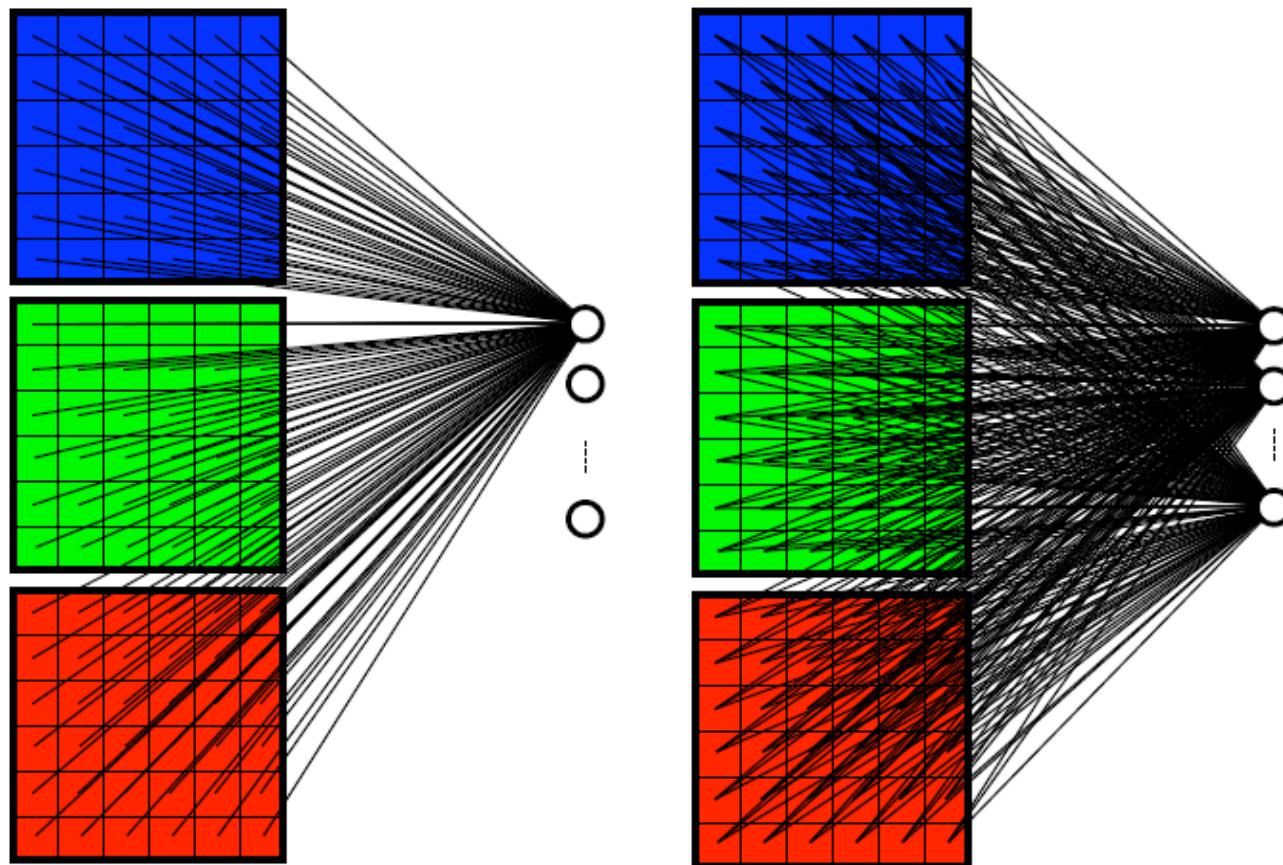
What's to Be Covered Today...

- General Framework for Visual Classification
 - Bag-of-Words Models as Image Features
- Intro to Neural Networks & CNN
 - Linear Classification
 - Neural Network for Machine Vision
 - Multi-Layer Perceptron
 - Convolutional Neural Networks



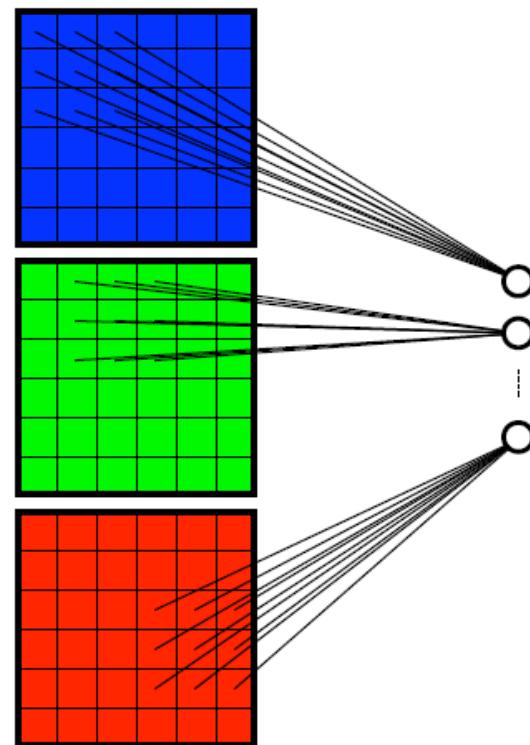
Convolutional Neural Networks

- How many weights for MLPs for images?



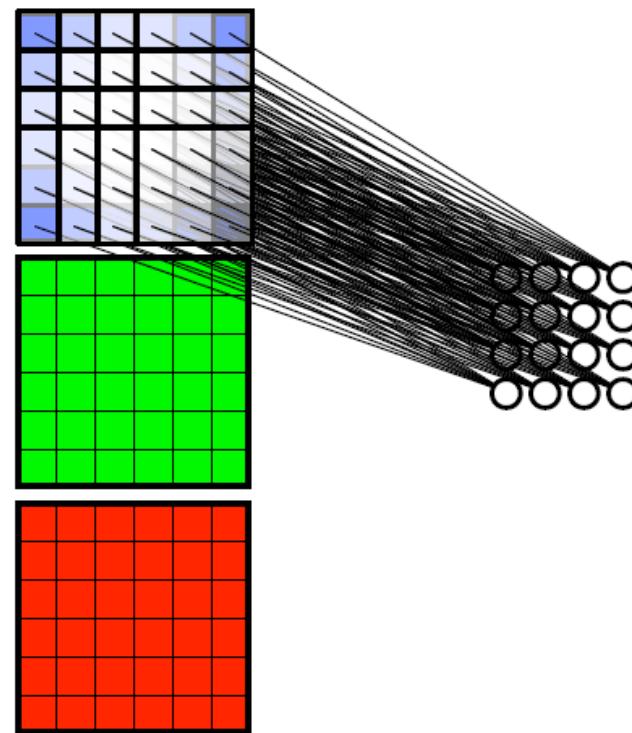
Convolutional Neural Networks

- Property I of CNN: Local Connectivity
 - Each neuron takes info only from a **neighborhood** of pixels.

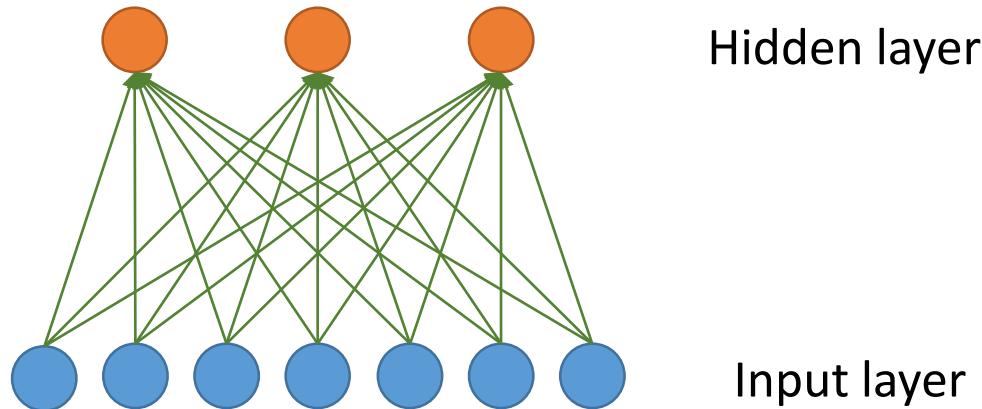


Convolutional Neural Networks

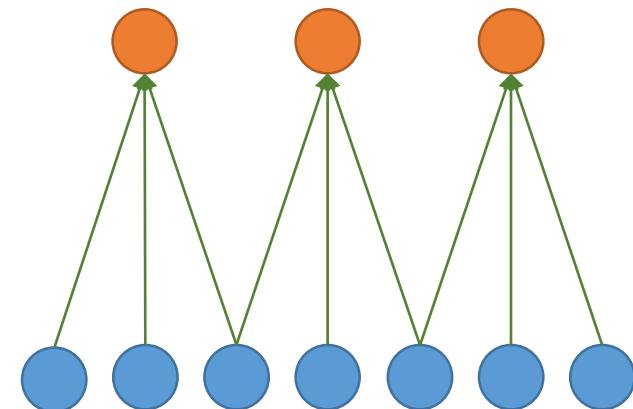
- Property II of CNN: Weight Sharing
 - Neurons connecting all neighborhoods have **identical** weights.



CNN: Local Connectivity



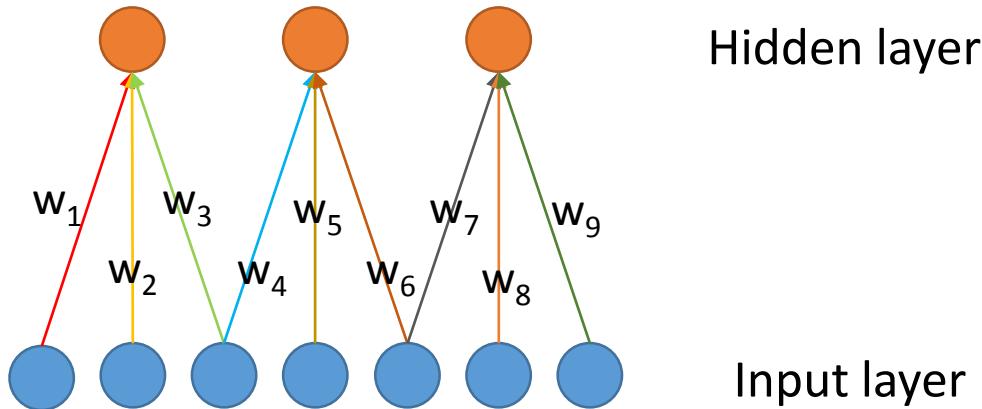
Global connectivity



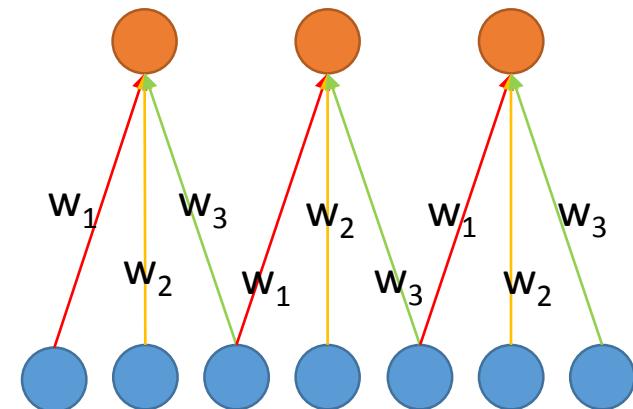
Local connectivity

- # input units (neurons): 7
- # hidden units: 3
- Number of parameters
 - Global connectivity:
 - Local connectivity:

CNN: Weight Sharing



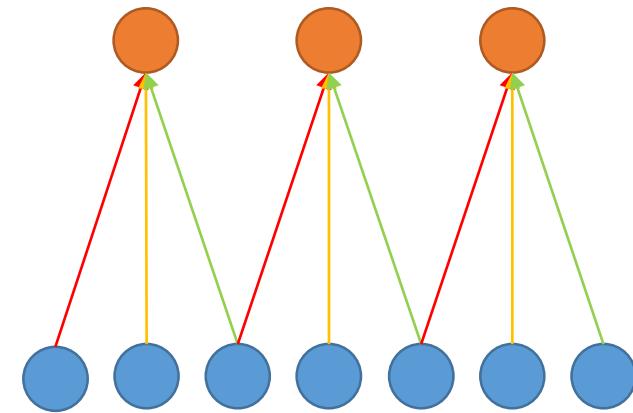
Without weight sharing



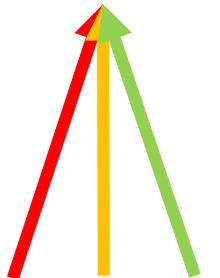
With weight sharing

- # input units (neurons): 7
- # hidden units: 3
- Number of parameters
 - Without weight sharing:
 - With weight sharing :

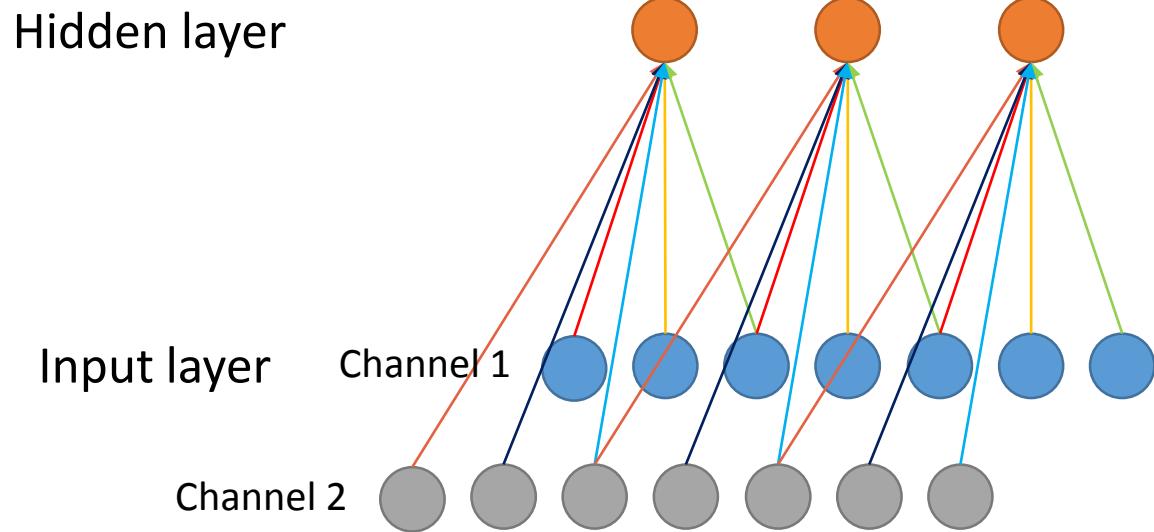
CNN with Multiple Input Channels



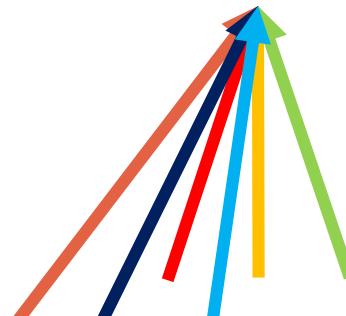
Single input channel



Filter weights

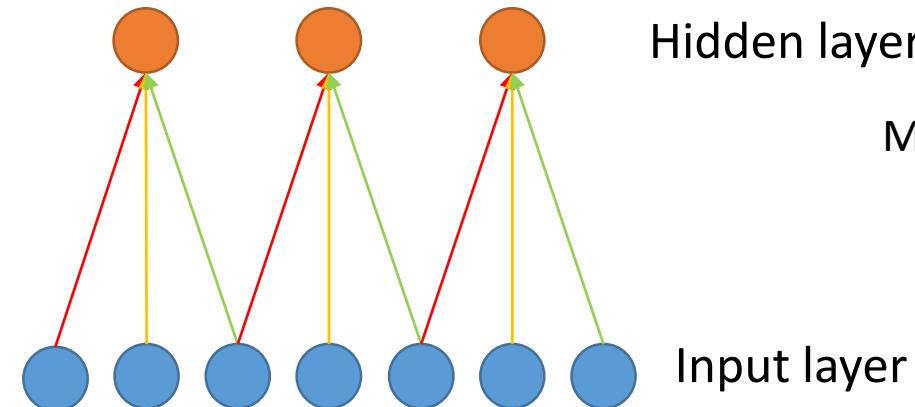


Multiple input channels

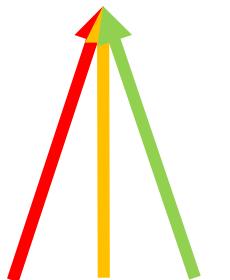


Filter weights

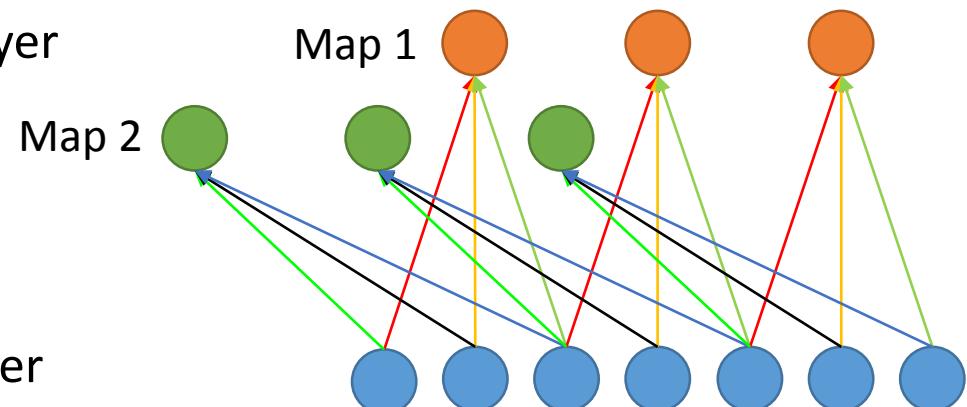
CNN with Multiple Output Maps



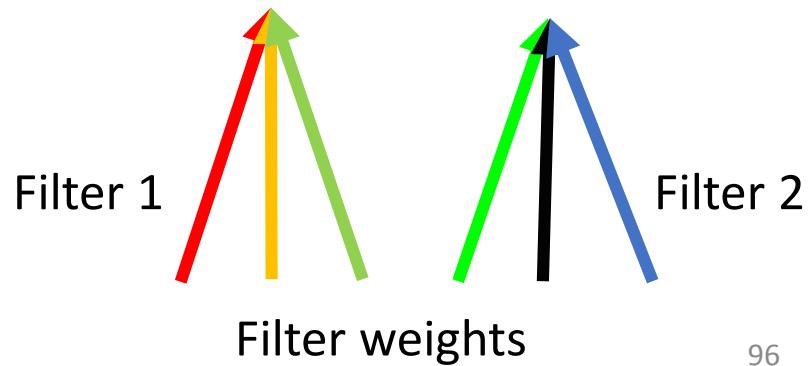
Single output map



Filter weights

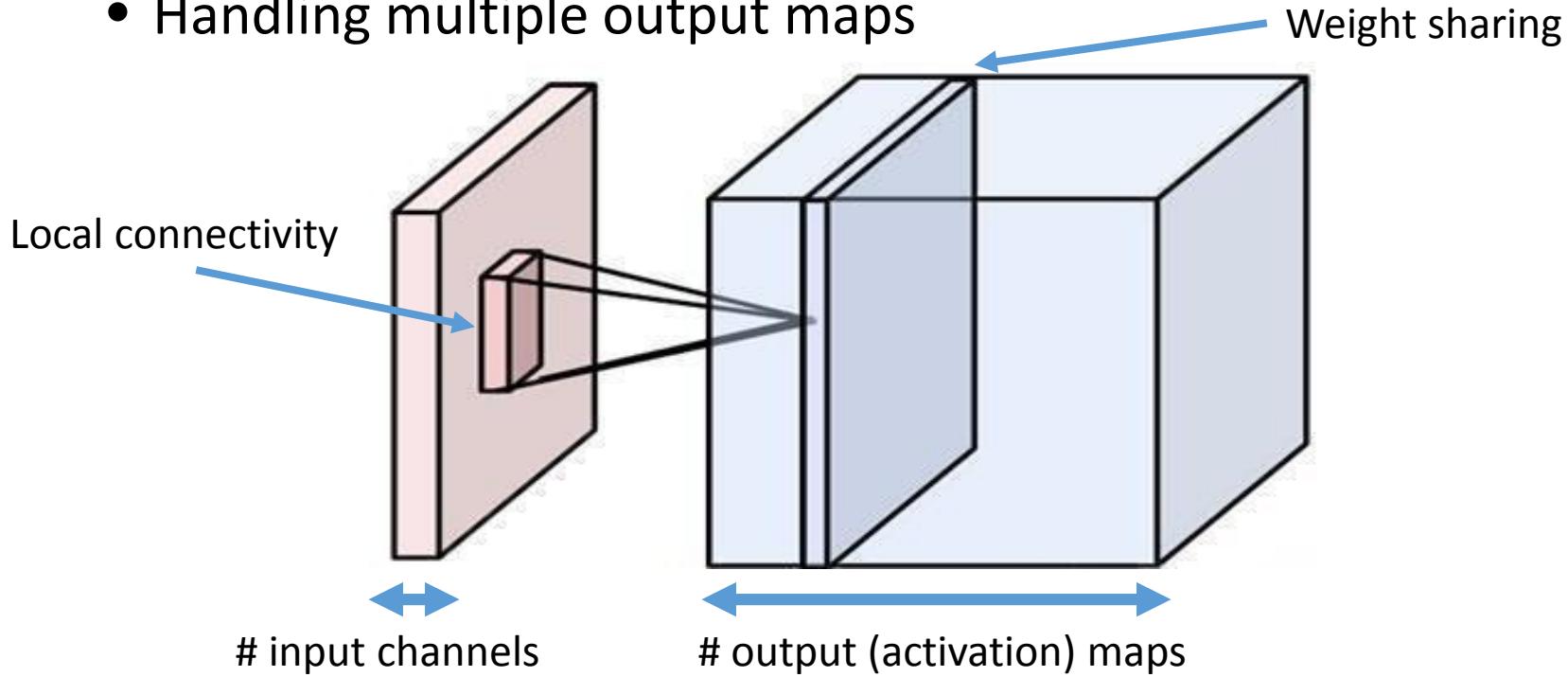


Multiple output maps

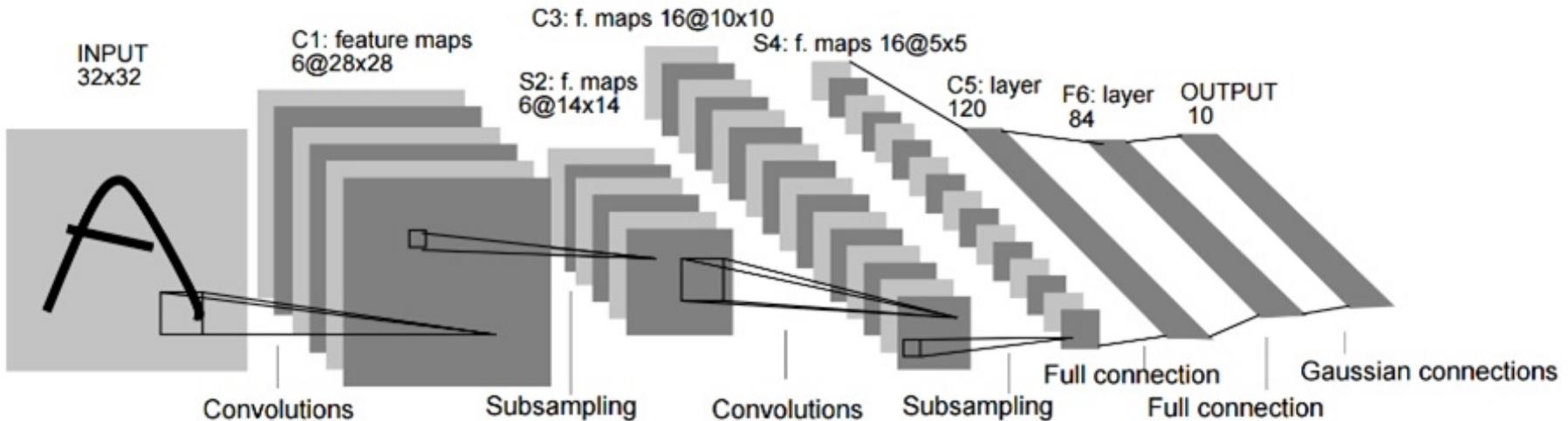


Putting them together

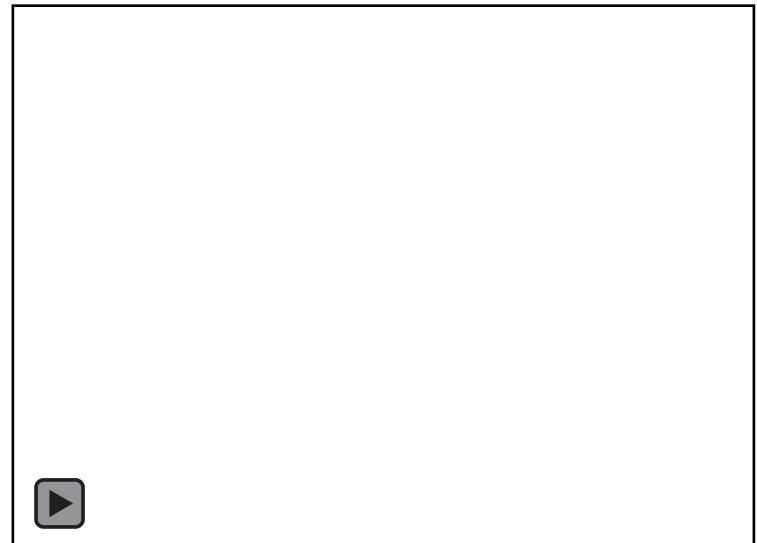
- Local connectivity
- Weight sharing
- Handling multiple input channels
- Handling multiple output maps



LeNet [LeCun et al. 1998]

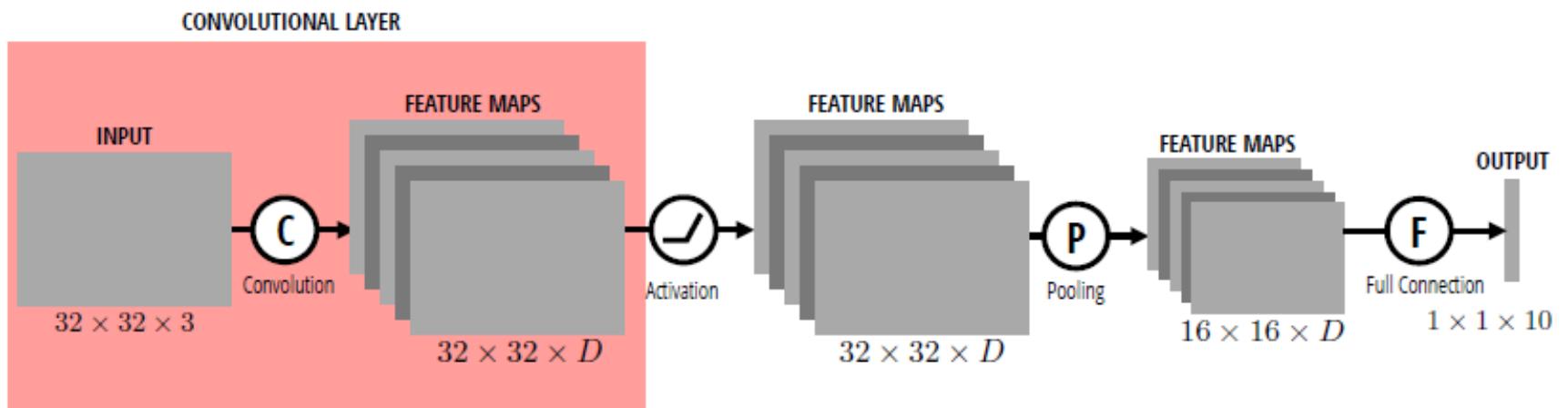


Gradient-based learning applied to document recognition
[[LeCun, Bottou, Bengio, Haffner 1998](#)]



LeNet-1 from 1993

Convolution Layer in CNN

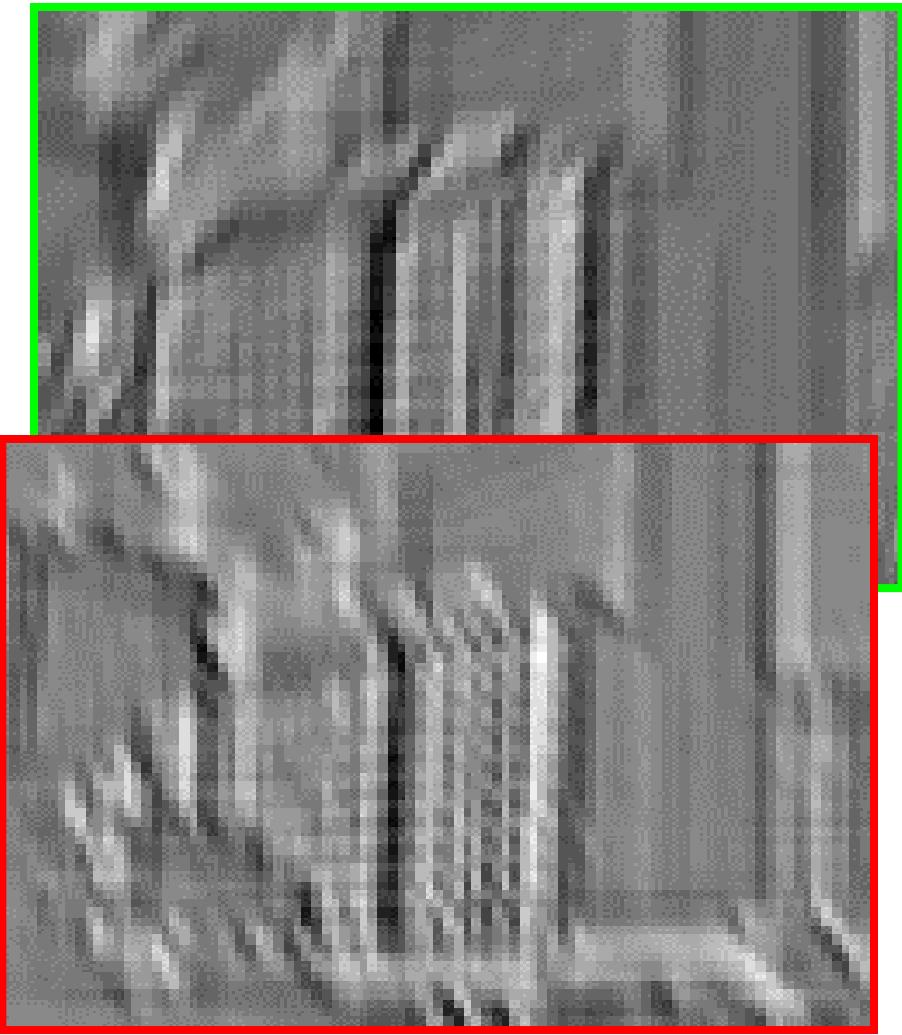


What is a Convolution?

- Weighted moving sum



Input



Feature Activation Map 100

What is a Convolution?

5	9	2	6	7	9	8
---	---	---	---	---	---	---

Signal

1	0	-1
---	---	----

Filter

-1	0	1
----	---	---

5	9	2	6	7	9	8
---	---	---	---	---	---	---

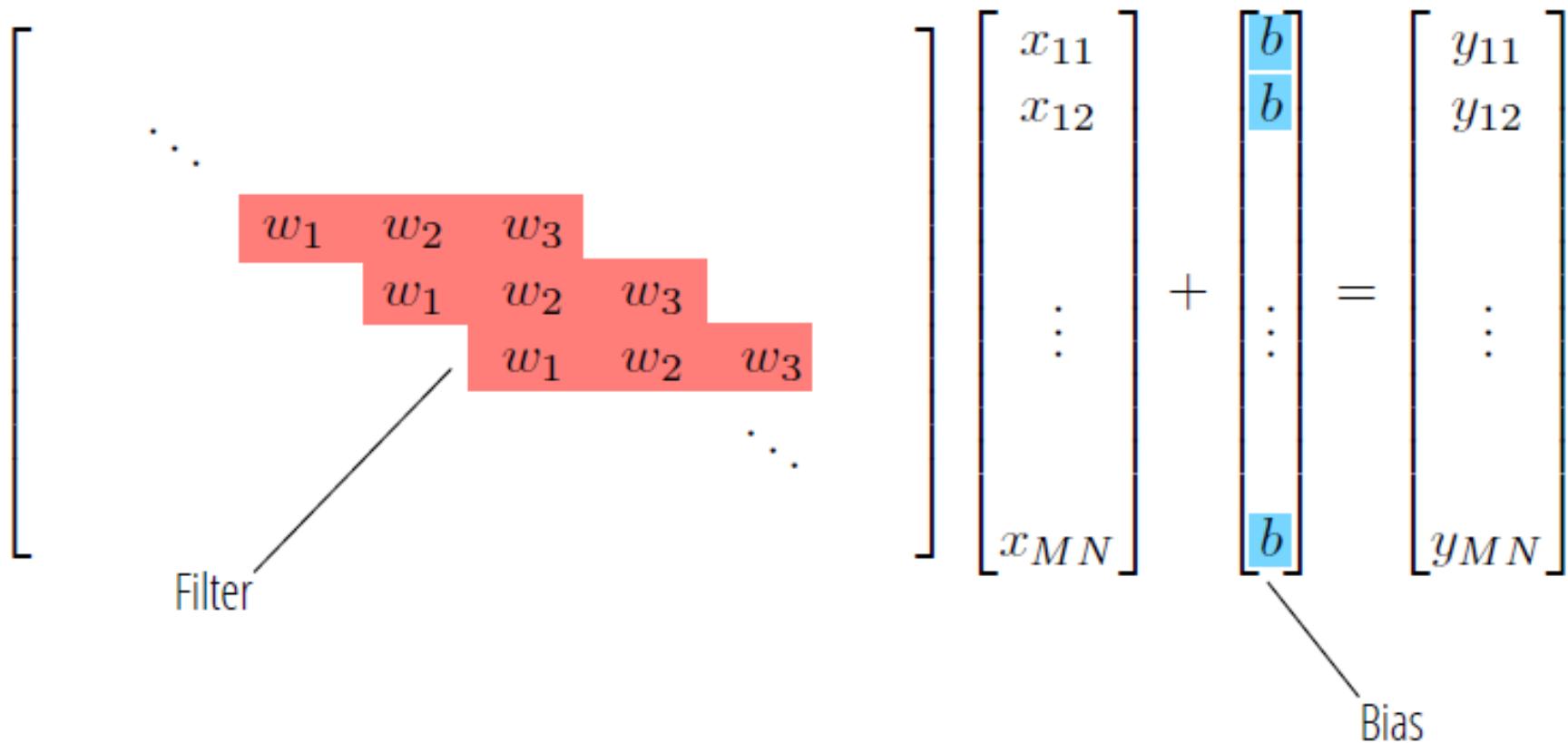
Output

-3	-3	5	3	1
----	----	---	---	---

Convolution is a local linear operator

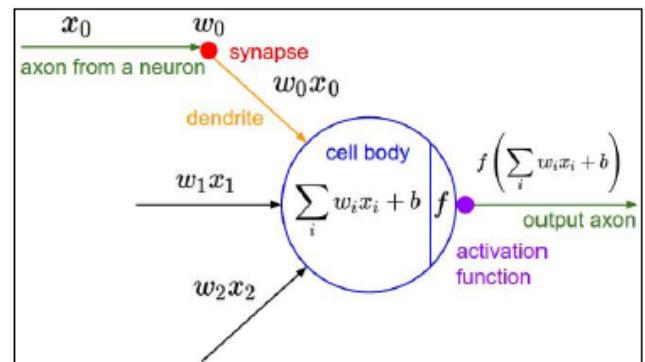
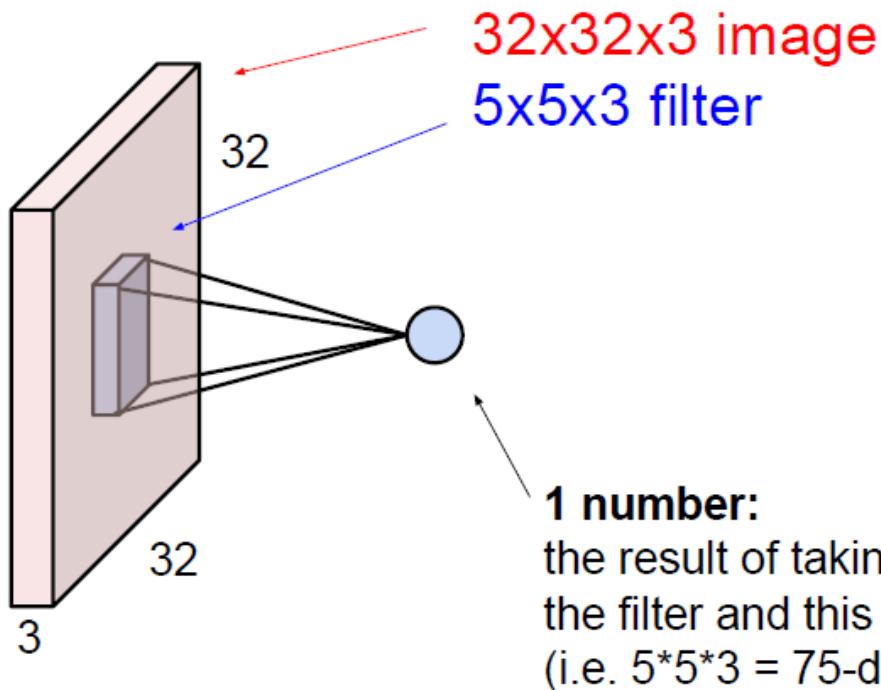
What is a Convolution?

- Toeplitz Matrix Form



Putting them together (cont'd)

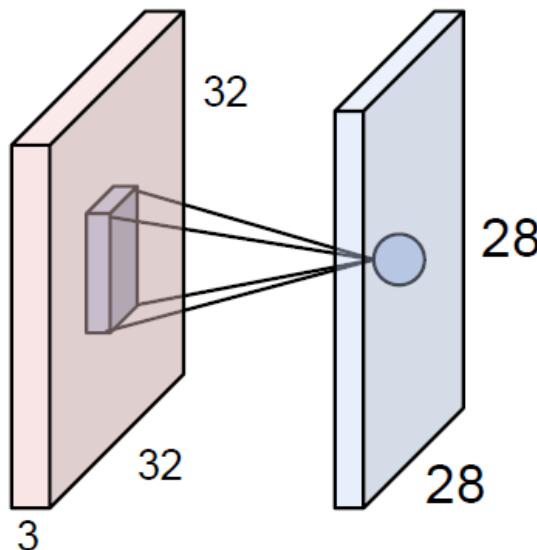
- The brain/neuron view of CONV layer



It's just a neuron with local connectivity...

Putting them together (cont'd)

- The brain/neuron view of CONV layer



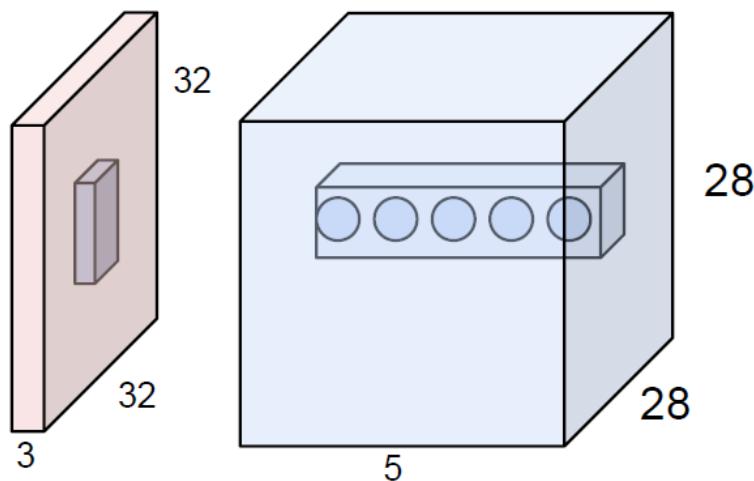
An activation map is a 28×28 sheet of neuron outputs:

- Each is connected to a small region in the input
- All of them share parameters

“5x5 filter” -> “5x5 receptive field for each neuron”

Putting them together (cont'd)

- The brain/neuron view of CONV layer

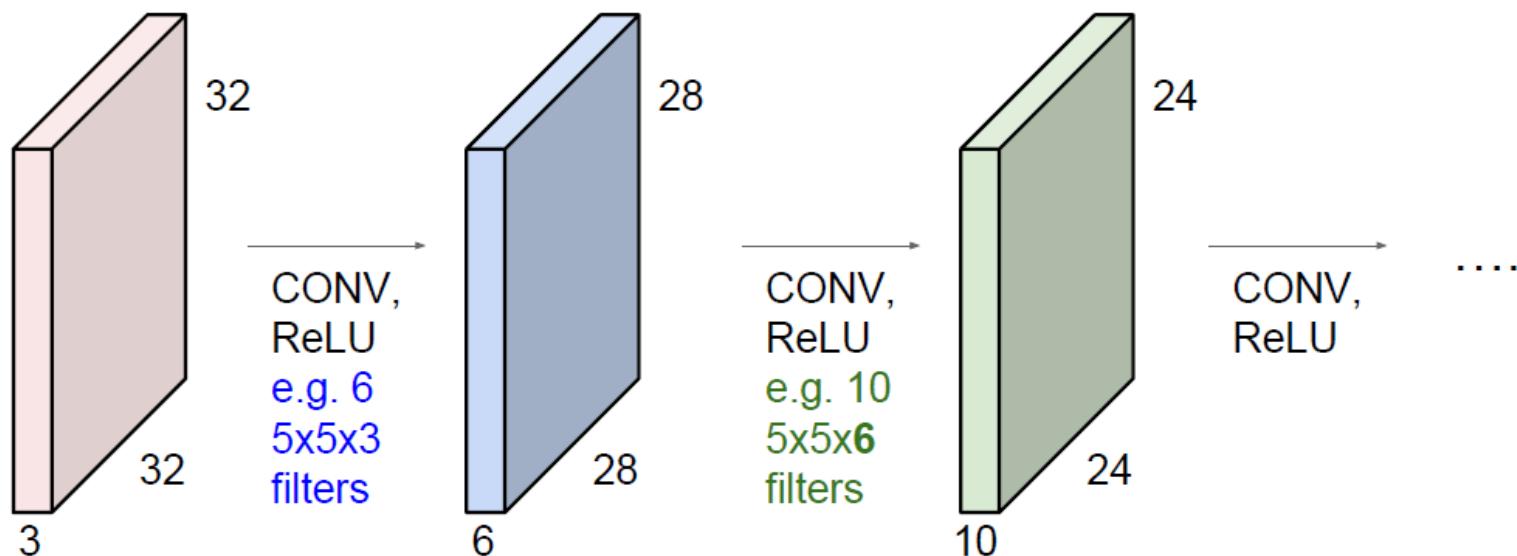


E.g. with 5 filters,
CONV layer consists of
neurons arranged in a 3D grid
($28 \times 28 \times 5$)

There will be 5 different
neurons all looking at the same
region in the input volume

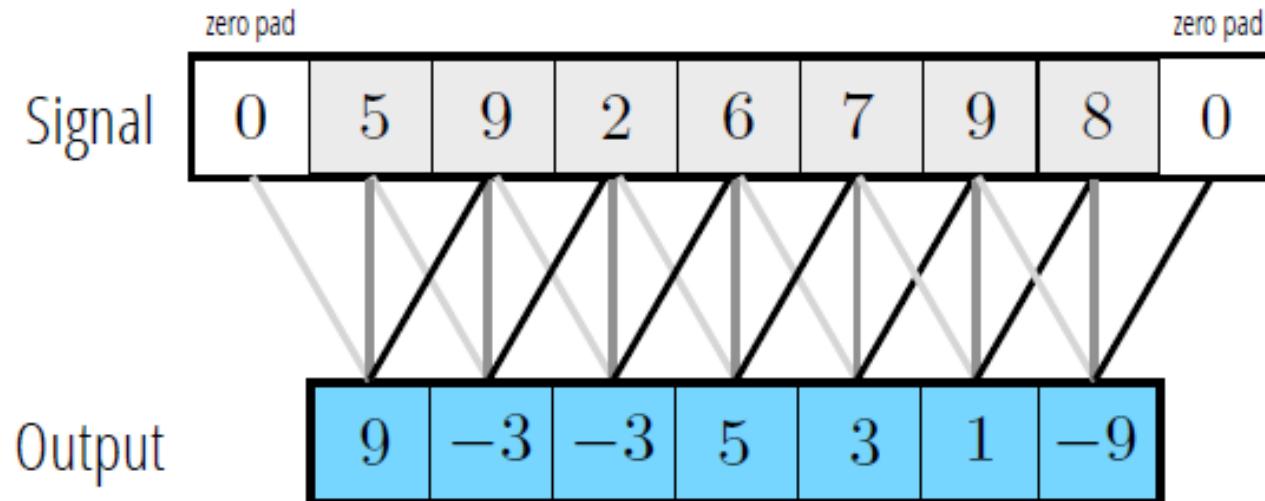
Putting them together (cont'd)

- Image input with 32×32 pixels convolved repeatedly with $5 \times 5 \times 3$ filters shrinks volumes spatially ($32 \rightarrow 28 \rightarrow 24 \rightarrow \dots$).



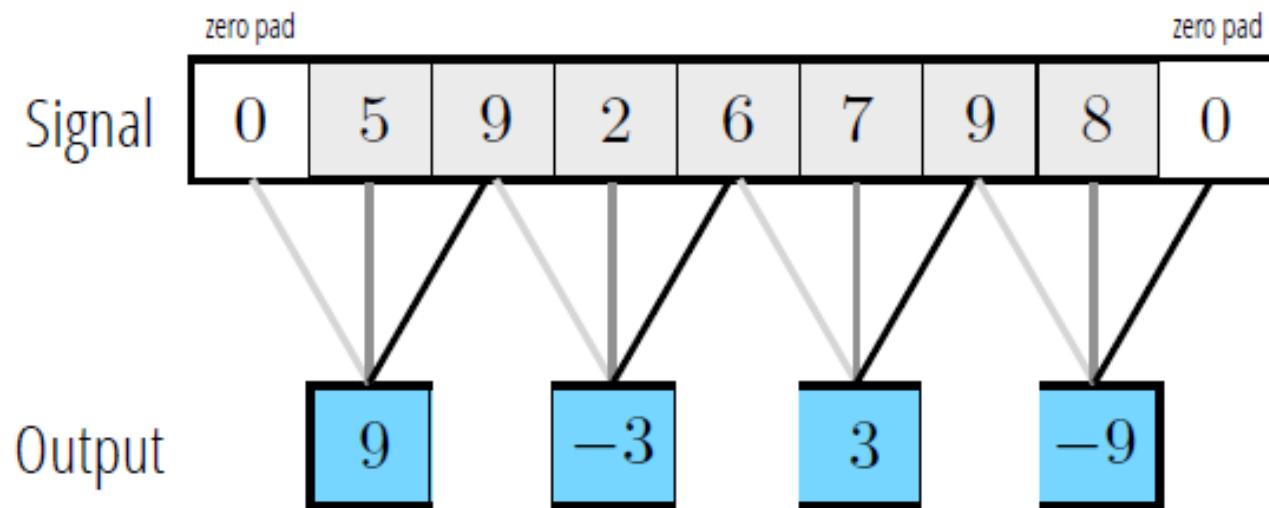
What is a Convolution?

- Zero Padding
 - Output is the same size as input (doesn't shrink as the network gets deeper).



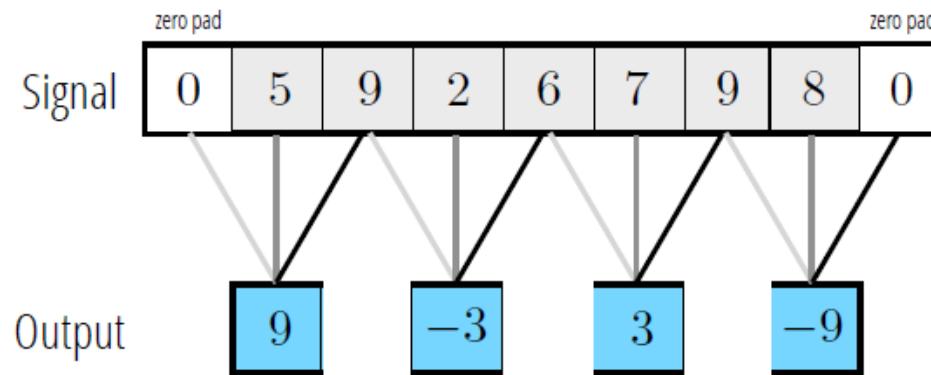
What is a Convolution?

- Stride
 - Step size across signals



What is a Convolution?

- Stride
 - Step size across signals

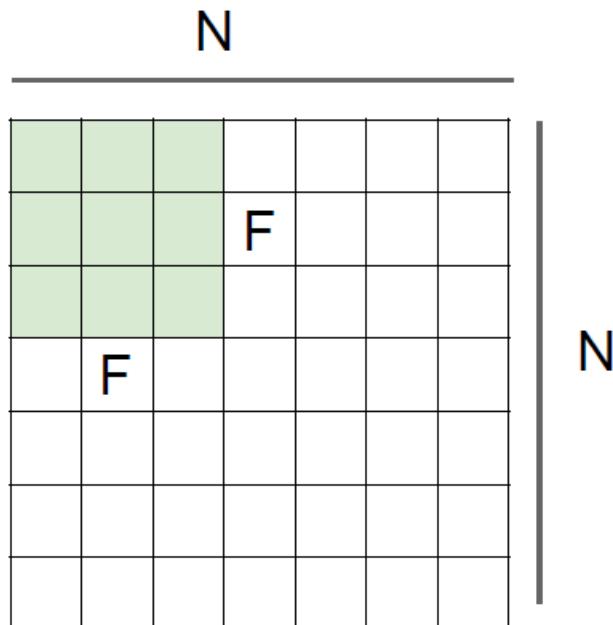


$$\frac{N - c}{s} + 1$$

Input Size Filter Size
Output Size
Stride: step size across the signal

What is a Convolution?

- Stride
 - Step size across signals



Output size:
 $(N - F) / \text{stride} + 1$

e.g. $N = 7$, $F = 3$:
stride 1 => $(7 - 3)/1 + 1 = 5$
stride 2 => $(7 - 3)/2 + 1 = 3$
stride 3 => $(7 - 3)/3 + 1 = 2.33 \vdots$

What is a Convolution?

- Zero Padding + Stride

0	0	0	0	0	0		
0							
0							
0							
0							

e.g. input 7x7

3x3 filter, applied with **stride 1**

pad with 1 pixel border => what is the output?

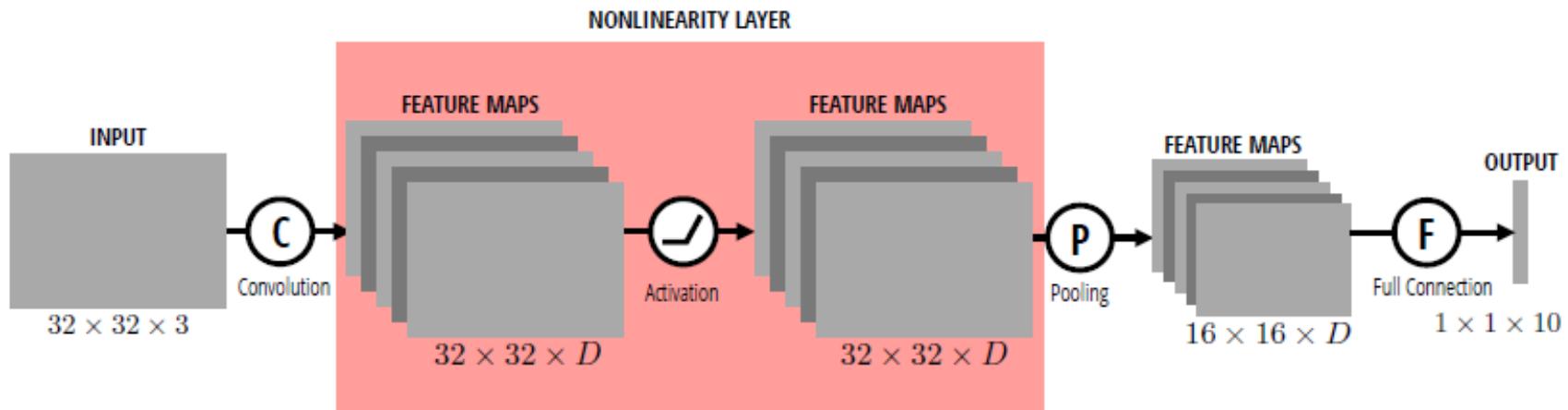
in general, common to see CONV layers with stride 1, filters of size FxF, and zero-padding with $(F-1)/2$. (will preserve size spatially)

e.g. $F = 3 \Rightarrow$ zero pad with 1

$F = 5 \Rightarrow$ zero pad with 2

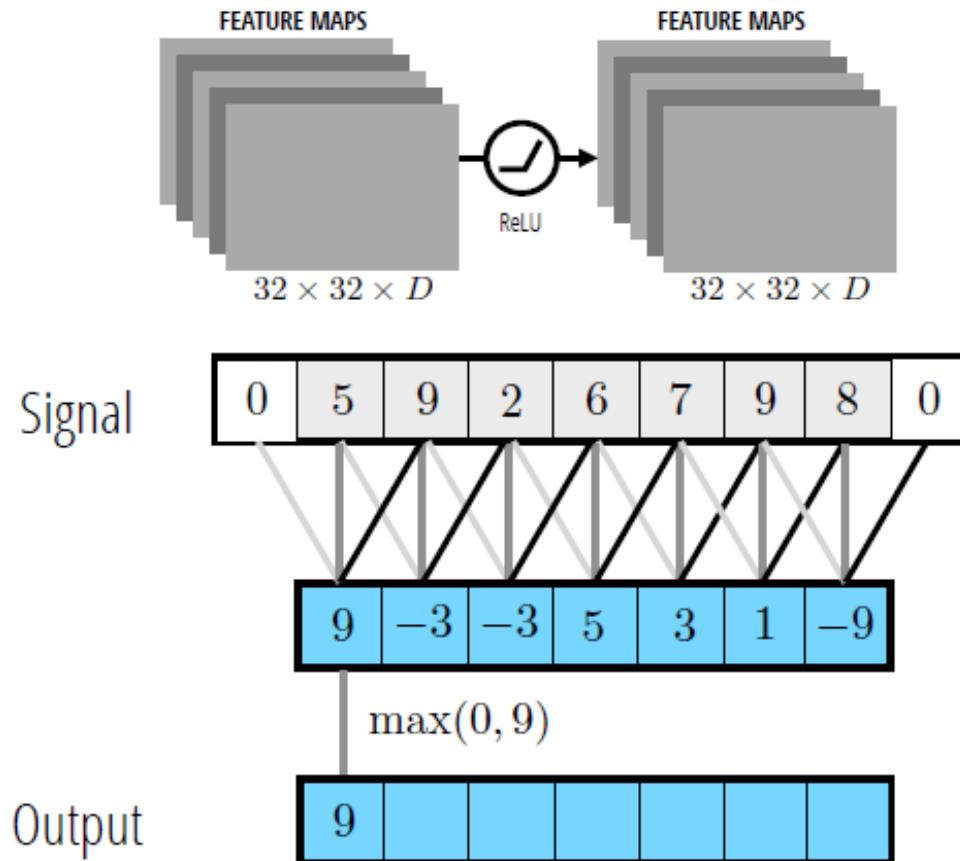
$F = 7 \Rightarrow$ zero pad with 3

Nonlinearity Layer in CNN



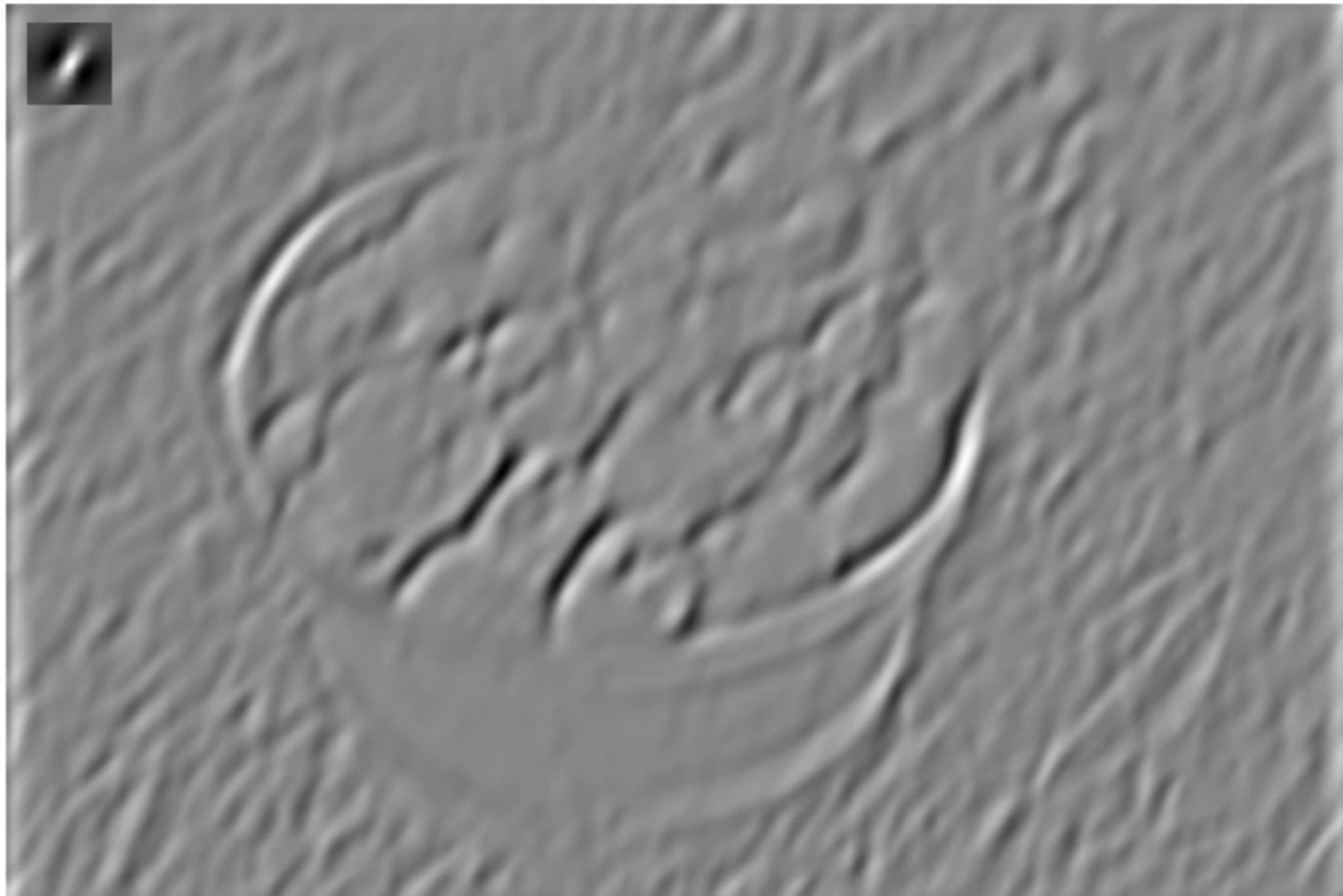
Nonlinearity Layer

- E.g., ReLU (Rectified Linear Unit)
 - Pixel by pixel computation of $\max(0, x)$



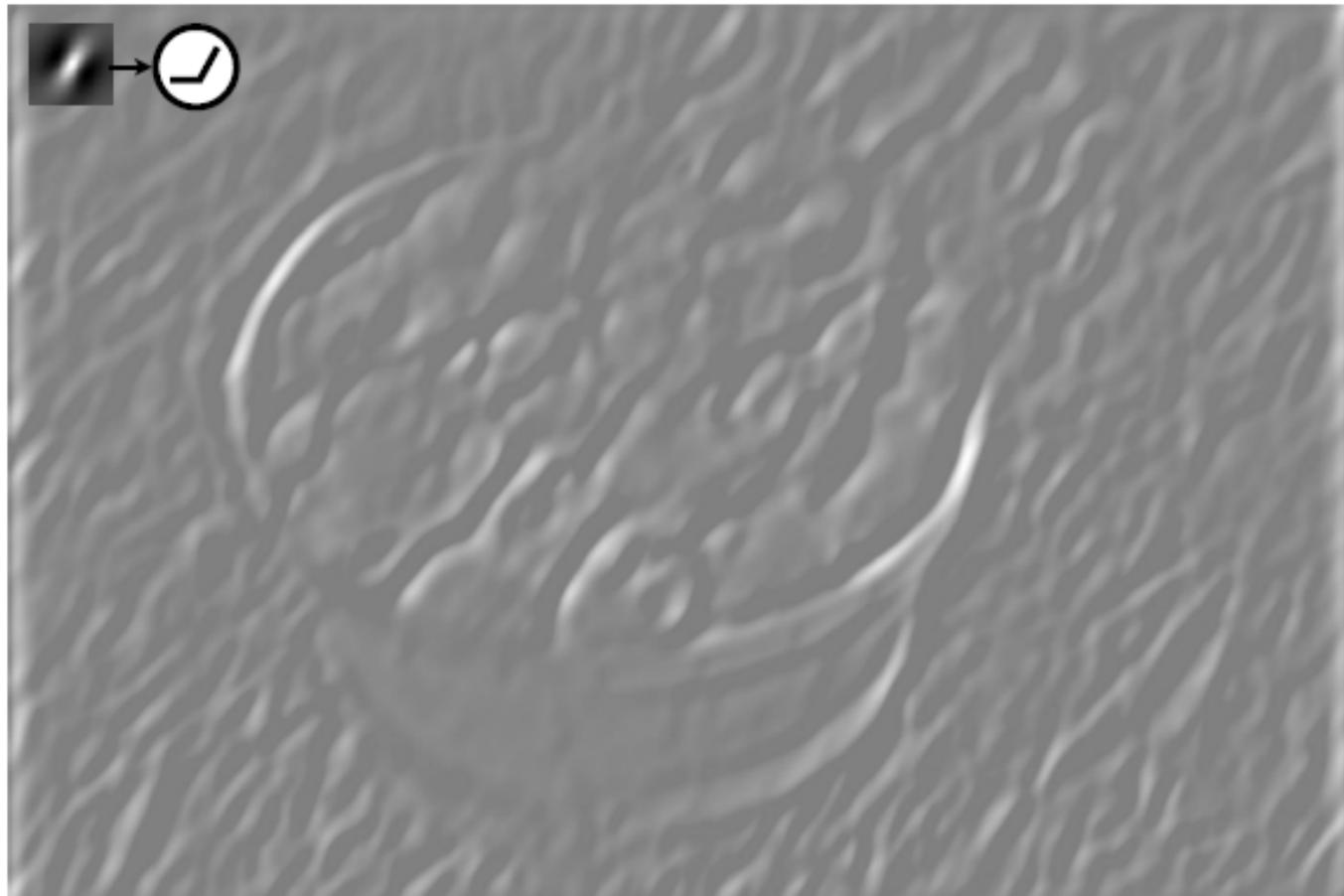
Nonlinearity Layer

- E.g., ReLU (Rectified Linear Unit)
 - Pixel by pixel computation of $\max(0, x)$

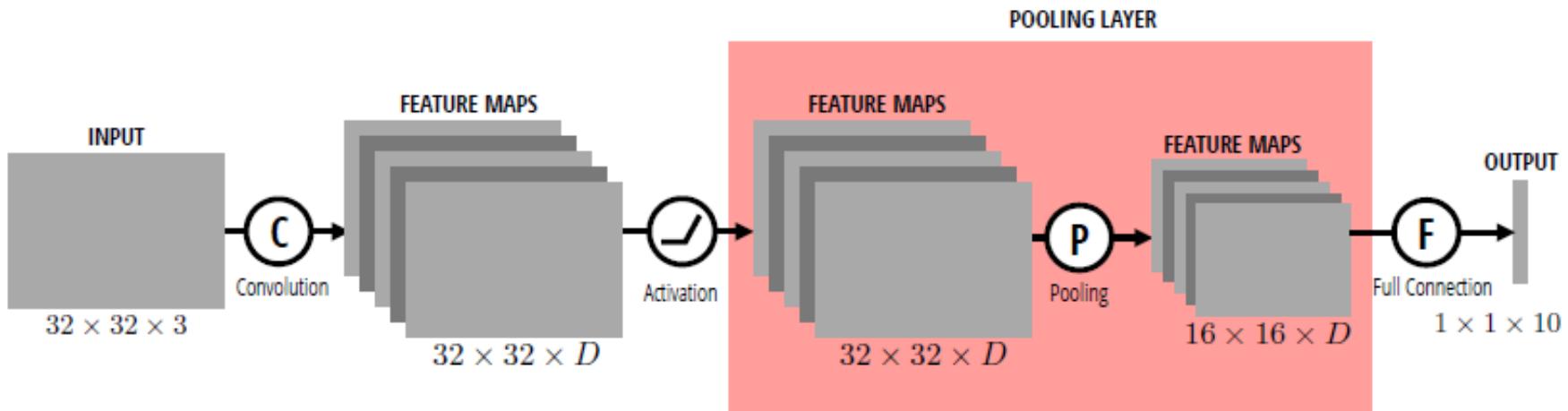


Nonlinearity Layer

- E.g., ReLU (Rectified Linear Unit)
 - Pixel by pixel computation of $\max(0, x)$

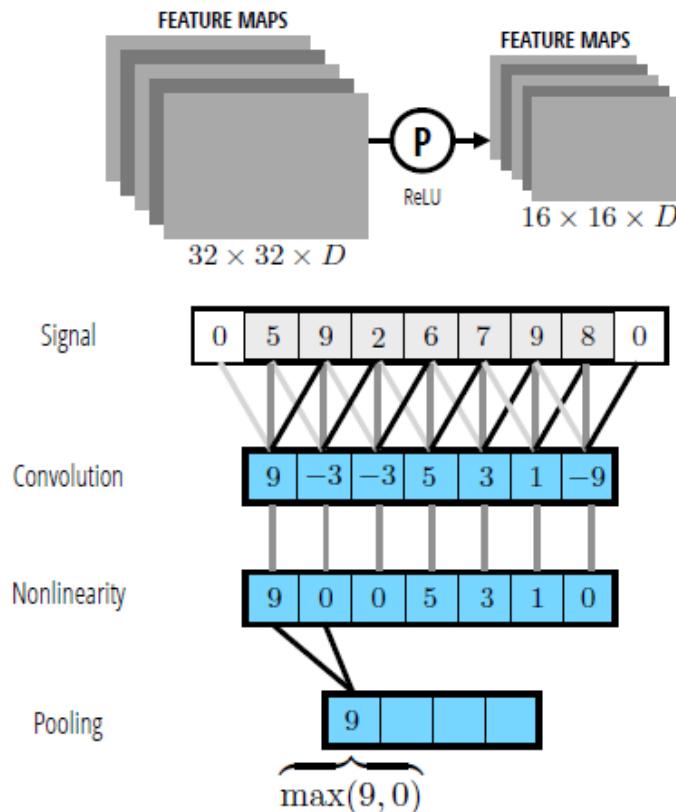


Pooling Layer in CNN



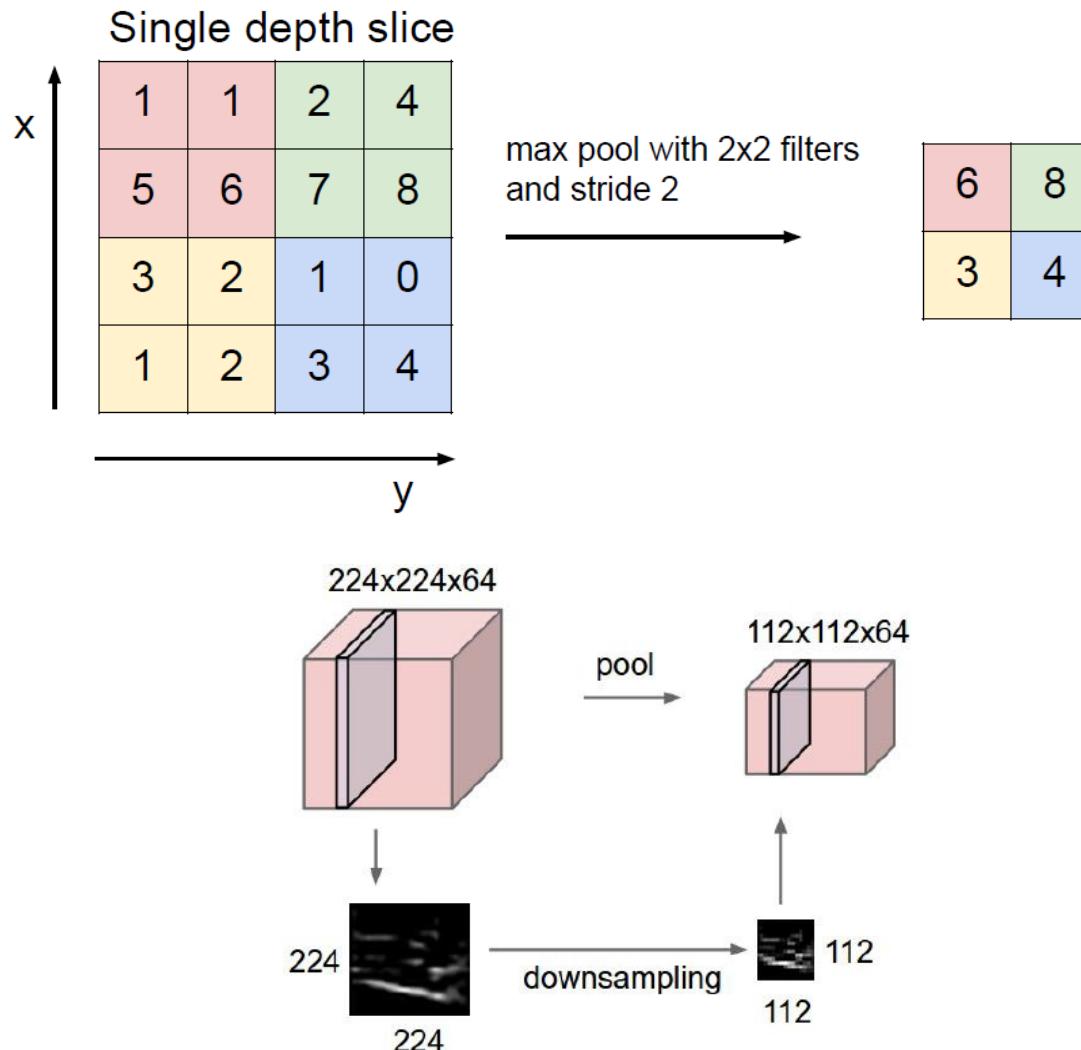
Pooling Layer

- Makes the representations smaller and more manageable
- Operates over each activation map independently
- E.g., Max Pooling

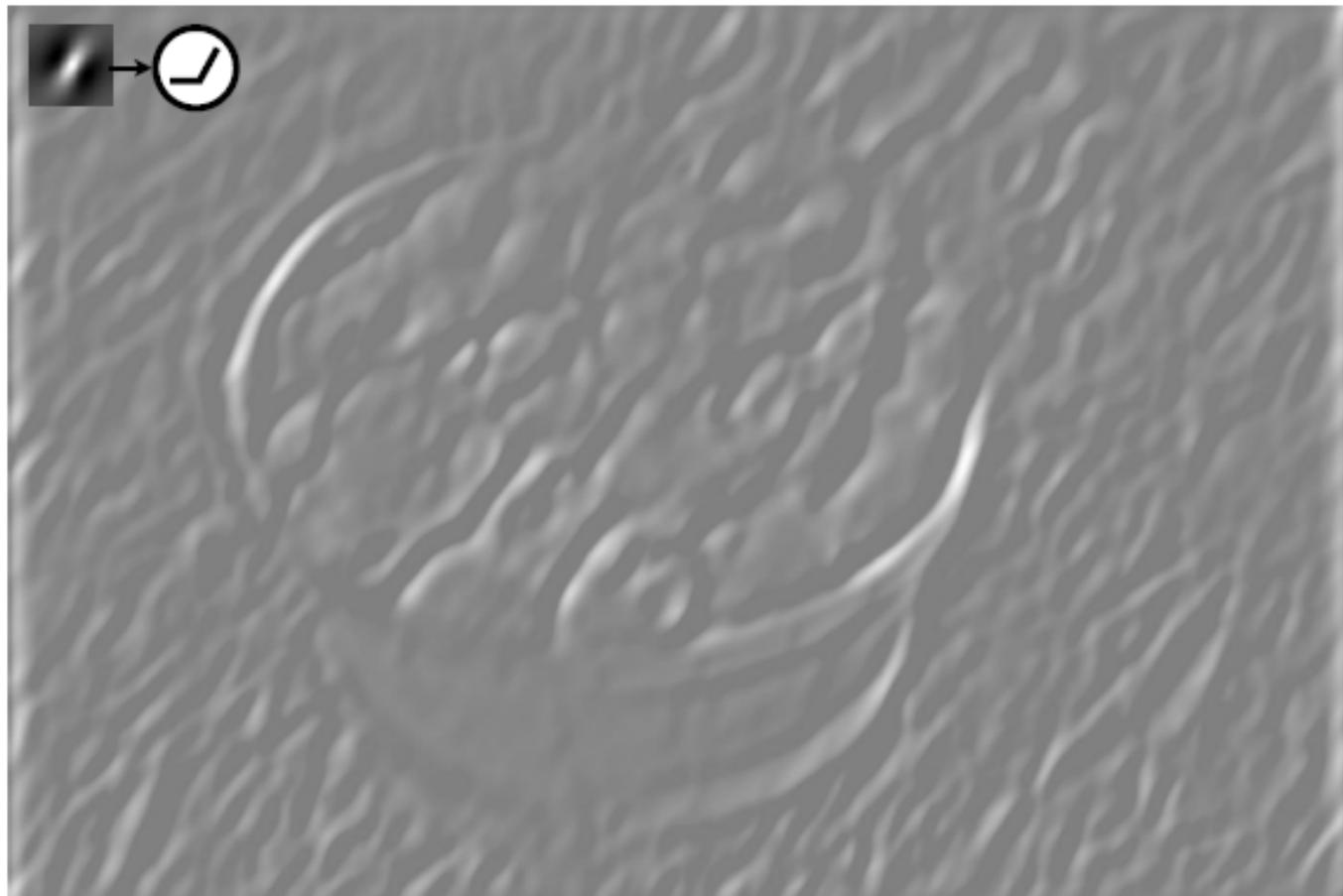


Pooling Layer

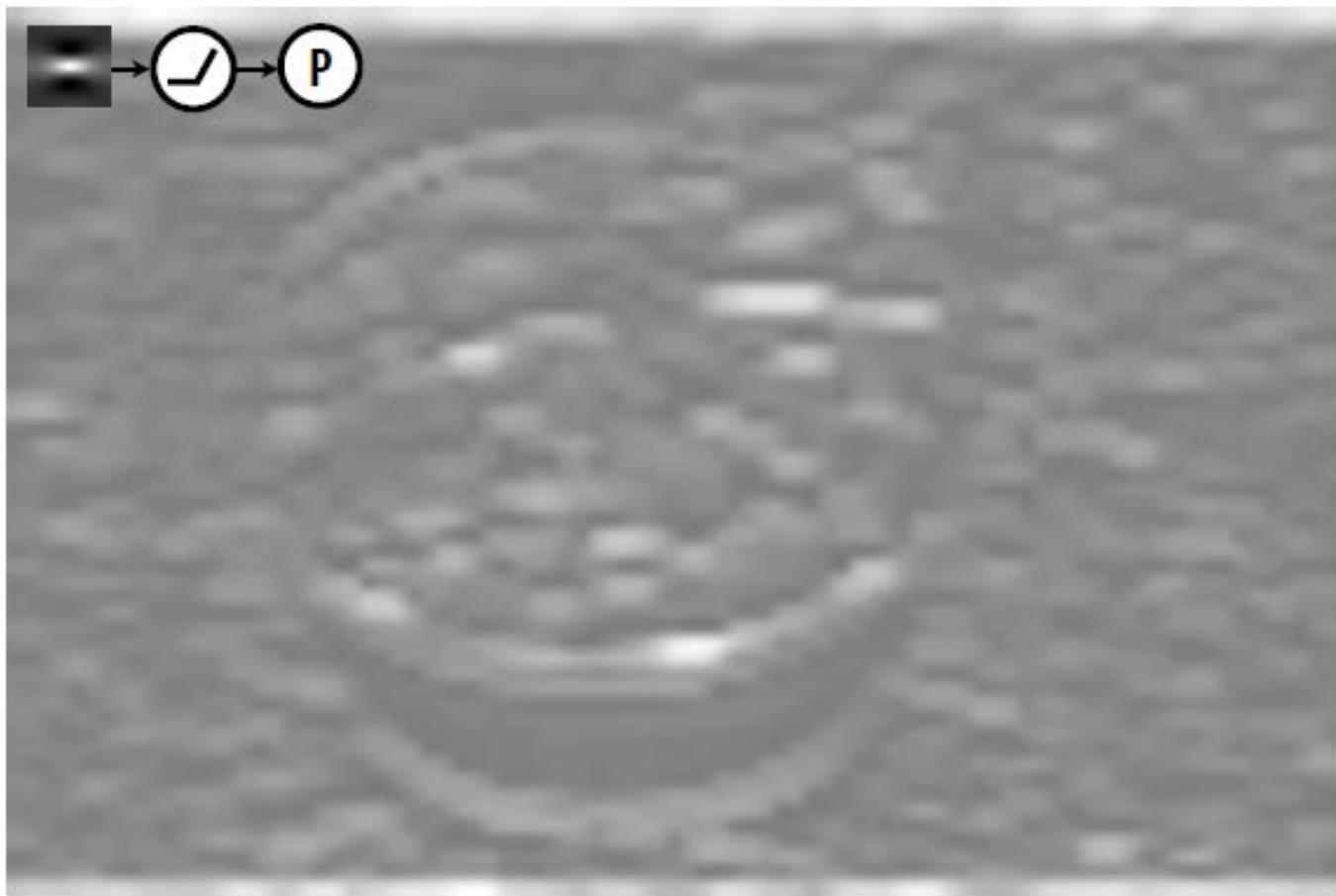
- Reduces the spatial size and provides spatial invariance



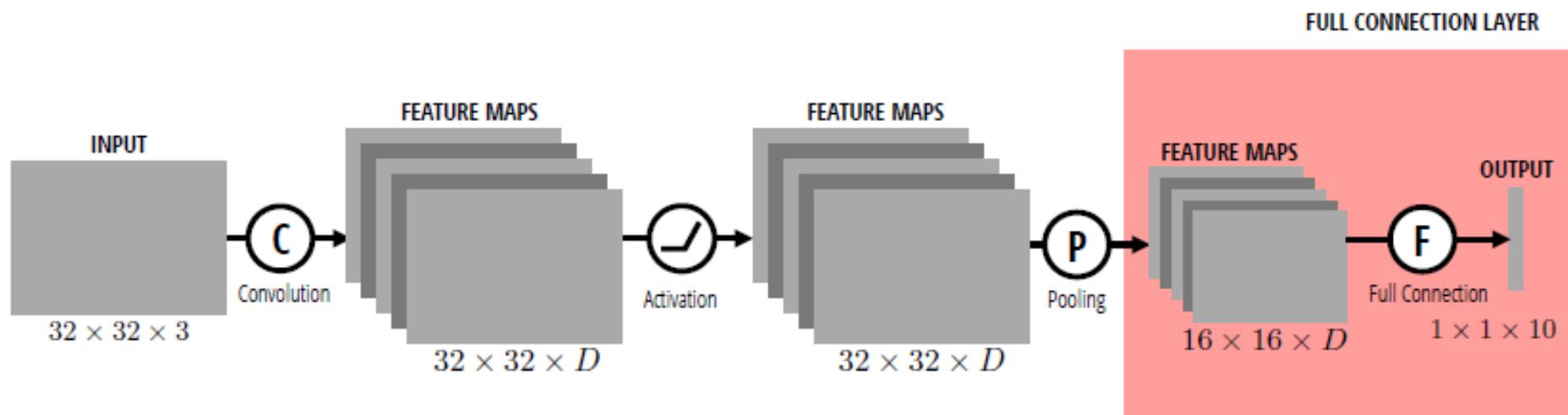
- Example
 - Nonlinearity by ReLU



- Example
 - Max pooling

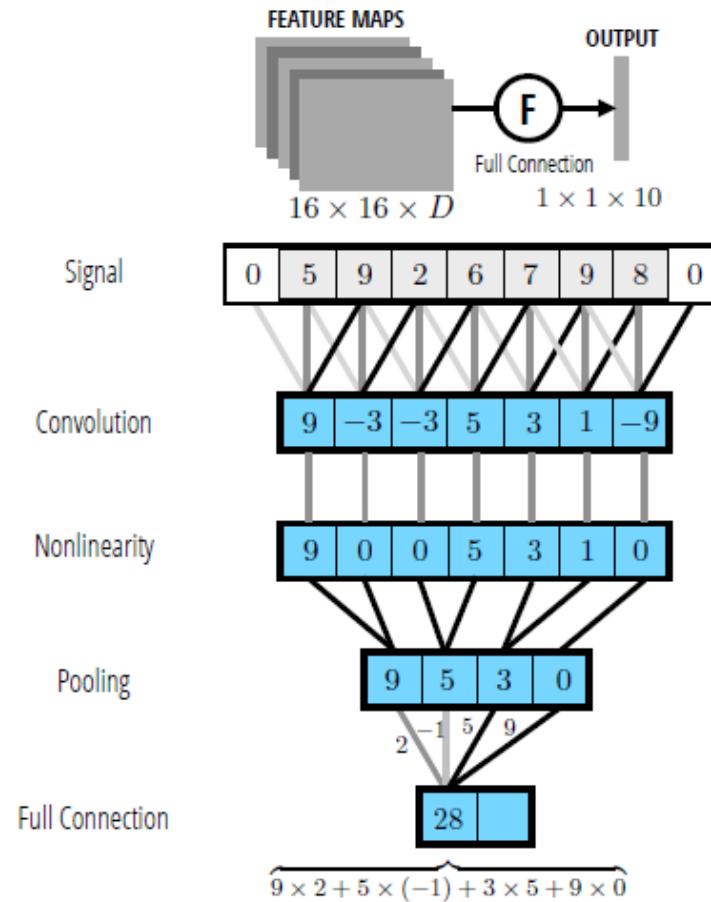


Fully Connected (FC) Layer in CNN



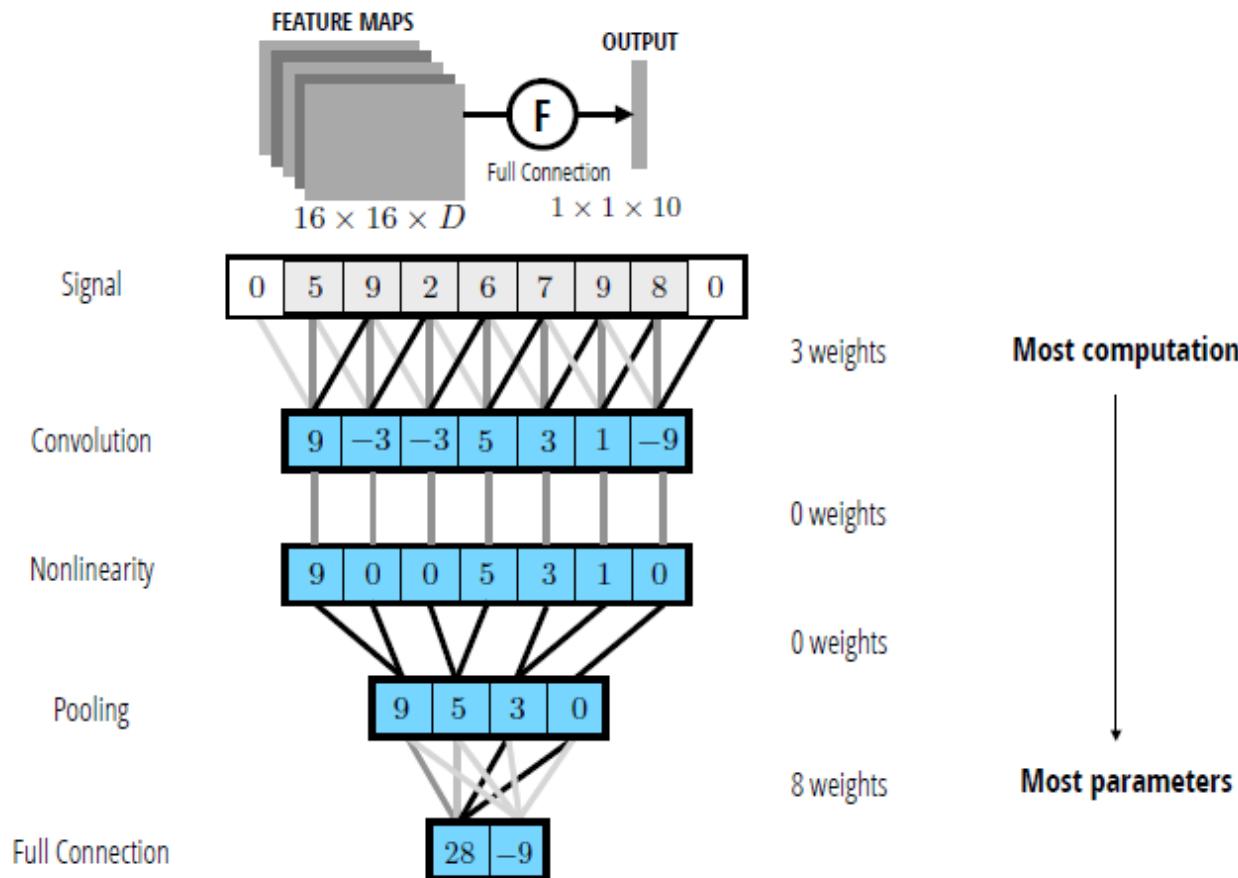
FC Layer

- Contains neurons that connect to the entire input volume, as in ordinary neural networks



FC Layer

- Contains neurons that connect to the entire input volume, as in ordinary neural networks



CNN

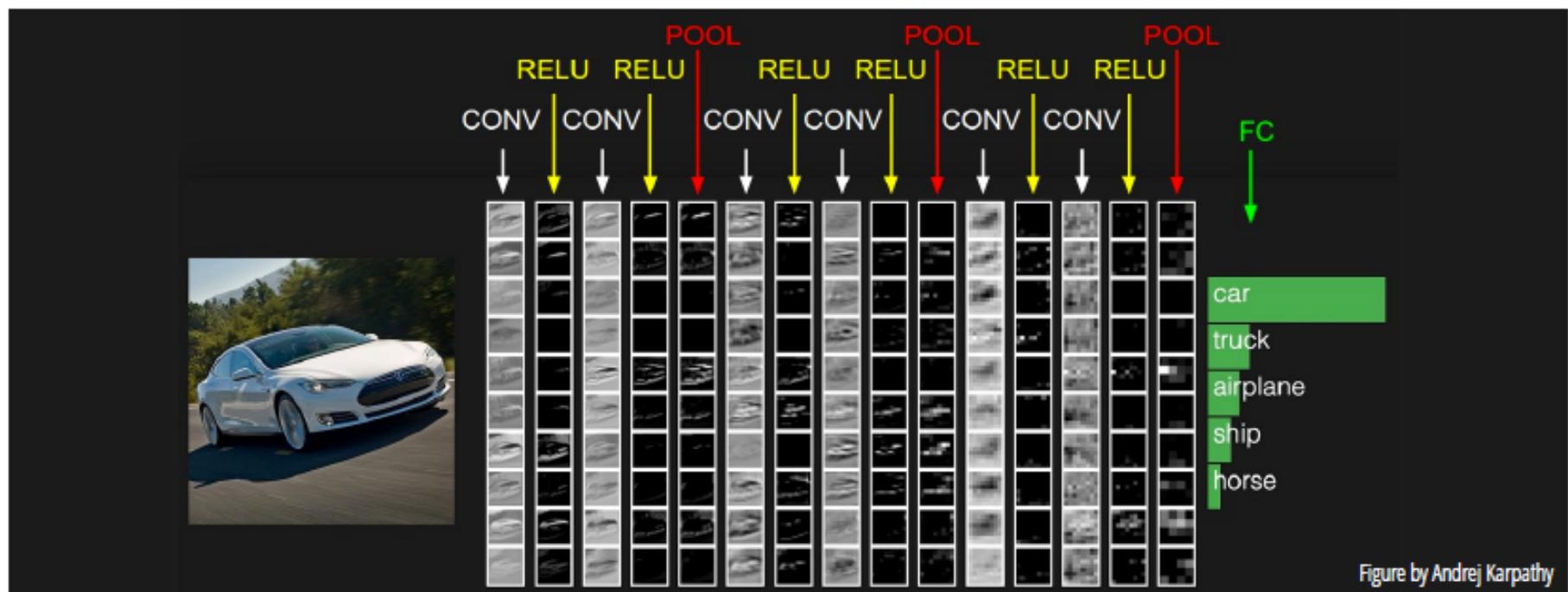
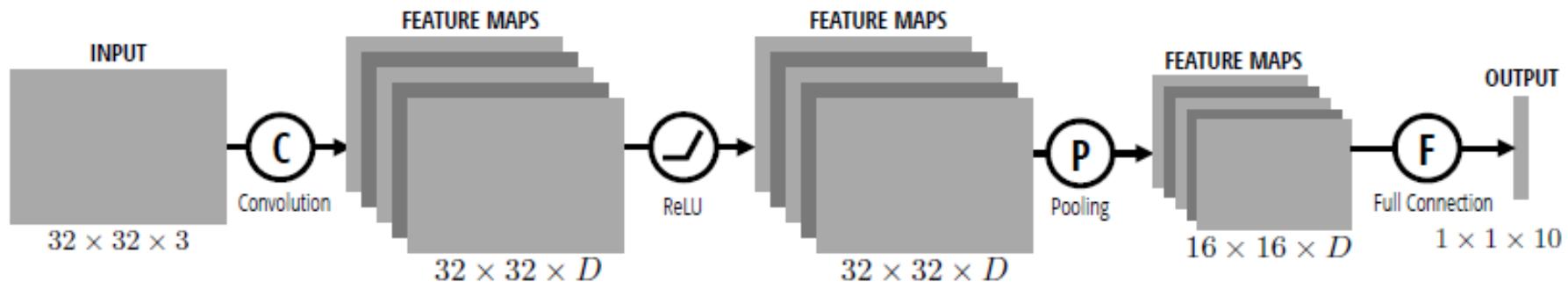
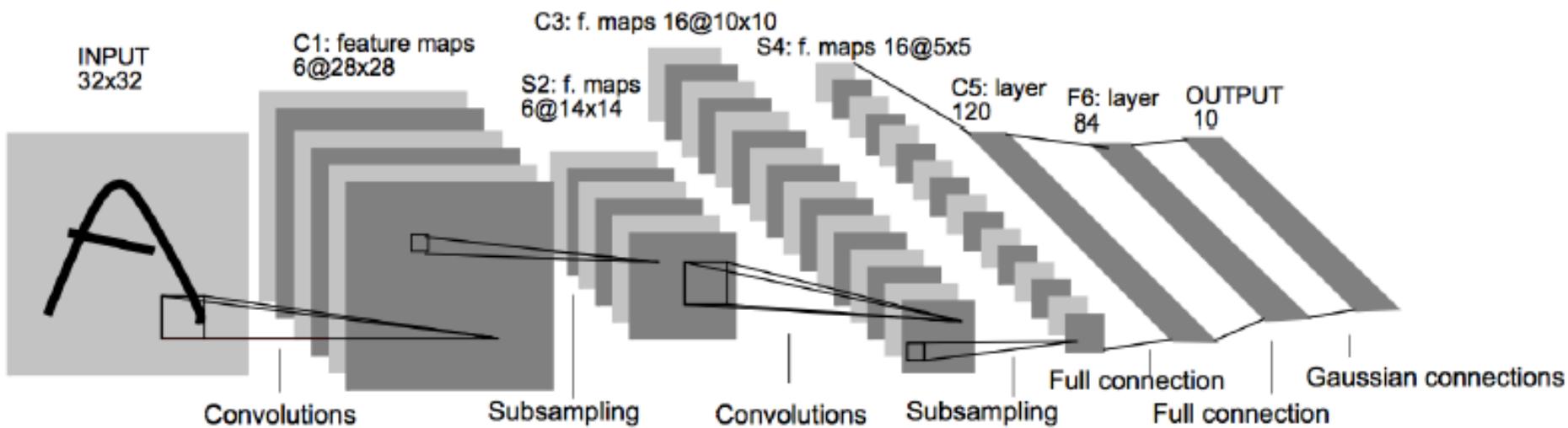


Figure by Andrej Karpathy

LeNet

- Presented by Yann LeCun during the 1990s for reading digits
- Has the elements of modern architectures

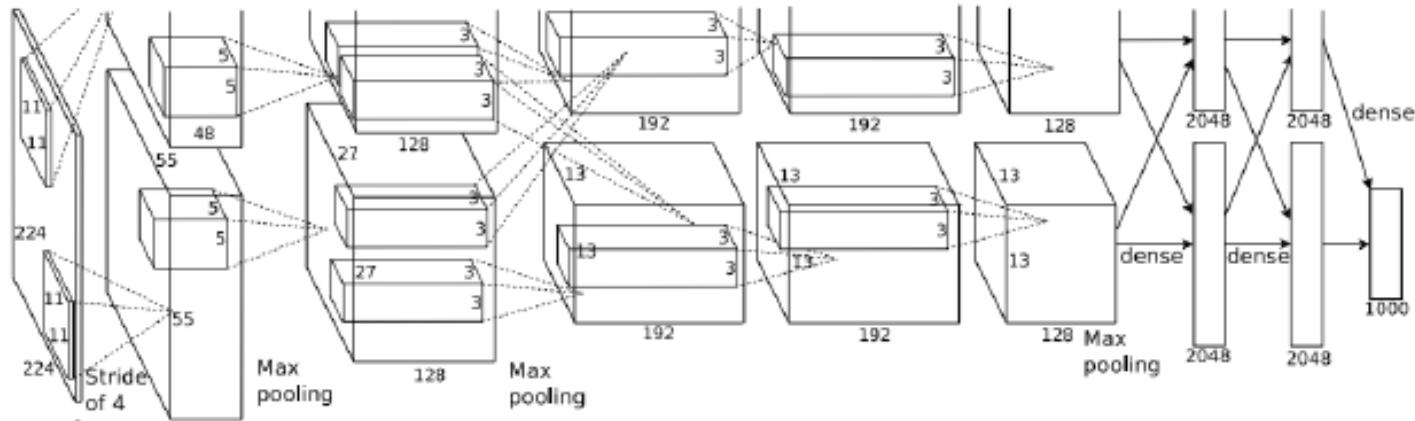
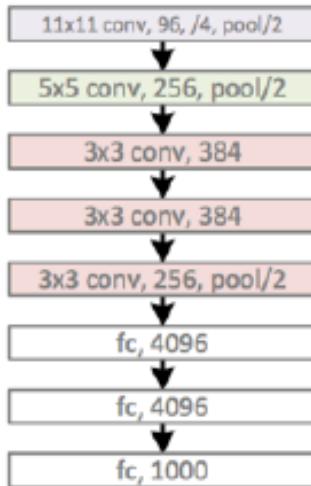


AlexNet [Krizhevsky et al., 2012]

- Repopularized CNN by winning the ImageNet Challenge 2012
- 7 hidden layers, 650,000 neurons, 60M parameters
- Error rate of 16% vs. 26% for 2nd place.

Full (simplified) AlexNet architecture:

- [227x227x3] INPUT
- [55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0
- [27x27x96] MAX POOL1: 3x3 filters at stride 2
- [27x27x96] NORM1: Normalization layer
- [27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2
- [13x13x256] MAX POOL2: 3x3 filters at stride 2
- [13x13x256] NORM2: Normalization layer
- [13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1
- [13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1
- [13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1
- [6x6x256] MAX POOL3: 3x3 filters at stride 2
- [4096] FC6: 4096 neurons
- [4096] FC7: 4096 neurons
- [1000] FC8: 1000 neurons (class scores)



Deep or Not?

- Depth of the network is critical for performance.



AlexNet: 8 Layers with 18.2% top-5 error

Removing Layer 7 reduces 16 million parameters, but only 1.1% drop in performance!

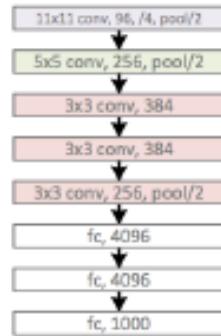
Removing Layer 6 and 7 reduces 50 million parameters, but only 5.7% drop in performance

Removing middle conv layers reduces 1 million parameters, but only 3% drop in performance

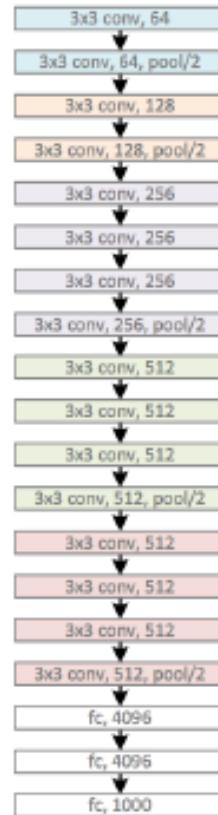
Removing feature & conv layers produces a **33% drop** in performance

CNN: A Revolution of Depth

AlexNet, 8 layers
(ILSVRC 2012)



VGG, 19 layers
(ILSVRC 2014)

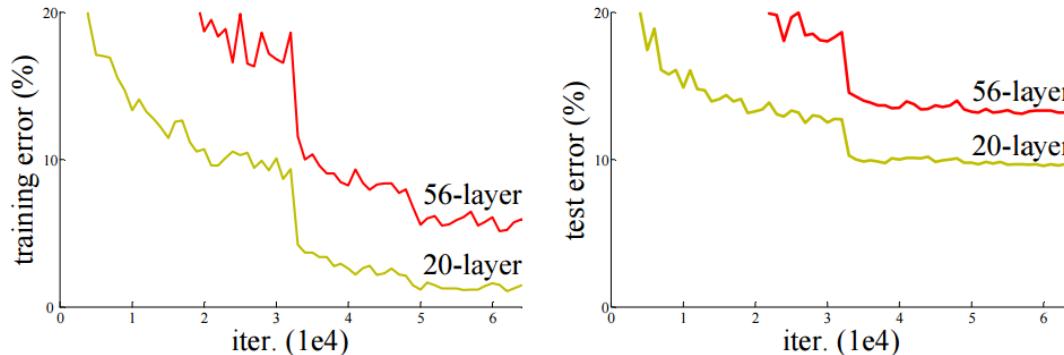


GoogleNet, 22 layers
(ILSVRC 2014)

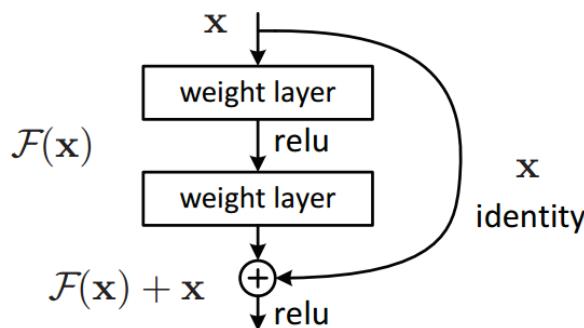


ResNet

- Can we just increase the #layer?



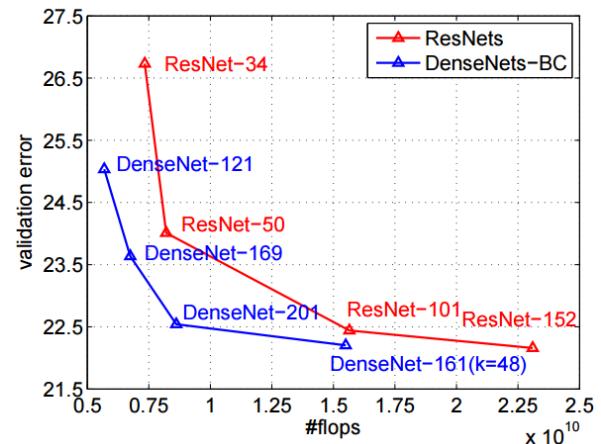
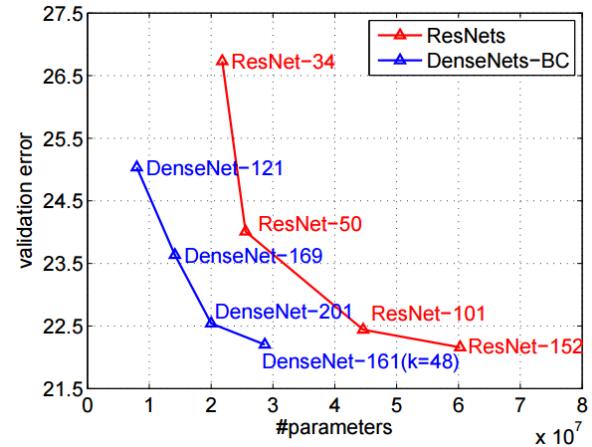
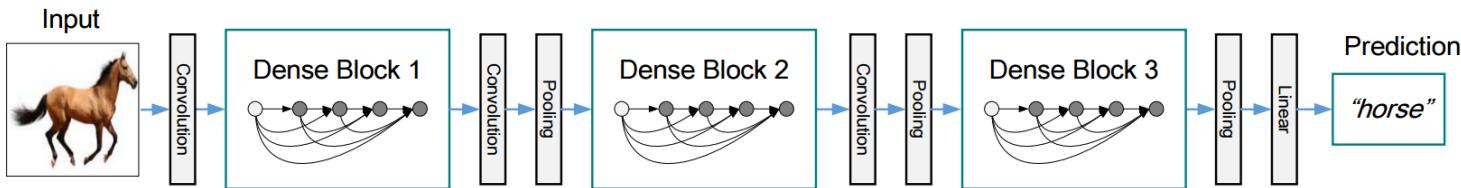
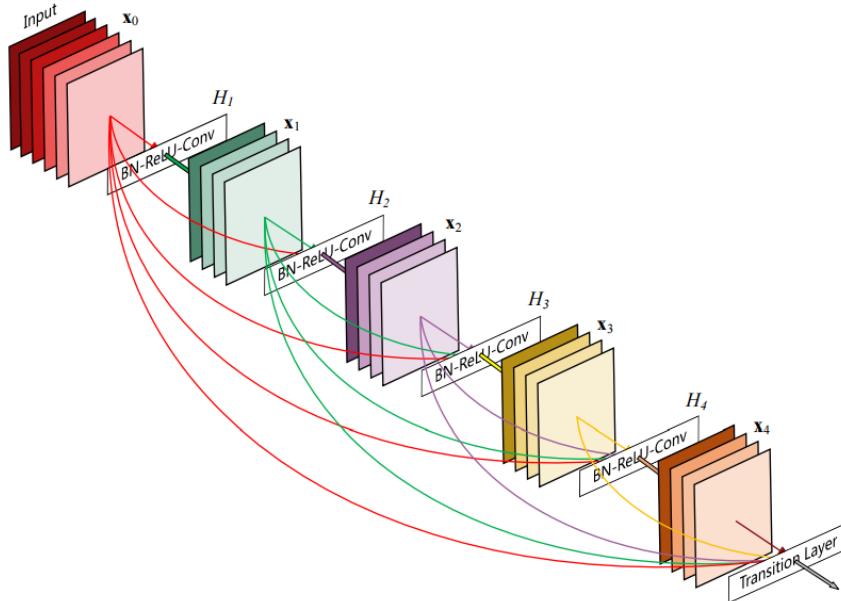
- How can we train very deep network?
 - Residual learning



method	top-5 err. (test)
VGG [41] (ILSVRC'14)	7.32
GoogLeNet [44] (ILSVRC'14)	6.66
VGG [41] (v5)	6.8
PReLU-net [13]	4.94
BN-inception [16]	4.82
ResNet (ILSVRC'15)	3.57

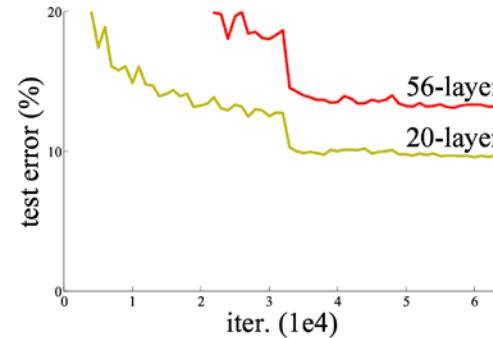
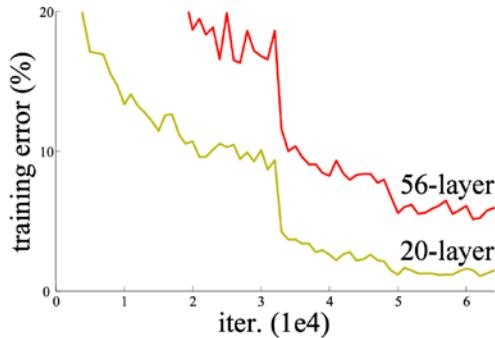
DenseNet

- Shorter connections (like ResNet) help
- Why not just connect them all?

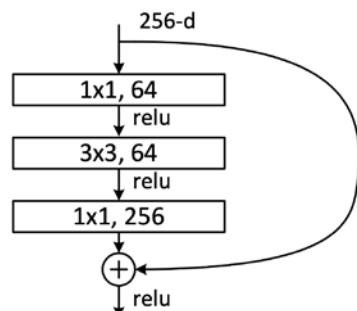


ResNet (cont'd)

- Can we just increase # of layers?



- How to train very deep networks?
 - Residual learning



Non-Bottleneck
(ResNet-18, 34)

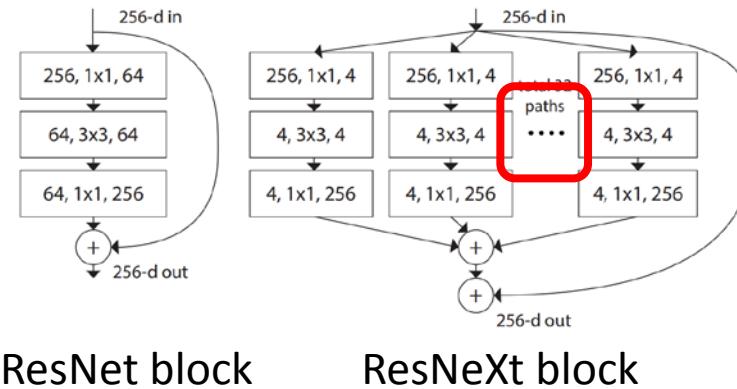
Bottleneck
(ResNet-50, 101, 152)

method	top-5 err. (test)
VGG [41] (ILSVRC'14)	7.32
GoogLeNet [44] (ILSVRC'14)	6.66
VGG [41] (v5)	6.8
PReLU-net [13]	4.94
BN-inception [16]	4.82
ResNet (ILSVRC'15)	3.57

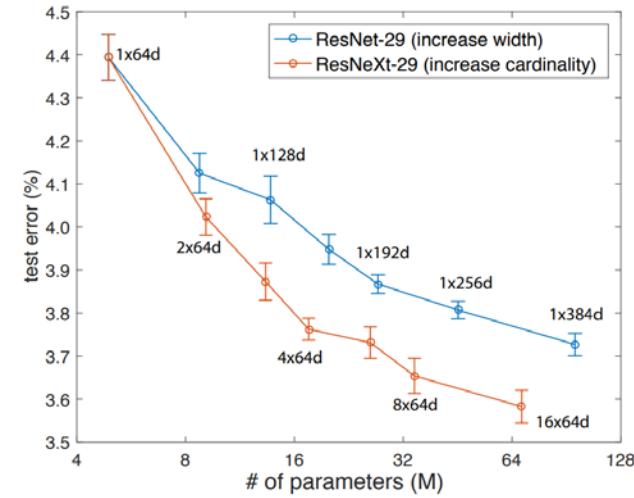
Ref: He, Kaiming, et al. "Deep residual learning for image recognition." CVPR, 2016.

ResNeXT

- Deeper and wider → better...what else?
 - Increase cardinality

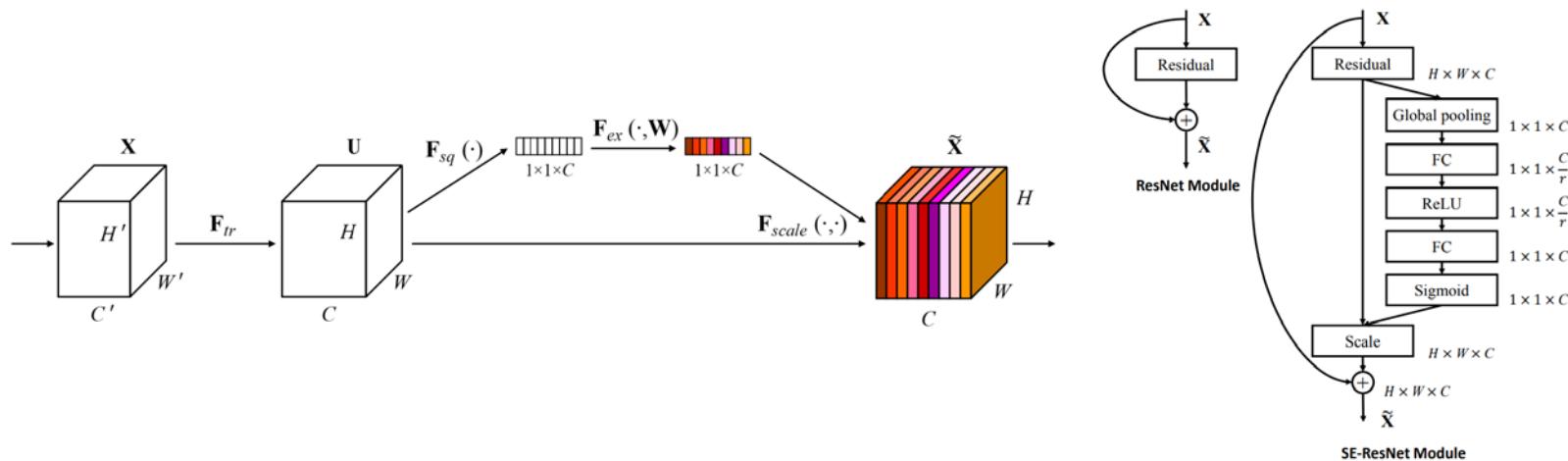


	setting	top-1 error (%)
ResNet-50	1 × 64d	23.9
ResNeXt-50	2 × 40d	23.0
ResNeXt-50	4 × 24d	22.6
ResNeXt-50	8 × 14d	22.3
ResNeXt-50	32 × 4d	22.2
ResNet-101	1 × 64d	22.0
ResNeXt-101	2 × 40d	21.7
ResNeXt-101	4 × 24d	21.4
ResNeXt-101	8 × 14d	21.3
ResNeXt-101	32 × 4d	21.2



Squeeze-and-Excitation Net (SENet)

- How to improve acc. without much overhead?
 - Feature recalibration (channel attention)



	original		re-implementation			SENet		
	top-1 err.	top-5 err.	top-1 err.	top-5 err.	GFLOPs	top-1 err.	top-5 err.	GFLOPs
ResNet-50 [13]	24.7	7.8	24.80	7.48	3.86	23.29 _(1.51)	6.62 _(0.86)	3.87
ResNet-101 [13]	23.6	7.1	23.17	6.52	7.58	22.38 _(0.79)	6.07 _(0.45)	7.60
ResNet-152 [13]	23.0	6.7	22.42	6.34	11.30	21.57 _(0.85)	5.73 _(0.61)	11.32
ResNeXt-50 [19]	22.2	-	22.11	5.90	4.24	21.10 _(1.01)	5.49 _(0.41)	4.25
ResNeXt-101 [19]	21.2	5.6	21.18	5.57	7.99	20.70 _(0.48)	5.01 _(0.56)	8.00
VGG-16 [11]	-	-	27.02	8.81	15.47	25.22 _(1.80)	7.70 _(1.11)	15.48
BN-Inception [6]	25.2	7.82	25.38	7.89	2.03	24.23 _(1.15)	7.14 _(0.75)	2.04
Inception-ResNet-v2 [21]	19.9 [†]	4.9 [†]	20.37	5.21	11.75	19.80 _(0.57)	4.79 _(0.42)	11.76

Remarks

- CNN:
 - Reduce the number of parameters
 - Reduce the memory requirements
 - Make computation independent of the size of the image
- Neuroscience provides strong inspiration on the NN design, but little guidance on **how to train CNNs**.
- Few structures discussed: convolution, nonlinearity, pooling

Training Convolutional Neural Networks

- Backpropagation +
stochastic gradient descent with momentum
 - [Neural Networks: Tricks of the Trade](#)
- Dropout
- Data augmentation
- Batch normalization

What Will We Cover Next Week?

- Pytorch Framework Tutorial
 - Introduction to Pytorch
 - Basis: Simple Classification Example
 - Advance : nn Module Class
 - Advance : DataSet & DataLoader Class
 - Advance : Finetuning with pretrained model
- HW #1 is due **3/23 Sat 3AM** & **no late submission!!**