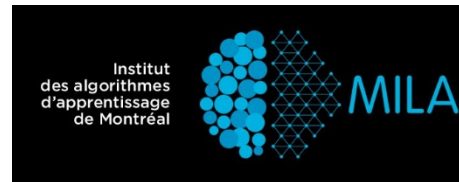


# Feedforward Neural Networks

**Jian Tang**

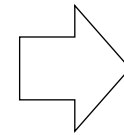
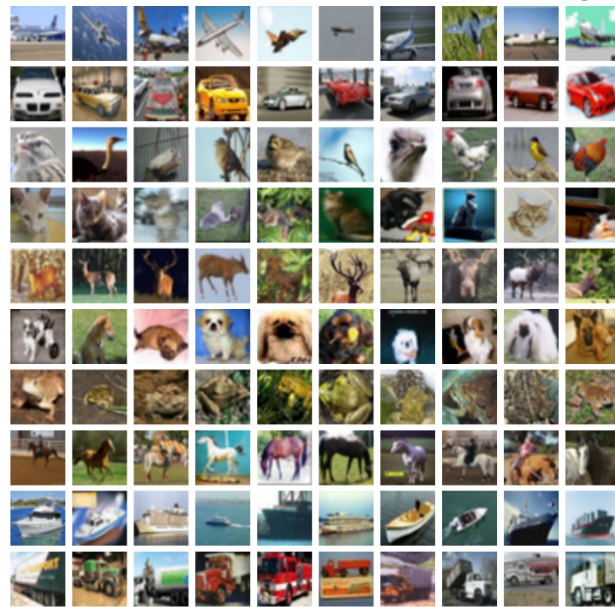
tangjianpku@gmail.com

**HEC MONTREAL**



# The task

- The goal is to learn a mapping function  $y = f(x; \theta)$  (e.g., for classification  $f: R^d \rightarrow C$ ).



**airplane**

**automobile**

**bird**

**cat**

**deer**

**dog**

**frog**

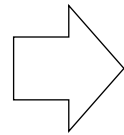
**horse**

**ship**

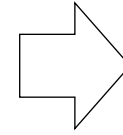
**truck**

Example: image classification

# Traditional Machine Learning



**Hand-crafted**  
Feature Extractor

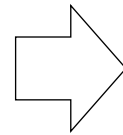


Simple Trainable Classifier  
e.g., SVM, LR

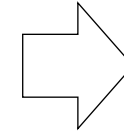


Domain experts

# Deep Learning= End-to-end Learning/Feature Learning



**Trainable**  
Feature Extractor

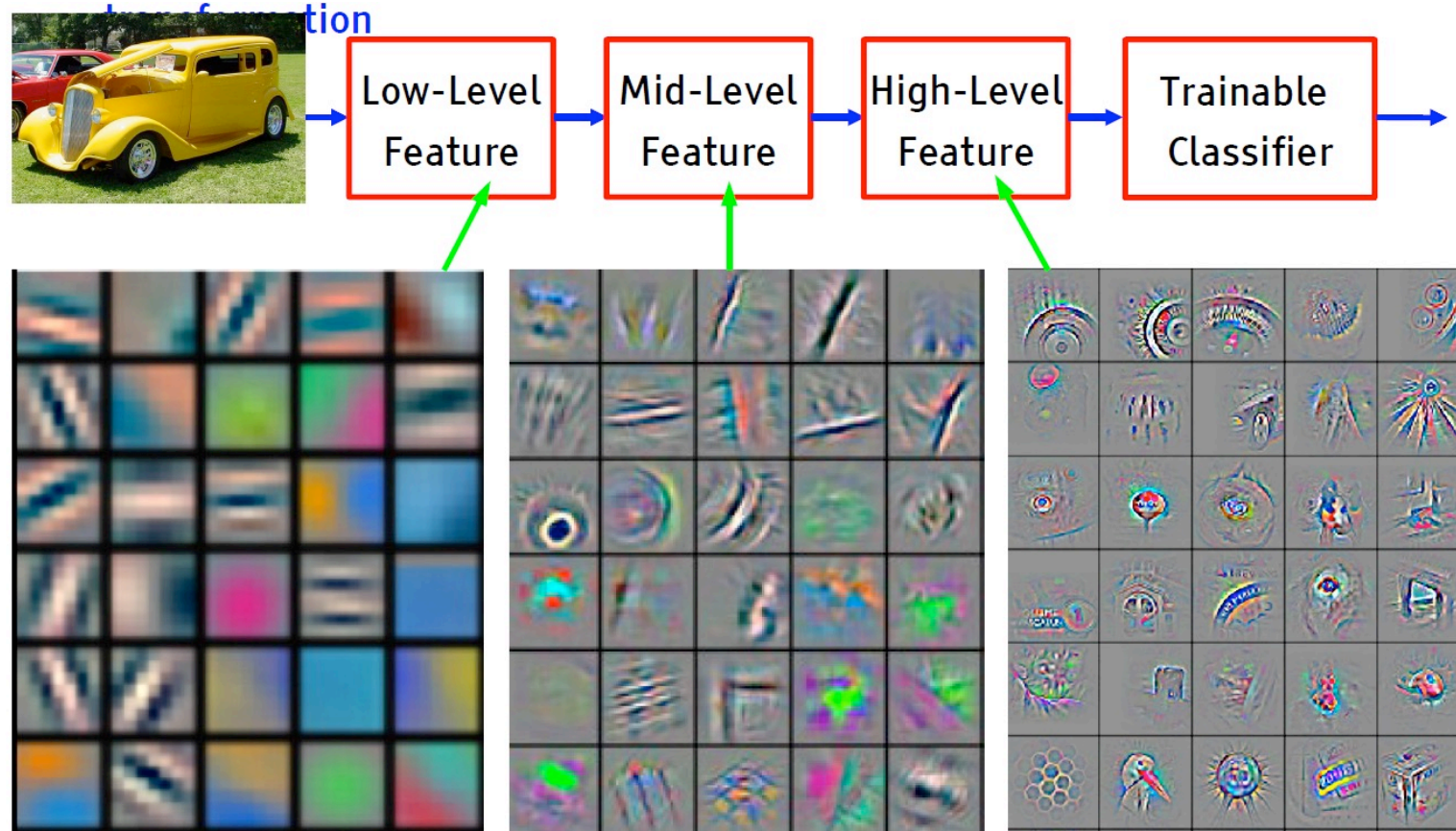


Trainable Classifier  
e.g., SVM, LR



Domain expert

# Deep Learning= Learning Hierarchical representations



(Figure from LeCun)

# Hierarchical representations with increasing level of abstraction

- Image recognition
  - Pixel -> edge -> texture-> motif -> part -> object
- Speech
  - Sample -> spectral band -> sound -> phone -> word...
- Text
  - Character -> word -> phrase->clause-> sentence  
->paragraph-> document

# Outline

- Network Components
  - Neurons (Hidden Units)
  - Output units
  - Cost functions
- Architecture design
  - Capacity of neural networks
- Training
  - Backpropagation with stochastic gradient descent

# Neuron: Nonlinear Functions

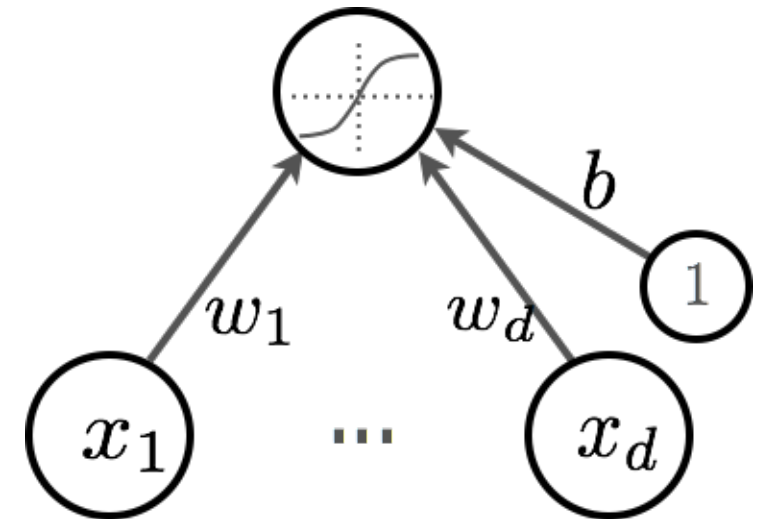
- Input: linear combination:

$$a(\mathbf{x}) = b + \sum_i w_i x_i = \mathbf{w}^T \mathbf{x} + b$$

- Output: nonlinear transformation:

$$h(\mathbf{x}) = g(a(\mathbf{x})) = g(\mathbf{w}^T \mathbf{x} + b)$$

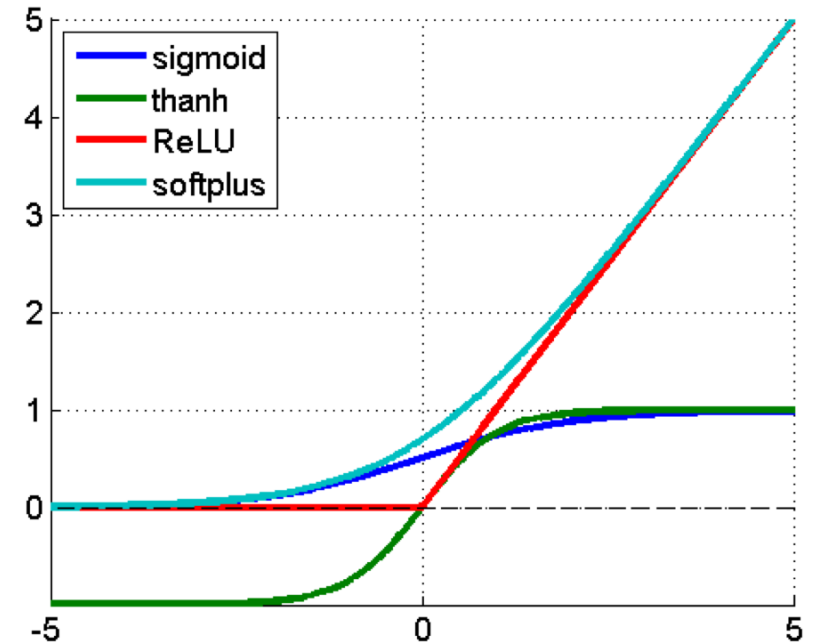
- $\mathbf{w}$ : are the weights (parameters)
- $b$  is the bias term
- $g(\cdot)$  is called the activation function





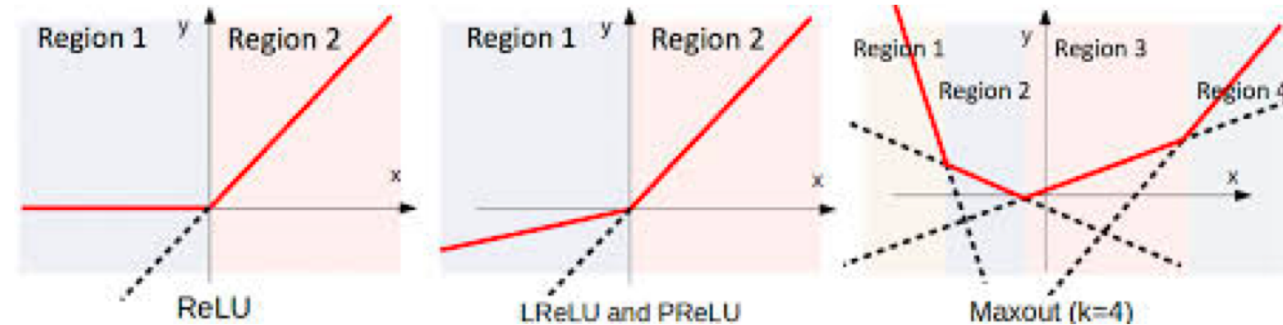
# Activation functions/Hidden Units

- Sigmoid function
  - $g(x) = 1/(1+\exp(-x))$
  - Map the input to (0,1)
- Tanh function
  - $g(x) = (1-\exp(-2x))/(1+\exp(-2x))$
  - Map the input to (-1,1)
- Rectified linear (ReLU) function
  - $g(x) = \max(0,x)$
  - No upper bounded



# Other activation functions

- Leaky ReLU (Maas et al. 2013)
  - $g(x) = \max(0, x) + \alpha \min(0, x)$
  - Fix  $\alpha$  to a small value, e.g., 0.01
- Parametric ReLU (He et al. 2015)
  - Treat  $\alpha$  as a parameter to learn
- Maxout units (Goodfellow et al. ,2013)
  - Generalize rectified linear units
  - Divide the output units into groups of  $k$  values, and output the maximum value in each group
  - Provides a way of learning a piecewise linear function that responds to multiple directions in the input  $x$  space.



# One Hidden layer Neural Networks

- Input of the hidden layer:

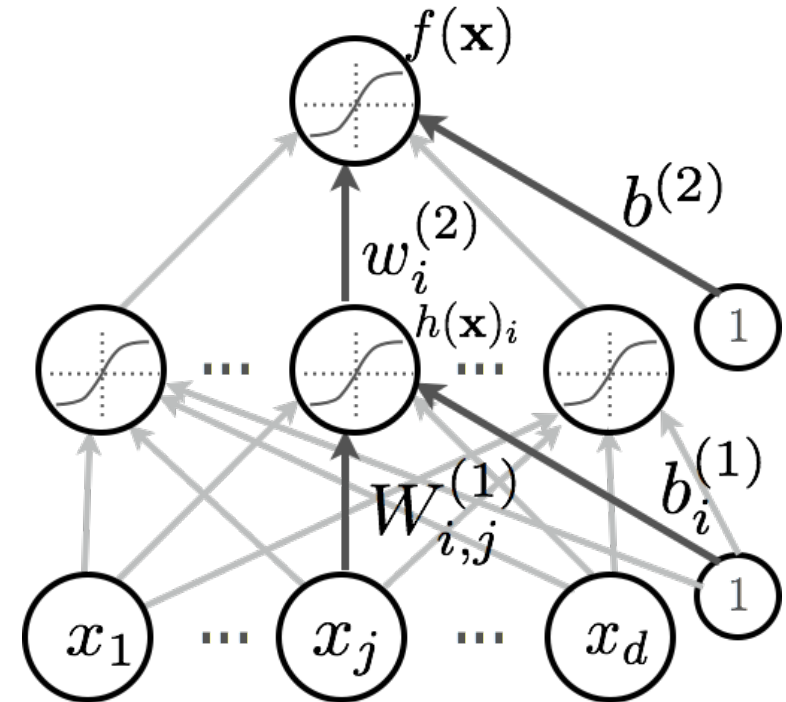
$$a(\mathbf{x}) = \mathbf{W}^T \mathbf{x} + \mathbf{b}$$

- Nonlinear transformation:

$$h(\mathbf{x}) = g_1(a(\mathbf{x}))$$

- Output layer

$$f(\mathbf{x}) = o(h(\mathbf{x}))$$



# Outline

- Network Components
  - Neurons (Hidden Units)
  - Output units
  - Cost functions
- Architecture design
  - Capacity of neural networks
- Training
  - Backpropagation with stochastic gradient descent

# Linear Units for Gaussian Output Distributions

- Given the hidden units  $\mathbf{h}$ , a layer of linear output units produces  $\hat{\mathbf{y}} = \mathbf{W}^T \mathbf{h} + \mathbf{b}$
- Linear output layers are often used to produce the mean of a conditional Gaussian distribution

$$p(\mathbf{y}|\mathbf{x}) = \text{N}(\mathbf{y}|\hat{\mathbf{y}}, \mathbf{I})$$

# Sigmoid Units for Bernoulli Output Distributions

- Bernoulli output distributions: binary classification
- The goal is to define  $p(y = 1|\mathbf{x})$ , which can be defined as follows:

$$p(y = 1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{h} + b)$$

# Softmax Units for Multinomial Output Distributions

- Bernoulli output distributions: multi-class classification
- First, define a linear layer to predict the unnormalized log probabilities of softmax:

$$\mathbf{z} = \mathbf{W}^T \mathbf{h} + \mathbf{b}$$

- where  $z_i = \log p(y = i | \mathbf{x})$ . Formally, the softmax function is given by
- 

$$p(y = i | \mathbf{x}) = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$

# Multilayer Neural Networks

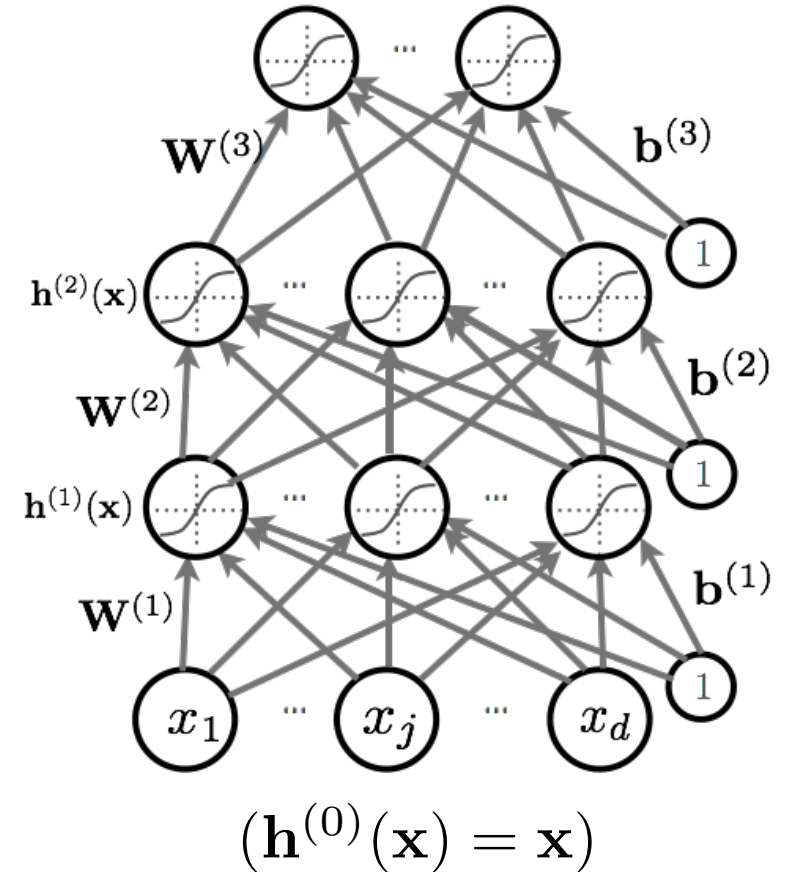
- Neural network with multiple hidden layers
- The output of previous layer as the input of next layer: ( $k=1\dots, L$ )

$$\mathbf{a}^{(k)}(\mathbf{x}) = \mathbf{b}^{(k)} + \mathbf{W}^{(k)}\mathbf{h}^{(k-1)}(\mathbf{x})$$

$$\mathbf{h}^{(k)}(\mathbf{x}) = \mathbf{g}(\mathbf{a}^{(k)}(\mathbf{x}))$$

- Final output layer

$$\mathbf{h}^{(L+1)}(\mathbf{x}) = \mathbf{o}(\mathbf{a}^{(L+1)}(\mathbf{x})) = \mathbf{f}(\mathbf{x})$$





# Outline

- Network Components
  - Neurons (Hidden Units)
  - Output units
  - Cost function
- Architecture design
  - Capacity of neural networks
- Training
  - Backpropagation with stochastic gradient descent

# Maximum Likelihood

- Most of the time, neural networks are used to define a distribution  $p(y^t | \mathbf{x}^t; \boldsymbol{\theta})$ . Therefore, the overall objective is defined as:

$$\operatorname{argmax}_{\boldsymbol{\theta}} \frac{1}{T} \sum_t \log p(y^t | \mathbf{x}^t; \boldsymbol{\theta}) + \lambda \Omega(\boldsymbol{\theta})$$

# Outline

- Network Components
  - Neurons (Hidden Units)
  - Output units
  - Cost functions
- Architecture design
  - Capacity of neural networks
- Training
  - Backpropagation with stochastic gradient descent

# Universal Approximation

- Universal Approximation Theorem (Hornik, 1991)
  - “a single hidden layer neural network with a linear output unit can approximate any continuous function arbitrary well, given enough hidden units”
- However, we may not be able to find the right parameters ....
  - The layer may be infeasibly large
  - Optimizing neural networks is difficult ...

# Deeper Networks are preferred

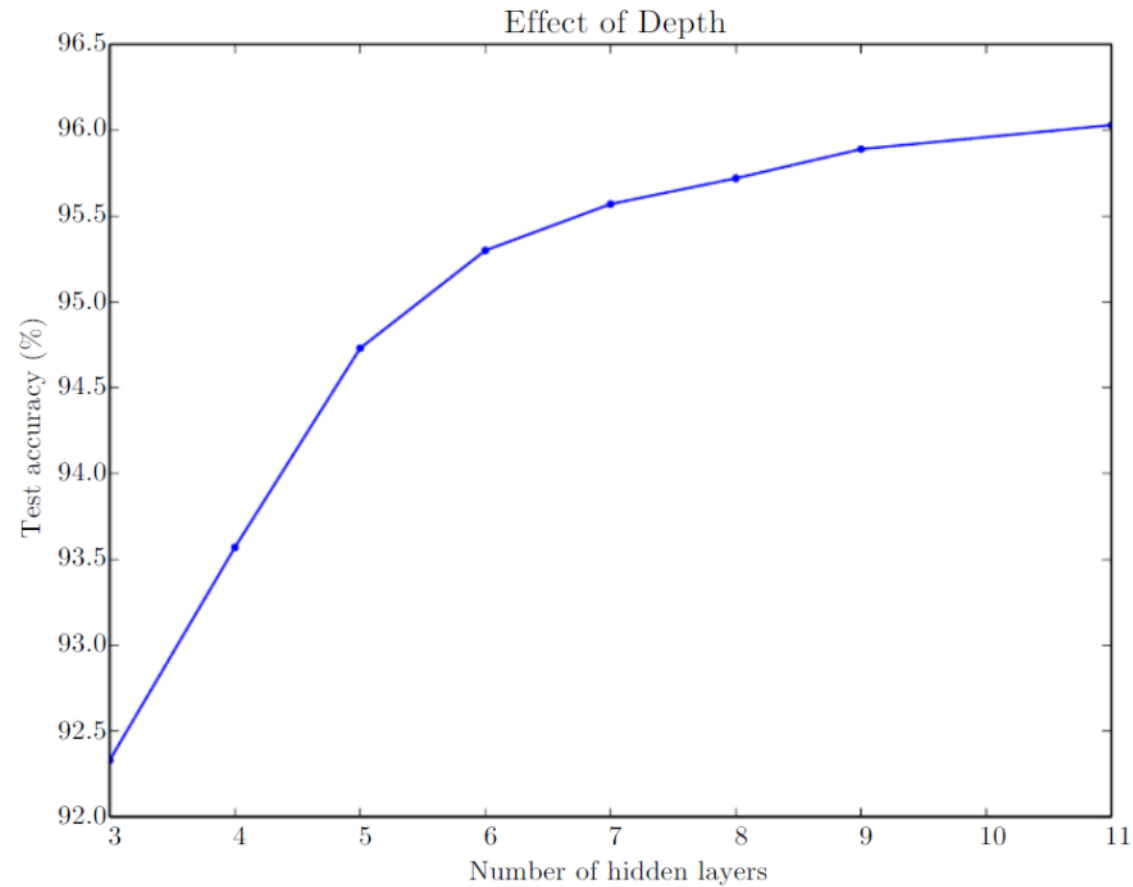
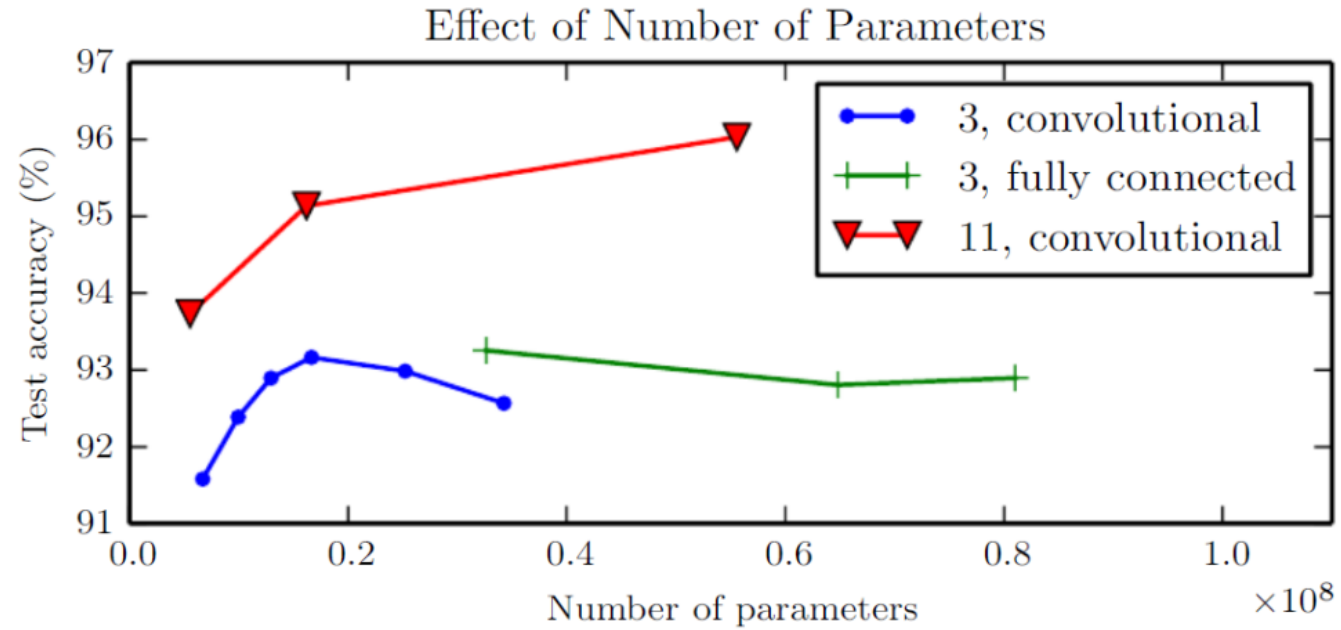


Figure: Empirical results showing that deeper networks generalize better

# Deeper Networks are preferred



**Figure:** Deeper models tend to perform better with the same number of parameter

# Deeper networks are preferred

- There exist families of functions which can be approximated efficiently with deep networks but require a much larger model for shallow networks
- Statistical reasons
  - a deep models encodes a very general belief that the function we want to learn should involve composition of several simple functions
  - Or we believe the learning problem consists of discovering different levels of variations, with the high-level ones defined on the low-level (simple) ones (e.g., Pixel -> edge -> texture -> motif -> part -> object).

# Outline

- Network Components
  - Neurons (Hidden Units)
  - Output units
  - Cost functions
- Architecture design
  - Capacity of neural networks
- Training
  - Backpropagation with stochastic gradient descent

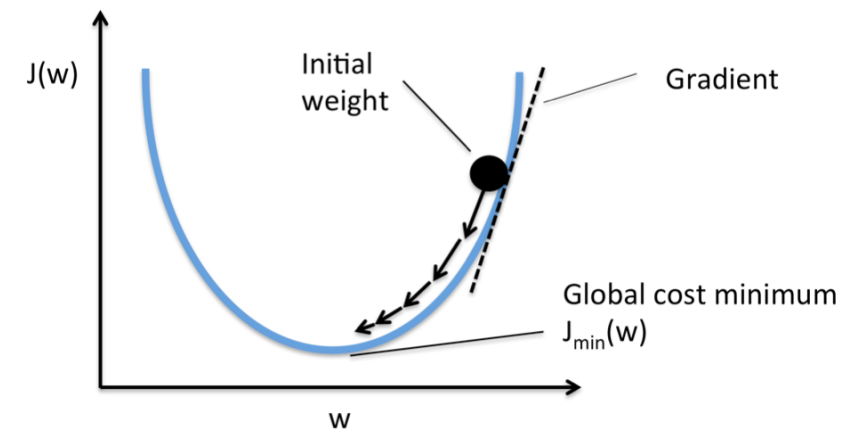


# Backpropagation with Stochastic Gradient Descent

- Gradient descent:
  - Update the parameters in the direction of gradients
  - Need to iterate over all the examples for every update
- Stochastic gradient descent
  - Perform updates after seeing each example
  - Initialize:  $\theta \equiv \{\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \dots, \mathbf{W}^{(L+1)}, \mathbf{b}^{(L+1)}\}$
  - For  $t=1:T$ 
    - for each training example  $(\mathbf{x}^{(t)}, y^{(t)})$

$$\Delta = -\nabla_{\theta} l(f(\mathbf{x}^{(t)}; \theta), y^{(t)}) - \lambda \nabla_{\theta} \Omega(\theta)$$

$$\theta \leftarrow \theta + \alpha \Delta$$

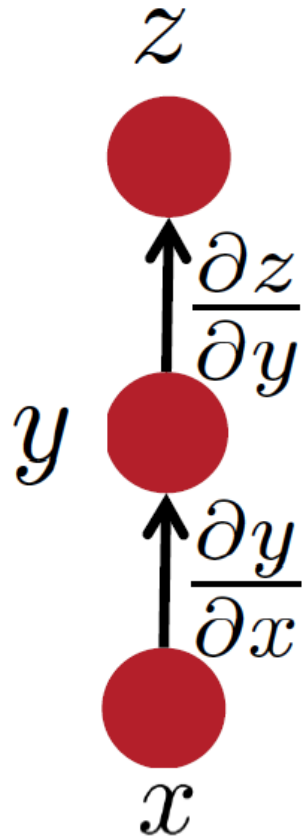


Training epoch

=

Iteration of all examples

# BackPropagation: Simple Chain Rule



$$\Delta z = \frac{\partial z}{\partial y} \Delta y$$

$$\Delta y = \frac{\partial y}{\partial x} \Delta x$$

$$\Delta z = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x} \Delta x$$

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$$

# Forward Propagation

---

**Require:** Network depth,  $l$

**Require:**  $\mathbf{W}^{(i)}, i \in \{1, \dots, l\}$ , the weight matrices of the model

**Require:**  $\mathbf{b}^{(i)}, i \in \{1, \dots, l\}$ , the bias parameters of the model

**Require:**  $\mathbf{x}$ , the input to process

**Require:**  $\mathbf{y}$ , the target output

$$\mathbf{h}^{(0)} = \mathbf{x}$$

**for**  $k = 1, \dots, l$  **do**

$$\mathbf{a}^{(k)} = \mathbf{b}^{(k)} + \mathbf{W}^{(k)} \mathbf{h}^{(k-1)}$$

$$\mathbf{h}^{(k)} = f(\mathbf{a}^{(k)})$$

**end for**

$$\hat{\mathbf{y}} = \mathbf{h}^{(l)}$$

$$J = L(\hat{\mathbf{y}}, \mathbf{y}) + \lambda \Omega(\theta)$$

---

# Backward Propagation

---

After the forward computation, compute the gradient on the output layer:

$$\mathbf{g} \leftarrow \nabla_{\hat{\mathbf{y}}} J = \nabla_{\hat{\mathbf{y}}} L(\hat{\mathbf{y}}, y)$$

**for**  $k = l, l - 1, \dots, 1$  **do**

Convert the gradient on the layer's output into a gradient into the pre-nonlinearity activation (element-wise multiplication if  $f$  is element-wise):

$$\mathbf{g} \leftarrow \nabla_{\mathbf{a}^{(k)}} J = \mathbf{g} \odot f'(\mathbf{a}^{(k)})$$

Compute gradients on weights and biases (including the regularization term, where needed):

$$\nabla_{\mathbf{b}^{(k)}} J = \mathbf{g} + \lambda \nabla_{\mathbf{b}^{(k)}} \Omega(\theta)$$

$$\nabla_{\mathbf{W}^{(k)}} J = \mathbf{g} \mathbf{h}^{(k-1)\top} + \lambda \nabla_{\mathbf{W}^{(k)}} \Omega(\theta)$$

Propagate the gradients w.r.t. the next lower-level hidden layer's activations:

$$\mathbf{g} \leftarrow \nabla_{\mathbf{h}^{(k-1)}} J = \mathbf{W}^{(k)\top} \mathbf{g}$$

**end for**

---

# References

- Deep Learning book, chap 6.