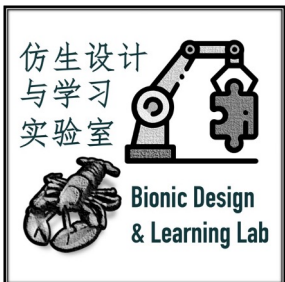


# Lecture 06

## Deep Networks I

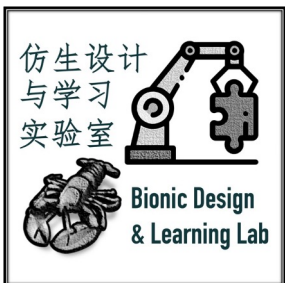


AncoraSIR.com

[Please refer to the course website for copyright credits]



# Neural Network



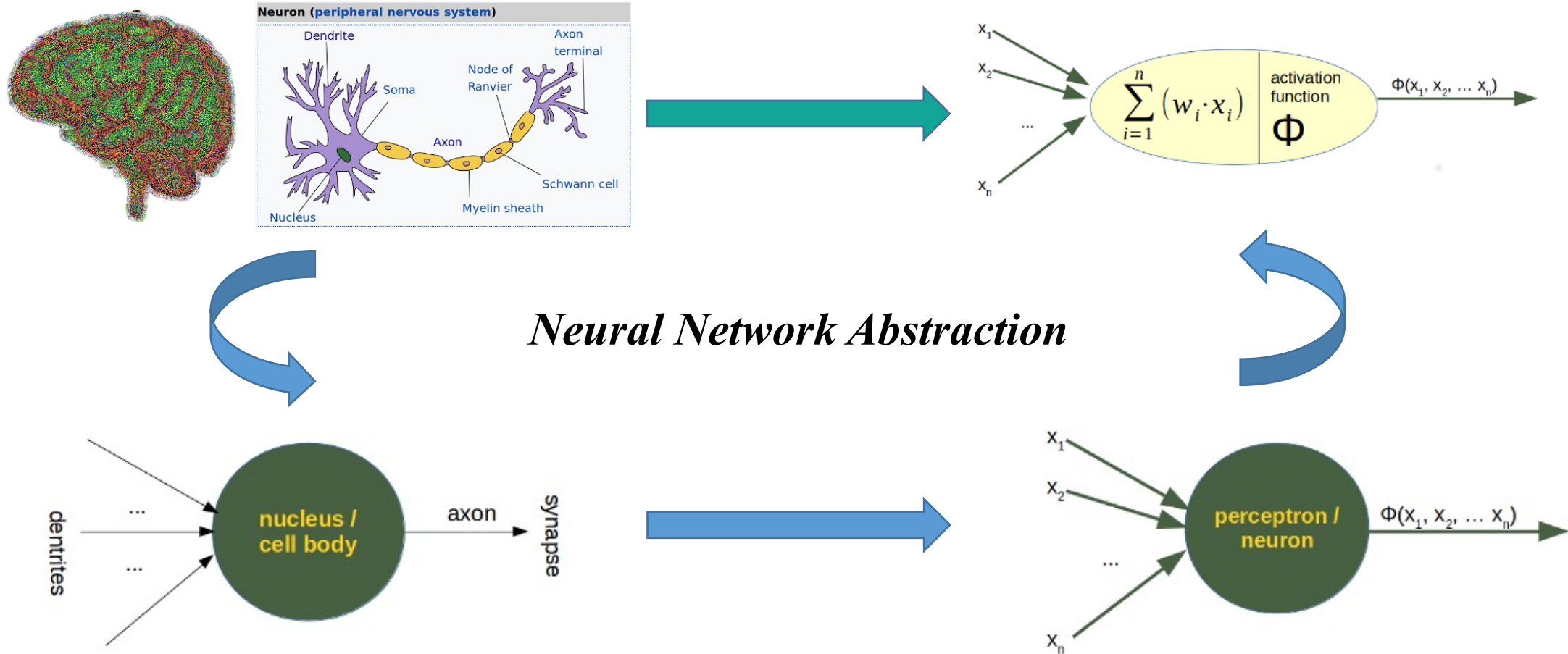
[AncoraSIR.com](http://AncoraSIR.com)



**SUSTech**  
Southern University  
of Science and Technology

# What is a Neural Network?

*From biological inspiration to mathematical modeling*



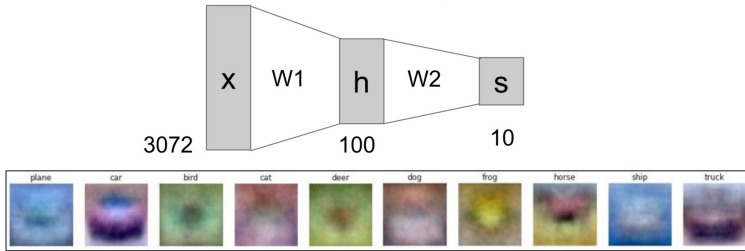
# A Perceptron as an Artificial Neuron

Linear score function:

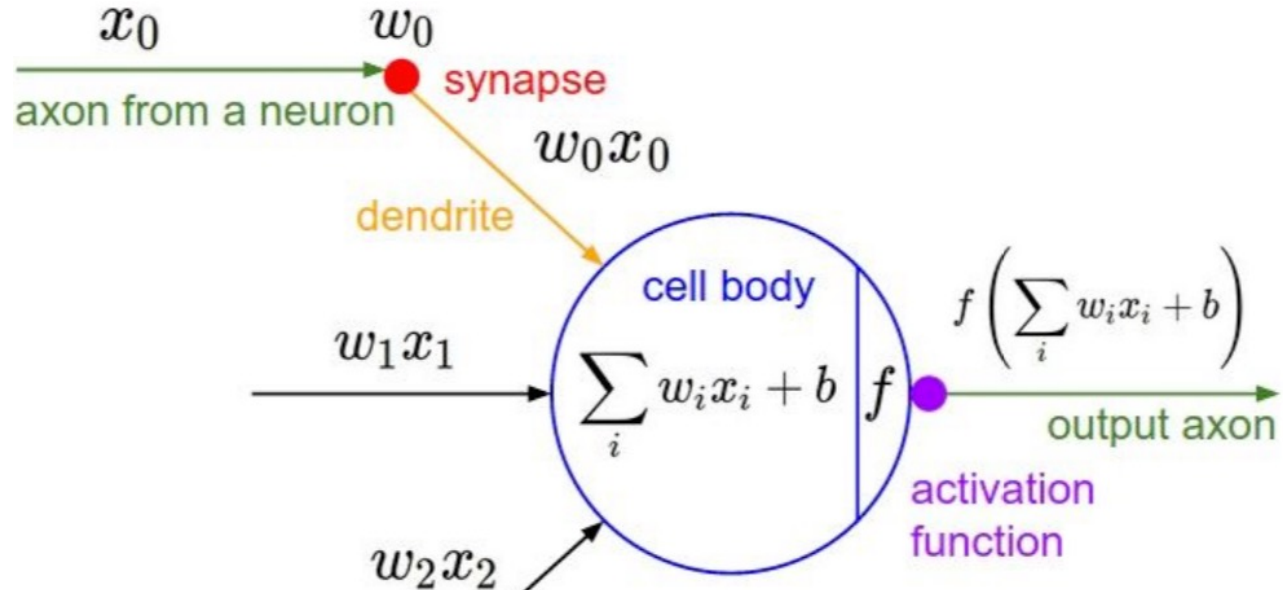
$$f = Wx$$

2-layer Neural Network

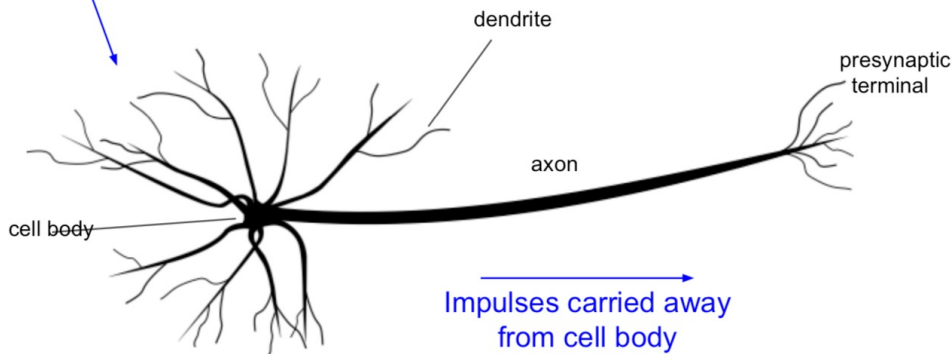
$$f = W_2 \max(0, W_1 x)$$



## Biological Inspiration



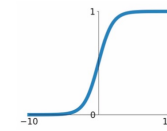
Impulses carried toward cell body



This image by Felipe Peruchio is licensed under [CC-BY 3.0](https://creativecommons.org/licenses/by/3.0/)

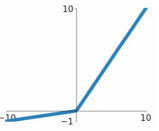
**Sigmoid**

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



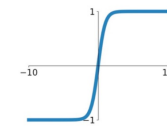
**Leaky ReLU**

$$\max(0.1x, x)$$



**tanh**

$$\tanh(x)$$

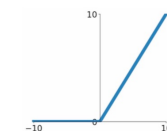


**Maxout**

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

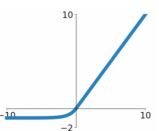
**ReLU**

$$\max(0, x)$$



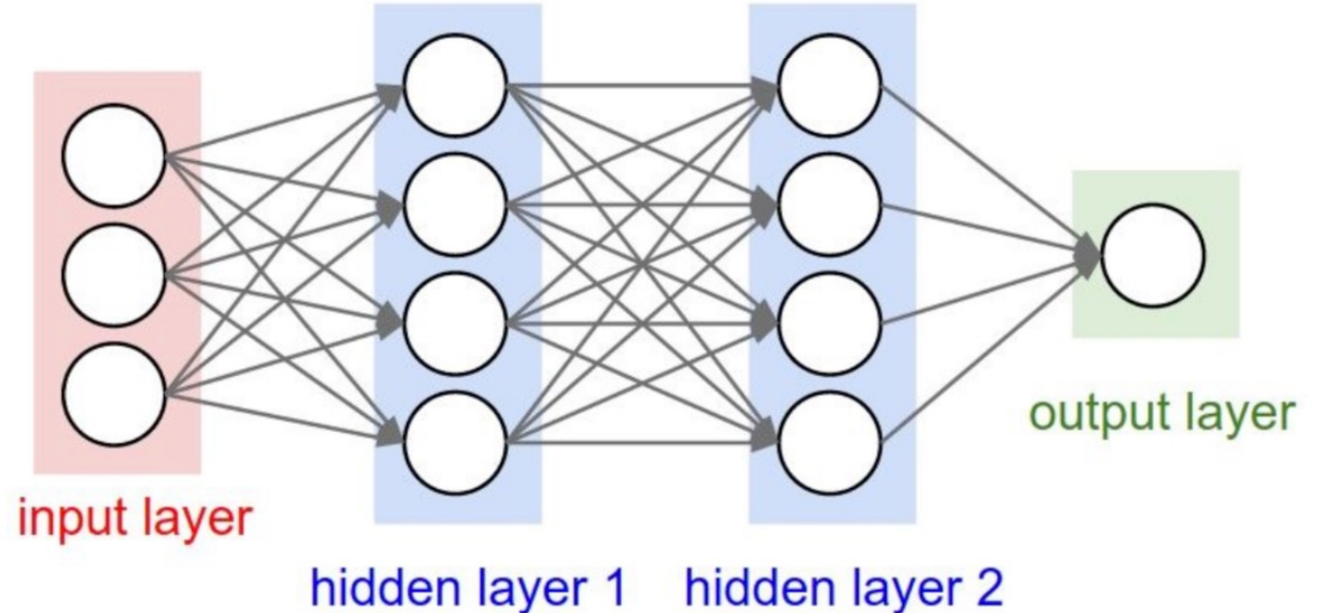
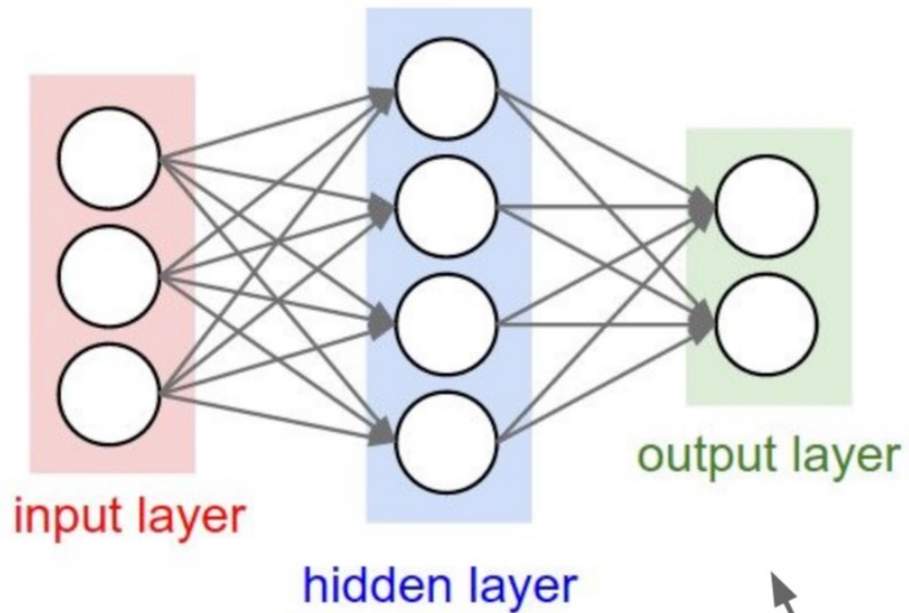
**ELU**

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



# Multi-Layer Perceptrons

## *Artificial Neural Networks*



“2-layer Neural Net”, or  
“1-hidden-layer Neural Net”

“3-layer Neural Net”, or  
“2-hidden-layer Neural Net”

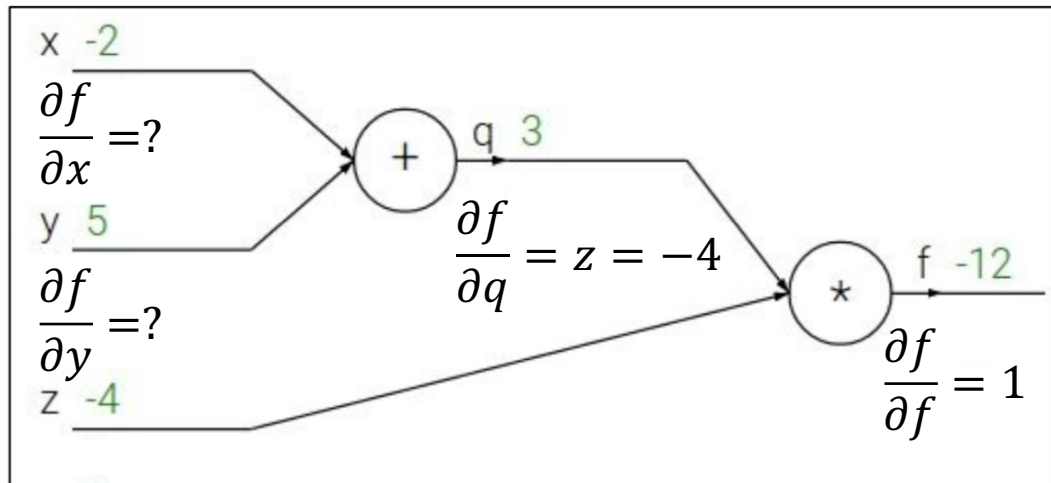
“Fully-connected” layers

# Computation Graph

*A simple example with backpropagation*

$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$



$$\frac{\partial f}{\partial z} = q = 3$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \quad \frac{\partial f}{\partial z} = q$$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \quad \frac{\partial q}{\partial y} = 1$$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

Chain rule:

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

Upstream gradient

Local gradient

Chain rule:

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y}$$

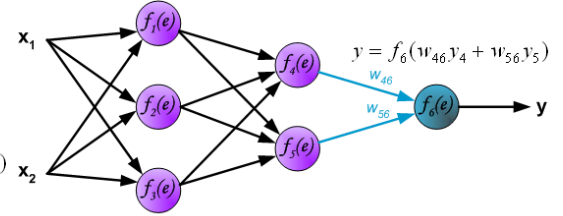
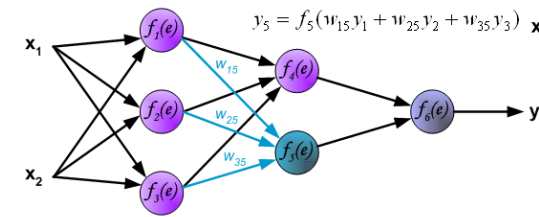
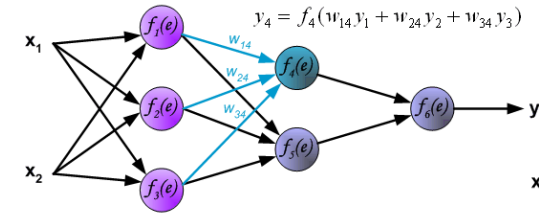
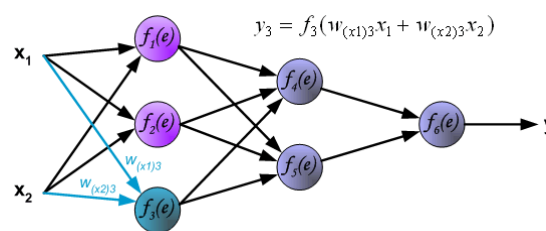
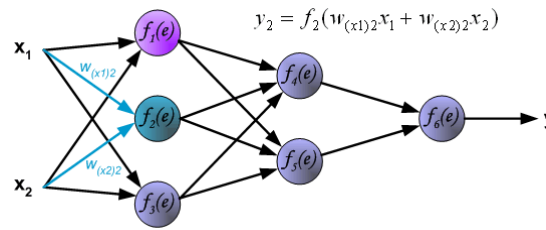
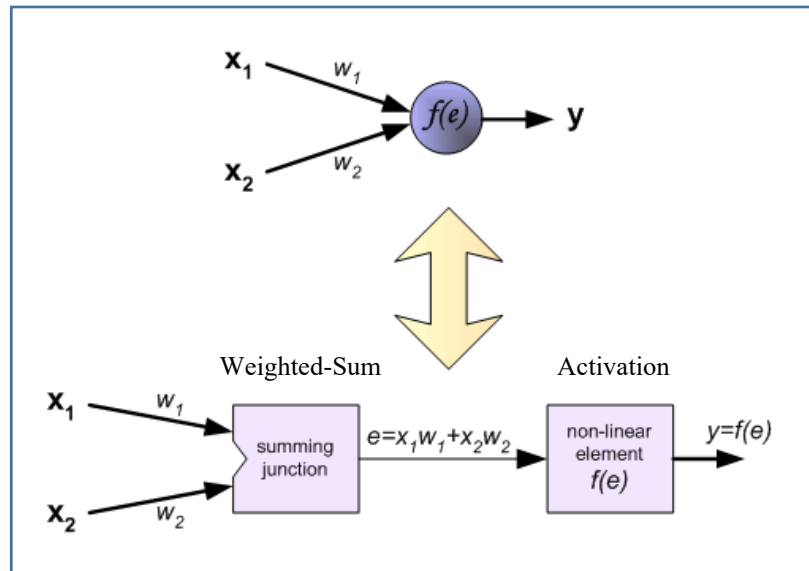
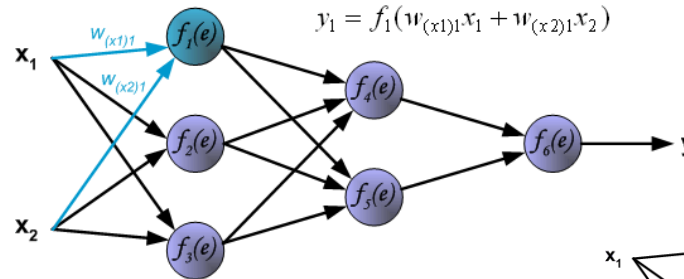
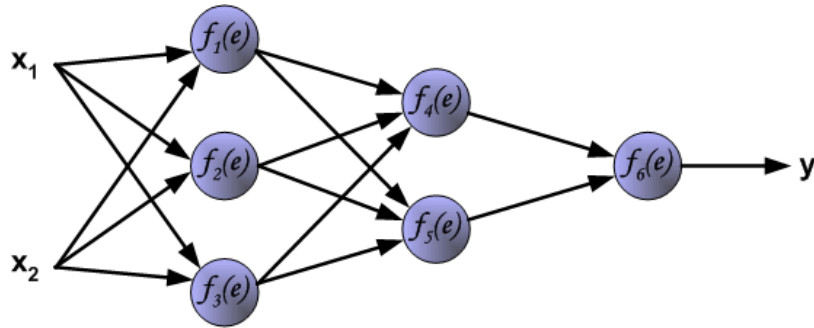
Upstream gradient

Local gradient



# Forward Propagation

*Accept inputs to train a Multi-layer Neural Network*



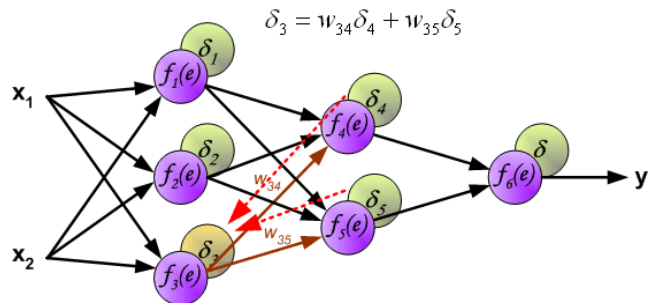
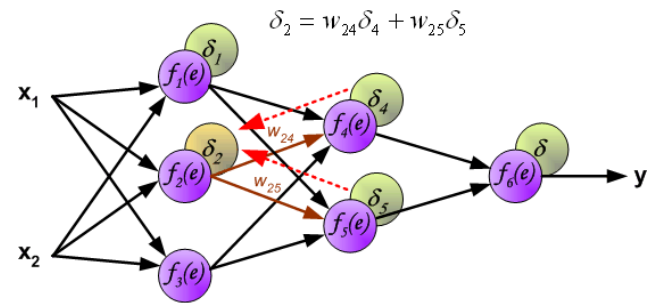
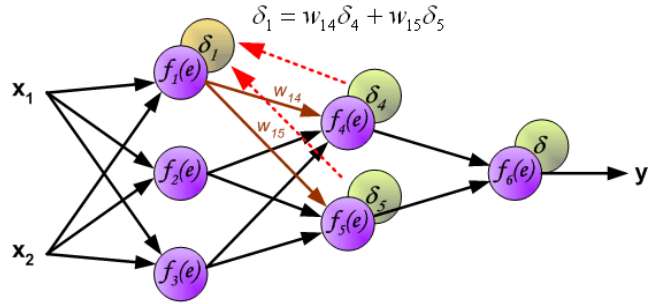
[http://galaxy.agh.edu.pl/%7Evlsi/AI/backp\\_t\\_en/backprop.html](http://galaxy.agh.edu.pl/%7Evlsi/AI/backp_t_en/backprop.html)



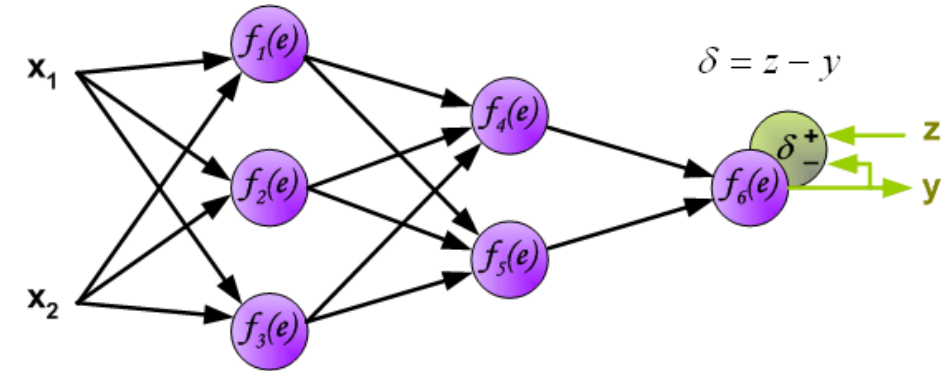
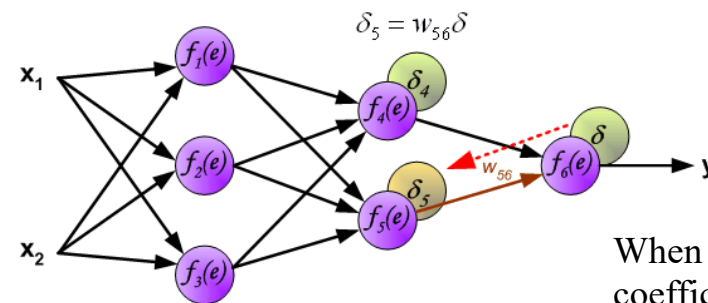
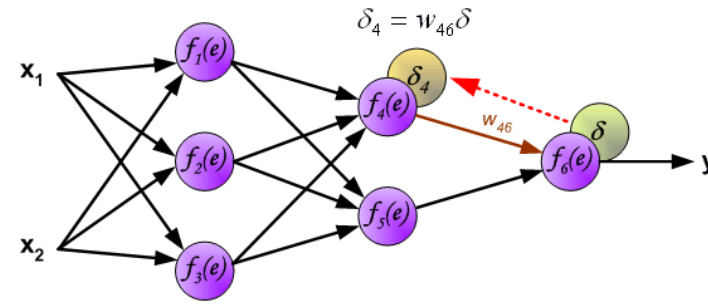
SUSTech  
Southern University  
of Science and Technology

# Backward Propagation

*Calculate the prediction error node-by-node*



The idea is to propagate error signal  $d$  (computed in single teaching step) back to all neurons, which output signals were input for discussed neuron.



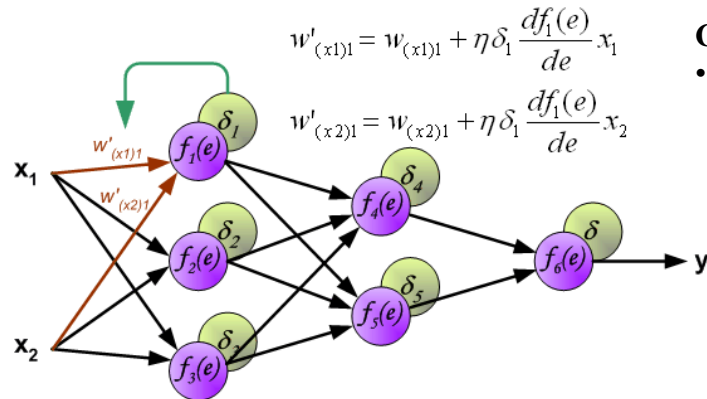
When the error signal for each neuron is computed, the weights coefficients of each neuron input node may be modified.

In formulas on the right,  $df(e)/de$  represents derivative of neuron activation function (which weights are modified).



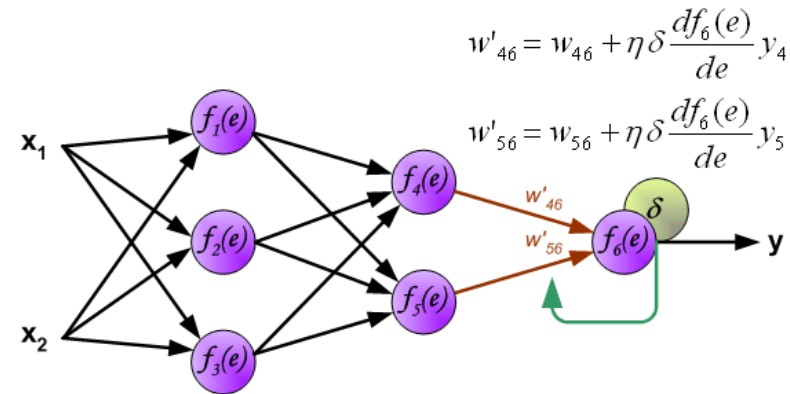
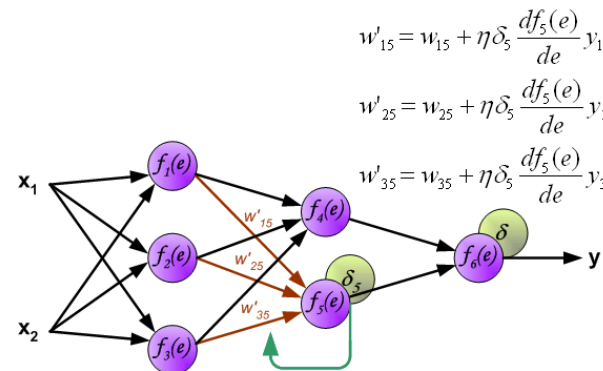
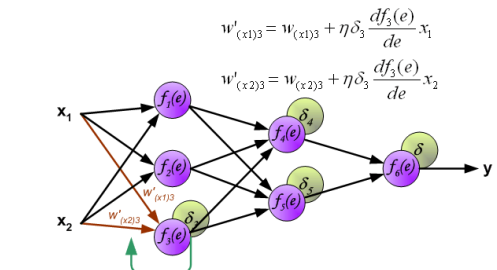
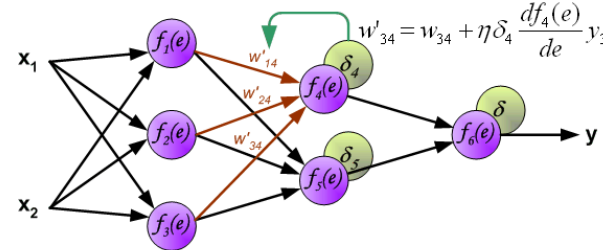
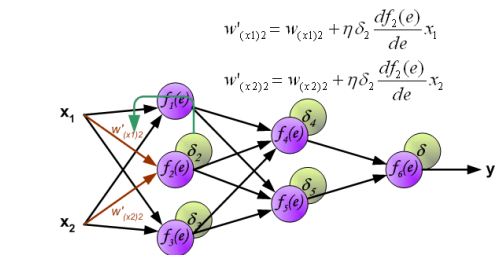
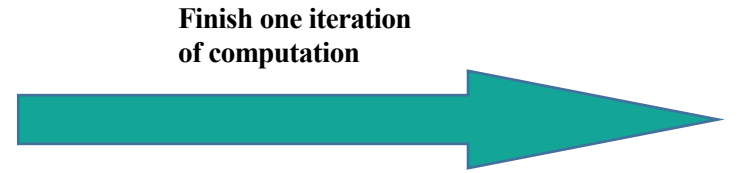
# Weight Update

*Update the weights to finish one iteration of computation, then repeat.*



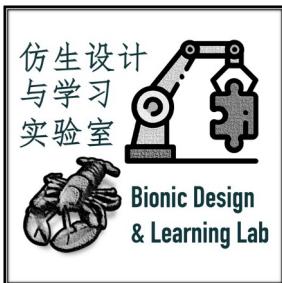
**Coefficient  $\eta$  affects network teaching speed**, to select this parameter:

- The **first** method is to start teaching process with large value of the parameter. While weights coefficients are being established the parameter is being decreased gradually.



- The **second**, more complicated, method starts teaching with small parameter value. During the teaching process the parameter is being increased when the teaching is advanced and then decreased again in the final stage. Starting teaching process with low parameter value enables to determine weights coefficients signs.

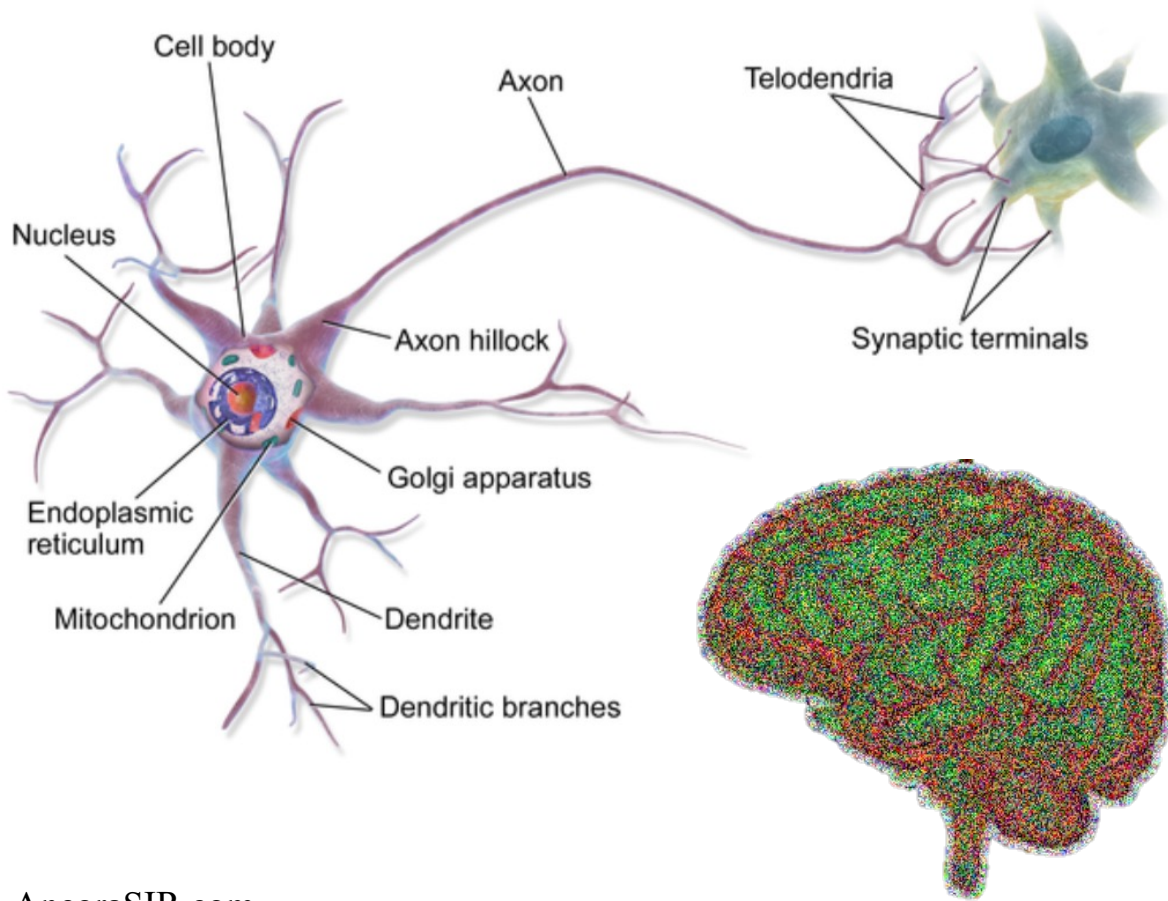
# Deep Forward Networks



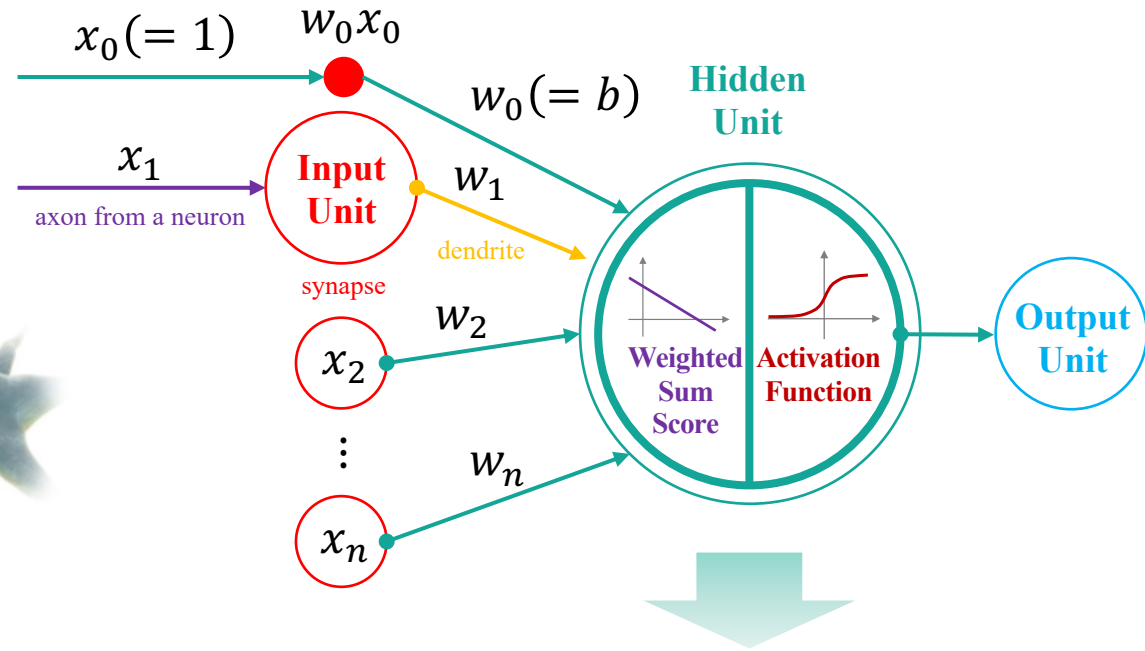
[AncoraSIR.com](http://AncoraSIR.com)

# From a Neuron to a Perceptron then a Neural Network

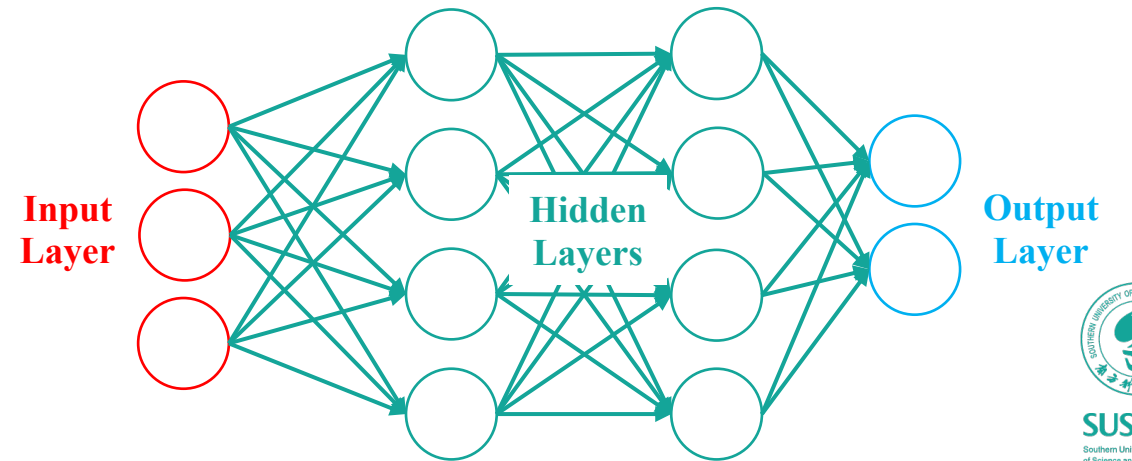
*Bio-inspired Architecture Design*



$$\hat{y} = g_{Activation}[f_{WeightedSum}(\mathbf{x})] = g(W^T \mathbf{x} + b)$$



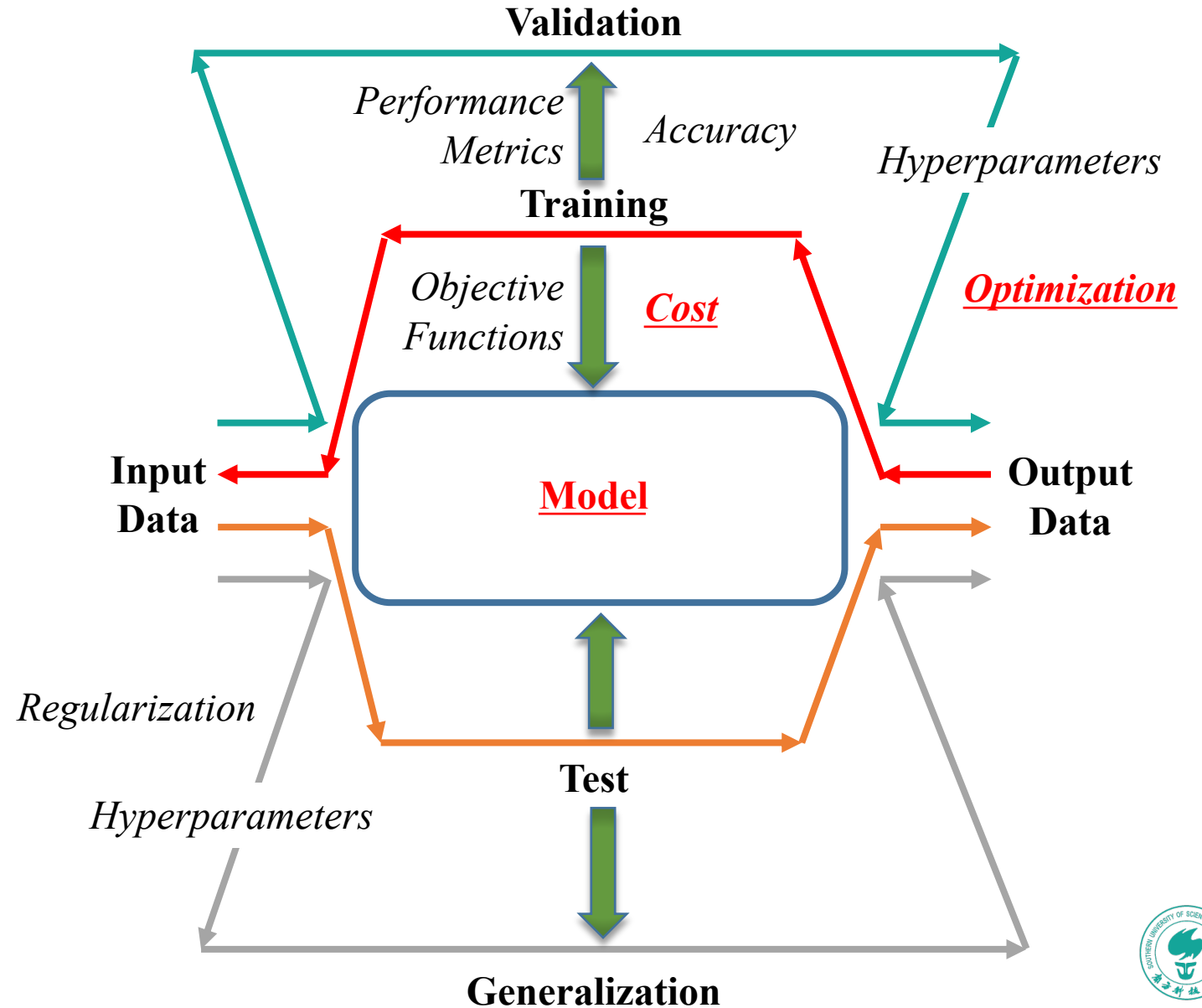
$$h^{(i)} = g^{(i)}(W^{(i)T} \mathbf{x} + b^{(i)}) \quad \hat{y} = g(W^T h^{(i)} + b)$$



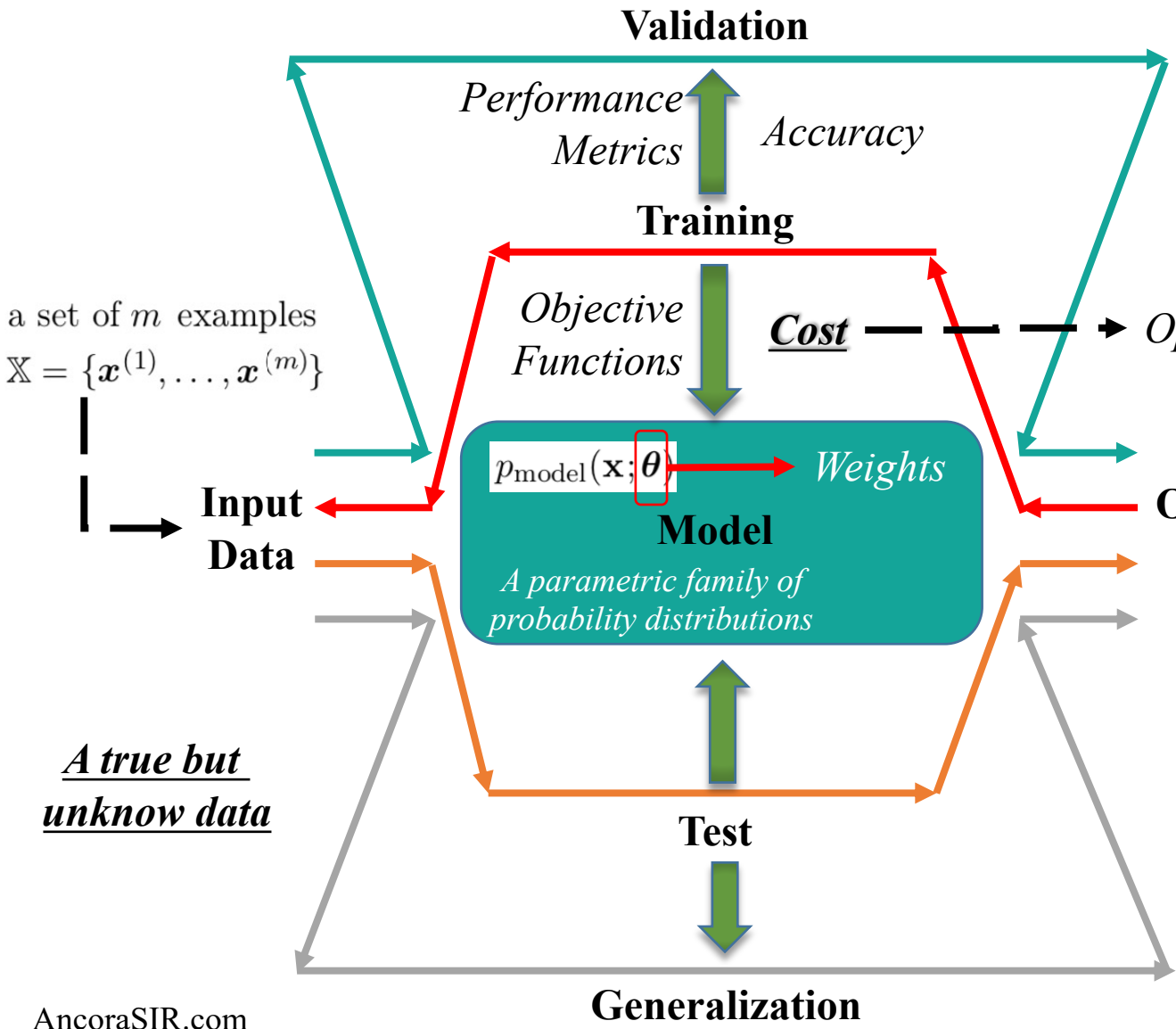
# Gradient-based Learning

## *ML vs NN*

- For supervised learning
  - NN can be viewed as ML with *gradient descent*
  - an optimization procedure
  - a cost function
  - a model family
- Difference
  - The *nonlinearity* of a neural network causes most interesting loss functions to become *non-convex*
  - Neural networks are usually trained by using *iterative, gradient-based optimizers* that merely drive the cost function to a very low value
- Next Steps
  - Choose a cost function
  - Choose model output



# Maximum Likelihood as a Cost Function



The maximum likelihood estimator for  $\theta$

$$\theta_{\text{ML}} = \arg \max_{\theta} p_{\text{model}}(\mathbb{X}; \theta)$$

$$= \arg \max_{\theta} \prod_{i=1}^m p_{\text{model}}(\mathbf{x}^{(i)}; \theta)$$

$\log_b(xy) = \log_b(x) + \log_b(y)$

$$\theta_{\text{ML}} = \arg \max_{\theta} \sum_{i=1}^m \log p_{\text{model}}(\mathbf{x}^{(i)}; \theta)$$

**To maximize**

$$\theta_{\text{ML}} = \arg \max_{\theta} \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(\mathbf{x}; \theta)$$

An empirical distribution

- Training Data, which consists of samples from  $p_{\text{data}}(\mathbf{x})$

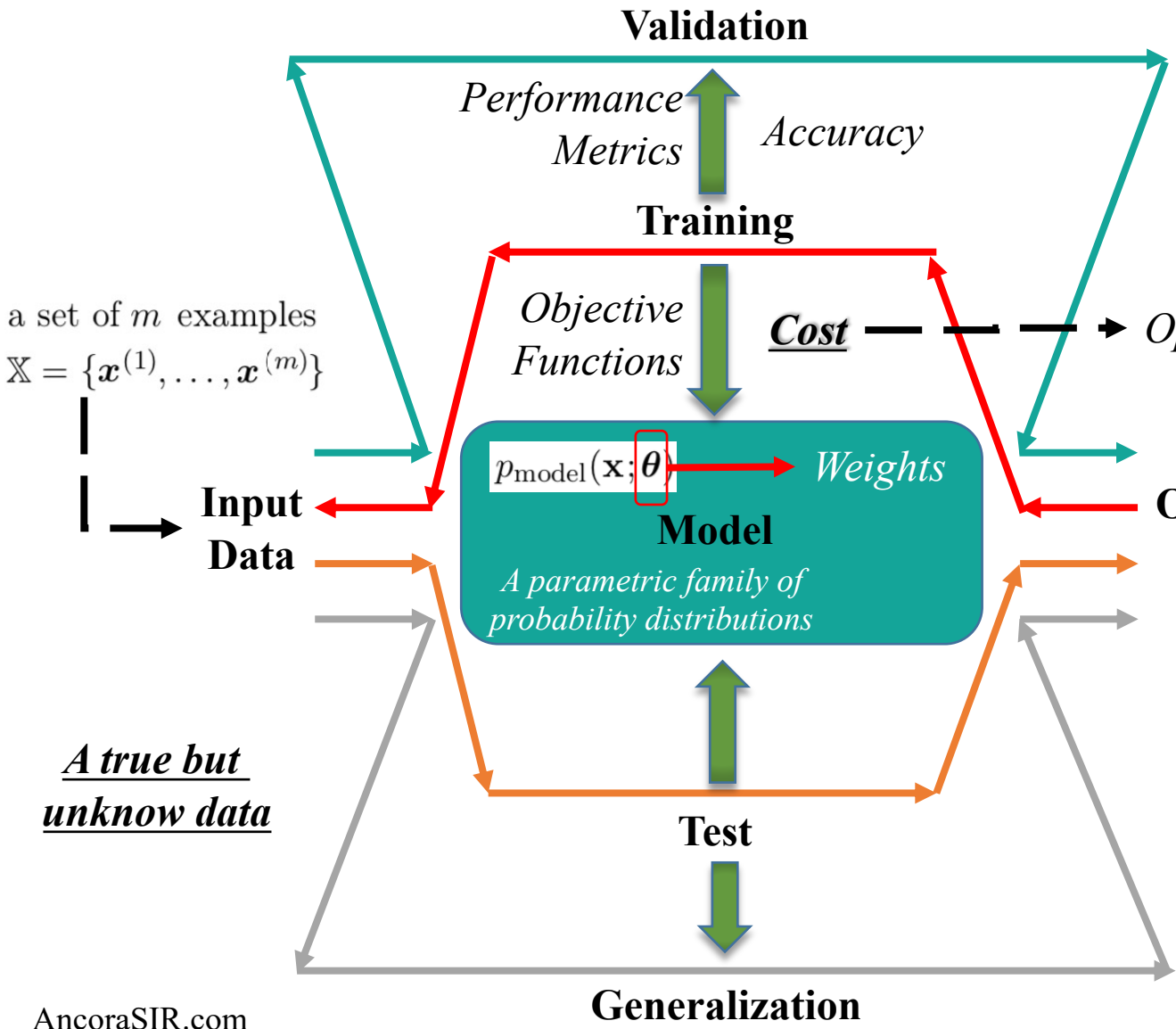
**To minimize**

KL Divergence as a Dissimilarity Metric

$$D_{\text{KL}}(\hat{p}_{\text{data}} || p_{\text{model}}) = \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log \hat{p}_{\text{data}}(\mathbf{x}) - \log p_{\text{model}}(\mathbf{x})]$$

Optimization  $\mathbb{E}[X] = \sum_{i=1}^k x_i p_i = x_1 p_1 + x_2 p_2 + \dots + x_k p_k$   
 $p_1 + p_2 + \dots + p_k = 1$

# Learning Conditional Distributions with Maximum Likelihood



$$\text{MSE}_{\text{train}} = \frac{1}{m} \sum_{i=1}^m \|\hat{y}^{(i)} - y^{(i)}\|^2$$

$$J(\theta) = \frac{1}{2} \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}} \|\mathbf{y} - f(\mathbf{x}; \theta)\|^2 + \text{const}$$

**Mean Squared Loss**

$$p_{\text{model}}(\mathbf{y} | \mathbf{x}) = \mathcal{N}(\mathbf{y}; f(\mathbf{x}; \theta), \mathbf{I})$$

**Cross-Entropy Loss**

$$J(\theta) = -\mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(\mathbf{y} | \mathbf{x})$$

*To minimize*

$$-\mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log p_{\text{model}}(\mathbf{x})]$$

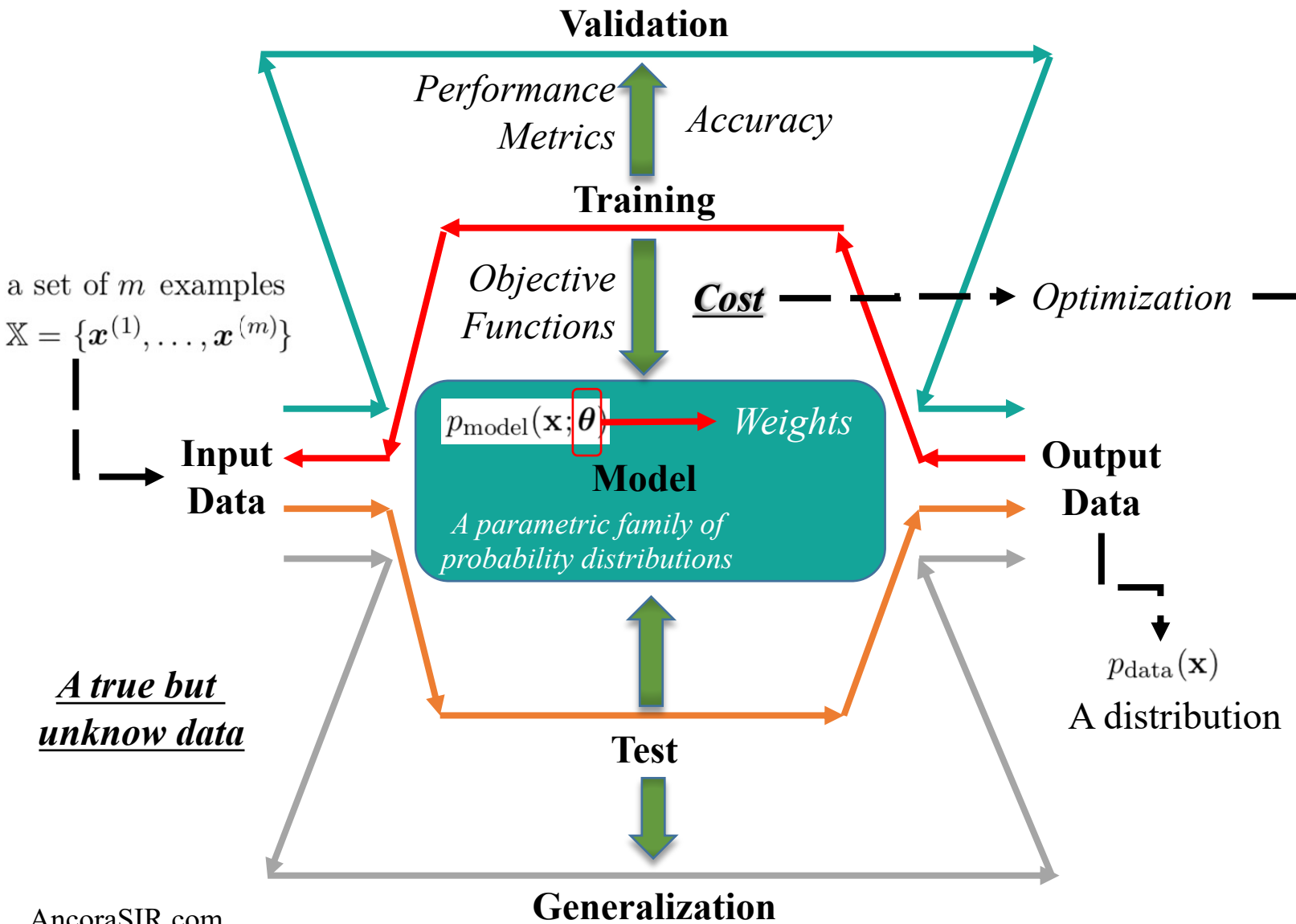
*KL Divergence as a Dissimilarity Metric*

$$D_{\text{KL}}(\hat{p}_{\text{data}} || p_{\text{model}}) = \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log \hat{p}_{\text{data}}(\mathbf{x}) - \log p_{\text{model}}(\mathbf{x})]$$



# Learning Conditional Statistics

just one conditional statistic of  $y$  given  $x$



Predicts mean value of  $y$  for each  $x$

$$f^* = \arg \min_f \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p_{\text{data}}} \|\mathbf{y} - f(\mathbf{x})\|^2$$

Cross-Entropy Loss

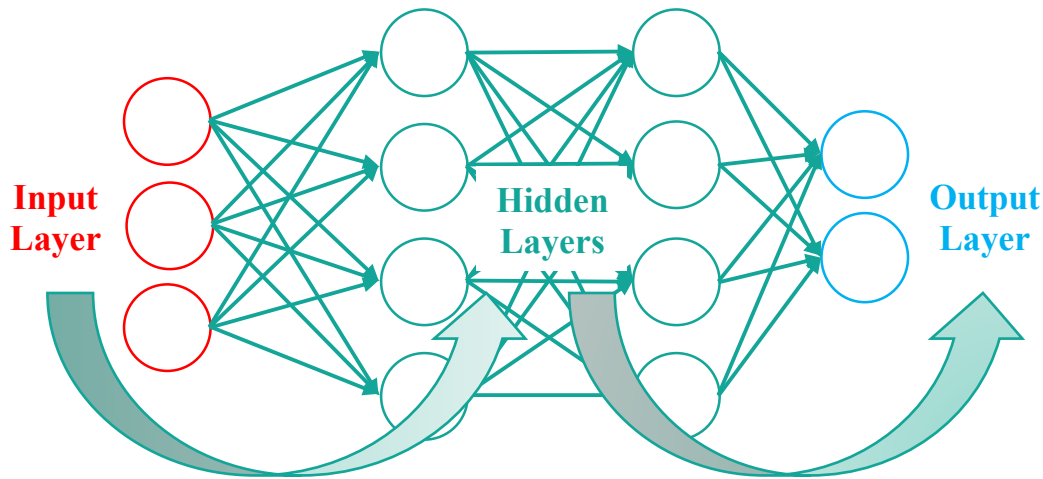
$$J(\theta) = -\mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(\mathbf{y} | \mathbf{x})$$

$$f^* = \arg \min_f \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p_{\text{data}}} \|\mathbf{y} - f(\mathbf{x})\|_1$$

Predicts median value of  $y$  for each  $x$

# Outputs Units from Hidden Layers

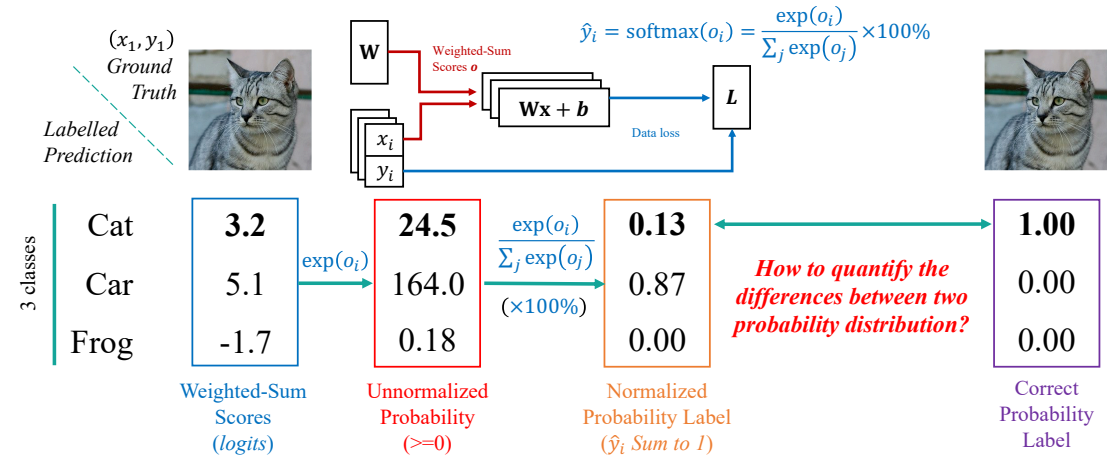
*Features (Inputs) of the Output Units provided by the Hidden Layers*



$$h = f(x; \theta)$$

$$\hat{y} = f(h)$$

Any kind of **neural network unit** that may be used as an **output** can also be used as a **hidden** unit.



Data Distribution

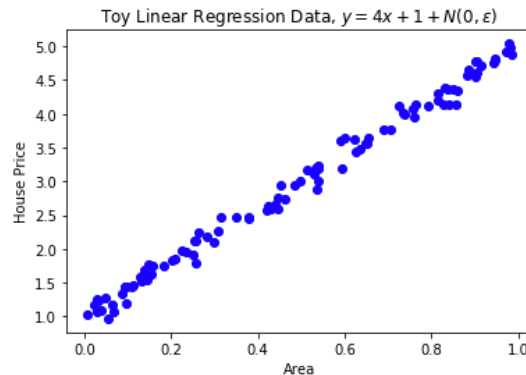
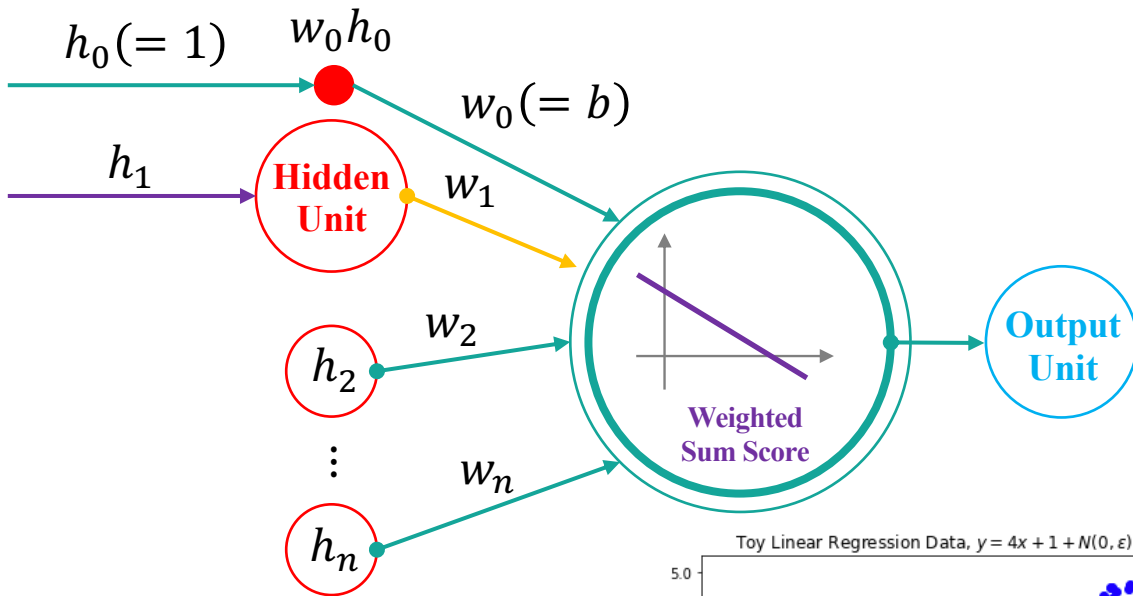
Model Distribution

**Cross-entropy** is a technique commonly used in deep neural networks

# Gaussian Output Distributions

Multiple Linear Regression as  $\hat{\mathbf{y}} = \mathbf{W}^T \mathbf{h} + \mathbf{b}$

$$\hat{\mathbf{y}} = f_{\text{weightedSum}}(\mathbf{h}) = \mathbf{W}^T \mathbf{h} + \mathbf{b}$$



- Linear Unit outputs the **mean** of a conditional Gaussian distribution
  - $p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}; \hat{\mathbf{y}}, \mathbf{I})$
- Cost Function
  - Loss function as the mean squared error

$$\frac{1}{2} \sum_{i=1}^n (y - \hat{y})^2$$

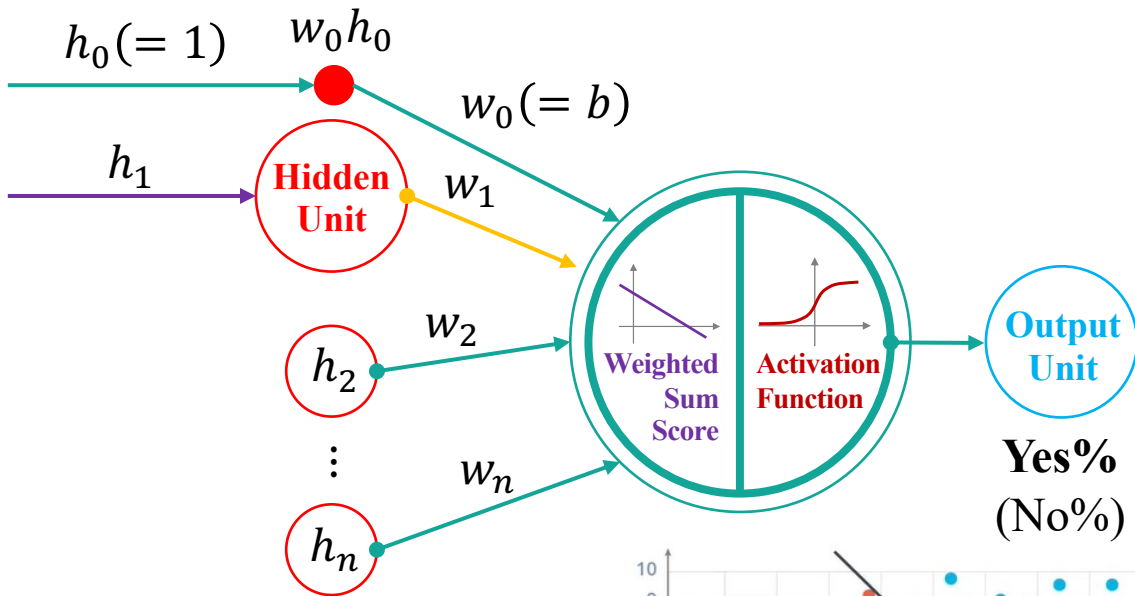
- Maximizing the log-likelihood

$$-\log p(\mathbf{y}|\mathbf{x}) = \frac{1}{2} \sum_{i=1}^n \left[ \log(2\pi\sigma^2) + \frac{1}{\sigma^2} (y - \hat{y})^2 \right]$$

# Bernoulli Output Distributions

Statistical Binary Classification as  $\hat{y} = \text{sigmoid}(\mathbf{w}^T \mathbf{h} + b)$

$$\hat{y} = g_{\text{Activation}}[f_{\text{WeightedSum}}(\mathbf{h})] = \text{sigmoid}(\mathbf{w}^T \mathbf{h} + b) \bullet \text{Outputs a Bernoulli distribution}$$

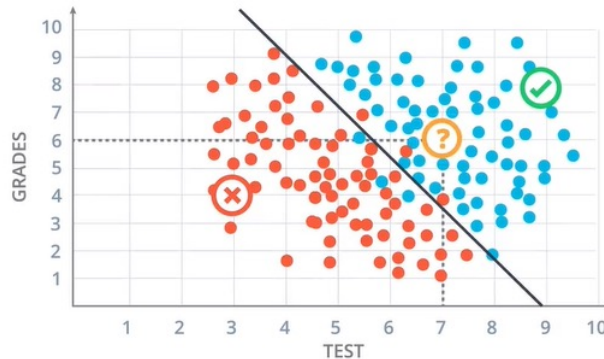


- Controlled by a sigmoidal transformation of the weighted-sum
- $P(y) = \text{sigmoid}[(2y - 1)(\mathbf{w}^T \mathbf{h} + b)]$

## • Cost Function

- Maximizing the log-likelihood

$$-\sum_{i=1}^n \sum_j y_j^{(i)} \log \hat{y}_j^{(i)}$$



# Multinoulli Output Distributions

Statistical Multi-class Classification as  $\hat{y} = \text{softmax}(W^T h + b)$

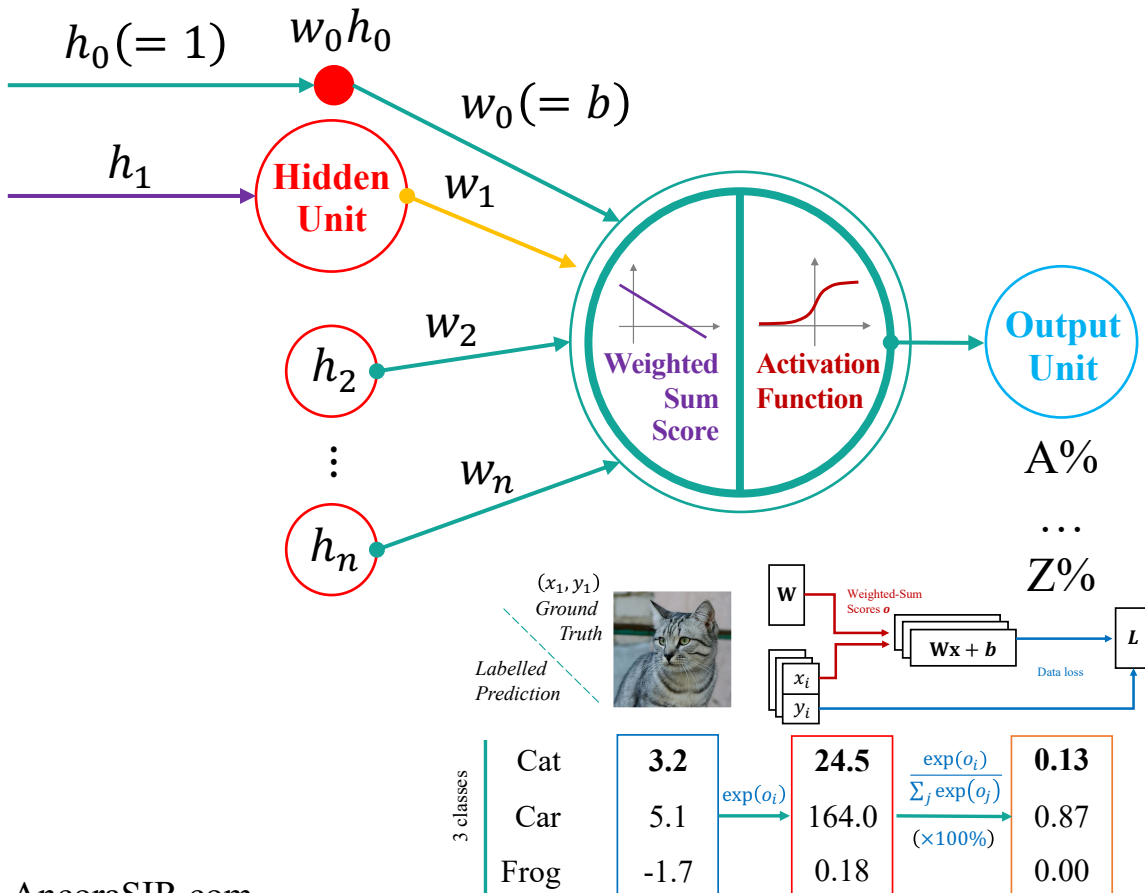
$\hat{y} = g_{\text{Activation}}[f_{\text{WeightedSum}}(h)] = \text{softmax}(W^T h + b)$  • Outputs a Multinoulli distribution

- Controlled by normalized exponentials of the weighted-sums
- $\hat{y} = \text{softmax}(W^T h + b)$

• Cost Function

- Averaged cross-entropy loss

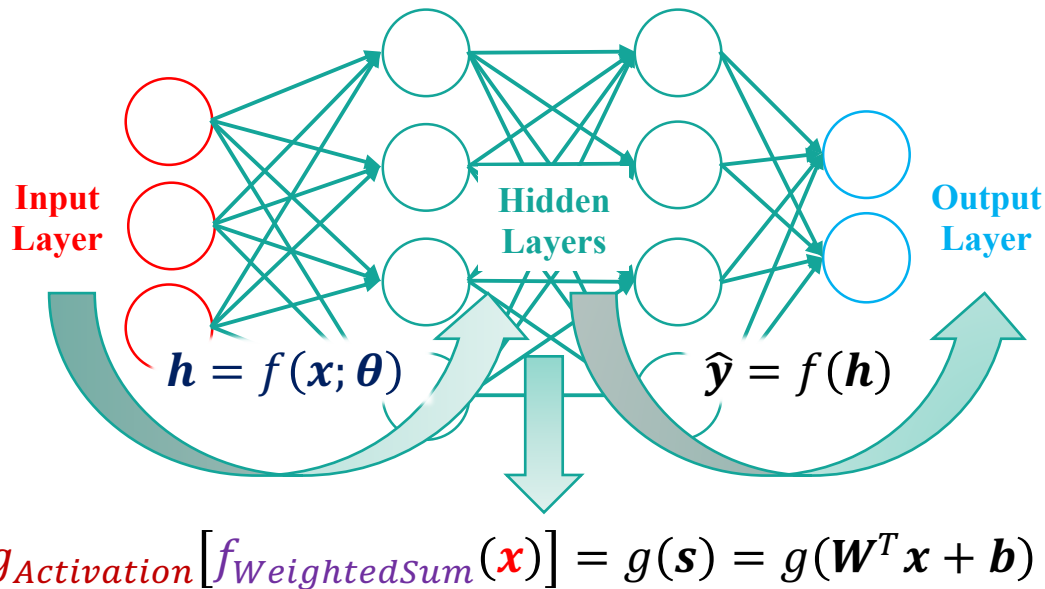
$$-\frac{1}{n} \sum_{i=1}^n y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)$$



# Hidden Units within the Hidden Layers

*A problem unique to deep neural networks (as they have hidden layers)*

- The activation design of hidden units
  - An extremely active area of research
  - Does not yet have many definitive guiding theoretical principles.



**Evaluating**  
its  
performance  
on a  
validation set

**Trial and  
Error**

*Usually impossible  
to predict in advance  
which will work best*

**Training a  
network with  
that kind of  
hidden unit**

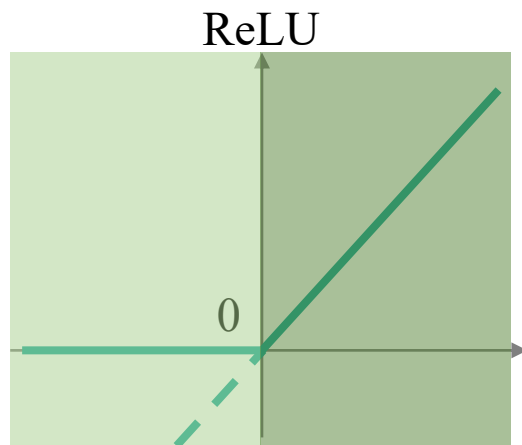
**Intuiting  
that a kind of  
hidden unit  
may work  
well**



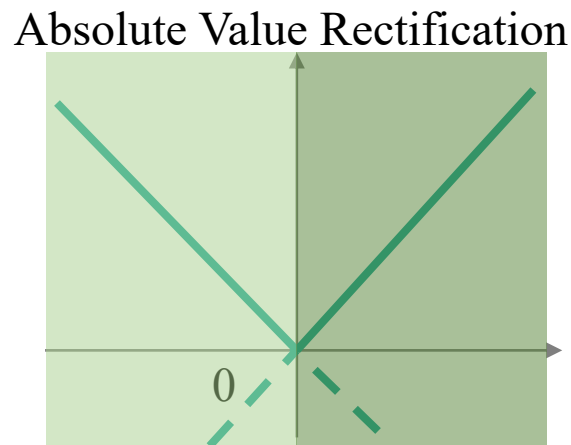
# Rectified Linear Units and Their Generalizations

$$g(s) = \max\{0, s\} \ \& \ g(s, \alpha) = \max\{0, s\} + \alpha \min\{0, s\}$$

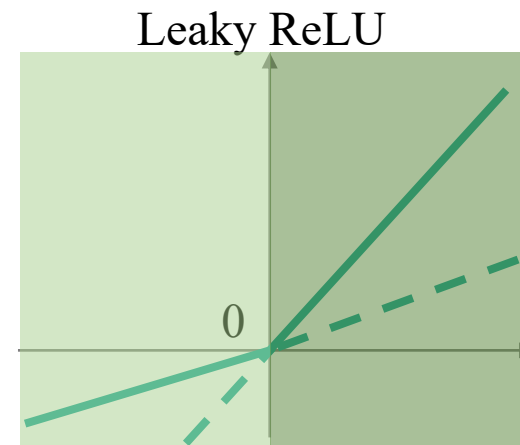
- Like a linear unit, easy to optimize
  - Output zero across half its domain => Large derivative whenever the unit is active
  - The 1st derivative is 1 whenever the unit is active
  - The 2nd derivative is 0 *almost* everywhere (not differentiable at  $z = 0$ )
  - A good practice to initialize the parameters with a small bias, such as 0.01



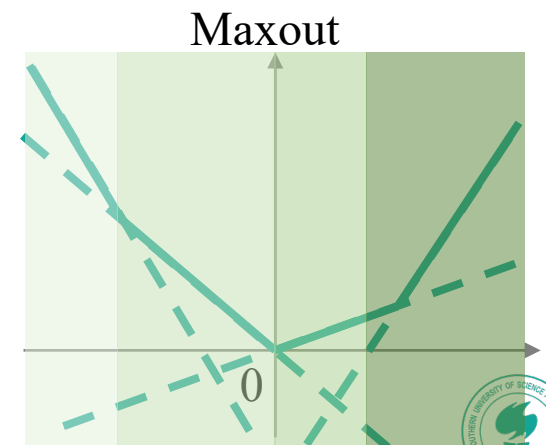
$$g(s) = \max\{0, s\}$$



$$g(s) = |s|$$



$$g(s) = \max\{0, s\} + 0.01 \times \min\{0, s\}$$



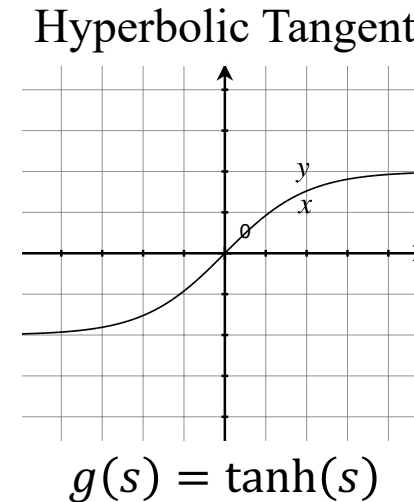
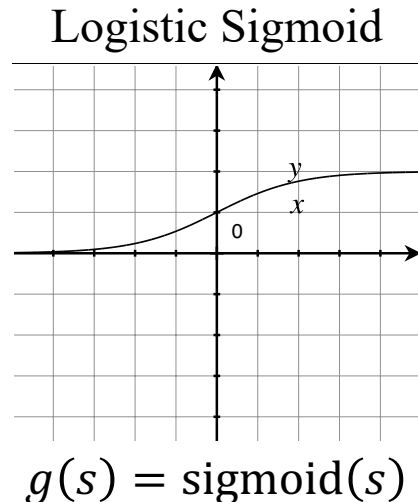
$$g(s) = \max_{k=1,2,3,4} \{s_k, s\}$$

# Logistic Sigmoid & Hyperbolic Tangent

$$g(s) = \text{sigmoid}(s) \text{ \& } g(s) = \text{tanh}(s)$$

- Popular before rectified linear units, used to predict classification probability
  - Closely related as  $\text{tanh}(s) = 2 \text{sigmoid}(2s) - 1$
- Widespread saturation
  - Approaching 1 when very positive, or approaching 0/-1 when very negative

- Difficult for gradient-based learning
- Discouraged for as hidden units for feedforward network
- Acceptable as output unit with appropriate cost function



- Typically performs better than the logistic sigmoid
- Resembles the identity function more closely
  - $\text{tanh}(0) = 0$
- Resembles a linear model more closely
  - Nearly linear with small activations

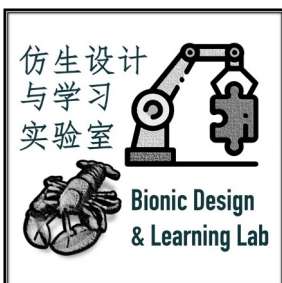
Bionic Design & Learning Lab  
@ SIR Group 仿生设计与学习实验室



Room 606  
7 Innovation Park  
南科创园7栋606室

Thank you~

[songcy@sustech.edu.cn](mailto:songcy@sustech.edu.cn)



AncoraSIR.com

