



# Week 06 | Lecture 07

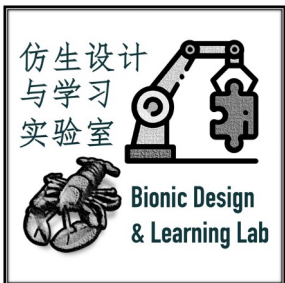
## From Logistic Regression to Neural Networks

Wan Fang

Southern University of Science and Technology

- **Statistical Binary Classification**
  - Logistic Regression
  - Stochastic Gradient Descent
- **Multi-class Classification**
  - Loss Function
  - Softmax Regression
  - Summary of Linear Classifications
- **Neural Networks**
  - Multi-Layer Perceptron
  - Forward & Backward Propagation
  - Exercises

# Statistical Binary Classification



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



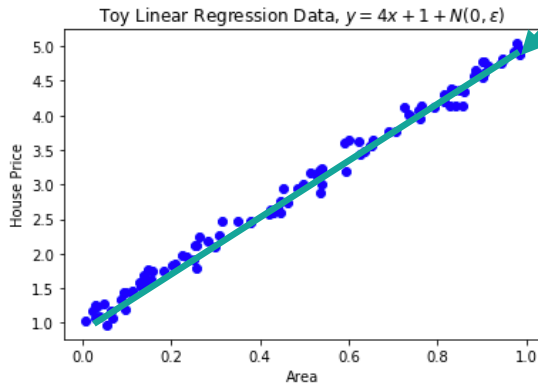
# Statistical Binary Classification

## Recall on linear regression and classification

- Linear Regression

- A basic linear model for line-fitting
- $\hat{y} = f_{weightedSum}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$

Estimate a line to describe a relationship



Predict University Acceptance based on Test and Grades (only two categories)

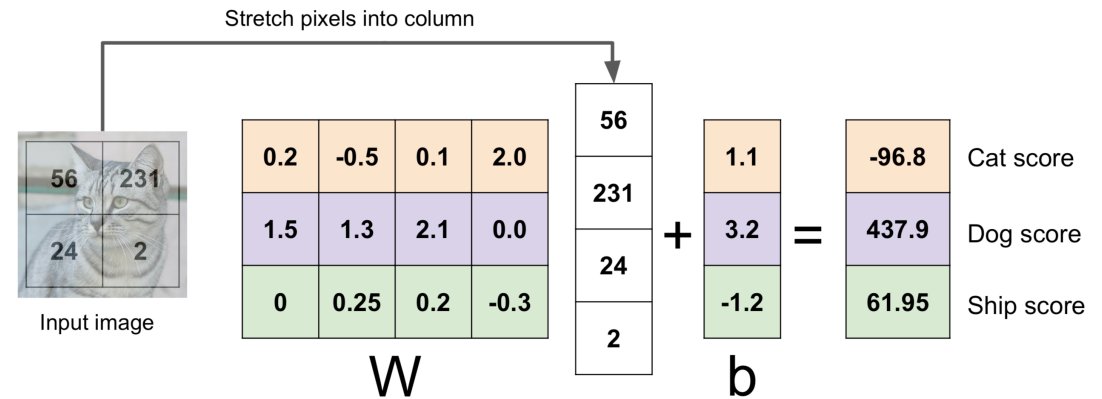


$$\hat{y} = g_{Activation}(s) = \begin{cases} 1 & \text{if } s \geq 0 \\ 0 & \text{if } s < 0 \end{cases}$$

- Linear Classification

- Vectorized weights for two or multiple classes
- $\mathbf{s} = f_{weightedSum}(\mathbf{x}) = \mathbf{W}\mathbf{x} + \mathbf{b}$

Is this picture a cat, a dog, or a ship?  
(Can we make a decision based on the results on the right?)



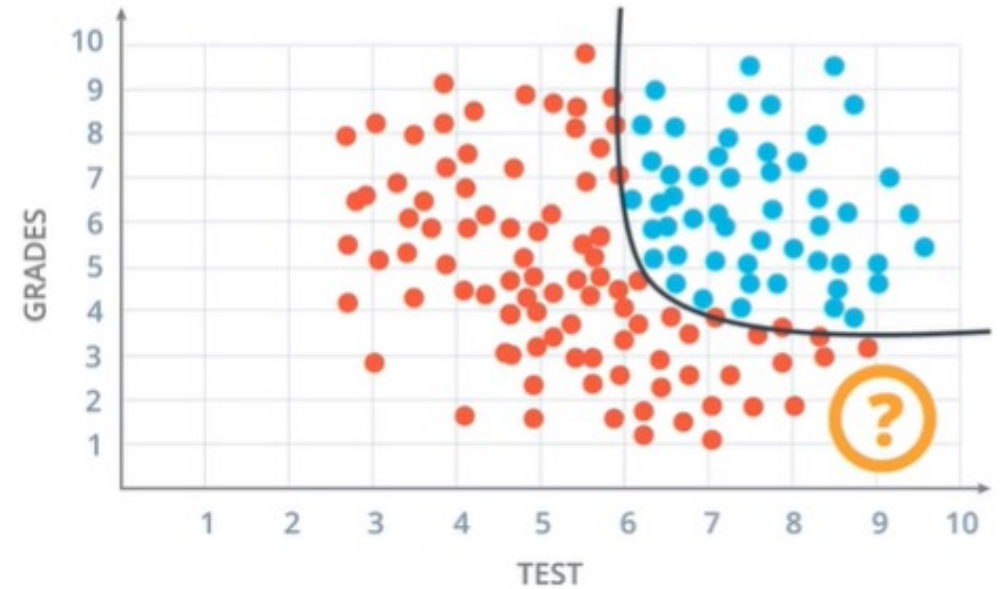
# Perceptron Algorithm

*If data is separable by a large margin, then Perceptron is a good algorithm to use.*



- Information lost about the distance to the cutoff value
- Uncertain about the final decision

*What if the boundary line is non-linear?*



- Unable to classify nonlinear scenarios

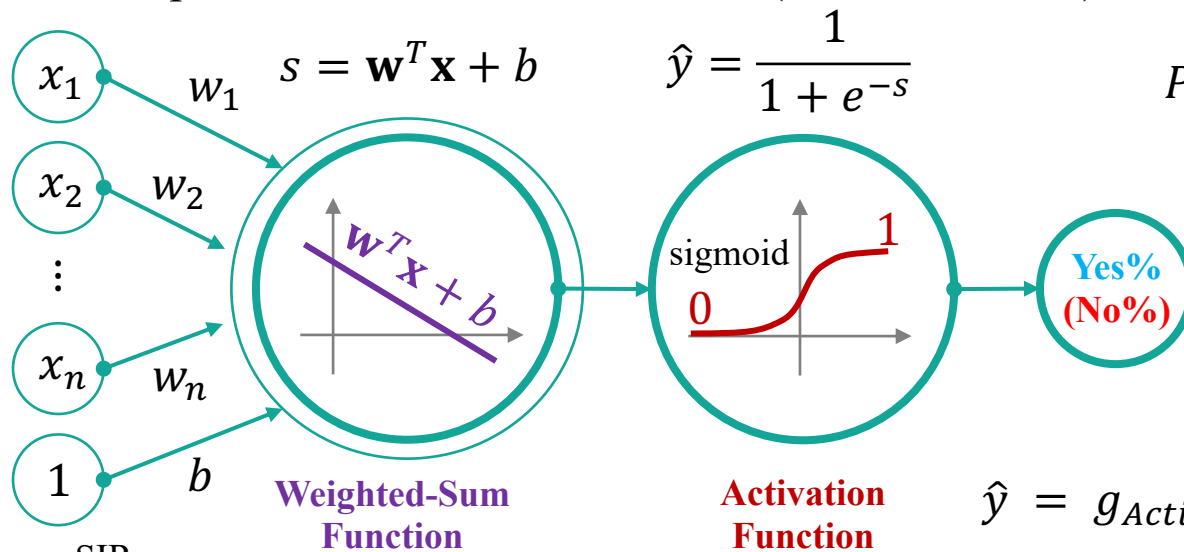
# By what chances will I get accepted to a University?

*Based on my Test and Grade scores ...*

- **Weighted-sum node**
  - Unchanged as the input data remains the same
- **Activation node**
  - Can be changed as we want a new expression of the output as a probability of prediction
  - Sigmoid function as a natural choice that transforms the output to a value between 0 and 1 (or within 100%)



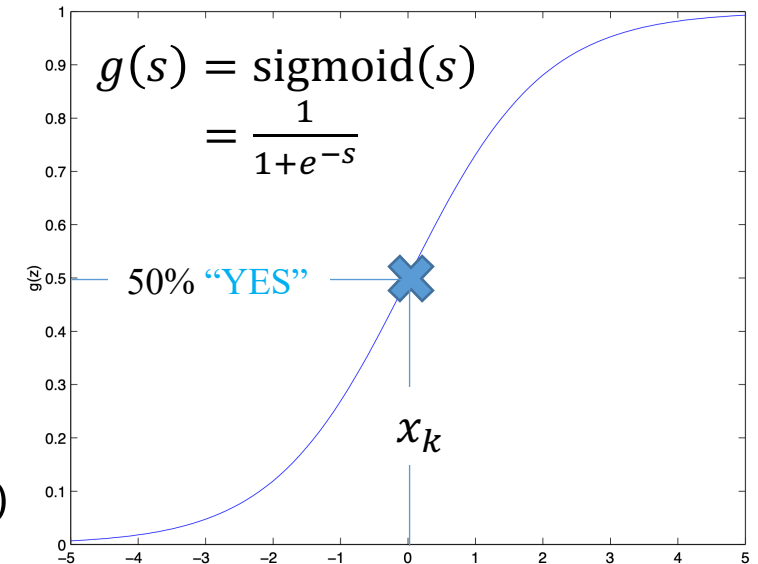
An example of acceptance at a University



$$P(\mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + b)}}$$

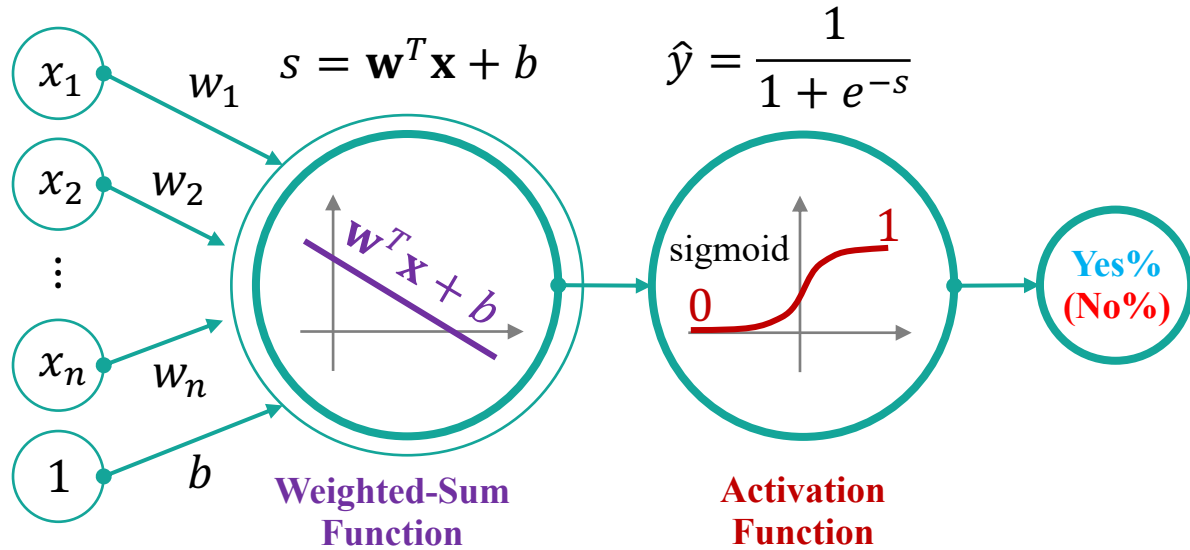
Can be viewed as a probability

$$\hat{y} = g_{Activation}(s) = \text{sigmoid}(s)$$



# Logistic Regression

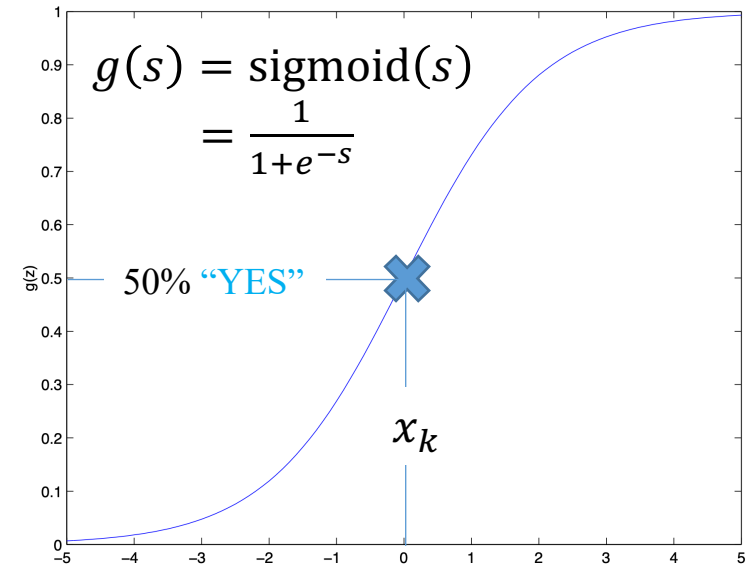
$$\hat{y} = g_{Activation}[f_{WeightedSum}(\mathbf{x})] = \text{sigmoid}(\mathbf{w}^T \mathbf{x} + b)$$



## Problem statement

- Assume  $\hat{y} = g_{Activation}[f_{WeightedSum}(\mathbf{x})] = \frac{1}{1+e^{-(\mathbf{w}^T \mathbf{x} + b)}}$
- How to minimize the **prediction error/loss** on a single training sample (with a maximum likelihood set of  $\mathbf{w}$ ) ?

- Hypothesis Function:  $h_{\mathbf{w}}(x) = \text{sigmoid}(\mathbf{w}^T \mathbf{x} + b)$
- Model output with a probability:  $P(y | x; \mathbf{w}) = [h_{\mathbf{w}}(x)]^y [1 - h_{\mathbf{w}}(x)]^{1-y}$ 
  - Yes%:**  $P(y = 1 | x; \mathbf{w}) = h_{\mathbf{w}}(x)$
  - No%:**  $P(y = 0 | x; \mathbf{w}) = 1 - h_{\mathbf{w}}(x)$



# Loss Function for Logistic Regression

*It measures how well you are doing on a single training example*

- Assume that  $m$  training examples were generated independently  $h_w(x) = \text{sigmoid}(\mathbf{w}^T \mathbf{x} + b)$
- We can write the likelihood of the parameters
  - $L(w) = p(\vec{y} | X; w)$   
 $= \prod_{i=1}^m p(y^{(i)} | x^{(i)}; w)$   
 $= \prod_{i=1}^m [h_w(x^{(i)})]^{y^{(i)}} [1 - h_w(x^{(i)})]^{1-y^{(i)}}$
- Take the log expression, we have the **loss function**
  - $\ell(w) = \log L(w)$   
 $= \sum_{i=1}^m y^{(i)} \log h_w(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_w(x^{(i)}))$
  - Usually take a “-” sign to indicate loss



# Stochastic Gradient Descent

*Finding the maximum likelihood of estimation*

- Rewrite the weight parameters in vectorized form
  - $w := w + \alpha \cdot \nabla_w \cdot \ell(w)$
  - + sign here to **maximize** likelihood
- When working with a single training example  $(x, y)$ ,
  - $\frac{\partial}{\partial w_j} \ell(w) = \left( y \frac{1}{g(w^T x)} - (1 - y) \frac{1}{1 - g(w^T x)} \right) \frac{\partial}{\partial w_j} g(w^T x) = (y - h_w(x)) x_j$
- Therefore, we can derive the stochastic gradient ascent rule
  - $w_j := w_j + \alpha \left( y^{(i)} - h_w(x^{(i)}) \right) x_j^{(i)}$

# Cost Function

*It measures how well you are doing on an entire training set*

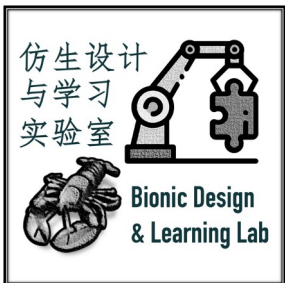
- We want the loss/error function to be as small as possible
  - If  $y^{(i)} = 1$ , then
    - $\text{LossFunc}(\hat{y}, y) = -\left[y^{(i)} \log h_w(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_w(x^{(i)}))\right] = -\log h_w(x^{(i)}) = -\log \hat{y}$
    - It means that we want  $\log \hat{y}$  to be as big as possible, but remember that it is bounded by 1
  - If  $y^{(i)} = 0$ , then
    - $\text{LossFunc}(\hat{y}, y) = -\left[y^{(i)} \log h_w(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_w(x^{(i)}))\right] = -\log(1 - \hat{y})$
    - It means that we want  $\log \hat{y}$  to be as small as possible, or close to 0
- Cost Function
  - The average of the loss functions of the entire training set, which is to be minimized

$$J(w, b) = -\frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})]$$

# Summary

	<b>Linear Regression</b>	<b>Perceptron</b>	<b>Logistic Regression</b>
<b>Problem</b>	Value Prediction	Binary Classification with a threshold	Binary Classification with a probability
<b>Weighted-Sum</b>	$w^T \mathbf{x} + b$	$w^T \mathbf{x} + b$	$w^T \mathbf{x} + b$
<b>Activation Function</b>	NA	Step Function	Sigmoid Function
<b>Prediction Outputs</b>	Continuous Value	Discrete Value $\{0, 1\}$	Continuous Probability $(0, 1)$
<b>Loss</b>	Squared Loss	Hinge Loss	Log-Loss

# Multi-class Classification

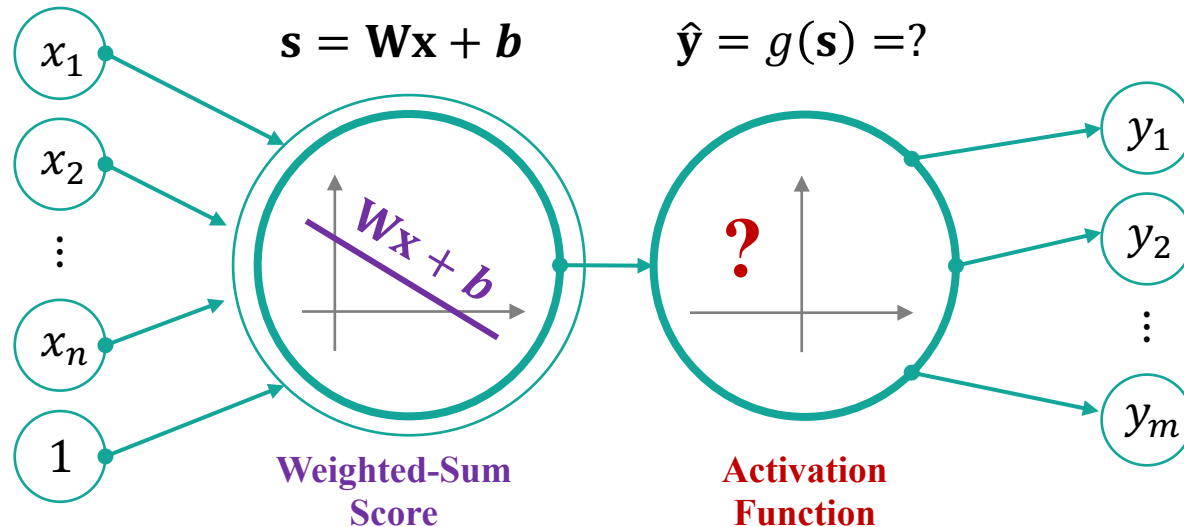
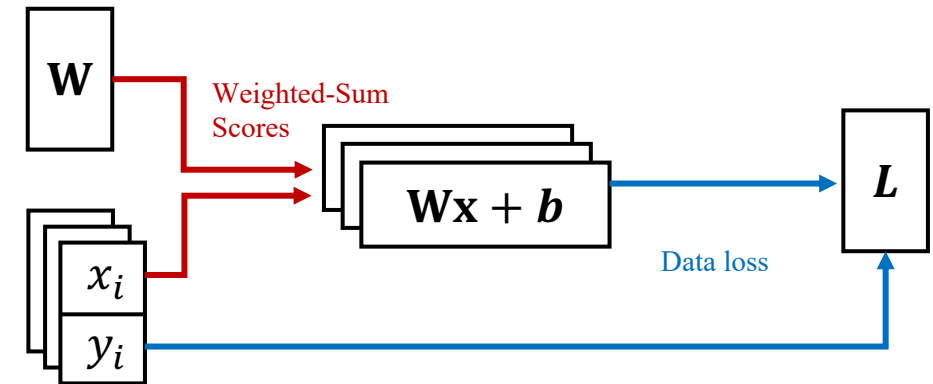
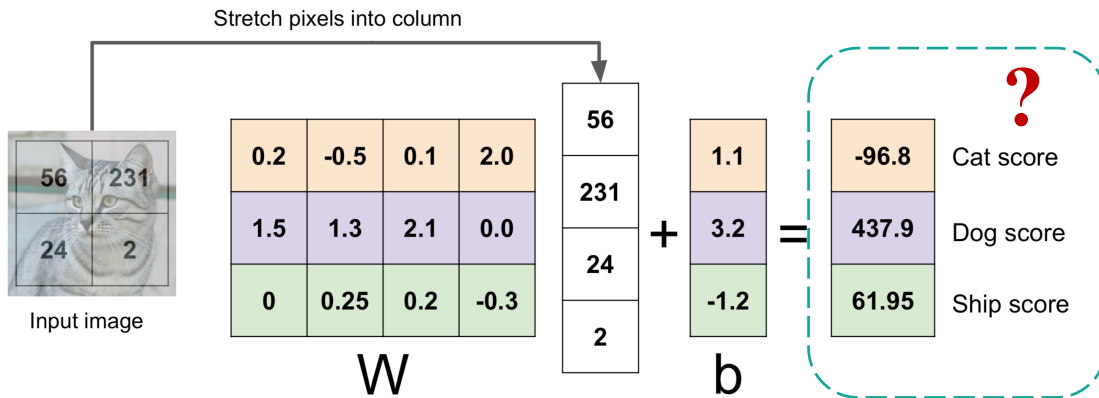


[AncoraSIR.com](http://AncoraSIR.com)



# Multi-class Classification

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



1. **Define a loss function** that quantifies our unhappiness with the scores across the training data.
2. Come up with a way of efficiently finding the parameters that minimize the loss function. (*optimization*)

# Define a Loss Function


*Quantify how good our current classifier is*

3 training samples

$\{(x_i, y_i)\}_{i=1}^3$

Ground Truth

Labelled Prediction  $\hat{y}_i$



$x_i$  image

$y_i$  label

3 classes	Cat	<b>3.2</b>	1.3	2.2
	Car	5.1	<b>4.9</b>	2.5
	Frog	-1.7	2.0	<b>-3.1</b>

Loss over the dataset is a sum of loss over examples

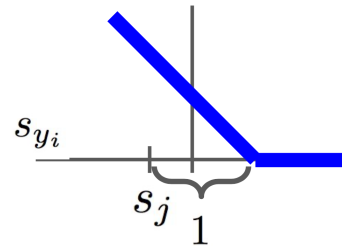
$$L = \frac{1}{N} \sum_i L_i(\hat{y}_i, y_i)$$

Denote Weighted-Sum score vector as  $\mathbf{s} = f_{WeightedSum}(\mathbf{x})$

Let's try with the hinge loss:

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_j + 1 \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases}$$

$$= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$



# Define a Loss Function

*Quantify how good our current classifier is*

3 training samples

$\{(x_i, y_i)\}_{i=1}^3$

Ground Truth

Labelled Prediction



3 classes

Cat	<b>3.2</b>	1.3	2.2
Car	5.1	<b>4.9</b>	2.5
Frog	-1.7	2.0	<b>-3.1</b>
<b>Loss</b>	<b>2.9</b>	<b>0</b>	<b>12.9</b>

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_j + 1 \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases}$$

$$\begin{aligned} L_1 &= \max(0, 5.1 - 3.2 + 1) + \max(0, -1.7 - 3.2 + 1) \\ &= \max(0, 2.9) + \max(0, -3.9) \\ &= 2.9 + 0 = 2.9 \end{aligned}$$

$$\begin{aligned} L_2 &= \max(0, 1.3 - 4.9 + 1) + \max(0, 2.0 - 4.9 + 1) \\ &= \max(0, -2.6) + \max(0, -1.9) \\ &= 0 + 0 = 0 \end{aligned}$$

$$\begin{aligned} L_3 &= \max(0, 2.2 + 3.1 - 1) + \max(0, 2.5 + 3.1 - 1) \\ &= \max(0, 6.3) + \max(0, 6.6) \\ &= 6.3 + 6.6 = 12.9 \end{aligned}$$

# Define a Loss Function

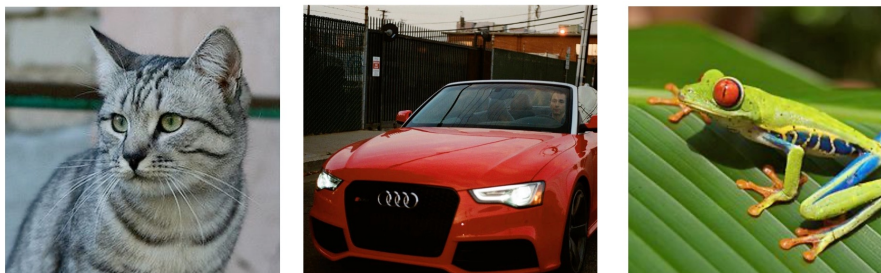
*Quantify how good our current classifier is*

3 training samples

$\{(x_i, y_i)\}_{i=1}^3$

Ground Truth

Labelled Prediction



3 classes

Cat	3.2	1.3	2.2
Car	5.1	4.9	2.5
Frog	-1.7	2.0	-3.1
<b>Loss</b>	<b>2.9</b>	<b>0</b>	<b>12.9</b>

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_j + 1 \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases}$$

Loss over full dataset is average:

$$\begin{aligned} L &= \frac{1}{N} \sum_i L_i(\hat{y}_i, y_i) \\ &= \frac{1}{3} (2.9 + 0 + 12.9) \\ &= 5.27 \end{aligned}$$

Recall that our goal is to find a set of  $\mathbf{W}$  with minimum loss over full dataset, i.e. the cost = 0

- Suppose that we found a  $\mathbf{W}$  such that  $L = 0$ . Is this  $\mathbf{W}$  unique?
  - $L$  is still 0 with  $2\mathbf{W}$
- **How do we choose between  $\mathbf{W}$  and  $2\mathbf{W}$ ?**
  - Let's try regularization



# Regularization

*Prevent the model from doing too well on training data*

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(\hat{y}_i, y_i) + \lambda R(W)$$

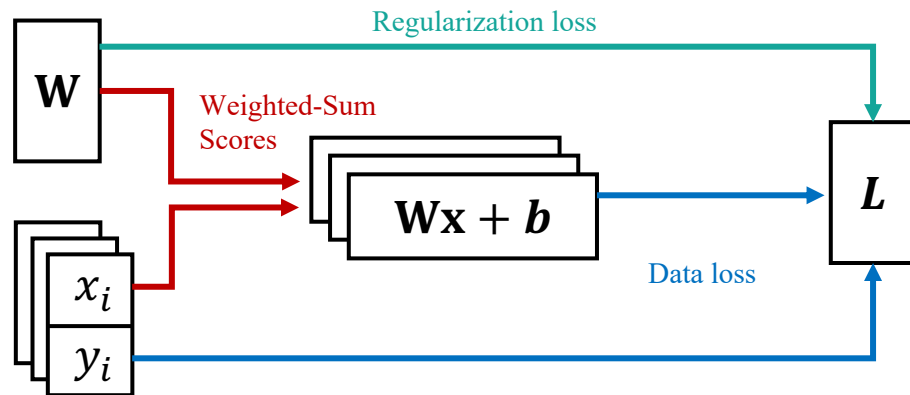
$\lambda$  as strength of  
Regularization  
(hyperparameter)

**Data loss**

**Regularization**

Model predictions should  
match training data

Prevent the model from doing  
too well on training data



## Simple examples

L2 regularization:  $R(W) = \sum_k \sum_l W_{k,l}^2$

L1 regularization:  $R(W) = \sum_k \sum_l |W_{k,l}|$

Elastic net (L1 + L2):  $R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$

## More complex:

Dropout

Batch normalization

Stochastic depth, fractional pooling, etc

## Why regularize?

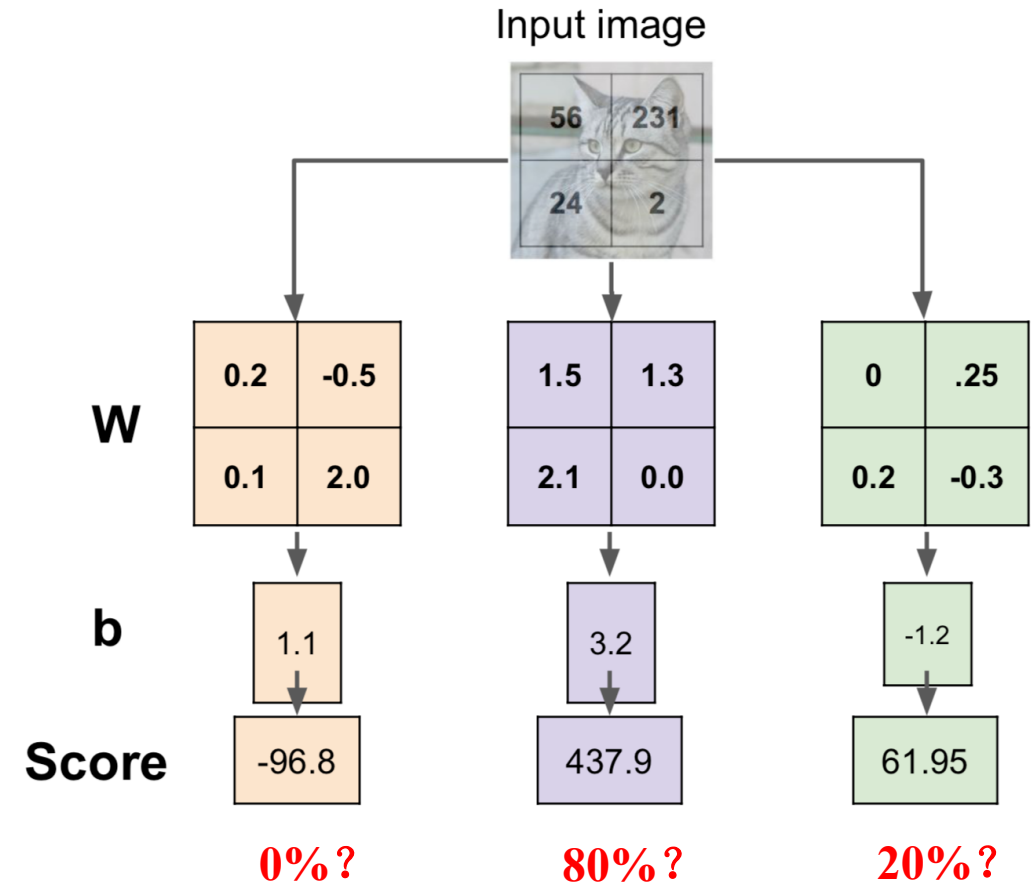
- Express preferences over weights
- Make the model simple so it works on test data
- Improve optimization by adding curvature

# Softmax Operation

*Interpret the outputs of our model as probabilities*

$$\hat{y}_i = \text{softmax}(o_i) = \frac{\exp(o_i)}{\sum_j \exp(o_j)} \times 100\%$$

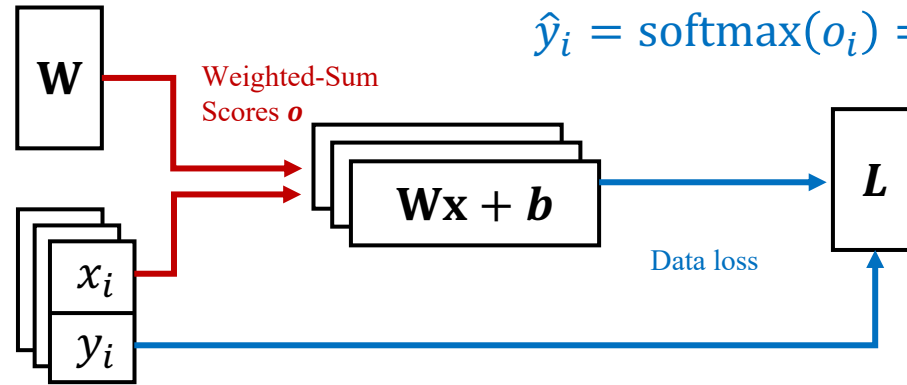
- One can interpret outputs  $\hat{y}_i$  as the probability that a given item belongs to class  $i$ .
- Then we can choose the class with the largest output value as our prediction
  - Why using  $o_i$  directly, instead of a probability?
  - What if the sum of probability is not 100%?
  - What if when  $o_i$  becomes negative?



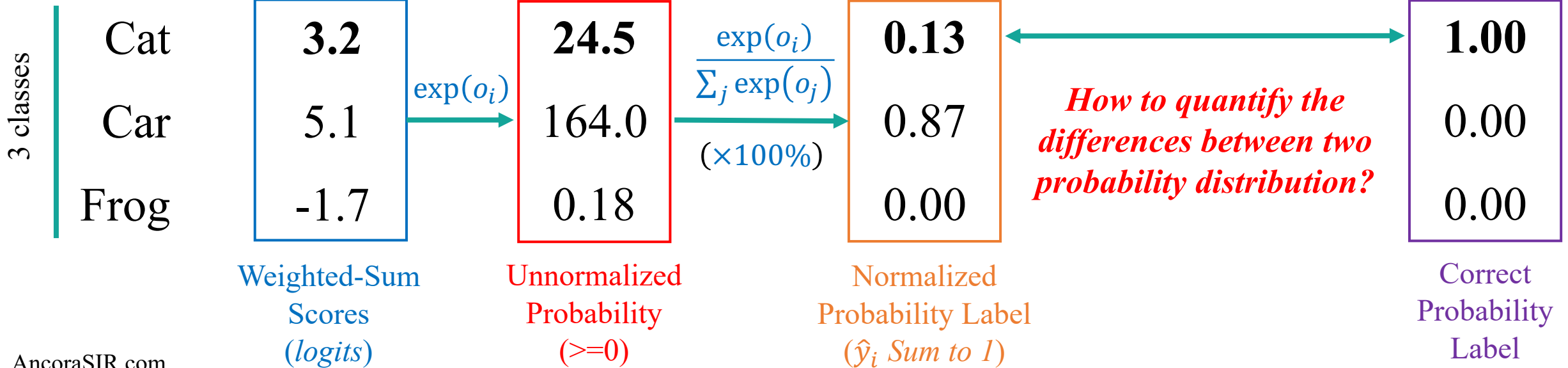
# Softmax Classifier

*Want to interpret raw classifier scores as probabilities*

$(x_1, y_1)$   
Ground Truth  
Labelled Prediction



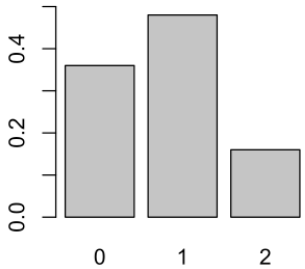
$$\hat{y}_i = \text{softmax}(o_i) = \frac{\exp(o_i)}{\sum_j \exp(o_j)} \times 100\%$$



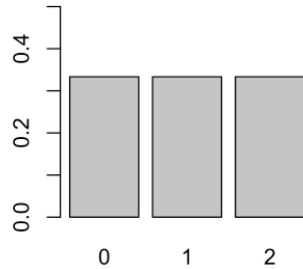
# Kullback–Leibler Divergence

*How to quantify the differences between two probability distribution?*

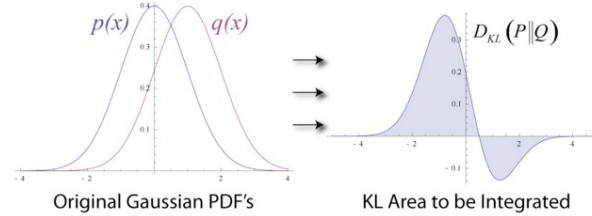
**Distribution P**  
Binomial with  $p = 0.4$ ,  $N = 2$



**Distribution Q**  
Uniform with  $p = 1/3$



<b>x</b>	<b>0</b>	<b>1</b>	<b>2</b>
Distribution $P(x)$	0.36	0.48	0.16
Distribution $Q(x)$	0.333	0.333	0.333



$$D_{KL}(P \parallel Q) = \sum_{y \in \mathcal{Y}} P(y) \log \frac{P(y)}{Q(y)}$$

$$= \sum_{y \in \mathcal{Y}} P(y) \log P(y) - \sum_{y \in \mathcal{Y}} P(y) \log Q(y)$$

$$= \left[ - \sum_{y \in \mathcal{Y}} P(y) \log Q(y) \right] - \left[ - \sum_{y \in \mathcal{Y}} P(y) \log P(y) \right]$$

$$= H(P, Q) - H(P)$$

**A good candidate of loss function for softmax**

*Can be minimized to update the weights*

$$\begin{aligned} D_{KL}(P \parallel Q) &= \sum_{x \in \mathcal{X}} P(x) \ln \left( \frac{P(x)}{Q(x)} \right) \\ &= 0.36 \ln \left( \frac{0.36}{0.333} \right) + 0.48 \ln \left( \frac{0.48}{0.333} \right) + 0.16 \ln \left( \frac{0.16}{0.333} \right) \\ &= 0.0852996 \end{aligned}$$

$$H(P, Q) = - \sum_{y \in \mathcal{Y}} P(y) \log Q(y) \quad H(P) = - \sum_{y \in \mathcal{Y}} P(y) \log P(y)$$

the cross-entropy of  $P$  and  $Q$

the cross-entropy of  $P$  with itself (or the entropy of  $P$ )

# Loss Function

*Log-Likelihood expressed in cross-entropy*

- The **likelihood** of the actual classes according to our model is

$$P(Y | X) = \prod_{i=1}^n P(y^{(i)} | x^{(i)}) \quad \longrightarrow \quad -\log P(Y | X) = \sum_{i=1}^n \boxed{-\log P(y^{(i)} | x^{(i)})}$$

- Maximizing the likelihood is equivalent to minimizing the log-likelihood.
- **Cross-entropy** loss for a single example (dropped superscript  $i$ )

$$l = -\log P(y | x) = -\sum_j y_j \log \hat{y}_j$$

- As  $\hat{y}$  is a discrete probability distribution and  $y$  is a one-hot vector, the sum over all  $j$  vanishes for all but one term.

# Cross-Entropy Loss and its Derivative

*Also called softmax loss*

- Plugging  $\mathbf{o}$  into the definition of the cross-entropy loss, we obtain:

$$l = - \sum_j y_j \log \hat{y}_j = \sum_j y_j \log \sum_k \exp(o_k) - \sum_j y_j o_j = \log \sum_k \exp(o_k) - \sum_j y_j o_j$$

- The derivative with respect to  $\mathbf{o}$  is

$$\partial_{o_j} l = \frac{\exp(o_j)}{\sum_k \exp(o_k)} - y_j = \text{softmax}(\mathbf{o})_j - y_j = P(y = j | x) - y_j$$

- The gradient is  $P(y = j | x) - y_j$ 
  - The difference between the probability predicted by our model  $P(y = j | x)$  and the true label  $y$ .
- Similar to regression where the gradient is  $\hat{y} - y$ 
  - The difference between the true label  $y$  and the estimation  $\hat{y}$

# Vectorization for Minibatches

We typically carry out vector calculations for minibatches of data for efficiency

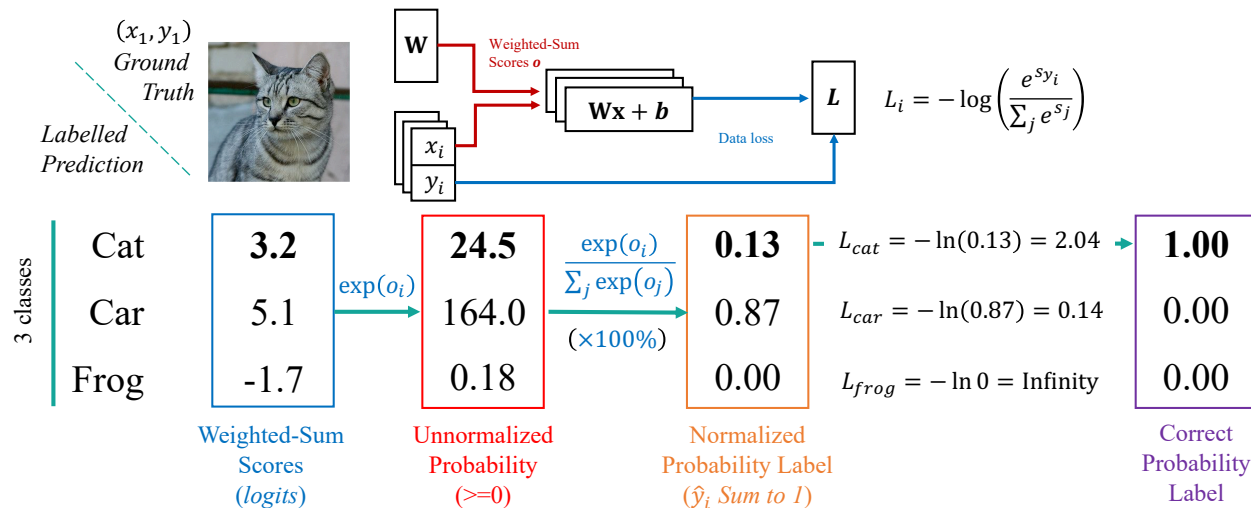
$$\hat{y}_i = \text{softmax}(o_i) = \frac{\exp(o_i)}{\sum_j \exp(o_j)} \times 100\%$$

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{o}) \text{ where } \hat{y}_i = \frac{\exp(o_i)}{\sum_j \exp(o_j)} \times 100\%$$

minibatch features  $\mathbf{X}$  are in  $\mathbb{R}^{n \times d}$ ,  
weights  $\mathbf{W} \in \mathbb{R}^{d \times q}$ , and  
the bias satisfies  $\mathbf{b} \in \mathbb{R}^q$ .

$$\mathbf{O} = \mathbf{XW} + \mathbf{b},$$

$$\hat{\mathbf{Y}} = \text{softmax}(\mathbf{O})$$



A minibatch  $\mathbf{X}$  of examples

- dimensionality  $d$  and batch size  $n$

Assume that we have  $q$  categories (outputs)

More efficient matrix-matrix computation  $\mathbf{XW}$   
Exponentiating all entries in  $\mathbf{O}$  then sum

# Understanding of Softmax Regression

- When there are two classes, softmax regression reduces to logistic regression.

Softmax

Binary Classes

Logistic

$$\hat{y}_j = \frac{\exp(o_j)}{\sum_j \exp(o_j)}$$

Activation

$$\hat{y} = \frac{\exp(o)}{\exp(o) + 1}$$

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

Loss

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

- Softmax when  $j=2$ 
$$\hat{y}_0 = \frac{\exp(o_0)}{\exp(o_0) + \exp(o_1)}$$
$$= \frac{\exp(o_0 - o_1)}{\exp(o_0 - o_1) + 1}$$

- The cross-entropy classification can be thought in two ways
  - As maximizing the likelihood of the observed data.
  - As minimizing out surprise required to communicate the labels.



# Summary & Comparison

## *Linear Neural Network*

	<b>Linear Regression</b>	<b>Perceptron</b>	<b>Logistic Regression</b>	<b>Softmax Regression</b>
<b>Problem</b>	Value Prediction	Binary Classification	Binary Classification	Multi-Class Classification
<b>Weights</b>	$wx + b$	$wx + b$	$wx + b$	$Wx + B$
<b>Activation Function</b>	NA	Step Function	Sigmoid Function	Softmax
<b>Prediction Outputs</b>	Continuous Value	Discrete Value 0, 1	Continuous Probability in (0,1)	A vector of Continuous Probabilities
<b>Loss</b>	Squared Loss	Hinge Loss	Log Loss (Binary cross entropy)	Cross Entropy
<b>Decision Boundary</b>			Linear	

# Neural Network



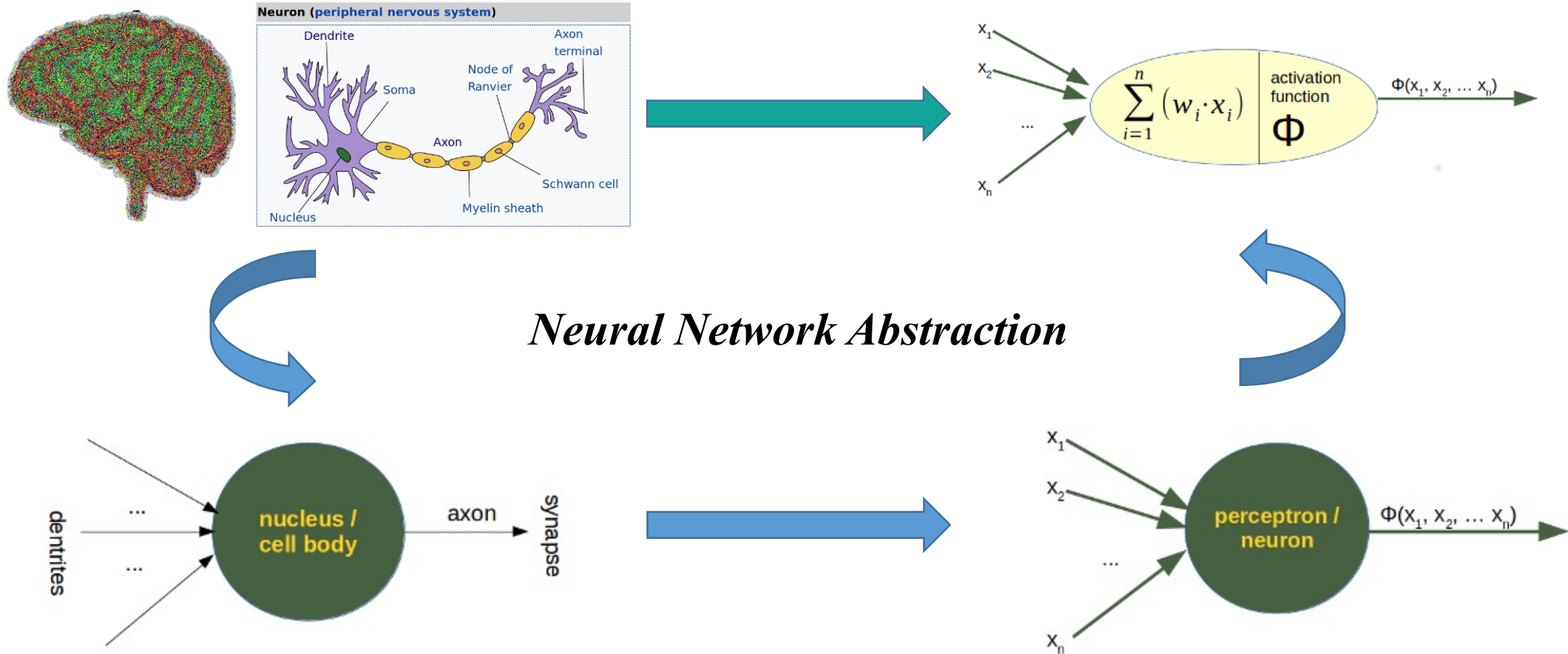
[AncoraSIR.com](http://AncoraSIR.com)



**SUSTech**  
Southern University  
of Science and Technology

# What is a Neural Network?

*From biological inspiration to mathematical modeling*



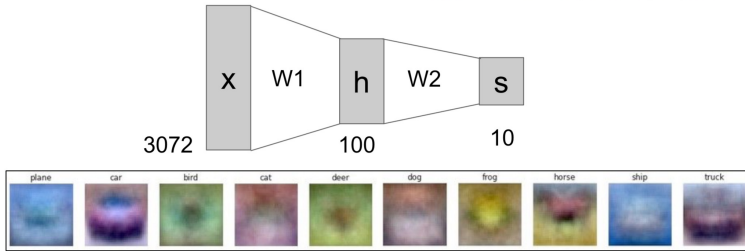
# A Perceptron as an Artificial Neuron

Linear score function:

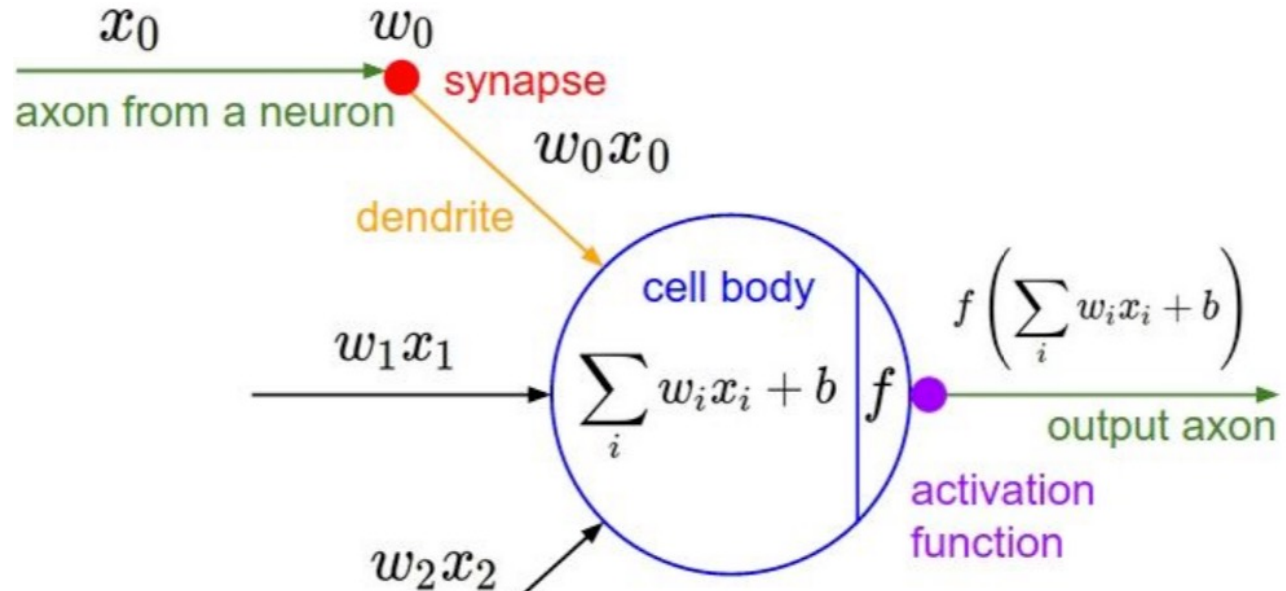
$$f = Wx$$

2-layer Neural Network

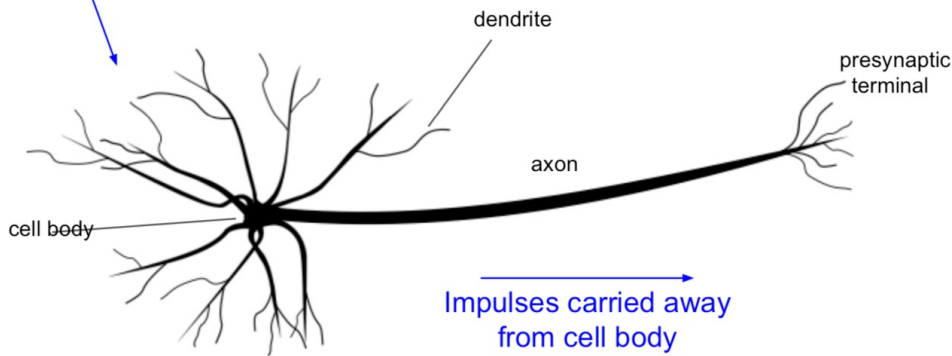
$$f = W_2 \max(0, W_1 x)$$



## Biological Inspiration



Impulses carried toward cell body

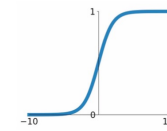


Impulses carried away from cell body

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

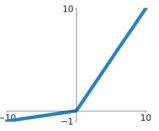
**Sigmoid**

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



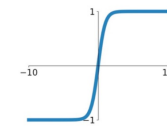
**Leaky ReLU**

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



**tanh**

$$\tanh(x)$$

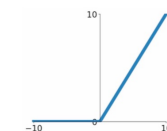


**Maxout**

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

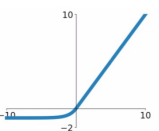
**ReLU**

$$\max(0, x)$$



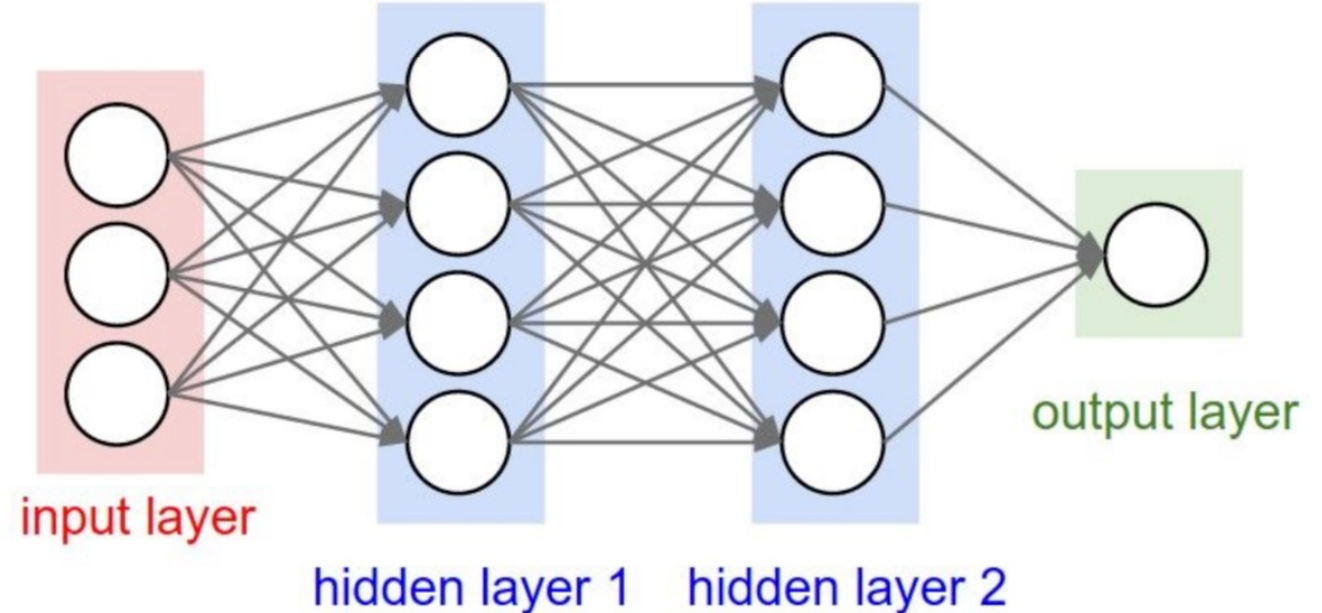
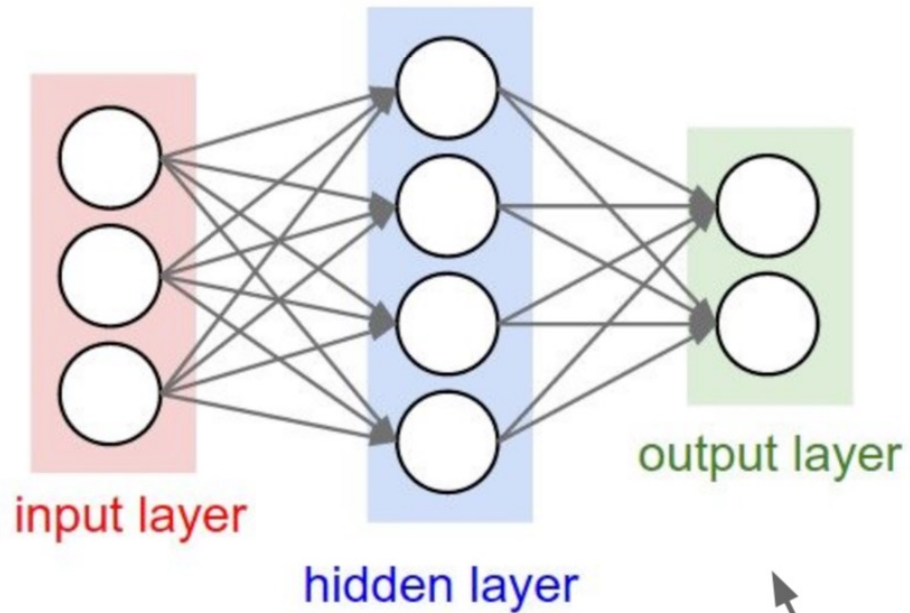
**ELU**

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



# Multi-Layer Perceptrons

## *Artificial Neural Networks*



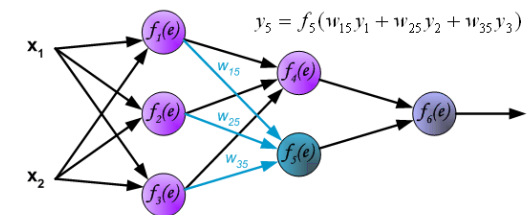
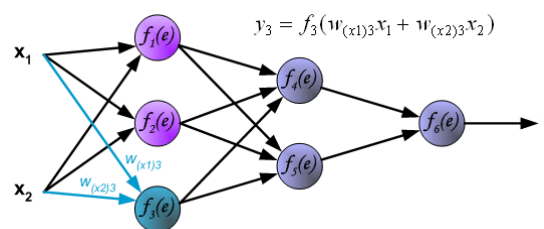
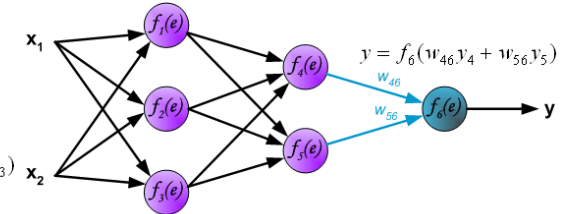
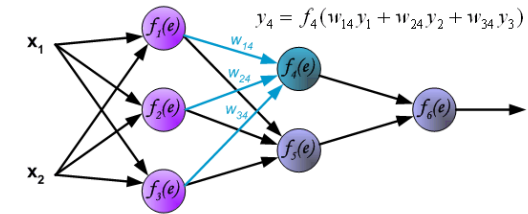
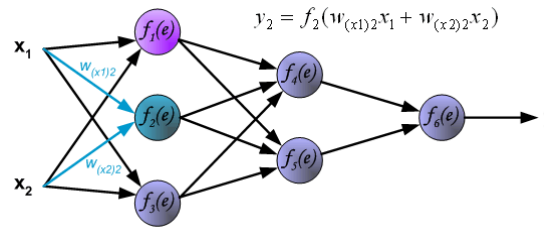
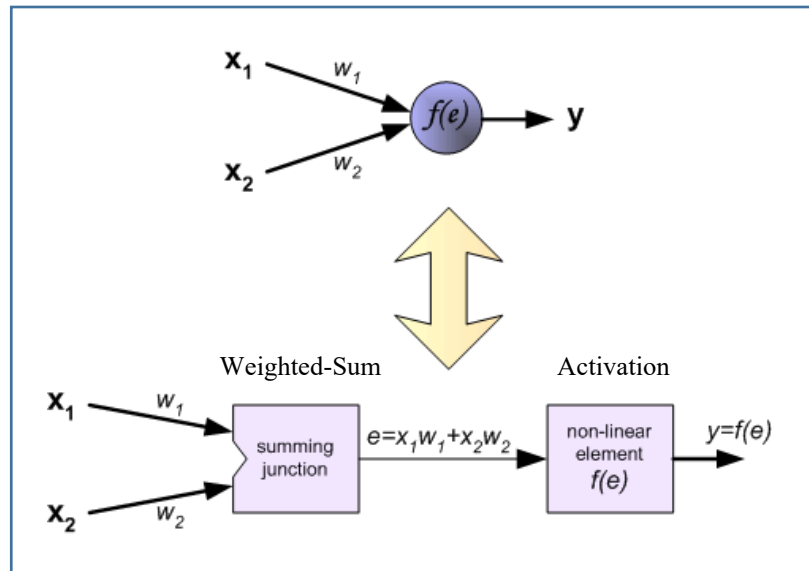
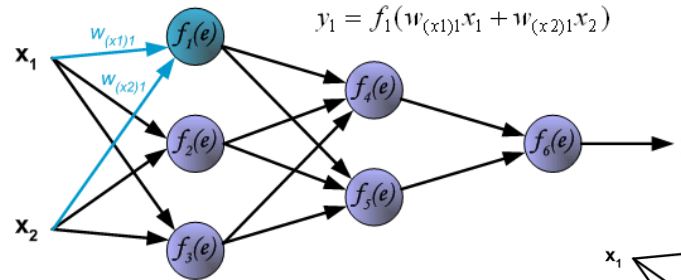
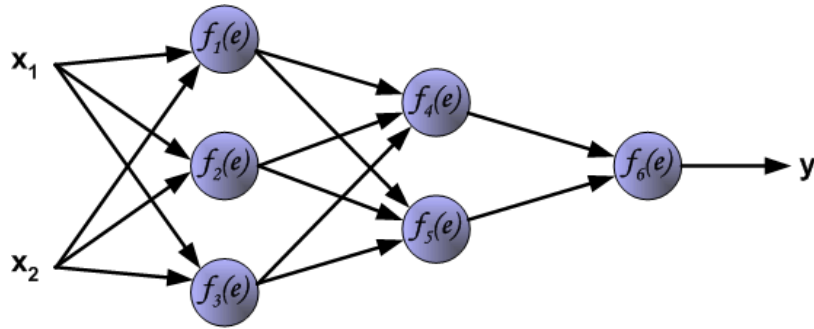
“2-layer Neural Net”, or  
“1-hidden-layer Neural Net”

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

**“Fully-connected” layers**

# Forward Propagation

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



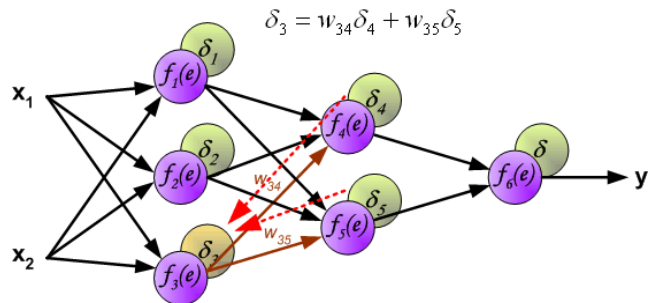
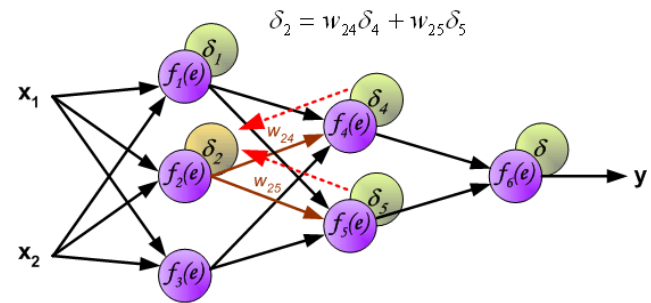
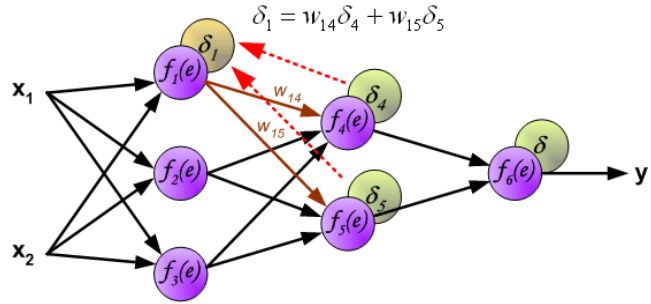
[http://galaxy.agh.edu.pl/%7Evlsi/AI/backp\\_t\\_en/backprop.html](http://galaxy.agh.edu.pl/%7Evlsi/AI/backp_t_en/backprop.html)



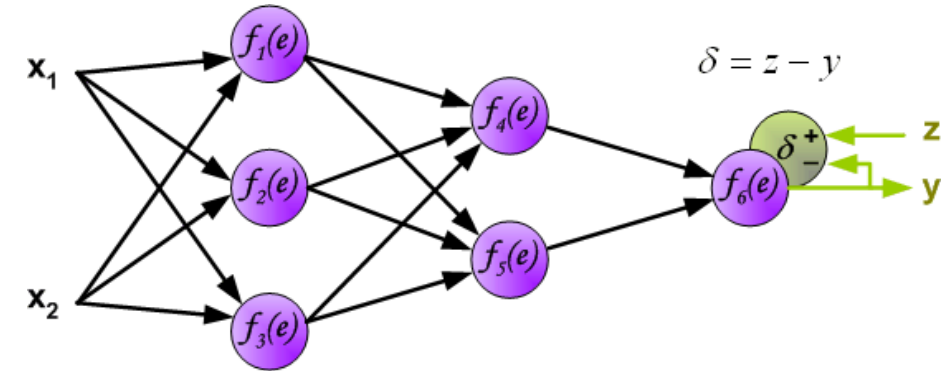
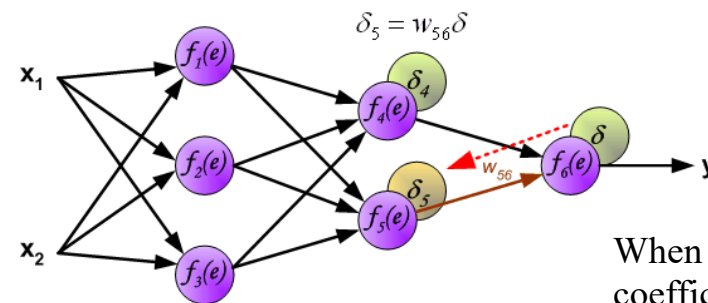
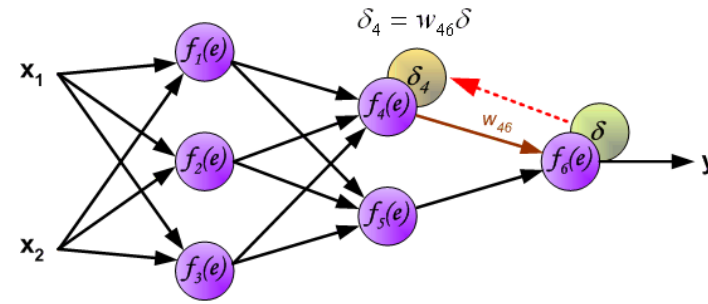
SUSTech  
Southern University  
of Science and Technology

# Backward Propagation

*Calculate the prediction error node-by-node*



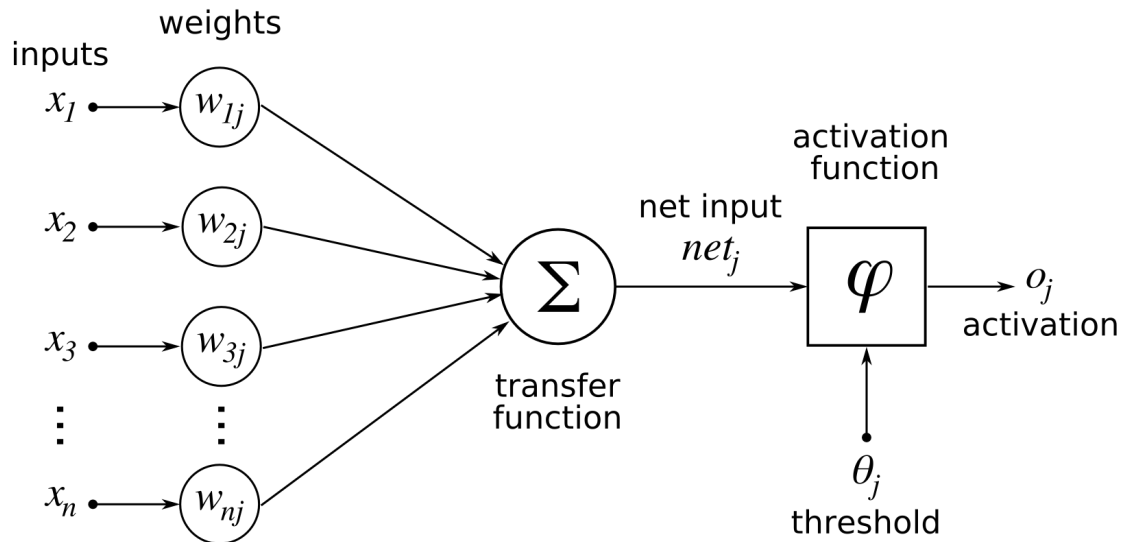
The idea is to propagate error signal  $d$  (computed in single teaching step) back to all neurons, which output signals were input for discussed neuron.



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).

# Backward Propagation



$$\frac{\partial E}{\partial w_{ij}} = o_i \delta_j$$

with

$$\delta_j = \frac{\partial E}{\partial o_j} \frac{\partial o_j}{\partial net_j} = \begin{cases} \frac{\partial L(o_j, t)}{\partial o_j} \frac{d\varphi(net_j)}{dnet_j} & \text{if } j \text{ is an output neuron,} \\ (\sum_{\ell \in L} w_{j\ell} \delta_\ell) \frac{d\varphi(net_j)}{dnet_j} & \text{if } j \text{ is an inner neuron.} \end{cases}$$

$$\Delta w_{ij} = -\eta \frac{\partial E}{\partial w_{ij}} = -\eta o_i \delta_j$$



# Exercise I

---

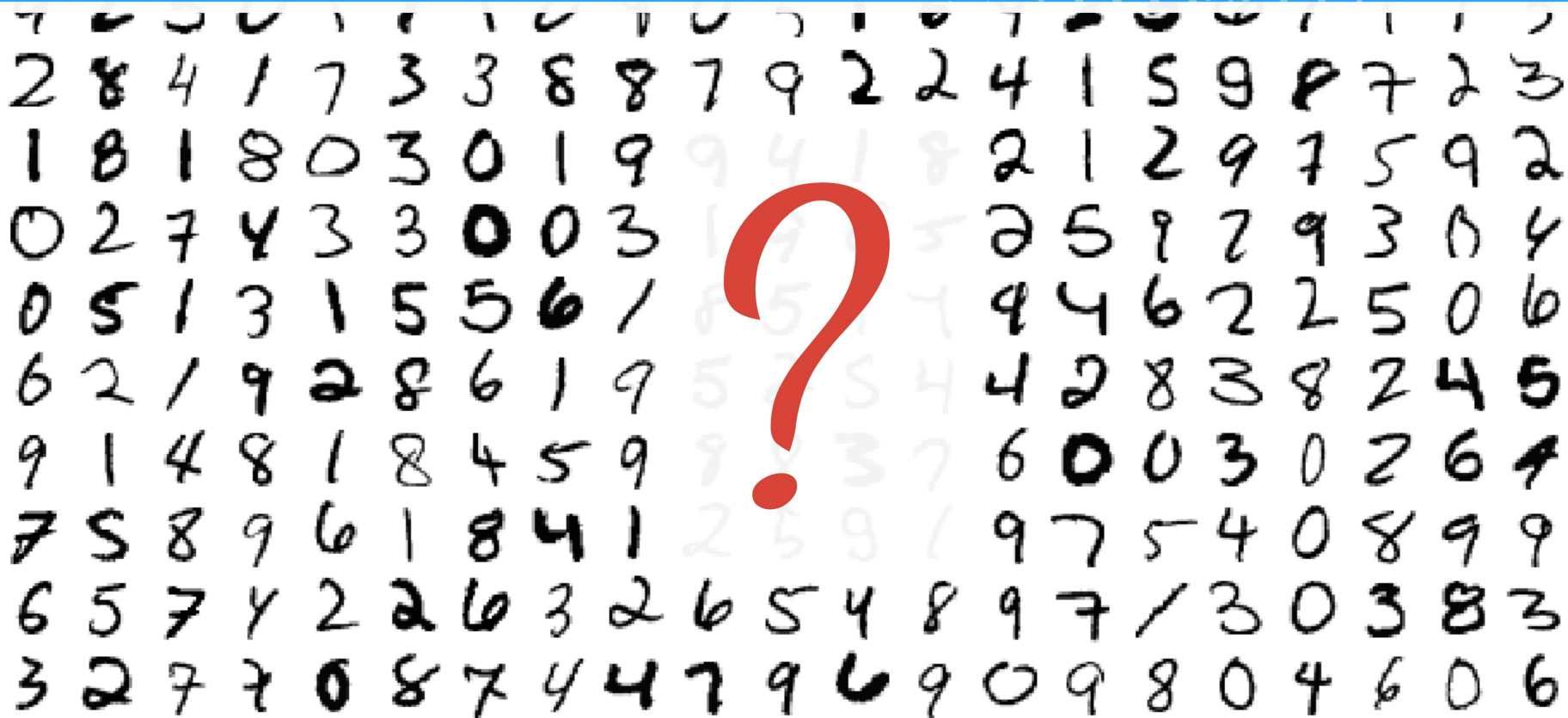
## *Neural Networks Playground from Google*

- Build your first neural networks on the website
  - <https://playground.tensorflow.org/>
- Play with different data types, features, network structures. Can Neural networks separate nonlinear features?
- Try different hyperparameters.

# Exercise II

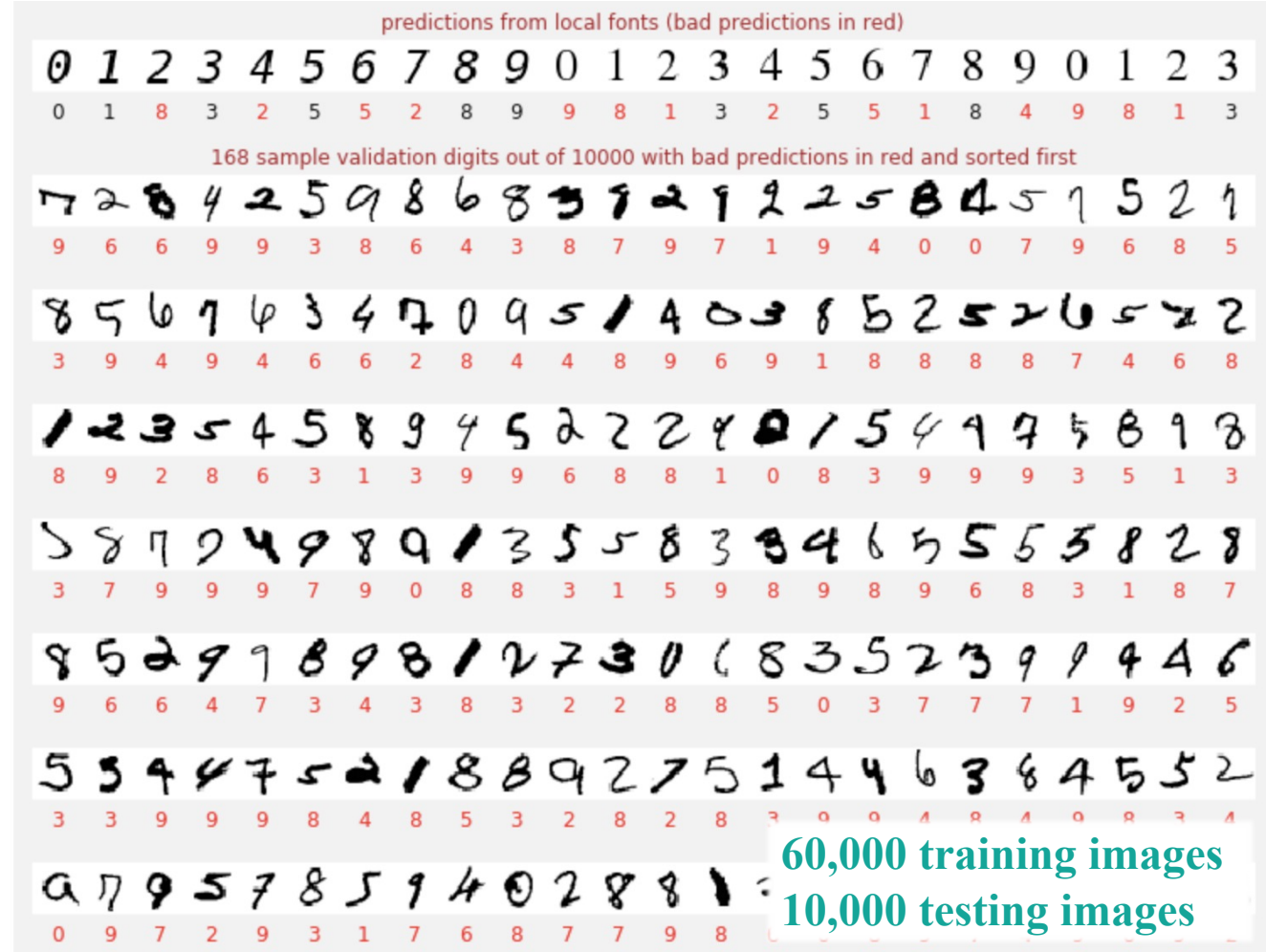
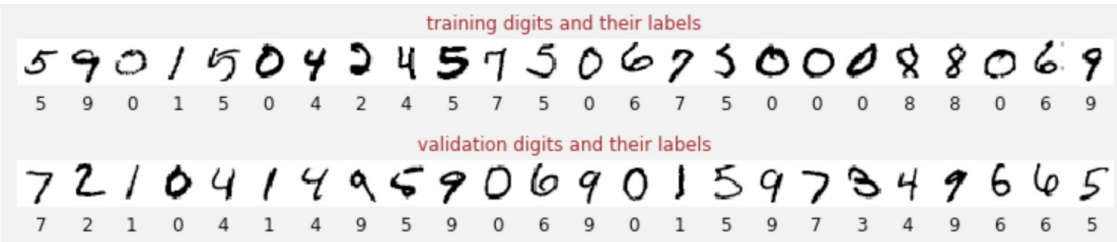
*“Hello, MNIST!”*

Hello World: handwritten digits classification - MNIST



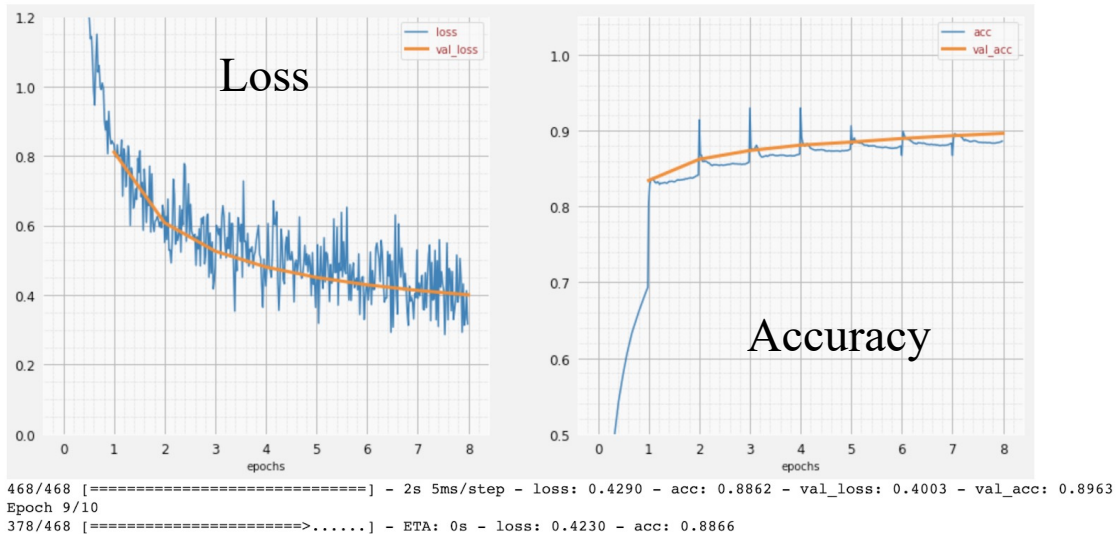
# A Toy Example of Training a Neural Network

*Training => Validation => Regularization => Testing => Optimization*



**Tensors** as matrices storage of data

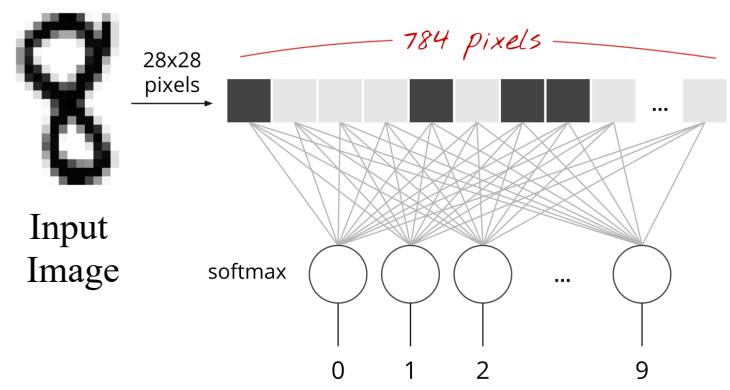
- [28, 28, 1]: A 28x28 pixel grayscale image (Gray)
- [128, 28, 28, 3]: A batch of 128 color images (RGB)



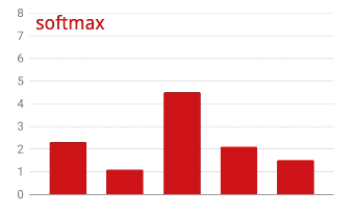
# A Single-layer Network of Image Classification

*Modified National Institute of Standards and Technology database*

Straightened into an array of features (x)



Through Softmax operation to predict 10 labels (y) with a normalized distribution



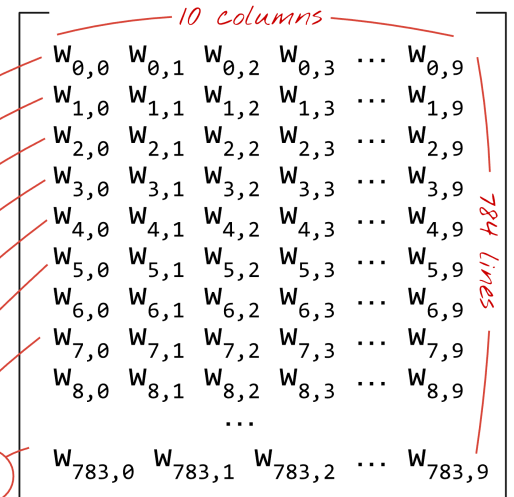
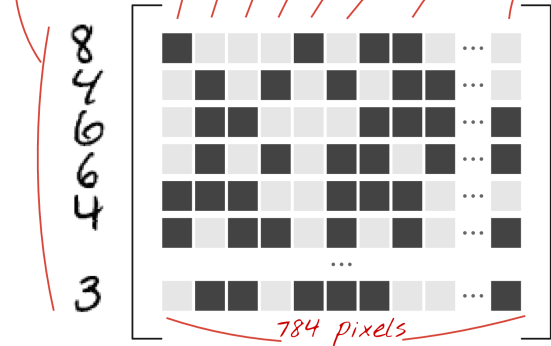
*weighted sum of all pixels + bias*

$$\text{softmax}(L_n) = \frac{e^{L_n}}{\|e^L\|}$$

*neuron outputs*

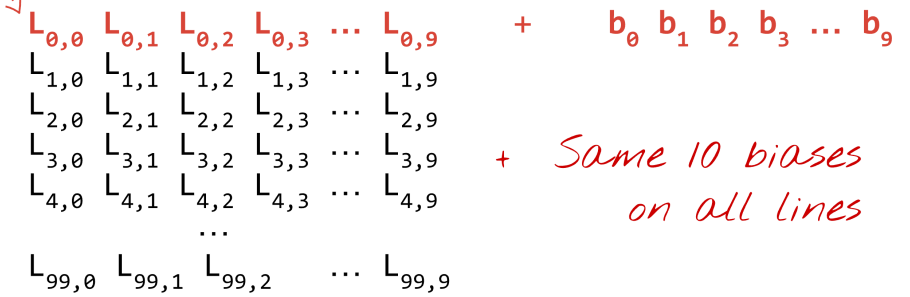
Vectorized dataset into a **batch** (of 100 images) for efficient operation

*X: 100 images, one per line, flattened*



*broadcast*

$$L = X \cdot W + b$$



# Softmax on a Batch of Images

*Output with a distribution of computed probabilities*

Predictions  
 $Y[100, 10]$

Images      Weights      Biases  
 $X[100, 784]$     $W[784,10]$     $b[10]$

$$Y = \text{softmax}(X \cdot W + b)$$

applied line by line

matrix multiply

broadcast on all lines

tensor shapes in [ ]

tensor shapes:  $X[100, 784]$     $W[784,10]$     $b[10]$

$$Y = \text{tf.nn.softmax}(\text{tf.matmul}(X, W) + b)$$

matrix multiply

broadcast on all lines

0	1	2	3	4	5	6	7	8	9
0	0	0	0	0	0	1	0	0	0

actual probabilities, "one-hot" encoded

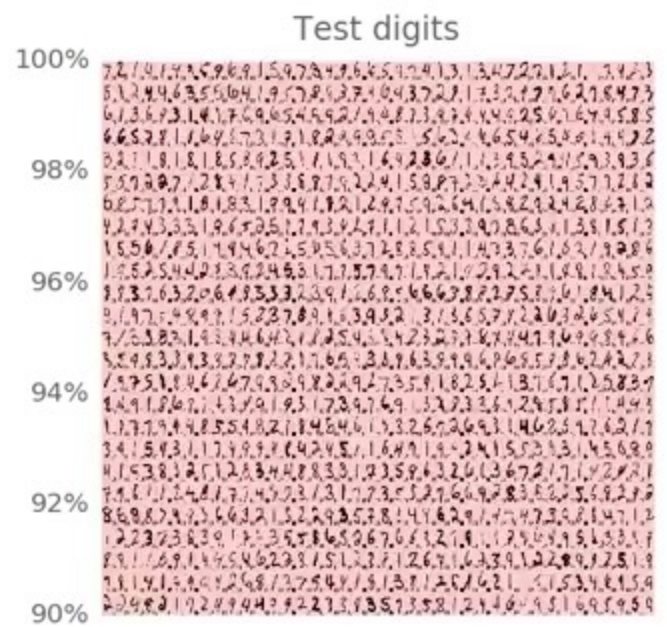
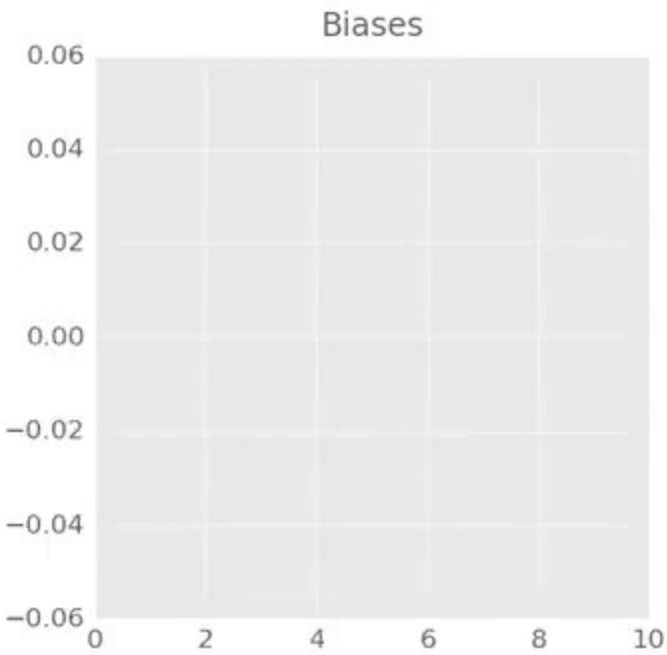
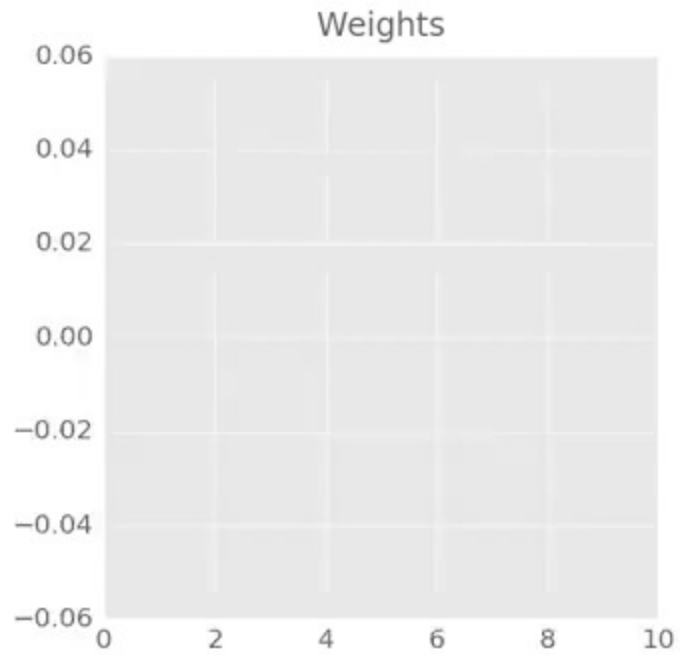
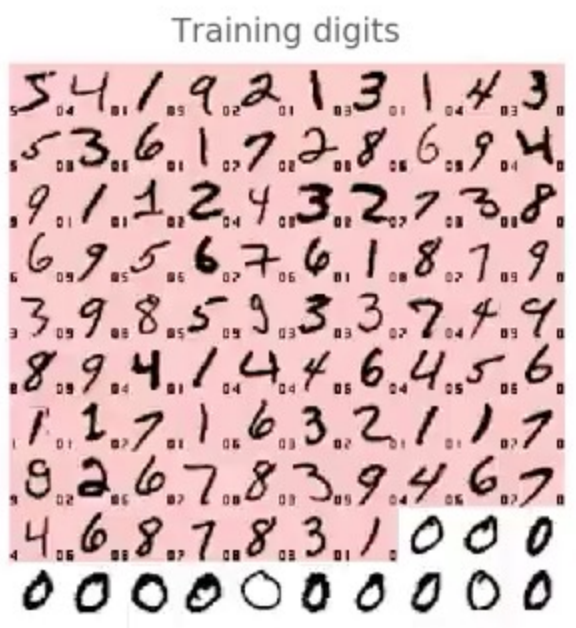
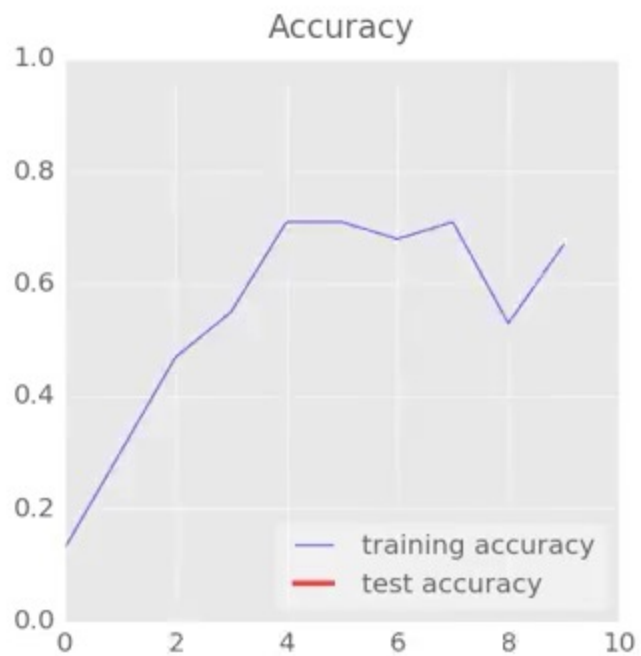
$$-\sum Y_i' \cdot \log(Y_i) \text{ Cross entropy}$$

computed probabilities

this is a "6"

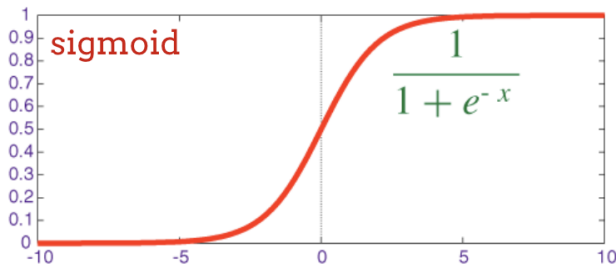
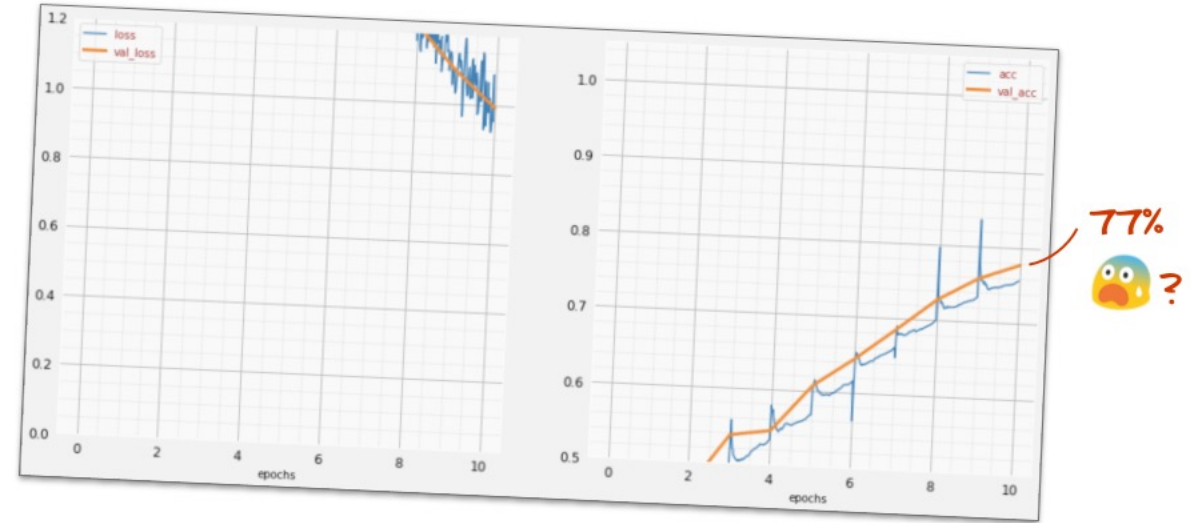
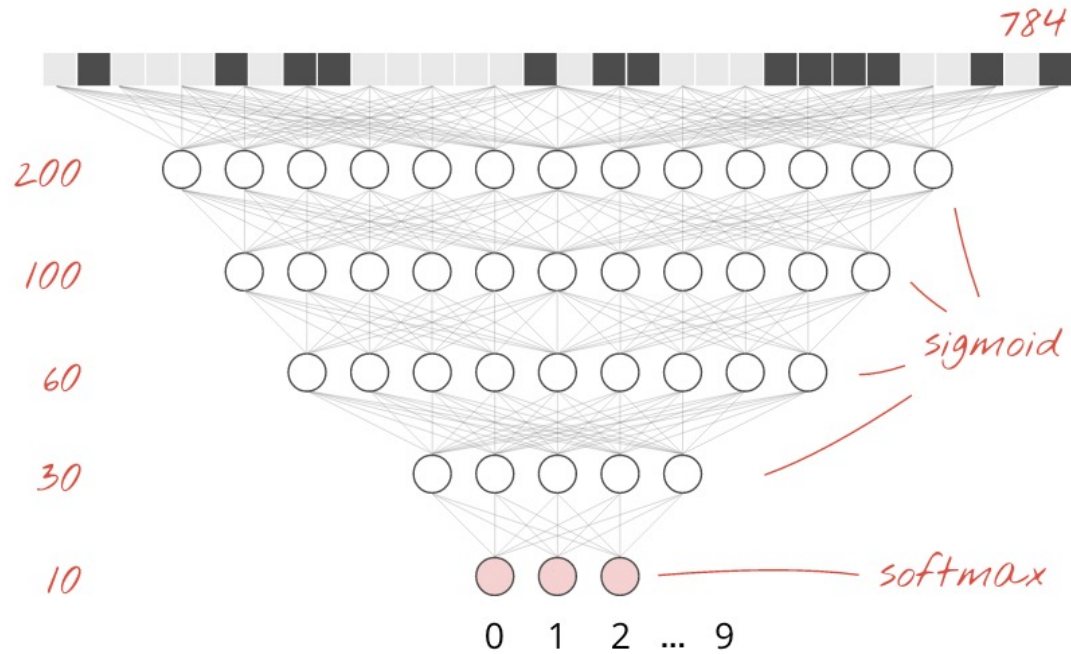
.01	.03	.00	.04	.03	.05	0.8	.02	.01	.01
0	1	2	3	4	5	6	7	8	9

[https://github.com/GoogleCloudPlatform/tensorflow-without-a-phd/blob/master/tensorflow-mnist-tutorial/keras\\_01\\_mnist.ipynb](https://github.com/GoogleCloudPlatform/tensorflow-without-a-phd/blob/master/tensorflow-mnist-tutorial/keras_01_mnist.ipynb)



# Adding Layers

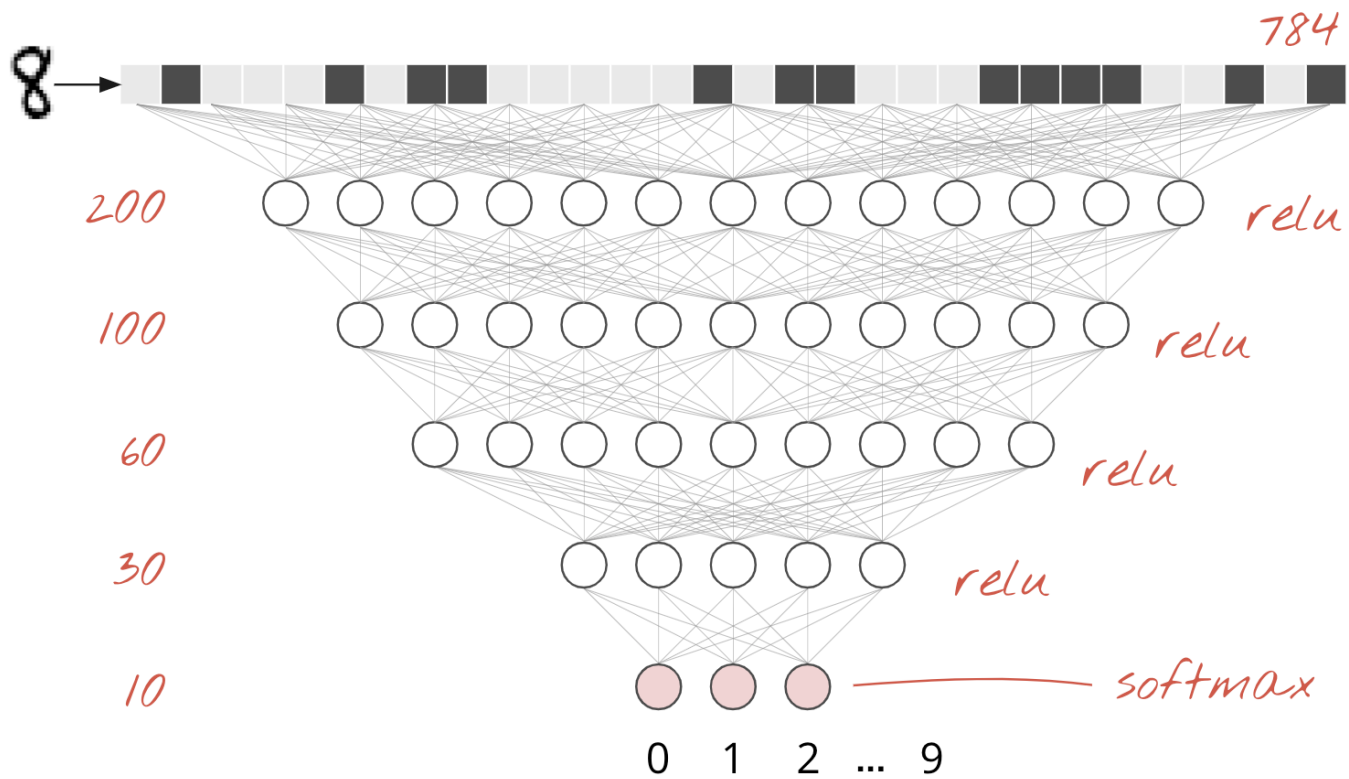
*Simply adding more layers with sigmoid activations does not give us the expected results ...*



Getting flat

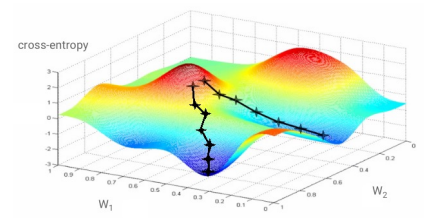
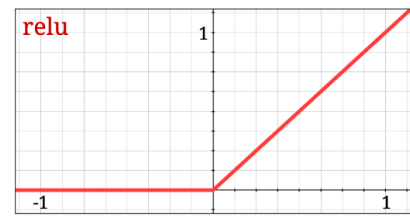
- The gradient can become very small and training get slower and slower.

# Special Care for Deep Networks



$$Y = \text{relu}\left(\sum_i W_i X_i + b\right)$$

*activation* (pointing to 'relu')  
*weights* (pointing to  $W_i$ )  
*inputs* (pointing to  $X_i$ )  
*bias* (pointing to  $b$ )

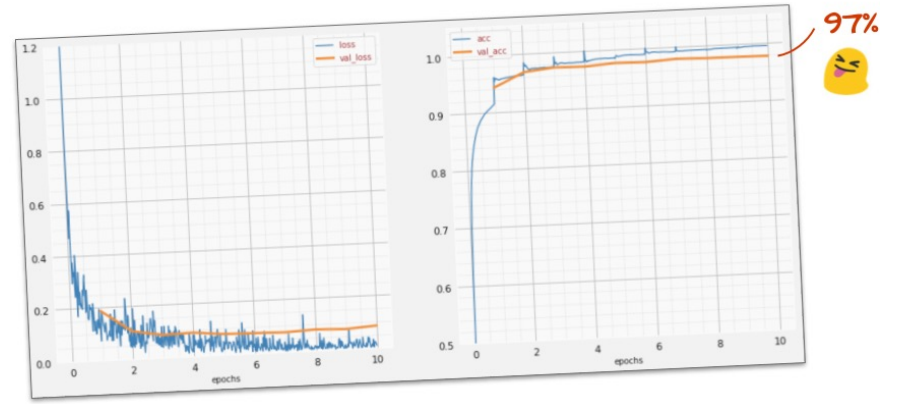


## HANDS ON:

Replace the 'sgd' optimizer with a better one, for example 'adam' and train again.

## HANDS ON:

Replace all `activation='sigmoid'` with `activation='relu'` in your layers and train again.







# DES 5002: Designing Robots for Social Good

Autumn 2022

**Thank you~**

Wan Fang

Southern University of Science and Technology