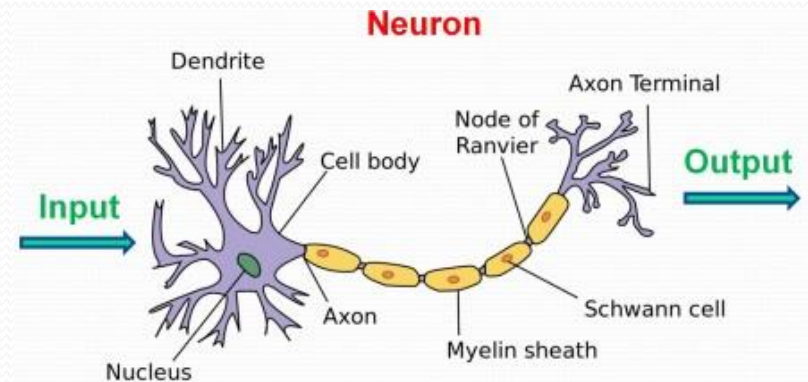# ESS2222

# Lecture 6 – Neural Networks

*Hosein Shahnas*

*University of Toronto, Department of Earth Sciences,*

# Outline

❑ **Combining Perceptrons**
❑ **Optimization**
❑ **Neural Networks**
❑ **Applying SGD & Recursion Relation**
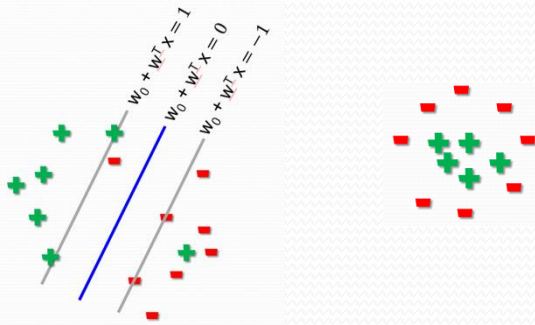❑ **Back Propagation Algorithm**

# Review of Lecture 5

**Soft margin SVM** for **slightly** nonlinear problems
**Kernel method** for **seriously** nonlinear problems

**Minimize** $\dfrac{1}{2}\|w\|^2 \; {}^{+}C\sum_i \xi^i \; , \; \xi^i \geq 0$

**Subject to:** $y^i(\,w_0 + w^T x^i\,) \geq 1 - \xi^i \quad \forall\, i \quad$ where $\xi^i \geq 0$

## One-vs.-All (OvA)



## Image Recognition



2D: 28x28 pixels
2D to 1D array
1D: 784
784
n images, one per line

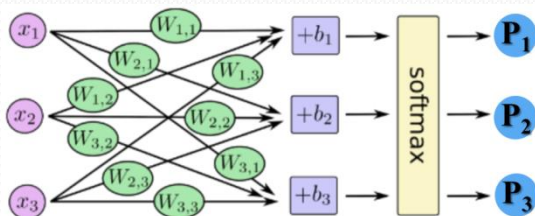## Softmax

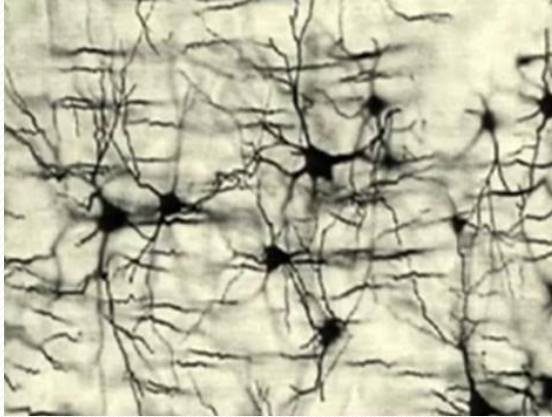**Generalization of the logistic function to multi-class settings**

$\Phi(z_j) = \dfrac{1}{1+e^{-z_j}} \longrightarrow \text{softmax}(z_j) = \dfrac{e^{z_j}}{\sum_{k=1}^{K} e^{z_k}}$

# Biological Neural Structure

**Biology as inspiration**

**Engineering**

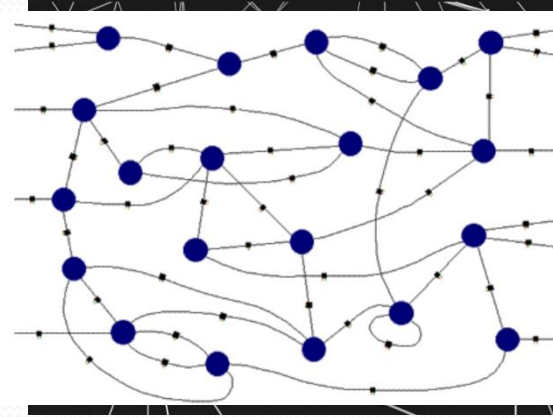**Bilogical function**          →          **Biological structure**





**Perceptrons** are the building blocks of the neural networks connected by synapses. So we may get the human intelligence by combining these building blocks.

**Imitating not exact:** Imitating biology has a limit. The airplane flies but doesn't flap wings! Our engineering does not depend on the details.
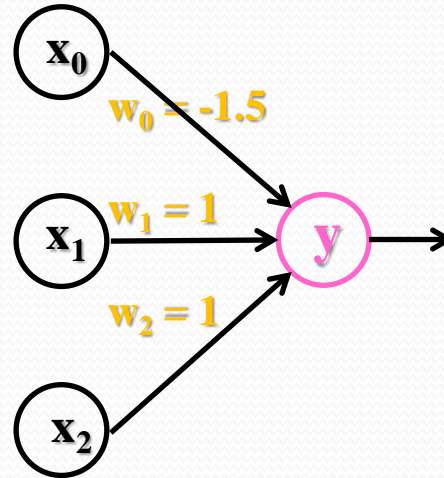
# Combining Perceptrons

Let's **explore** what we can do with combinations of perceptrons rather than single ones.

Let's consider the classification problem for which the **perceptron algorithm failed.**



This problem **cannot** be classified by a single perceptron. But what about with two perceptrons?

# Combining Perceptrons



|  | AND | |
|---|---|---|
| $x_1$ | $x_2$ | Output |
| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

$w_0 = -1.5$

$w_1 = 1$

$w_2 = 1$

🟥 **Class 0**

➕ **Class 1**

$$\sigma \equiv \phi(z) = \begin{cases} 0, & z < 0 \\ 1, & z \geq 0 \end{cases} \quad z = \sum_{i=0}^{n} w_i \, x_i$$

$$\sigma = x_1 \, w_1 + x_2 \, w_2 + w_0$$

$$\sigma = \phi(1 * 0 + 1 * 0 - 1.5) = 0$$
$$\sigma = \phi(1 * 0 + 1 * 1 - 1.5) = 0$$
$$\sigma = \phi(1 * 1 + 1 * 0 - 1.5) = 0$$
$$\sigma = \phi(1 * 1 + 1 * 1 - 1.5) = 1$$

**y**

**And**

$$x_1 \, w_1 + x_2 \, w_2 + w_0 = 0 \qquad x_2 = \frac{-w_1}{w_2} \, x_1 - \frac{w_0}{w_2}$$

# Logic Gates



| OR | | |
|:---:|:---:|:---:|
| $x_1$ | $x_2$ | Output |
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

| NAND | | |
|:---:|:---:|:---:|
| $x_1$ | $x_2$ | Output |
| 0 | 0 | 1 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

| AND | | |
|:---:|:---:|:---:|
| $x_1$ | $x_2$ | Output |
| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

| XOR | | |
|:---:|:---:|:---:|
| $x_1$ | $x_2$ | Output |
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

$x_2$

$x_1$

$w_0 = -10$
$w_1 = 20$
$w_2 = 20$

$w_0 = -30$

$x_1$

$h_1$

$w_1 = 20$

$y$

$w_2 = 20$

$w_1 = -20$
$w_2 = -20$
$w_0 = 30$

$x_2$

$h_2$

$$\sigma \equiv \phi(z) = \begin{cases} 0, & z < 0 \\ 1, & z \geq 0 \end{cases} \quad z = \sum_i w_i x_i$$

$\sigma = \phi(20 * 0 + 20 * 0 - 10) = 0$
$\sigma = \phi(20 * 1 + 20 * 1 - 10) = 1$
$\sigma = \phi(20 * 0 + 20 * 1 - 10) = 1$
$\sigma = \phi(20 * 1 + 20 * 0 - 10) = 1$

$\sigma = \phi(-20 * 0 - 20 * 0 + 30) = 1$
$\sigma = \phi(-20 * 1 - 20 * 1 + 30) = 0$
$\sigma = \phi(-20 * 0 - 20 * 1 + 30) = 1$
$\sigma = \phi(-20 * 1 - 20 * 0 + 30) = 1$

$\sigma = \phi(20 * 0 + 20 * 1 - 30) = 0$
$\sigma = \phi(20 * 1 + 20 * 0 - 30) = 0$
$\sigma = \phi(20 * 1 + 20 * 1 - 30) = 1$
$\sigma = \phi(20 * 1 + 20 * 1 - 30) = 1$

$h_1$

$h_2$

$y$

*Or*

*Nand*

**And**

$$\sigma \equiv \phi(z) = \begin{cases} 0, & z < 0 \\ 1, & z \geq 0 \end{cases} \quad z = \sum_i w_i x_i$$

$\sigma = \phi(20 * 0 + 20 * 0 - 10) = 0$    $\sigma = \phi(-20 * 0 - 20 * 0 + 30) = 1$    $\sigma = \phi(20 * 0 + 20 * 1 - 30) = 0$

$\sigma = \phi(20 * 1 + 20 * 1 - 10) = 1$    $\sigma = \phi(-20 * 1 - 20 * 1 + 30) = 0$    $\sigma = \phi(20 * 1 + 20 * 0 - 30) = 0$

$\sigma = \phi(20 * 0 + 20 * 1 - 10) = 1$    $\sigma = \phi(-20 * 0 - 20 * 1 + 30) = 1$    $\sigma = \phi(20 * 1 + 20 * 1 - 30) = 1$

$\sigma = \phi(20 * 1 + 20 * 0 - 10) = 1$    $\sigma = \phi(-20 * 1 - 20 * 0 + 30) = 1$    $\sigma = \phi(20 * 1 + 20 * 1 - 30) = 1$

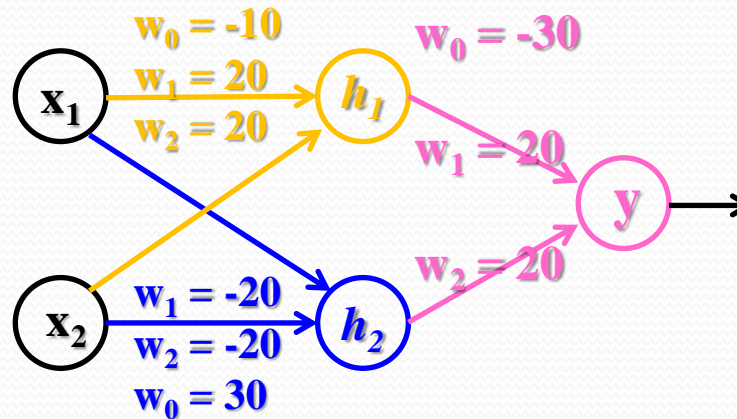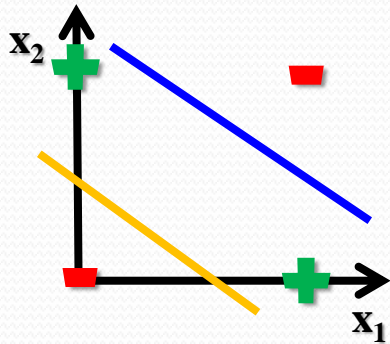$$h_1 \qquad\qquad\qquad h_2 \qquad\qquad\qquad y$$

$$Or \qquad\qquad\qquad Nand \qquad\qquad\qquad And$$

# Combining Perceptrons

# Combining Perceptrons
# Multilayer Perceptron

# Combining Perceptrons
# Multilayer Perceptron

**8-Perceptrons**          **16-Perceptrons**

**We solved this problem using <span style="color:red">feature-transformation</span> before.**

**We can use <span style="color:green">neural networks.</span>**

# Optimization

There are many perceptrons (and therefore many parameters), so optimization might be a problem (remember that for single perceptron when the data were nonlinear, we had convergence problem).

Solution:

1- We chose a soft threshold (tanh) rather than a hard threshold (step function).

2- We use SGD

3- And an efficient way to find weight factors w.

$$\phi(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\phi'(z) = 1 - \tanh(z)^2$$





——— **Step function**

——— **tanh**

**Input**     **Hidden layers**   **output**

**0**              **1≤ $l$ < $L$**              **L**

# Neural Networks

$$w_{ij}{}^l \begin{cases} 1 \leq l < L & \textbf{hidden layers} \\ l = 0 & \textbf{input layer} \\ l = L & \textbf{output layer} \\ 0 \leq i \leq d^{(l-1)} & \textbf{inputs} \\ 1 \leq j \leq d^{(l)} & \textbf{outputs} \end{cases}$$



$$x_j{}^l = \phi\left(z_j{}^{(l)}\right) = \phi\left(\sum_{i=0}^{d^{(l-1)}} w_{ij}^{(l)} \; x_i{}^{(l-1)}\right)$$

$$\begin{bmatrix} x_1{}^{(l)} \\ x_2{}^{(l)} \\ \\ \\ x_{d^{(l)}} \end{bmatrix}$$

$d^{(0)}$ **: dimension of the feature space**

$$\begin{bmatrix} x_0{}^{(0)} \\ x_1{}^{(0)} \\ x_2{}^{(0)} \end{bmatrix} \begin{bmatrix} x_0{}^{(l)} \\ x_1{}^{(l)} \\ x_2{}^{(l)} \\ x_3{}^{(l)} \\ x_4{}^{(l)} \\ x_5{}^{(l)} \end{bmatrix} \begin{bmatrix} x_0{}^{(2)} \\ x_1{}^{(2)} \\ x_2{}^{(2)} \\ x_3{}^{(2)} \\ x_4{}^{(2)} \\ x_5{}^{(2)} \end{bmatrix} \begin{bmatrix} x_1{}^{(3)} \end{bmatrix}$$

# Neural Networks

$$w_{ij}^l \begin{cases} 1 \le l < L & \textbf{\textit{hidden layers}} \\ l = 0 & \textbf{\textit{input layer}} \\ l = L & \textbf{\textit{output layer}} \\ 0 \le i \le d^{(l-1)} & \textbf{\textit{inputs}} \\ 1 \le j \le d^{(l)} & \textbf{\textit{outputs}} \end{cases}$$

$$x_j{}^1 = \phi\left(z_j{}^{(l)}\right) = \phi\left(\sum_{i=0}^{d^{(l-1)}} w_{ij}^{(l)} x_i{}^{(l-1)}\right)$$

$$\begin{bmatrix} x_1^{(l)} \\ x_2^{(l)} \\ \\ \\ x_{d^{(l)}} \end{bmatrix}$$

$w_{ij}{}^1$

$w_{ij}{}^2$

$w_{ij}{}^L$

output

$(ij)$   $(ij)$   $(ij)$   $(ij)$

$d^{(0)}$ : **dimension of the feature space**

$$\begin{bmatrix} x_0^{(0)} \\ x_1^{(0)} \\ x_2^{(0)} \end{bmatrix} \begin{bmatrix} x_0^{(l)} \\ x_1^{(l)} \\ x_2^{(l)} \\ x_3^{(l)} \\ x_4^{(l)} \\ x_5^{(l)} \end{bmatrix} \begin{bmatrix} x_0^{(2)} \\ x_1^{(2)} \\ x_2^{(2)} \\ x_3^{(2)} \\ x_4^{(2)} \\ x_5^{(2)} \end{bmatrix} \begin{bmatrix} x_1^{(3)} \end{bmatrix}$$

15

# Applying SGD

All the weights $\mathbf{w} = \left\{ w_{ij}^{(l)} \right\}$ determine $h(x)$.

Error on sample $(x_n, y_n)$:     $e(w) = e(h(x_n), y_n)$

To implement SGD, we need to calculate:  $\nabla e(\mathbf{w}) = \dfrac{\partial e(w)}{\partial w_{ij}^{(l)}}$  $\forall\, i, j, l$

Computing $\dfrac{\partial e(w)}{\partial w_{ij}^{(l)}}$ :

$\dfrac{\partial e(w)}{\partial w_{ij}^{(l)}}$ can be calculated by perturbing $w_{ij}^{(l)}$ and observing the variations on the error at the output and get numerical estimates for partial derivatives. The problem with this approach is that we have to do this for all $w_{ij}^{(l)}$ .

**But we can obtain a <span style="color:blue">recursion relation</span> and then get all coefficients using this formula.**

$$x_j^{(i)} = \phi\left(z_j^{(l)}\right)$$

$$x_j^{(l)}$$
$$z_j^{(l)}$$
$$w_{ij}^{(l)}$$
$$x_i^{(l-1)}$$

# Recursion Relation

$$\frac{\partial e(w)}{\partial w_{ij}^{(l)}} = \frac{\partial e(w)}{\partial z_j^{(l)}} \frac{\partial z_j^{(l)}}{\partial w_{ij}^{(l)}} \qquad \text{(chain rule)}$$

Let $\qquad \boldsymbol{\delta_j}^{(l)} = \frac{\partial e(w)}{\partial z_j^{(l)}} \qquad \longrightarrow \qquad \frac{\partial e(w)}{\partial w_{ij}^{(l)}} = \frac{\partial z_j^{(l)}}{\partial w_{ij}^{(l)}} \boldsymbol{\delta_j}^{(l)}$

But since $\quad \boldsymbol{z_j}^{(l)} = \sum_{i=0}^{d^{(l-1)}} \boldsymbol{w_{ij}^{(l)}} \boldsymbol{x_i}^{(l-1)} \qquad \longrightarrow \qquad \frac{\partial z_j^{(l)}}{\partial w_{ij}^{(l)}} = x_i^{(l-1)}$

Then: $\frac{\partial e(w)}{\partial w_{ij}^{(l)}} = \frac{\partial z_j^{(l)}}{\partial w_{ij}^{(l)}} \boldsymbol{\delta_j}^{(l)} = x_i^{(l-1)} \boldsymbol{\delta_j}^{(l)}$

The only thing we need is $\qquad \boldsymbol{\delta_j}^{(l)} = \frac{\partial e(w)}{\partial z_j^{(l)}}$

If we can find a recursion relation for $\boldsymbol{\delta_j}^{(l)}$, then we can compute all of them by knowing one of them.

We compute $\boldsymbol{\delta_j}^{(l)}$ <span style="color:red">for the final layer,</span> because if we know $\delta$ later we can obtain $\delta$ earlier <span style="color:red">(back propagation).</span>

# $\delta$ for Final Layer

For l = L,  j = 1,  $\delta_j{}^{(l)} = \dfrac{\partial e(w)}{\partial z_j^{(l)}}$  $\longrightarrow$  $\delta_1{}^{(L)} = \dfrac{\partial e(w)}{\partial z_1^{(L)}}$

We have  $e(w) = e(h(x_n), y_n)$, but for the final layer:  $h(x_n) = \phi(z_1{}^{(L)}) = x_1{}^{(L)}$

Then: $e(w) = e(x_1{}^{(L)}, y_n)$

If $e(w) = (h(x_n) - y_n)^2$  then $e(w) = (x_1{}^{(L)} - y_n)^2$

For tanh-activation function: $\phi(z) = \tanh(z) = \dfrac{e^z - e^{-z}}{e^z + e^{-z}},$  $\phi'(z) = 1 - \tanh(z)^2$

Now we want to calculate: $\delta_i^{(l-1)} = \frac{\partial e(w)}{\partial z_i^{(l-1)}}$

$\delta_j^{(l-1)} = \sum_j^{d(l)} \frac{\partial e(w)}{\partial z_j^{(l)}} \frac{\partial z_j^{(l)}}{\partial x_i^{(l-1)}} \frac{\partial x_i^{(l-1)}}{\partial z_i^{(l-1)}}$ (chain rule)

$\delta_j^{(l-1)} = \sum_j^{d(l)} = \delta_i^{(l)} \ w_{ij}^{(l)} \ \phi'(z_i^{(l-1)})$

$\delta_j^{(l-1)} = (1 - (x_i^{(l-1)})^2) \ \sum_j^{d(l)} \delta_i^{(l)} \ w_{ij}^{(l)}$ (For tanh-activation function)

$\frac{\partial e(w)}{\partial w_{ij}^{(l)}} = x_i^{(l-1)} \ \boldsymbol{\delta_j}^{(l)}$

$\Delta w^{(l)}_{ij} = -\eta \frac{\partial e(w)}{\partial w_{ij}^{(l)}}$

# Back Propogation Algorithm

For tanh-activation function

1 - Initialize $w_{ij}^{(l)}$ at random
2 - For t = 0, 1, 2 ………
3 - pick n $\in$ $\{1, 2, ….. N\}$ (random pickup, i.e. SGD)
4 - Forward compute all $x_j^{(l)}$

5 - Backward compute all $\delta_j^{(l)}$

6 - Update weights $w_{ij}^{(l)} = w_{ij}^{(l)} - x_i^{(l-1)} \delta_j^{(l)}$

7 - Iterate until the stopping criterion is achieved.

8 - Return the final weights $w_{ij}^{(l)}$

Be careful: Initialize $w_{ij}^{(l)}$ at random and not to zero

If we do so, either $x_j^{(l)}$ or $\delta_j^{(l)}$ will become zero and therefore not useful.

$\delta_j^{(l)}$

$w_{ij}^{(l)}$

$x_i^{(l-1)}$

# Remark

Neural networks can be thought as Learned Nonlinear Transform. Note that the nonlinear transformation of features (e.g., polynomial, RBF, etc.) are not learned transformation.

Since in the hidden layers the features are higher order features (leaned features), then we can implement a better learning. Indeed the network looks for weight factors for a proper transform the factors that fits data.



Raw input: features of dimension d

Output

Hidden layers : higher order features or learned features

# Dropout

**Dropout is a <span style="color:red">regularization technique</span> for neural network models <span style="color:blue">(Srivastava, et al. , 2014).</span> This is a simple way to prevent neural networks from overfitting.**

**Some key points:**
**1) Use 20%-50% dropout**
**2) Dropout with larger network in general provides better performance, giving the model more of an opportunity to learn independent representations.**
**3) Dropout can be used on visible (input) as well as hidden layers**



**Standard neural network**

**Neural network with dropout**

## Example 1: MNIST Database - Handwritten digits

A simple sequential deep learning model for handwritten digits recognition using Keras and TensorFlow,

**784**

**512**

**0.2**

**Input**

**10**

**28×28**

**Output**

tf.keras.layers.Flatten(input_shape=(28, 28)
tf.keras.layers.Dense(512, activation=tf.nn.relu)
tf.keras.layers.Dropout(0.2)
tf.keras.layers.Dense(10, activation=tf.nn.softmax)

# A Simple Network
# MNIST handwritten digits

```
1 #============================================================ import tensorflow and keras libraries
2 import tensorflow as tf
3 mnist = tf.keras.datasets.mnist
4 import sys
5 #============================================================ import tensorflow and keras libraries
6
7 #============================================================ import MNIST data
8 # MNIST (Modified National Institute of Standards and Technology database) a large database of handwritten digits
9 (x_train, y_train),(x_test, y_test) = mnist.load_data()
10 x_train, x_test = x_train / 255.0, x_test / 255.0      # scale
11 print('x_train.shape= ', x_train.shape)
12 print('x_test.shape = ', x_test.shape)
13 #sys.exit()
14 #============================================================ import MNIST data
15
16 #============================================================ define a sequential model (deep learning) with keras
17 model = tf.keras.models.Sequential([
18   tf.keras.layers.Flatten(input_shape=(28, 28)),          # flatten data 28x28 --> 784
19   tf.keras.layers.Dense(512, activation=tf.nn.relu),      # dense layer =512
20   tf.keras.layers.Dropout(0.2),                           # dropout 512*0.2
21   tf.keras.layers.Dense(10, activation=tf.nn.softmax)     # output layer = 10
22 ])
23 #============================================================ define a sequential model (deep learning) with keras
24
25 #============================================================ compile the model
26 '''
27 Cross entropy loss function: measures the dissimilarity between the distribution of observed class labels
28 and the predicted probabilities of class lables.
29 Categorical refers to the possibility of having more than two classes.
30 Sparse refers to using a single integer from zero to the number of classes minus one (e.g. { 0; 1; or 2 }
31 for a three-class problem), instead of a dense one-hot encoding of the class label (e.g. { 1,0,0; 0,1,0; or 0,0,1 }
32 for a class label for the same three-class problem).
33 '''
34
35 model.compile(optimizer='adam',                              # adam optimizer
36               loss='sparse_categorical_crossentropy',        # define the loss function
37               metrics=['acc'])                               # define the metric ('accuracy', 'mse', 'msle', 'mae')
38 #============================================================ compile the model
39
40 #============================================================ fit the model
41 model.fit(x_train, y_train, epochs=10)
42 #============================================================ fit the model
43
44 #============================================================ evaluate the model
45 print('=================================================test')
46 model.evaluate(x_test, y_test)
47 print('=================================================test')
48 #============================================================ evaluate the model
49 sys.exit()
50
```
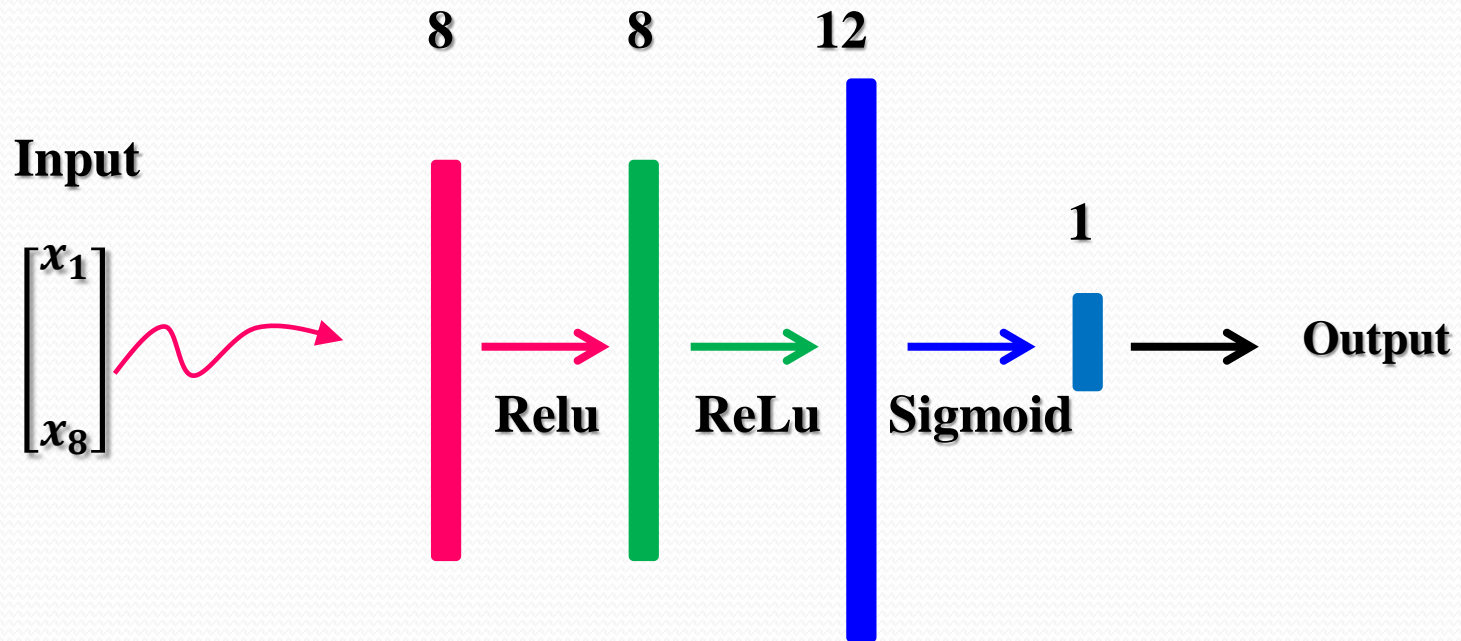
## Example 2: Pima Indians onset of diabetes dataset

A simple sequential deep learning model for predicting handwritten digits using Keras,



```
model.add(Dense(8, input_dim=8, activation='relu'))
model.add(Dense(12, activation='relu'))
model.add(Dense(1, activation='sigmoid'))
```

# A Simple Network
## Onset of Diabetes Dataset

```
 1 #https://machinelearningmastery.com/tutorial-first-neural-network-python-keras/
 2 # Create your first MLP in Keras
 3 #=========================================================
 4 from keras import optimizers
 5 from keras.models import Sequential
 6 from keras.layers import Dense
 7 import keras
 8 import numpy
 9 import sys
10 #=========================================================
11 # fix random seed for reproducibility
12 #numpy.random.seed(7)
13 #========================================================= import data
14 # load pima indians dataset
15 dataset = numpy.loadtxt("pima-indians-diabetes.csv", delimiter=",")
16 #========================================================= import data
17
18 #========================================================= split data inti test and sample data
19 # split into input (X) and output (Y) variables
20 X = dataset[:,0:8]
21 Y = dataset[:,8]
22 print('X.shape = ', X.shape)
23 print('Y.shape = ', Y.shape)
24 #========================================================= split data inti test and sample data
25
26 #========================================================= define a sequential model (deep learning) with keras
27 # create model
28 model = Sequential()
29 model.add(Dense(8, input_dim=8, activation='relu'))
30 model.add(Dense(12, activation='relu'))
31 #model.add(keras.layers.Dropout(0.1))
32 model.add(Dense(1, activation='sigmoid'))
33 #========================================================= define a sequential model (deep learning) with keras
34
```

```
34
35#============================================================ compile the model
36# Compile model                                      # set loss function, optimizer, metrics
37model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
38
39'''
40#model.compile(loss='binary_crossentropy', optimizer='sgd', metrics=['accuracy'])
41sgd = optimizers.SGD(lr=0.01, decay=1e-6, momentum=0.9, nesterov=True)
42model.compile(loss='binary_crossentropy', optimizer=sgd, metrics=['accuracy'])
43
44Adagrad = optimizers.Adagrad(lr=0.01, epsilon=None, decay=0.0)
45model.compile(loss='binary_crossentropy', optimizer=Adagrad, metrics=['accuracy'])
46'''
47#============================================================ compile the model
48
49#============================================================ fit the model
50# Fit the model
51model.fit(X, Y, epochs=1000, batch_size=10)                    # set batch size and epoch-number
52
53#model.fit(X, Y, epochs=500)
54#============================================================ fit the model
55
56#============================================================ evaluate the model
57# evaluate the model
58scores = model.evaluate(X, Y)
59print('===================================== test')
60print("\n%s: %.2f%%" % (model.metrics_names[1], scores[1]*100))
61print('===================================== test')
62#============================================================ evaluate the model
63sys.exit()
64
```