

Linear Algebra Preparation for Emergent Neural Network Architectures (LAPENNA) MLP Example

**Kwai Wong, Stan Tomov
Rocco Febbo, Julian Halloy**

University of Tennessee, Knoxville

Feb. 27, 2021

Acknowledgements:

- Support from NSF, UTK, JICS, ICL, NICS
- LAPENNA, www.jics.utk.edu/lapenna, NSF award #202409
- www.icl.utk.edu, cfdlab.utk.edu, www.xsede.org,
www.jics.utk.edu/recsem-reu,
- Autonomous vehicle, www.bitbucket.org/cfdl/opendnnwheel
- Parallel workflow framework, www.bitbucket.org/cfdl/opendiel
- magmaDNN, www.bitbucket.org/icl/magmadnn
- These projects come out from the RECSEM REU Summer program supported under NSF award #1659502



- ✓ **Introduction to LAPENNA Program**
- ✓ **Section 2 : Machine Learning and DNN**
 - Overview
 - MNIST Test case , TensorFlow example
 - Neural Network Computation, Forward path
 - Example of Forward path calculation
 - Backward path
 - Fundamental Mathematical Theory
 - Optimization
 - Complete the cycle

✓ **XSEDE user code, login, general operations**

Shared with me > LAPENNA material > References > XSEDE ▾ ⚙

Name ↑

Owner



Utilizing XSEDE Resources ⚙

Rocco Febbo

Ensure you have an XSEDE account and have setup 2 factor authentication.

Obtain access to an SSH terminal. If you're on Windows you can use MobAXterm.
(<https://mobaxterm.mobatek.net/>)

Type:

ssh login.xsede.org

Logging In

Enter your password and select 2 factor authentication option. (either phone Duo app)
Once logged in, you will be greeted by the Single Sign on Hub message.

\$> gsissh bridges

XSEDE Access : need to sign in to www.xsede.org

✓ Google drive share space, google colab, colab.research.google.com

Shared with me > LAPENNA material		Owner
Name	↑	
References		me
Google Colab		Julian Halloy
TensorFlow with GPU		Rocco Febbo

← → ⌂ ⌂ https://drive.google.com/drive/folders/19EhWRSgFbMnZoeA0ksc4lONYdiHMGDI9

Most Visited Getting Started

Drive Search in Drive

New Priority My Drive Shared drives Shared with me Recent Starred Trash

Shared with me > LAPENNA material > References

Name	↑
DNN	
General-HPC	
Linear Algebra	
Modeling	
Numerical LA	
XSEDE	

Shared with me > LAPENNA material

Name	↑
0-cs229-linalg-LA-review.pdf	
0-LA-book-Cheney-Math-UCDavis.pdf	
0-LA-vmls-slides.pdf	
0-linear-algebra-DNN.ppt	
0-LinearAlgebra_2016-DNN.ppt	
1-LA-book.pdf	
1-LinearAlgebra_Matlab_Review.ppt	
1-LinearAlgebra-CalTech.pdf	
Documents-LINKS	

“ To learn data science is to do data science”

Google Colab, jupyter notebook
python, R tutorials
upload file to /content
mount your google drive

<https://colab.research.google.com/notebooks/intro.ipynb>

<https://www.youtube.com/watch?v=inN8seMm7UI>

<https://github.com/dataprofessor>

<https://www.youtube.com/watch?v=huAWa0bqxtA>

<https://www.youtube.com/watch?v=Ri1MfaSISW0>

Use this link to run R code on Google COLAB

<https://colab.research.google.com/notebook#create=true&language=r>

✓ [colab.research.google.com,](https://colab.research.google.com)

The screenshot shows the Google Colab interface. At the top, there are two cards: "Google Colab" and "TensorFlow with GPU". Below them is a large button labeled "Open with Google Colaboratory". The main area displays the "TensorFlow with GPU" notebook. The notebook has a table of contents on the left with sections like "Tensorflow with GPU", "Enabling and testing the GPU", and "Observe TensorFlow speedup on GPU relative to CPU". The "Tensorflow with GPU" section contains code for checking GPU availability and performing a convolution operation. The "Enabling and testing the GPU" section provides instructions for enabling GPUs and includes a code cell showing the connection to a GPU. The "Observe TensorFlow speedup on GPU relative to CPU" section includes code for comparing CPU and GPU execution times.

```
#tensorflow_version 2.x
import tensorflow as tf
import timeit

device_name = tf.test.gpu_device_name()
if device_name != '/device:GPU:0':
    raise SystemError('GPU device not found')

print(
    '\n\nThis error most likely means that this notebook is not '
    'configured to use a GPU. Change this in Notebook Settings via the '
    'command palette (cmd/ctrl-shift-P) or the Edit menu.\n\n'
)
raise SystemError('GPU device not found')

def cpu():
    with tf.device('/cpu:0'):
        random_image_cpu = tf.random.normal((100, 100, 100, 3))
        net_cpu = tf.keras.layers.Conv2D(32, 7)(random_image_cpu)
    return tf.math.reduce_sum(net_cpu)

def gpu():
    with tf.device('/device:GPU:0'):
        random_image_gpu = tf.random.normal((100, 100, 100, 3))
        net_gpu = tf.keras.layers.Conv2D(32, 7)(random_image_gpu)
    return tf.math.reduce_sum(net_gpu)

# We run each op once to warm up; see: https://stackoverflow.com/a/45067900
cpu()
gpu()

# Run the op several times.
print('Time (s) to convolve 32x7x7x3 filter over random 100x100x100x3 images '
      '(batch x height x width x channel). Sum of ten runs.')
print('CPU (s):')
cpu_time = timeit.timeit('cpu()', number=10, setup="from __main__ import cpu")
print(cpu_time)
print('GPU (s):')
gpu_time = timeit.timeit('gpu()', number=10, setup="from __main__ import gpu")
print(gpu_time)
print('GPU speedup over CPU: {}'.format(int(cpu_time/gpu_time)))

[User]
Time (s) to convolve 32x7x7x3 filter over random 100x100x100x3 images (batch x height x width x channel). Sum of ten runs.
CPU (s):
3.862475891000031
GPU (s):
0.10837535100017703
GPU speedup over CPU: 35x
```

- ✓ COLAB Introduction <http://www.youtube.com/watch?v=vVe648dJOdl>
- ✓ Google COLAB : <https://www.youtube.com/watch?v=inN8seMm7UI>

List of references: Linear Algebra

- ✓ <https://ocw.mit.edu/courses/mathematics/18-06-linear-algebra-spring-2010/>
- ✓ Deep Learning UC Berkeley STAT-157 2019,
https://www.youtube.com/watch?v=Va8WWRfw7Og&list=PLZSO_6-bSqHQHBCoGaObUljoXAyyqhpFW
- ✓ Linear Algebra – Gilbert Strang:
https://www.youtube.com/watch?v=dtX8iQGQQkI&list=PLZSO_6-bSqHQHBCoGaObUljoXAyyqhpFW&index=4
- ✓ <http://www-math.mit.edu/~gs/>, video course
- ✓ Matrix Applied Linear Algebra-Jim Demmel
https://people.eecs.berkeley.edu/~demmel/ma221_Spr20/Lectures/index.html
- ✓ <https://www.3blue1brown.com/>
- ✓ https://www.youtube.com/playlist?list=PLZHQBObOWTQDPD3MizzM2xVFitgF8hE_ab
- ✓ <https://www.math.ucdavis.edu/~linear/linear-guest.pdf>
- ✓ <https://comp-think.github.io/>
- ✓ <https://edu.google.com/resources/programs/exploring-computational-thinking/>
- ✓ <https://www.youtube.com/watch?v=Nc-V948dXWI>
- ✓ Google drive → Linear Algebra
- ✓ Many more

Four Tiers – Computational Ecosystem

Advanced Computing Architectures

- Emergent Architectures
- Tactical Computing
- Next Generation Computing Systems
- High Performance Networking and Memory

High Performance

Predictive Simulation Sciences

- Computational Math & Algorithms
- Scientific Computing
- Verification, Validation & Uncertainty Quantification
- Applied Computer Modeling and Analysis

Computing Sciences

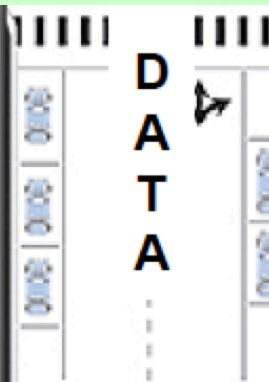
- Programming Environments
- Programming Languages
- Software Integration

Compute Ecosystem

Computing

Data Intensive Sciences

- Sciences of Large Data
- Computational Math for Data Analytics
- Real-time Data Access & Analytics



Modeling Cycle and Emergent Technology

Big Data → Deep Learning → Computing capacity

Combining the Learned the pattern and behavior hidden in sensor or
computed data to predict complicated phenomena

A growing field with evolving technology

Data- + model-based simulation --> Supercomputer

Fast enough to get approximate results (time)

Formulation, Algorithm, linear algebra

Enough memory to process the huge amount of data (length scale)

Length, size, reduction, representation, projection

Exercises

1. A desktop computer has the following specification: quad-core, 24GB RAM, 2.5 GHz, and each core has 8 floating units (perform 8 DP operations in one clock cycle). What is the theoretical double precision peak performance of the desktop in FLOPS?
2. Given A is a $M \times K$ matrix, B is a $K \times N$ matrix, what is the number of floating point operation (FLOP) of $A \times B$?
3. Given $M=K=N$ for A and B , what is the maximum dimension N that the computer can hold?
4. Based on the answer above, what is the time needed to compute $C=A \times B$ if 90% of the theoretical rate of the desktop can be attained for the BLAS3 operation?

Numbers : Lots of Them: bit, byte, FLOP (S)

- Core : computing unit : processor
- Dual core machine (Intel or AMD CPU) : a CPU with 2 cores, each core is a 2.4 GHz computing unit with 2GB of RAM (memory in the processor not disk space)
- Binary bits (b) : “0” or “1”, 1 Byte (B) = 8 bits
- Binary number : $11111111 = (2^7 + 2^6 + 2^5 + 2^4 + 2^3 + 2^2 + 2^1 + 2^0) = (2^8 - 1) = 255 !!$
- **32 bits** machine or operating system => largest integer (all positive) = $(2^{32} - 1) = (4,294,967,296 - 1)$ or range of integer = $-(2^{31})$ to $(2^{31} - 1)$
- **64 bits** machine or operating system => range of integer = $-(2^{63})$ to $(2^{63} - 1)$
- Kilo (K) = 10^3 (or 2^{10}) ; Mega (M) = 10^6 (or 2^{20}); **Giga (G) = 10^9 (or 2^{30}); Tera (T billion) = 10^{12} (or 2^{40}) ; Peta (P) = 10^{15} (or 2^{50})**
- GB = 1000^3 bytes, GiB = 1000^3 Bytes
- **FLoating Point Operation (+, -, / , *)** : $(10.1 + 0.1) * 1.0 / 2.0 = 5.1 \Rightarrow 3 \text{ FLOP}$
- FLOPS = FLOP per second :: 1 PetaFLOPS (kraken) = **10^{15} FLOP in one second**
- **FLOPS in a core = (clock rate) x (floating point operation in one clock cycle)**
- **Peak Rate = (FLOPS in one compute unit, core) x (no. of core)**

1. A desktop computer named MA6633DC has the following specification: quad-core, 24GB RAM, 2.5 GHz, and each core has 8 floating units (perform 8 DP operations in one clock cycle). What is the theoretical double precision peak performance of the desktop in FLOPS?

FLOPS in one core= (clock rate) x (floating point operation in one clock cycle)

Peak Rate = (FLOPS in one core) x (no. of core)

Peak FLOPS = 2.5 GHz * 8 * 4 cores = 80 GFLOPS

2. Given A is a M x K matrix, B is a K x N matrix, what is the number of floating point operation (FLOP) of A x B?

C=A x B ; A : A(M , K) ; B : B(K, N) , C = C(M, N) ; one element of C need (2K-1) DP

FLOP = (2K-1) x M x N

3. Given M=K=N for A and B, what is the maximum dimension N that MA6633DC can hold?

For double precision float, we need 8 bytes to store each element.

If M=N=K=N, in total we need $3N^2$ elements

$3N^2 * 8 = 24 \text{ GB}$, thus **N = floor(sqrt(1e9)) = 31622**

4. Based on the answer above, what is the time needed to compute C=A x B if 90% of the theoretical rate of MA6633DC can be attained for the BLAS3 operation?

$90\% * 80 \text{ GFLOPS} = 72 \text{ GFLOPS}$

The FLOP we need to compute A x B is $\sim 2N^3$; **Time = 2 (31622)³ FLOP / 72 GFLOPS = 878 seconds**

Question #1

- a) A supercomputer is performing a HPL test for submission to the top500 list. The computer has 5000 compute units. Each compute unit has 20 cores, 64 GB RAM, 2.5 GHz, and each core has 8 floating units (performing 8 DP operations in one clock cycle). What is the theoretical double precision peak performance of the supercomputer in TFLOPS?
- b) (Assuming solving $Ax=b$ requires $2N^3/3$ FLOP, where A is a $N \times N$ matrix, what is the maximum size of A (N) that the computer can solve? (only need to store A)
- c) Using the N in b), what is the time needed to solve $Ax=b$ based on the peak rate obtained in a)?

Question #1

- a) A supercomputer is performing a HPL test for submission to the top500 list. The computer has 5000 compute units. Each compute unit has 20 cores, 64 GB RAM, 2.5 GHz, and each core has 8 floating units (performing 8 DP operations in one clock cycle). What is the theoretical double precision peak performance of the supercomputer in TFLOPS?
- b) (Assuming solving $Ax=b$ requires $2N^3/3$ FLOP, where A is a $N \times N$ matrix, what is the maximum size of A (N) that the computer can solve? (only need to store A)
- c) Using the N in b), what is the time needed to solve $Ax=b$ based on the peak rate obtained in a)?

Linear Algebra



David Cherney, Tom Denton,
Rohit Thomas and Andrew Waldron

Caltech Division of the Humanities
and Social Sciences

Quick Review of Matrix and Real Linear Algebra

KC Border
Subject to constant revision
Last major revision: December 12, 2016
v. 2019.10.23:14.52

Introduction to Linear Algebra

Mark Goldman
Emily Mackevicius

Lecture slides for

Introduction to Applied Linear Algebra:
Vectors, Matrices, and Least Squares

Stephen Boyd Lieven Vandenberghe

Linear Algebra Review and Reference

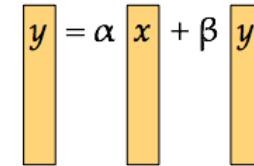
Zico Kolter (updated by Chuong Do)

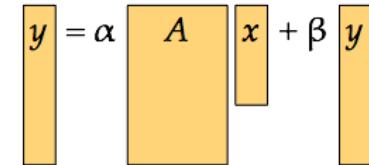
September 30, 2015

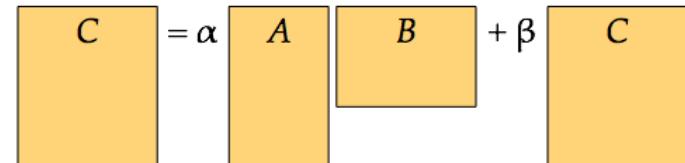
- ✓ **Level 1 BLAS — vector operations**
 - ✓ $O(n)$ data and flops (floating point operations)
 - ✓ Memory bound:
 $O(1)$ flops per memory access

- ✓ **Level 2 BLAS — matrix-vector operations**
 - ✓ $O(n^2)$ data and flops
 - ✓ Memory bound:
 $O(1)$ flops per memory access

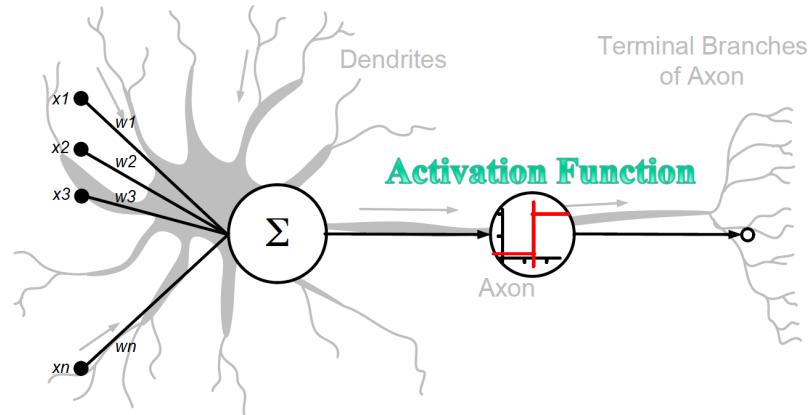
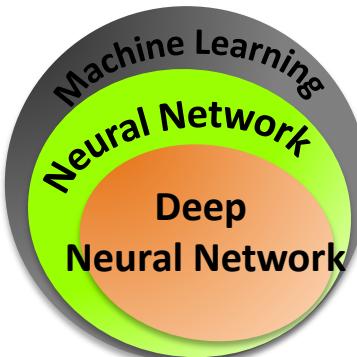
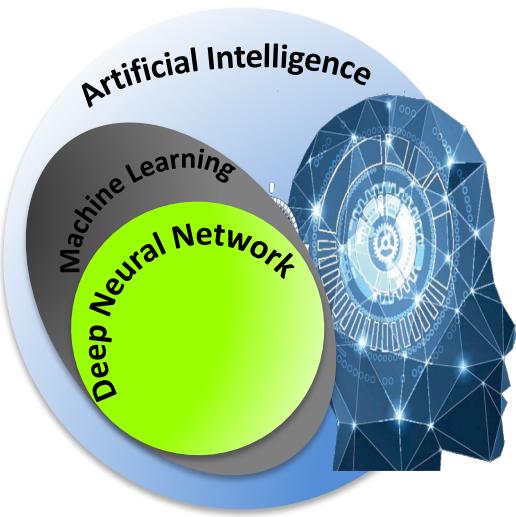
- ✓ **Level 3 BLAS — matrix-matrix operations**
 - ✓ $O(n^2)$ data, $O(n^3)$ flops
 - ✓ Surface-to-volume effect
 - ✓ Compute bound:
 $O(n)$ flops per memory access

$$y = \alpha x + \beta y$$


$$y = \alpha A x + \beta y$$


$$C = \alpha A B + \beta C$$


AI, ML, NN, DNN

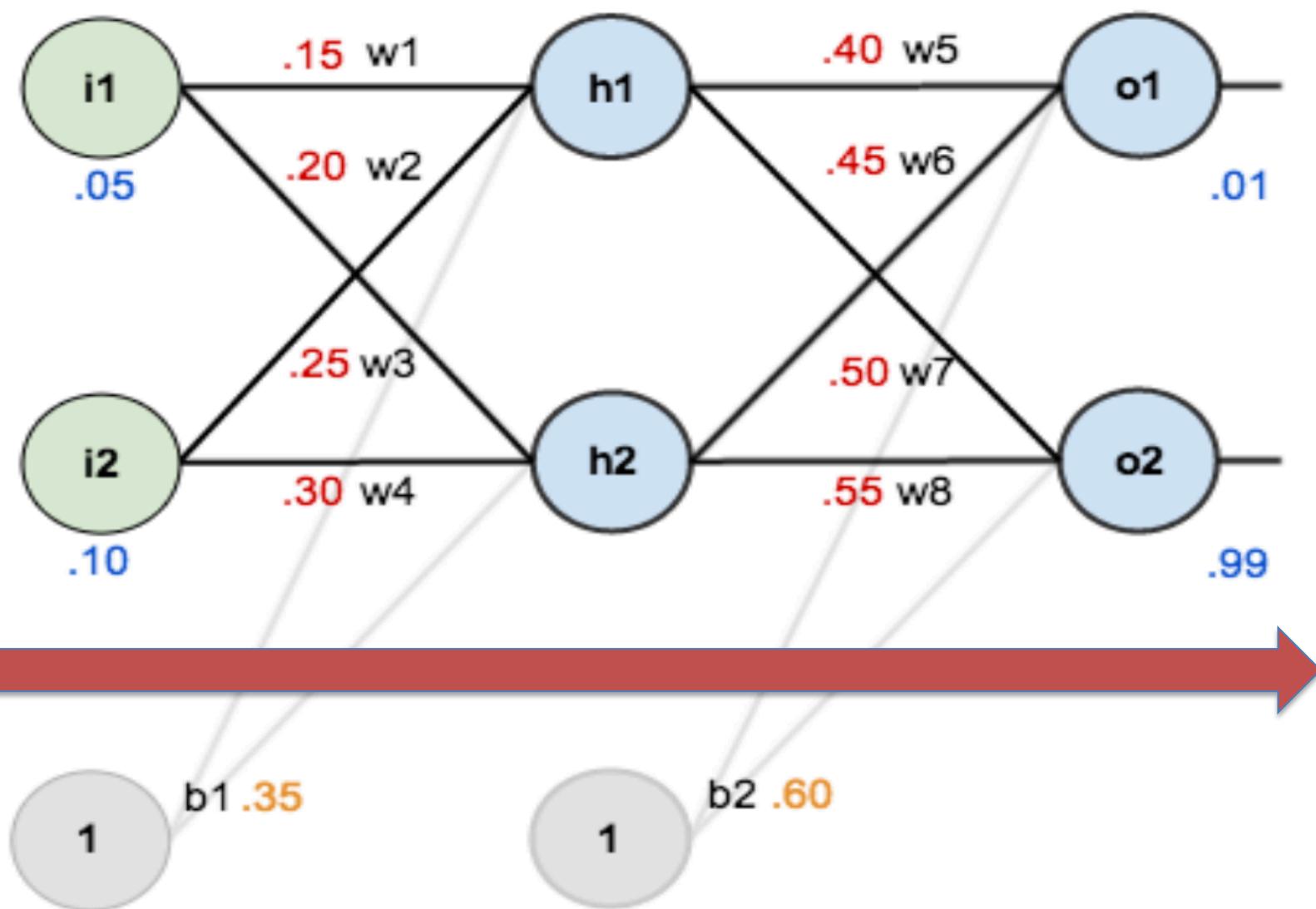


- ✓ **Artificial Intelligence (AI)**: science and engineering of making intelligent machines to perform the human tasks (John McCarthy, 1956). AI applications is ubiquitous.
- ✓ **Machine learning (ML)** : A field of study that gives computers the ability to learn without being explicitly programmed (Arthur Samuel, 1959). A computer program is said to learn from experience E with respect to some task T and some performance measure P, if its performance on T, as measured by P, improves with experience E (Tom Mitchell, 1998).
- ✓ **Neural Network (NN)** : Neural Network modeling, a subfield of **ML** is algorithm inspired by structure and functions of biological neural nets
- ✓ **Deep Neural Network (DNN)** : (aka deep learning): an extension of **NN** composed of many layers of functional neurons, is dominating the science of modern **AI** applications
- ✓ **Supervised Learning (SL)** : A class in ML, dataset has labeled values, use to predict output values associated with new input values.

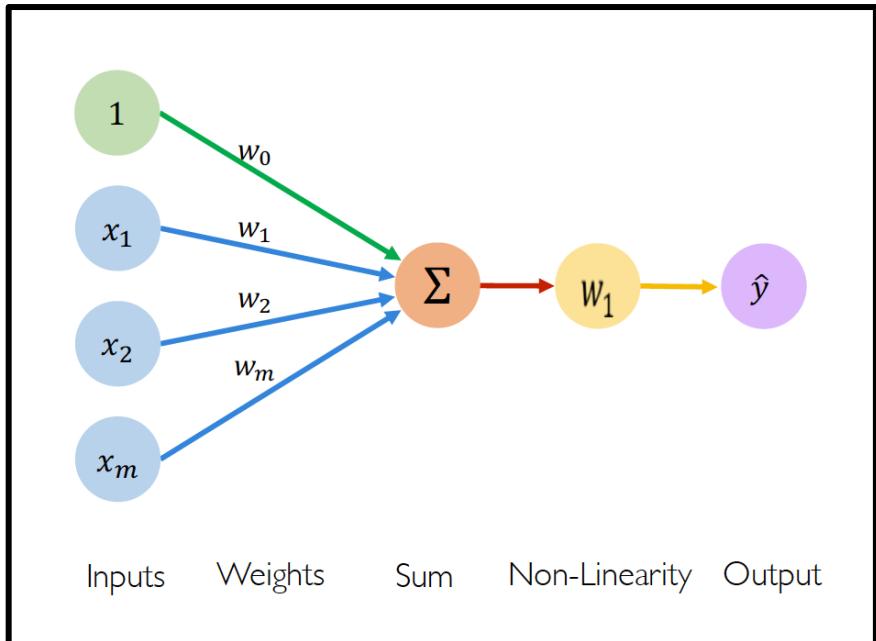
✓ **Section 2 : Machine Learning and DNN**

- **Neural Network Computation, Forward path**
- **Example of Forward path calculation**

Multilayer Perceptron : Forward Path Calculation

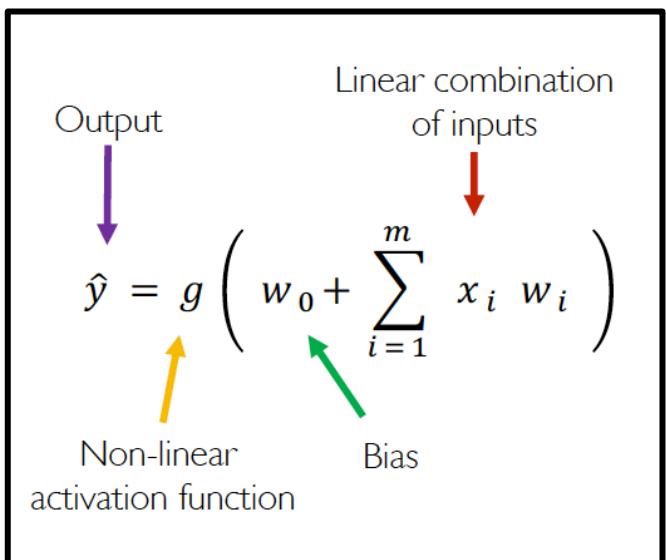


DNN MLP Forward Steps



$$\hat{y} = g(w_0 + X^T W)$$

where: $X = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$ and $W = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}$

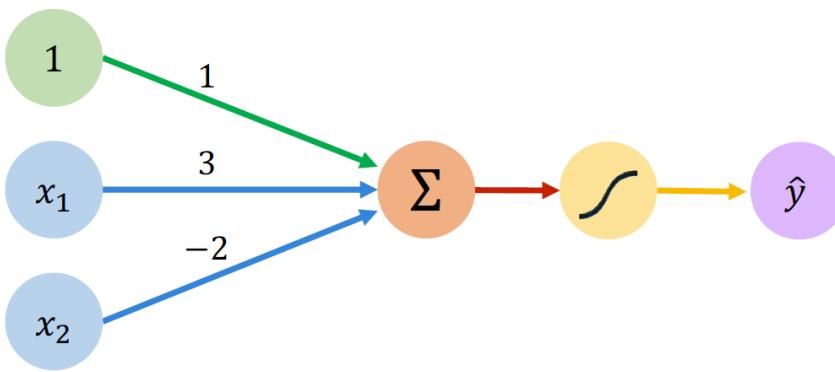


$$\hat{y} = g(w_0 + X^T W)$$

- Example: sigmoid function

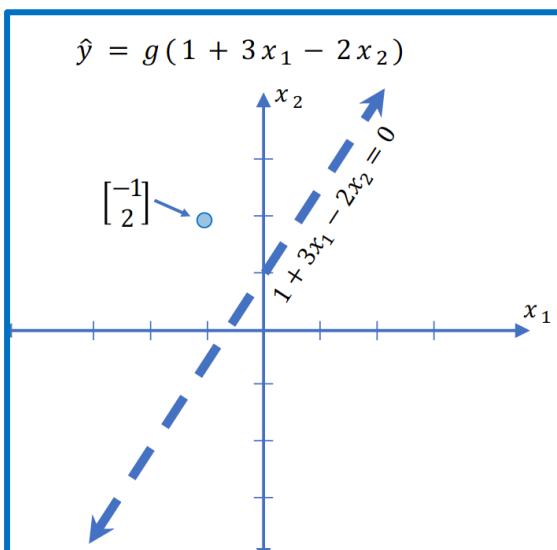
$$g(z) = \sigma(z) = \frac{1}{1 + e^{-z}}$$

Activation Function Example



Assume we have input: $\mathbf{X} = \begin{bmatrix} -1 \\ 2 \end{bmatrix}$

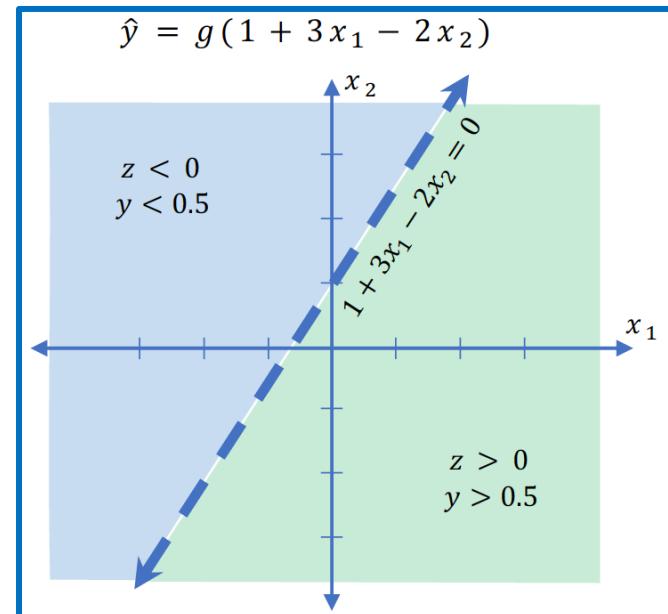
$$\begin{aligned}\hat{y} &= g(1 + (3 * -1) - (2 * 2)) \\ &= g(-6) \approx 0.002\end{aligned}$$



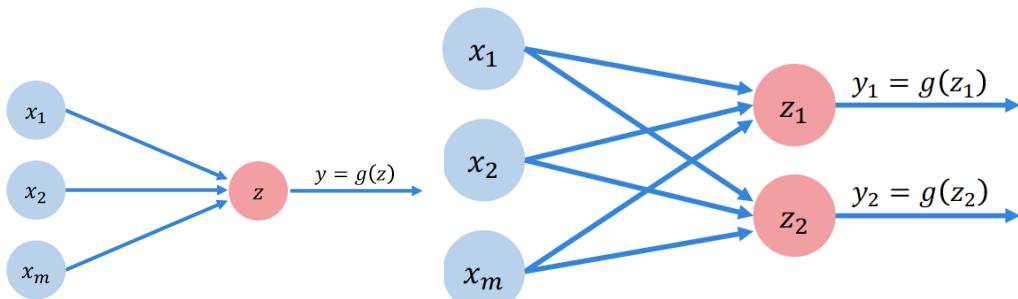
We have: $w_0 = 1$ and $\mathbf{W} = \begin{bmatrix} 3 \\ -2 \end{bmatrix}$

$$\begin{aligned}\hat{y} &= g(w_0 + \mathbf{X}^T \mathbf{W}) \\ &= g\left(1 + \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^T \begin{bmatrix} 3 \\ -2 \end{bmatrix}\right) \\ \hat{y} &= g(1 + 3x_1 - 2x_2)\end{aligned}$$

This is just a line in 2D!

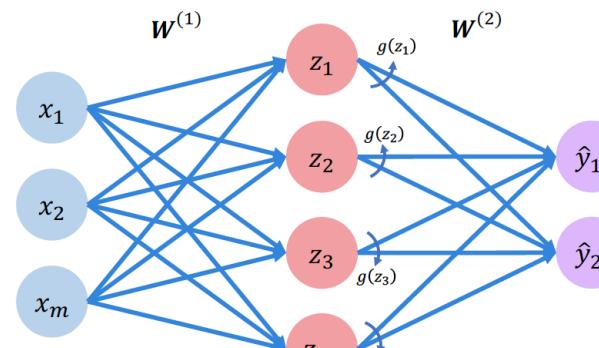


Multilayer Perceptron : Forward Path Calculation



$$z = w_0 + \sum_{j=1}^m x_j w_j$$

$$z_i = w_{0,i} + \sum_{j=1}^m x_j w_{j,i}$$

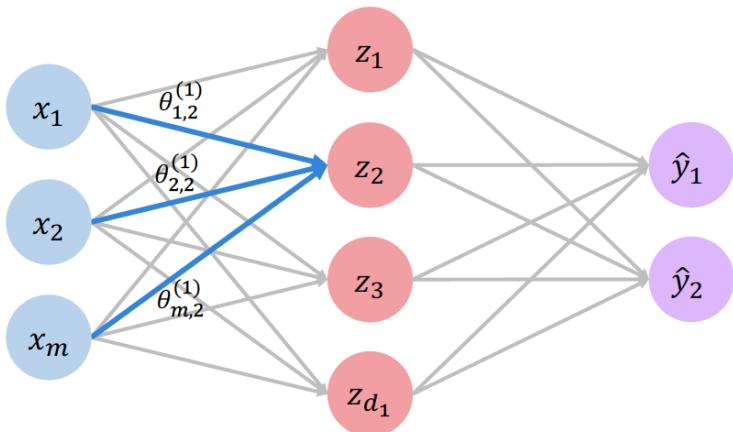


Inputs

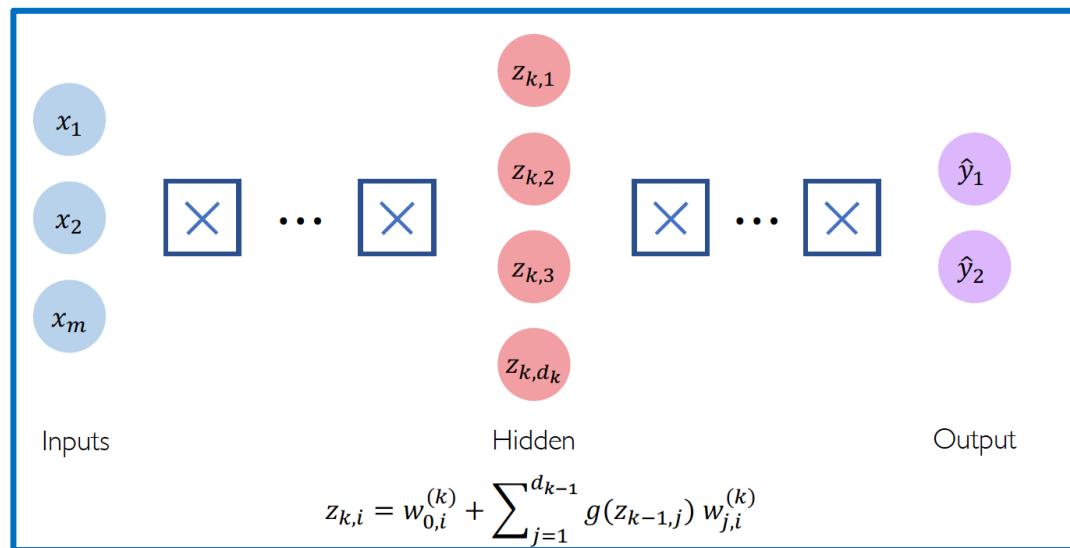
Hidden

Final Output

$$z_i = w_{0,i}^{(1)} + \sum_{j=1}^m x_j w_{j,i}^{(1)} \quad \hat{y}_i = g\left(w_{0,i}^{(2)} + \sum_{j=1}^{d_1} z_j w_{j,i}^{(2)}\right)$$



$$\begin{aligned} z_2 &= w_{0,2}^{(1)} + \sum_{j=1}^m x_j w_{j,2}^{(1)} \\ &= w_{0,2}^{(1)} + x_1 w_{1,2}^{(1)} + x_2 w_{2,2}^{(1)} + x_m w_{m,2}^{(1)} \end{aligned}$$



Inputs

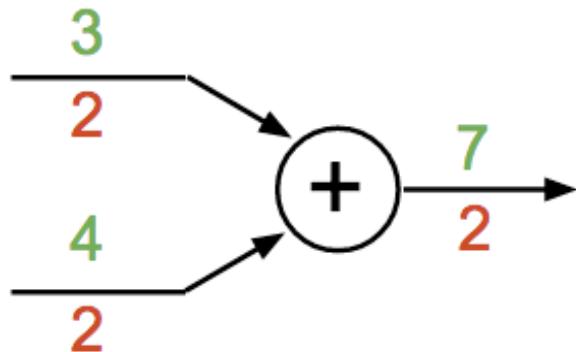
Hidden

Output

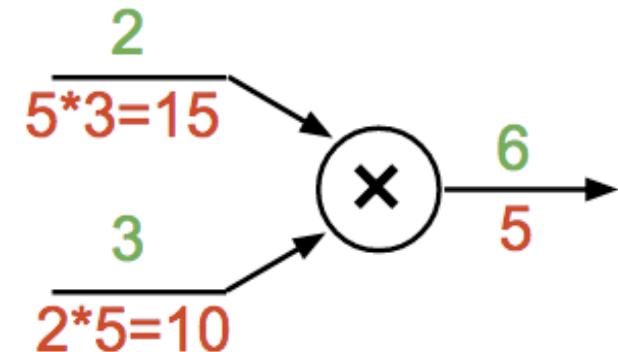
$$z_{k,i} = w_{0,i}^{(k)} + \sum_{j=1}^{d_{k-1}} g(z_{k-1,j}) w_{j,i}^{(k)}$$

Forward path : backward path

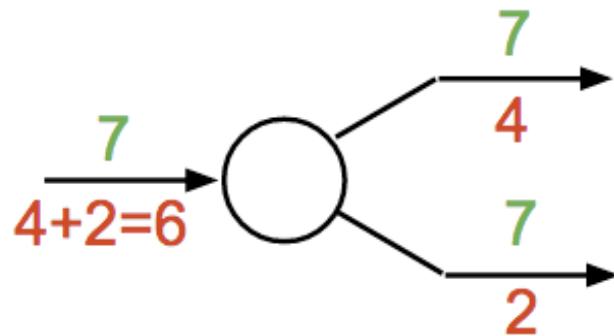
add gate: gradient distributor



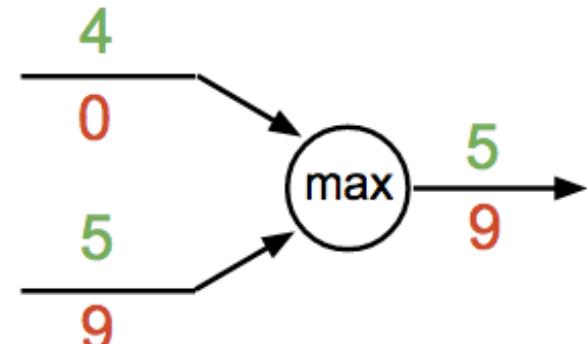
mul gate: “swap multiplier”



copy gate: gradient adder



max gate: gradient router

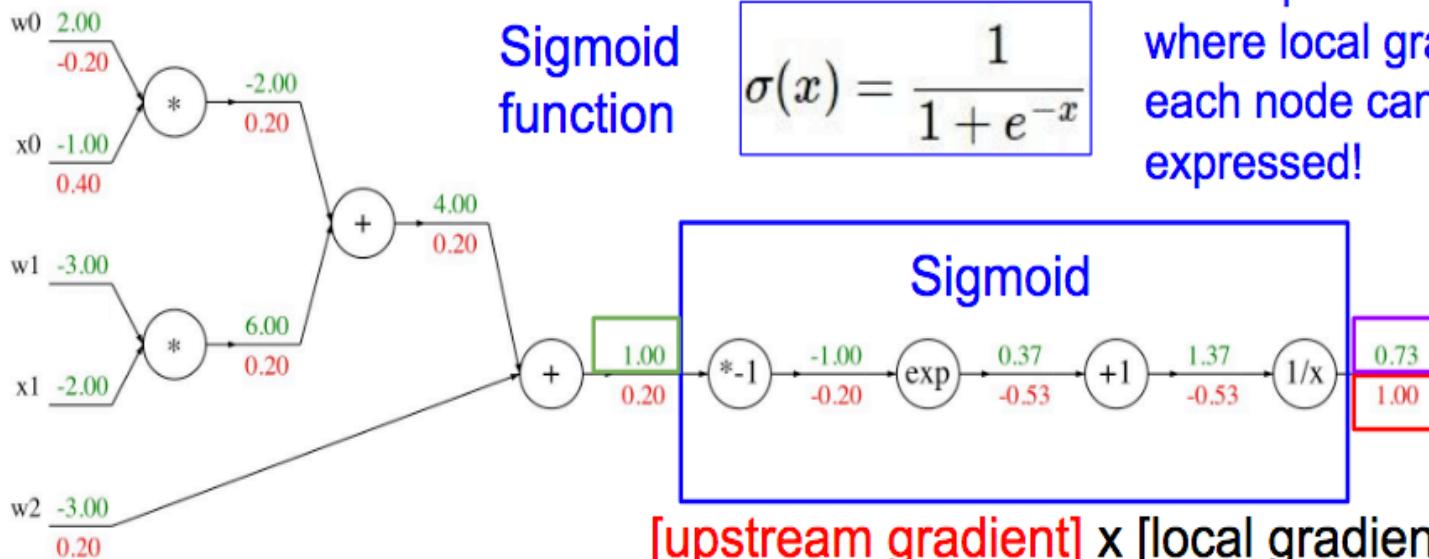


Computational Graph

<http://cs231n.stanford.edu/syllabus.html>

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

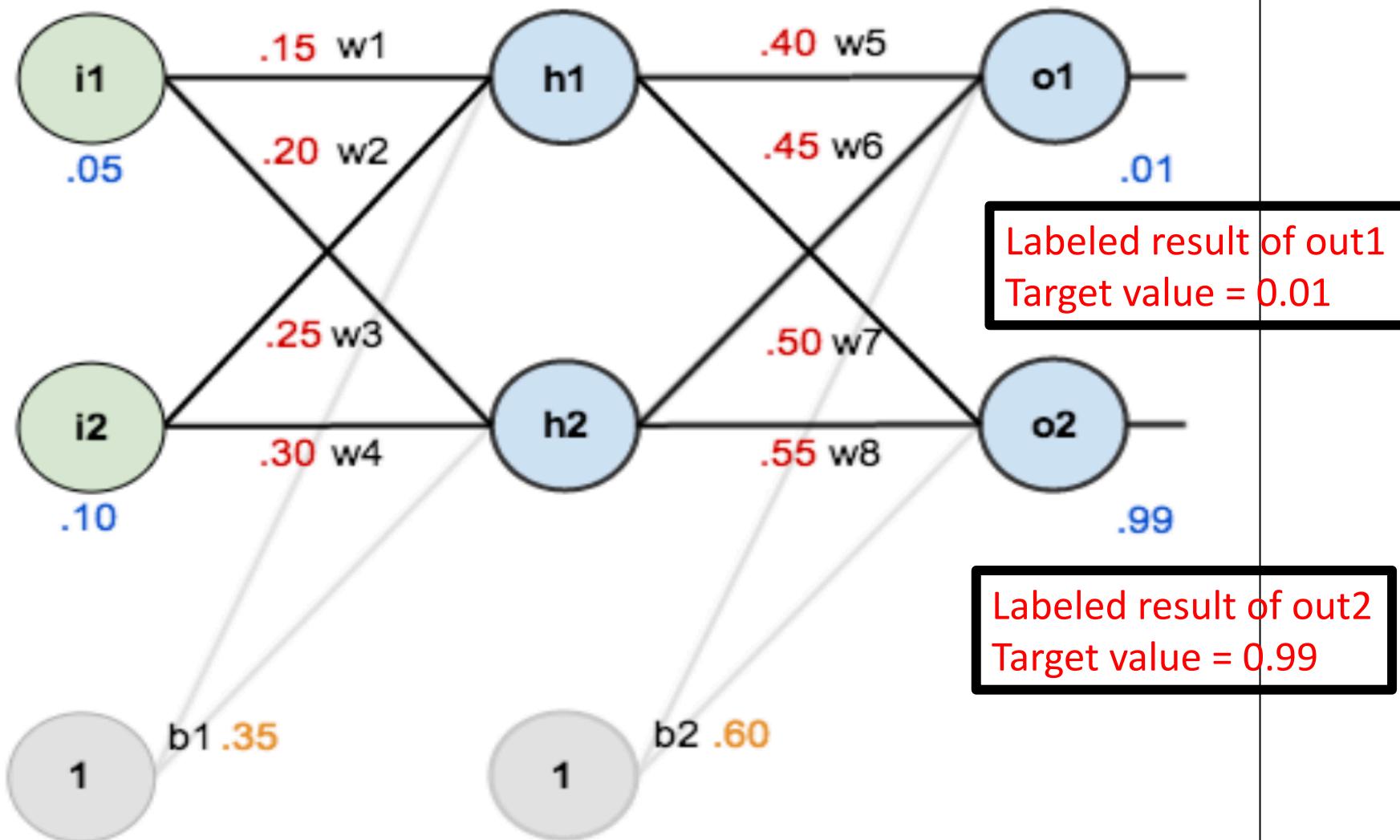


Sigmoid local
gradient:

$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left(\frac{1 + e^{-x} - 1}{1 + e^{-x}} \right) \left(\frac{1}{1 + e^{-x}} \right) = (1 - \sigma(x))\sigma(x)$$

Computational graph representation may not be unique. Choose one where local gradients at each node can be easily expressed!

DNN Example



DNN Example (1)

$w_1 = 0.15$



$i_1 = 0.05$



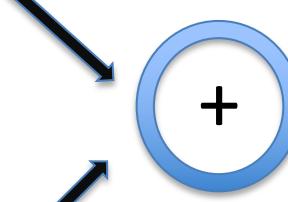
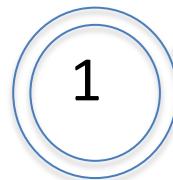
$w_2 = 0.2$



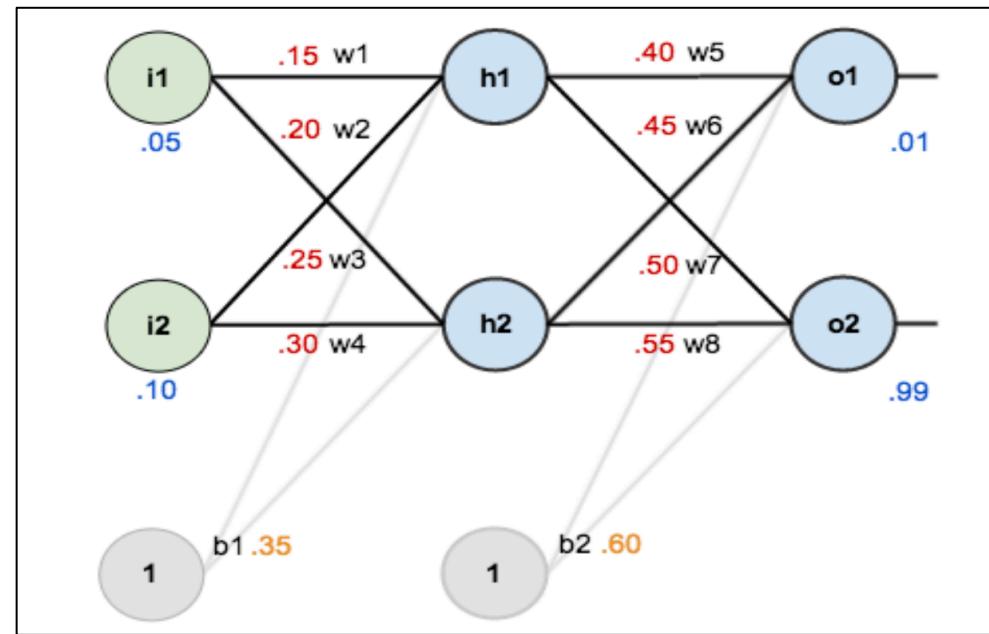
$i_2 = 0.1$



$b_1 = 0.35$



$h_1 = 0.59326992$



DNN Example (2)

$w3=0.25$



$i1 = 0.05$



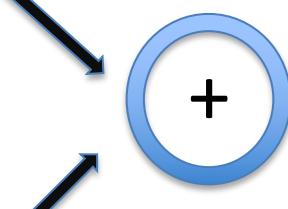
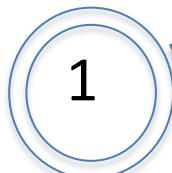
$w4=0.3$



$i2 = 0.1$



$b1=0.35$



$i1 \cdot w3 + i2 \cdot w4 + b1$

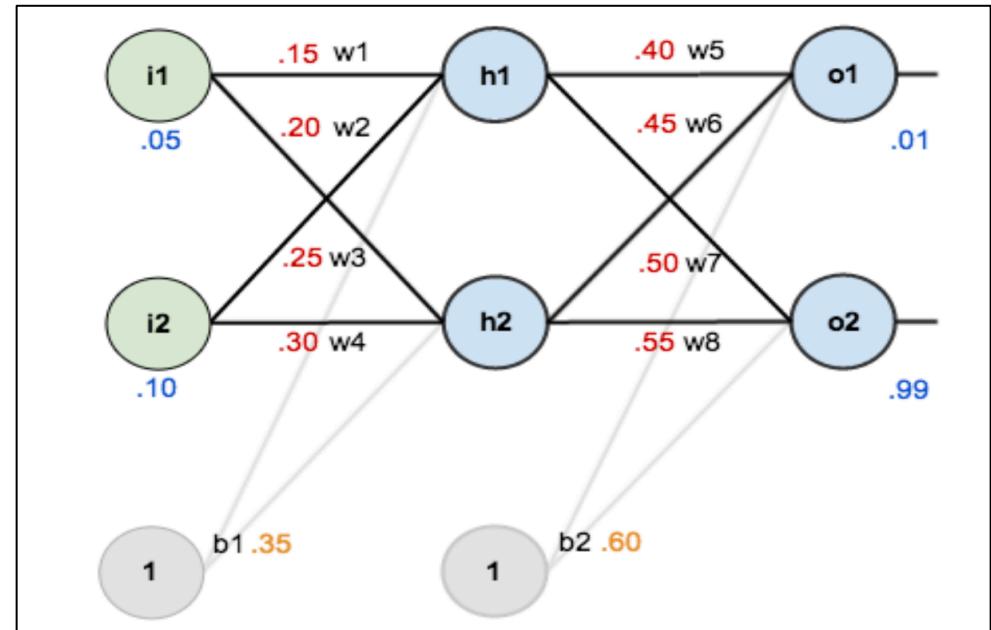
$+$

$+$

$+$



$h2 = 0.596884378$

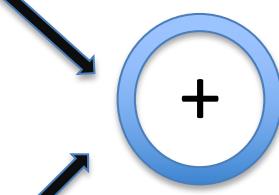


DNN Example (3)

$$w5=0.4$$



$$h1 = 0.59327$$



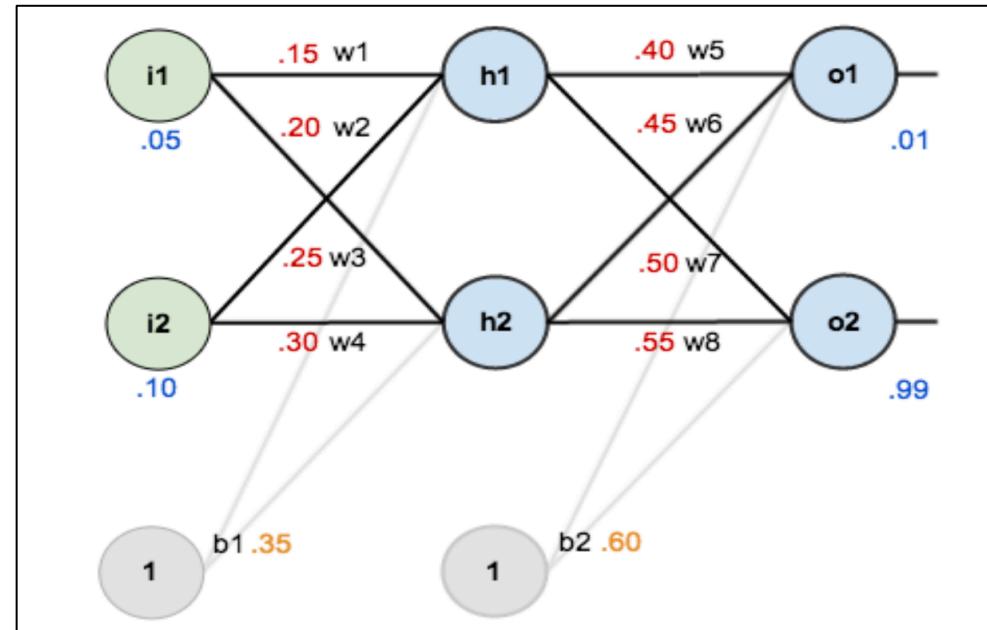
$$w6=0.45$$



$$h2 = 0.59688$$



$$b2=0.6$$



$$o1 = 0.75136$$



DNN Example (4)

$$w7=0.5$$



$$h1 = 0.59327$$



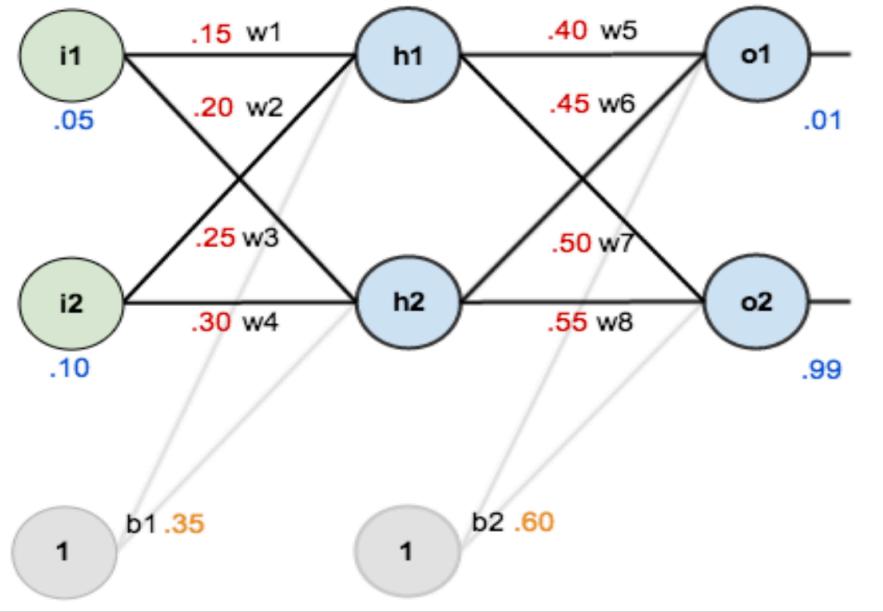
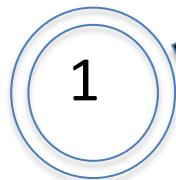
$$w8=0.55$$



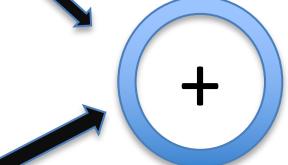
$$h2 = 0.59688$$



$$B2=0.6$$



$$+$$



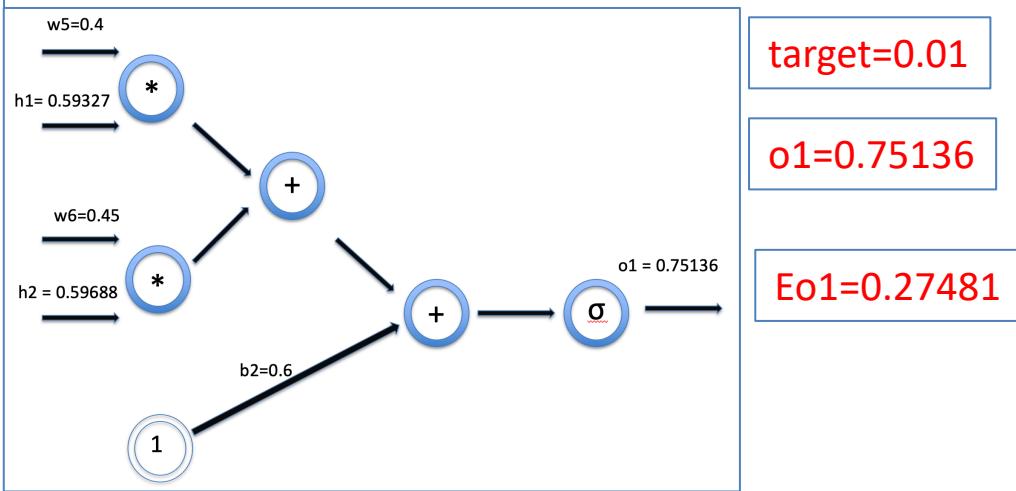
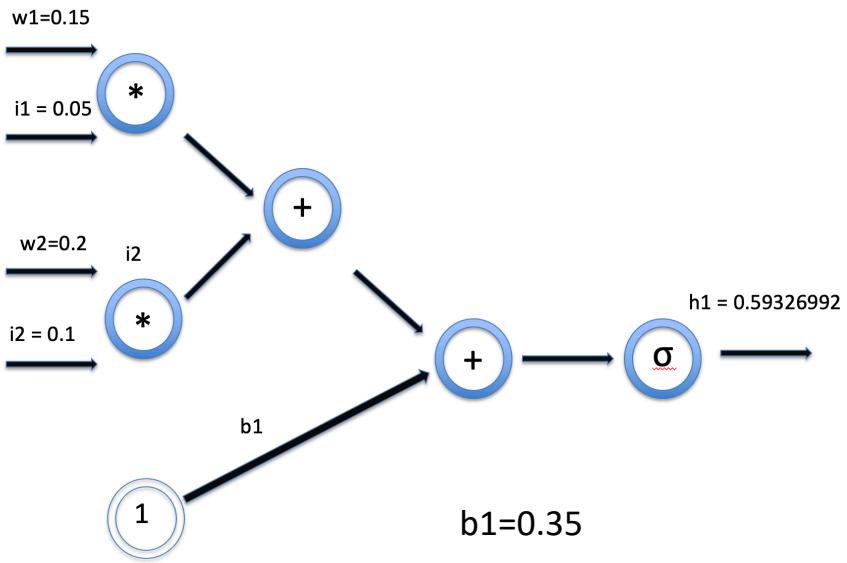
$$+$$



$$o2 = 0.77293$$

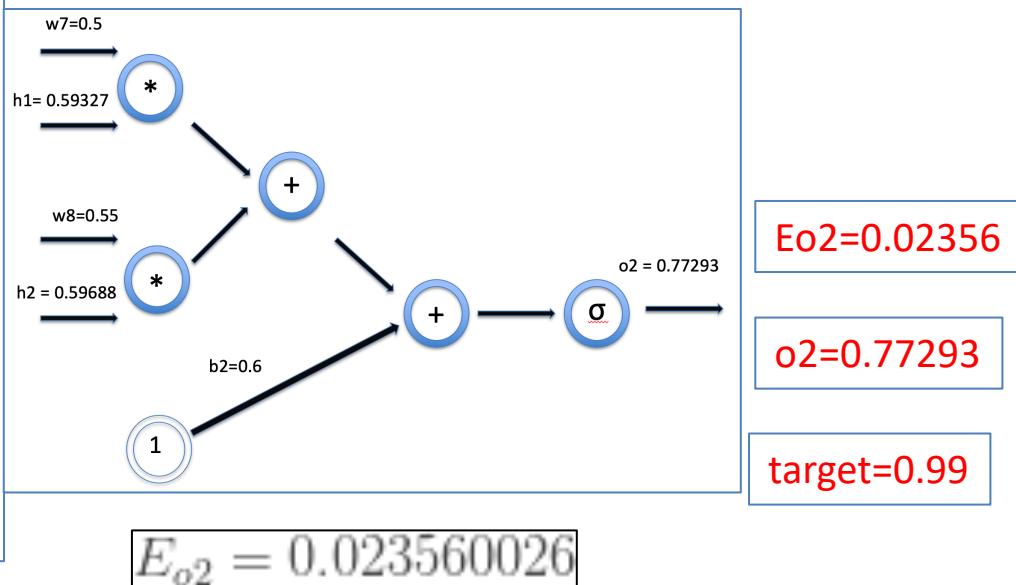
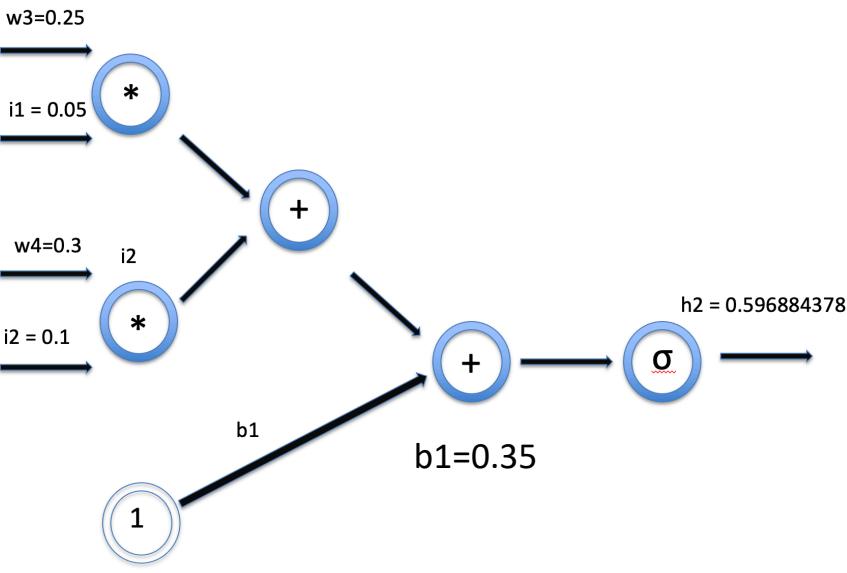


$$E_{o1} = \frac{1}{2}(target_{o1} - out_{o1})^2 = \frac{1}{2}(0.01 - 0.75136507)^2 = 0.274811083$$



$$E_{total} = \sum \frac{1}{2}(target - output)^2$$

$$E_{total} = E_{o1} + E_{o2} = 0.298371109$$



```

# e constant
e = 2.7182818284
# initial values
i1 = 0.05
i2 = 0.10
# initial weights
w1 = 0.15
w2 = 0.20
w3 = 0.25
w4 = 0.30
w5 = 0.40
w6 = 0.45
w7 = 0.50
w8 = 0.55
# bias
b1 = 0.35
b2 = 0.60
# targets
To1 = 0.01
To2 = 0.99

```

```

# forward propagation
h1 = 1/(1+e**(-(w1*i1 + w2*i2+b1)))
print("h1: " + str(h1))

h2 = 1/(1+e**(-(w3*i1 + w4*i2+b1)))
print("h2: " + str(h2))

o1 = 1/(1+e**(-(w5*h1 + w6*h2+b2)))
print("o1: " + str(o1))

o2 = 1/(1+e**(-(w7*h1 + w8*h2+b2)))
print("o2: " + str(o2))

```

h1: 0.5932699921052087 h2: 0.5968843782577157 o1:
 0.7513650695475076 o2: 0.772928465316421

```

# Error
Eo1 = 0.5*(To1-o1)**2
print("Error o1: " + str(Eo1))
Eo2 = 0.5*(To2-o2)**2
print("Error o2: " + str(Eo2))

E = Eo1 + Eo2
print("Total Error: " + str(E))

```

Error o1: 0.2748110831725904 Error o2:
 0.023560025584942117 Total Error: 0.29837110875753253

“Training a neural network”

What does it mean?

What does the network train?

What are the training parameters

How does it work?

“Training a neural network”

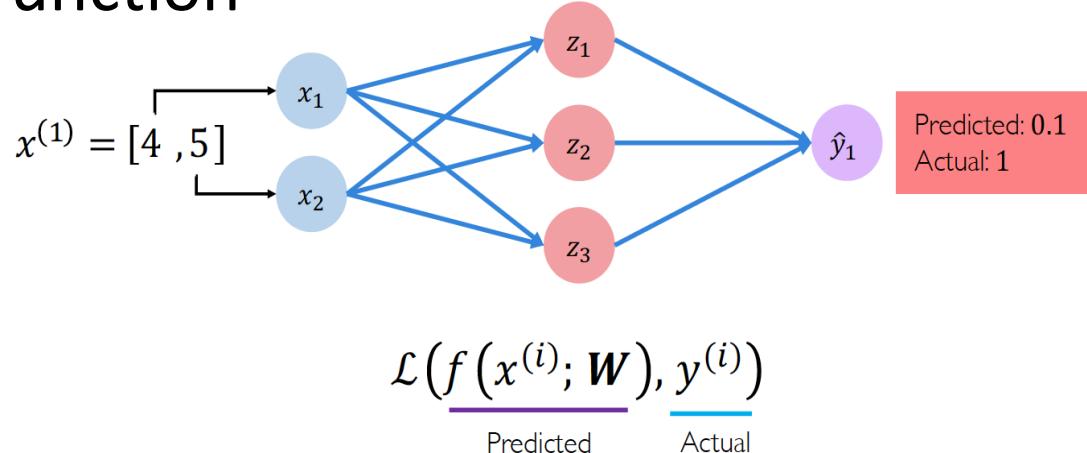
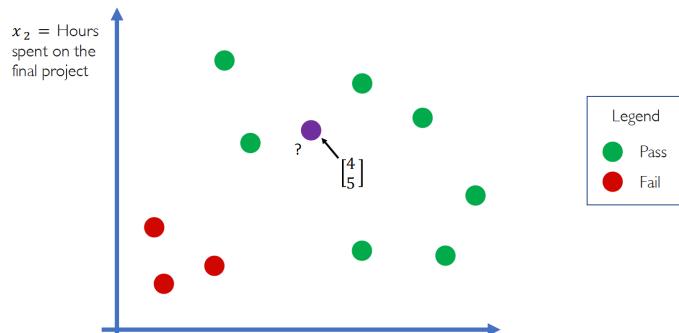
Supervised training

Compare the difference between the
computed and the labeled given
values

Loss values – loss function

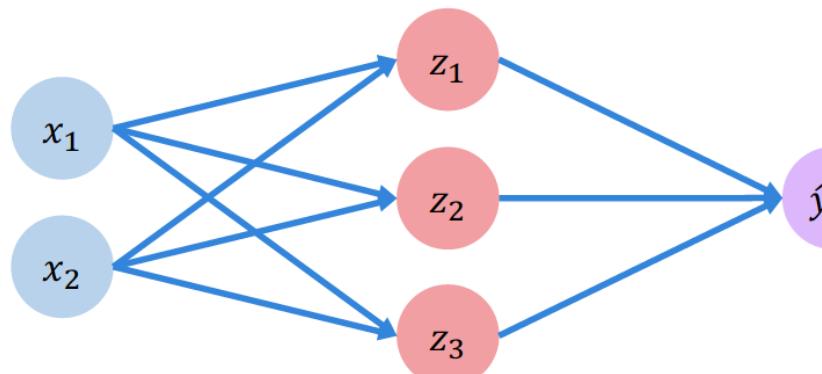
Loss Function

Example Problem: Will I pass this class?



The **empirical loss** measures the total loss over our entire dataset

$$X = \begin{bmatrix} 4, & 5 \\ 2, & 1 \\ 5, & 8 \\ \vdots & \vdots \end{bmatrix}$$



$f(x)$	y
0.1	1
0.8	0
0.6	1
\vdots	\vdots

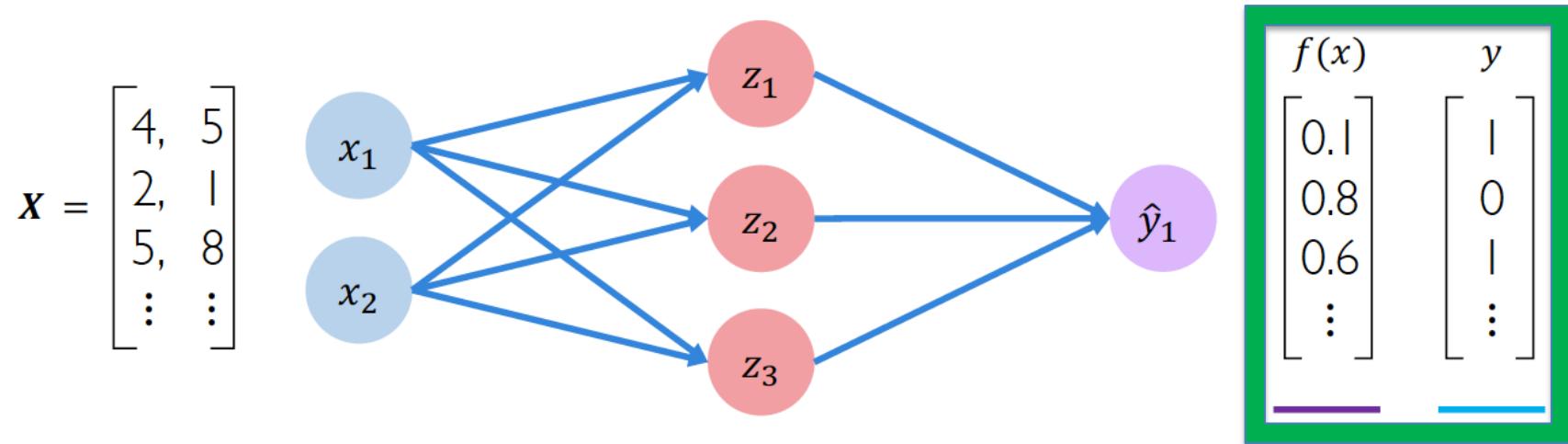
Also known as:

- Objective function
- Cost function
- Empirical Risk

$\curvearrowleft J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}\left(\underbrace{f\left(x^{(i)}; \mathbf{W}\right)}_{\text{Predicted}}, \underbrace{y^{(i)}}_{\text{Actual}}\right)$

Cross Entropy Loss

Cross entropy loss can be used with models that output a probability between 0 and 1



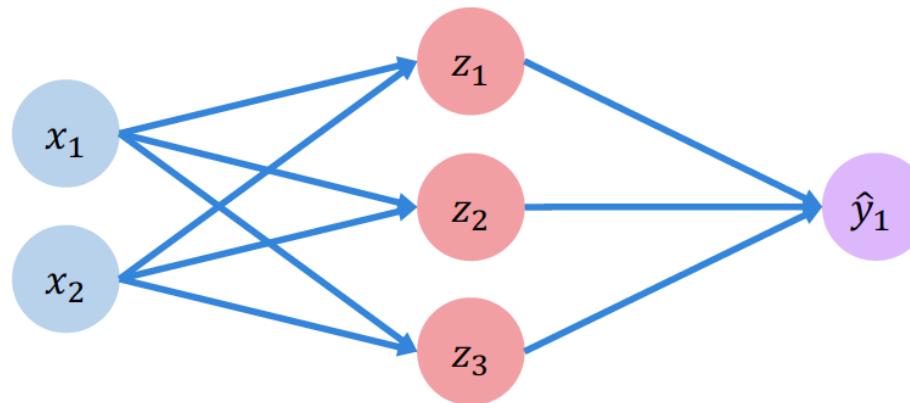
$$J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \underbrace{y^{(i)} \log(f(x^{(i)}; \mathbf{W}))}_{\text{Actual}} + \underbrace{(1 - y^{(i)}) \log(1 - f(x^{(i)}; \mathbf{W}))}_{\text{Actual}} \quad \underbrace{\text{Predicted}}_{\text{Predicted}}$$

loss = tf.reduce_mean(tf.nn.softmax_cross_entropy_with_logits(model.y, model.pred))

Mean Squared Error Loss

Mean squared error loss can be used with regression models that output continuous real numbers

$$X = \begin{bmatrix} 4, & 5 \\ 2, & 1 \\ 5, & 8 \\ \vdots & \vdots \end{bmatrix}$$



$f(x)$	y
30	90
80	20
85	95
\vdots	\vdots

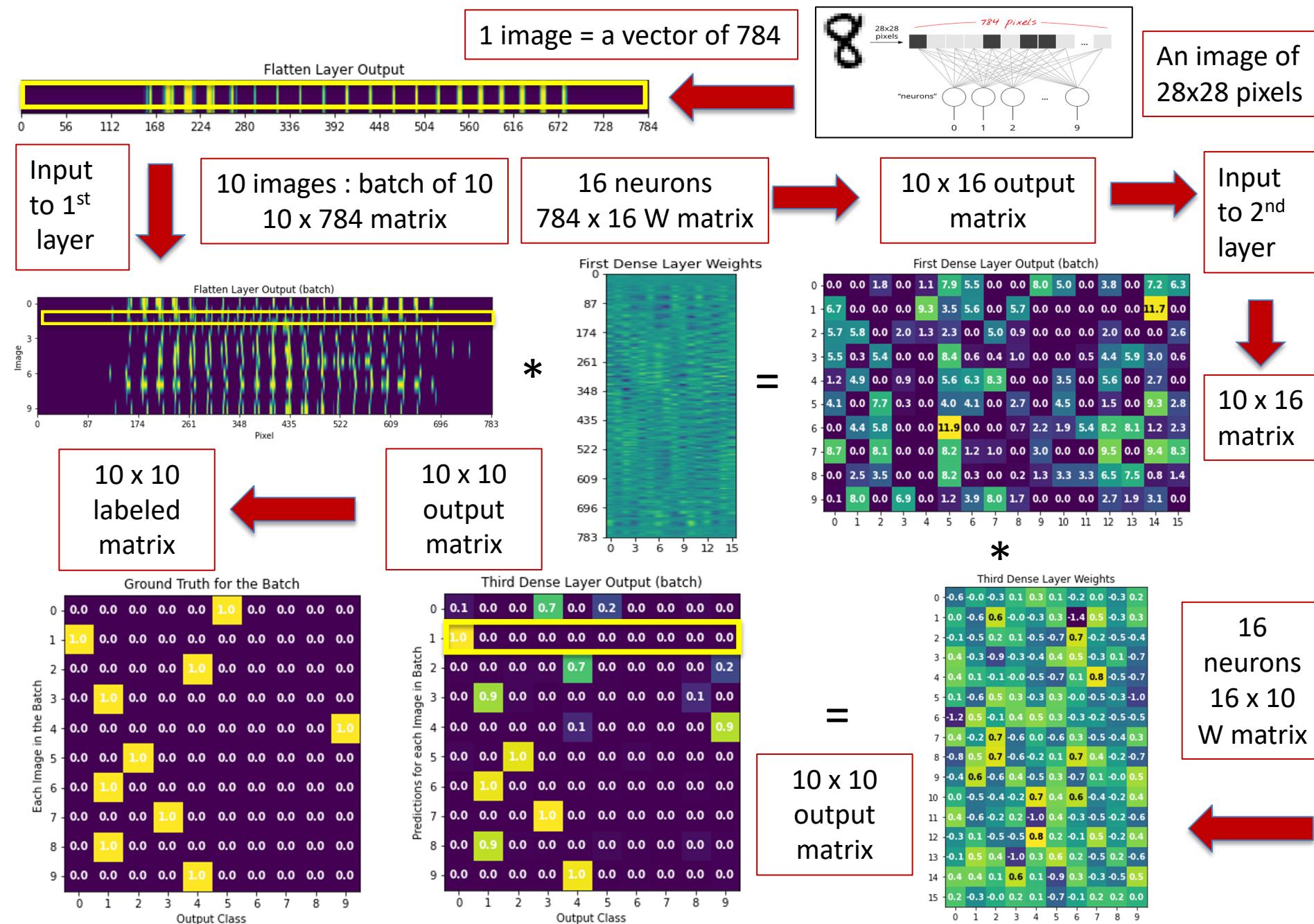
Final Grades
(percentage)

$$J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \underline{\text{Actual}} \underline{\text{Predicted}} \left(y^{(i)} - f(x^{(i)}; \mathbf{W}) \right)^2$$



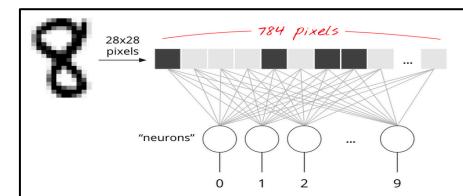
```
loss = tf.reduce_mean( tf.square(tf.subtract(model.y, model.pred) ) )
```

Summary : Flow of MLP Dense layer NN



Summary : Flow of MLP Dense layer NN

1 image = a vector of 784



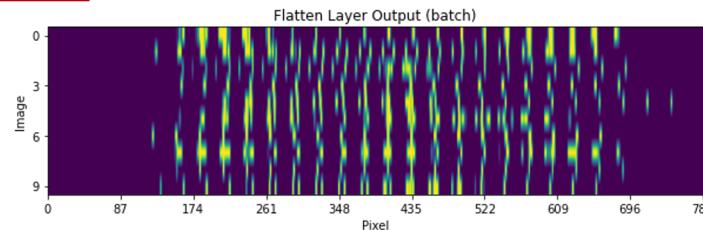
An image of 28x28 pixels

Input to 1st layer



10 images : batch of 10
10 x 784 matrix

10 x 10
labeled
matrix



10 x 10
output
matrix

Ground Truth for the Batch

Each Image in the Batch	0	1	2	3	4	5	6	7	8	9	Output Class
0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	4
1	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1
2	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	3
3	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0
4	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	9
5	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	2
6	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1
7	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	3
8	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0
9	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	4

Comparing
Computed NN
results
To
Labeled results
(target)

Each Image in the Batch	0	1	2	3	4	5	6	7	8	9	Output Class
0	0.1	0.0	0.0	0.7	0.0	0.2	0.0	0.0	0.0	0.0	4
1	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1
2	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.2	2
3	0.0	0.9	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.1	0
4	0.0	0.0	0.0	0.0	0.1	0.0	0.0	0.0	0.0	0.9	3
5	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1
6	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0
7	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	3
8	0.0	0.9	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0
9	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	4

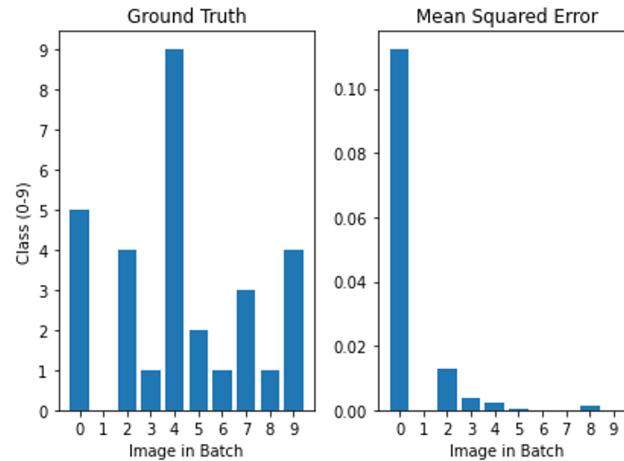
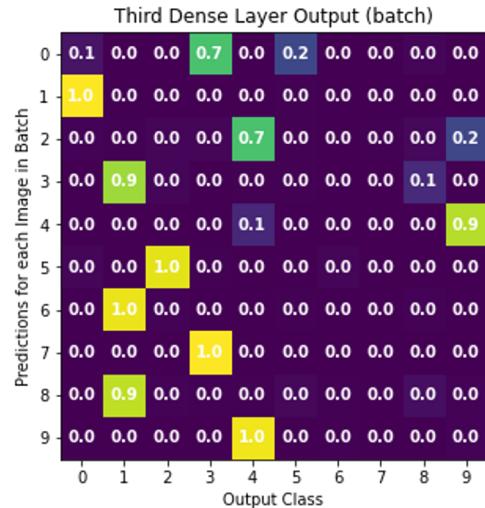
Predictions for each Image in Batch

TensorFlow 2.0 : Forward Pass

`tf.losses.mean_squared_error(y_true_tf, y_pred_tf)`

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

Evaluating the Error



Softmax Classifier (Multinomial Logistic Regression)



Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

Softmax Function

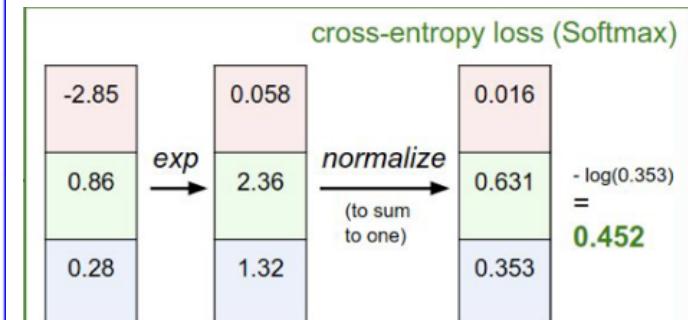
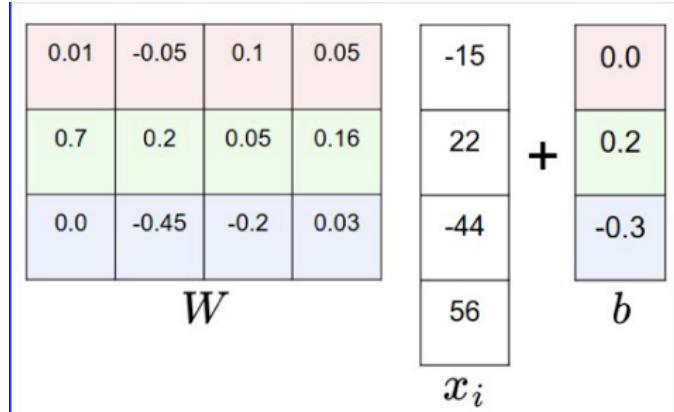
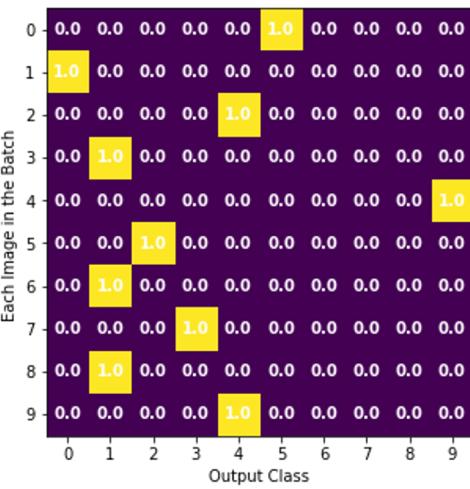
cat
car
frog

Unnormalized log-probabilities / logits	exp	normalize
3.2	5.1	24.5
-1.7	164.0	0.18
		probabilities

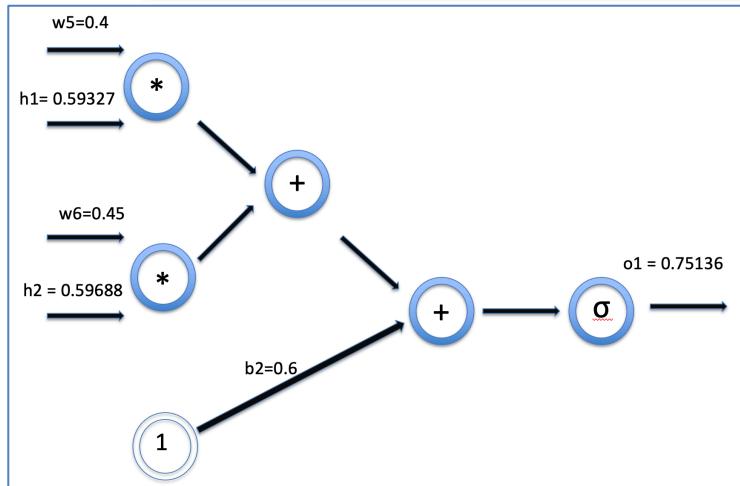
Probabilities must be >= 0	compare	1.00
L_i = -log P(Y = y_i X = x_i)		0.00
Cross Entropy		0.00
H(P, Q) = H(p) + D_KL(P Q)		Correct prob

$$L_i = -\log\left(\frac{e^{sy_i}}{\sum_j e^{sj}}\right)$$

Ground Truth for the Batch



$$E_{o1} = \frac{1}{2}(target_{o1} - out_{o1})^2 = \frac{1}{2}(0.01 - 0.75136507)^2 = 0.274811083$$



target=0.01 o1=0.75136 Eo1=0.27481

Forward Path to Objective Function

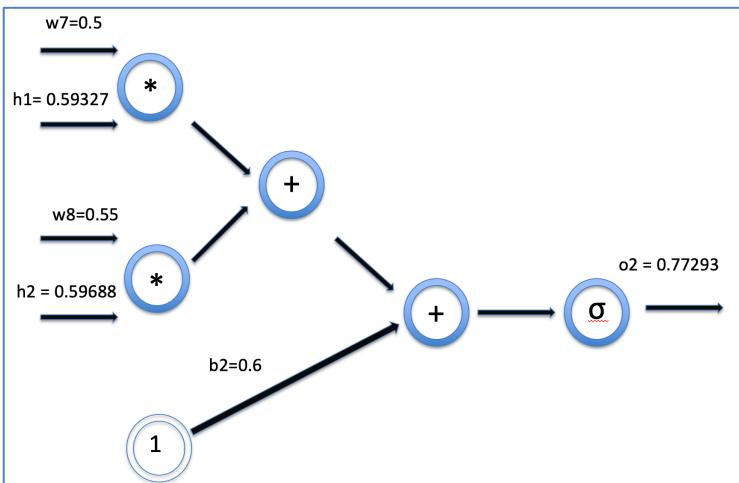
target1 = 0.01

$o_1 = 0.75136$

Cost f

$$E_{total} = \sum \frac{1}{2}(target - output)^2$$

$$E_{total} = E_{o1} + E_{o2} = 0.298371109$$



$o_2 = 0.77293$

Cost f

target2 = 0.99

target=0.99

$o_2=0.77293$

Eo2=0.02356



$E_{o1} = 0.27481$

$ET = 0.298371$

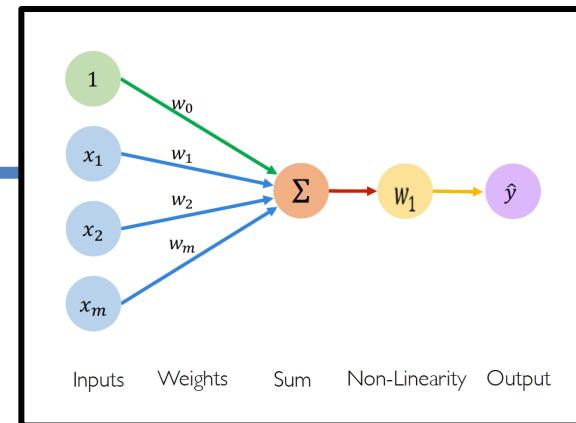
grad $ET = 1.0$

What are the training parameters How does it work?

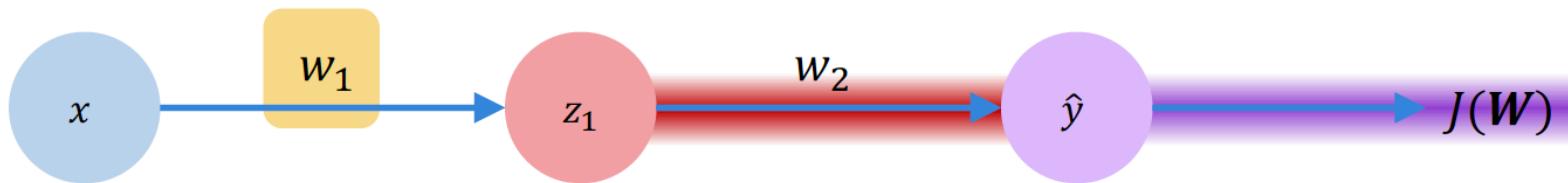
Weights and bias

Weights = $w_1, w_2, w_3, w_4, w_5, w_6, w_7, w_8$
= number of links connecting different layers of neurons

Bias = $b_1-h_1, b_1-h_2, b_2-o_1, b_2-o_2$ = number of neurons

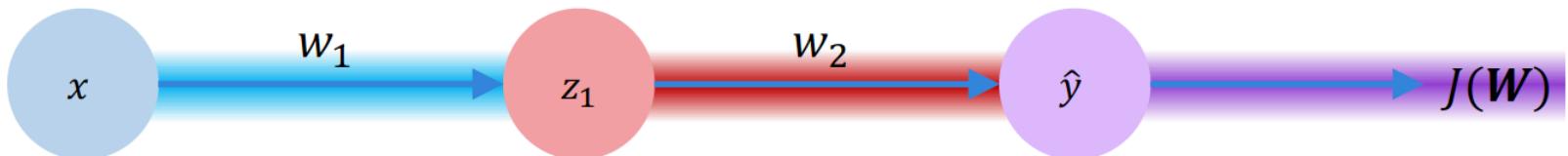


DNN MLP LA Kernels



$$\frac{\partial J(\mathbf{W})}{\partial w_1} = \frac{\partial J(\mathbf{W})}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w_1}$$

Apply chain rule! Apply chain rule!



$$\frac{\partial J(\mathbf{W})}{\partial w_1} = \frac{\partial J(\mathbf{W})}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1} * \frac{\partial z_1}{\partial w_1}$$

DNN Forward Backward Calculation : Weights and bias

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$



```
weights = tf.random_normal(shape, stddev=sigma)
```

2. Loop until convergence:

3. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$



```
grads = tf.gradients(ys=loss, xs=weights)
```

4. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$



```
weights_new = weights.assign(weights - lr * grads)
```

5. Return weights

In one epoch

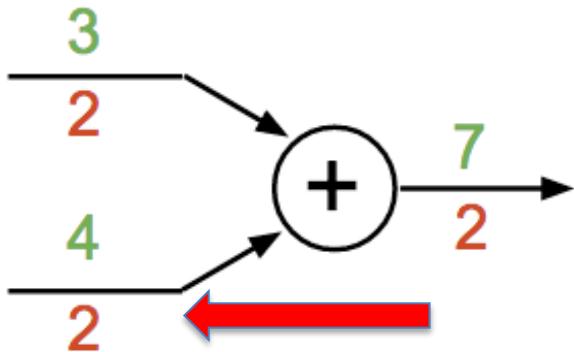
1. Pick a mini-batch
2. Feed it to Neural Network
3. Forward path calculation
4. Calculate the mean gradient of the mini-batch
5. Use the mean gradient calculated in step 4 to update the weights
6. Repeat steps 1–5 for the mini-batches we created

$$\frac{\sum_{j=1}^m \nabla C_{X_j}}{m} \approx \frac{\sum_x \nabla C_x}{n} = \nabla C,$$

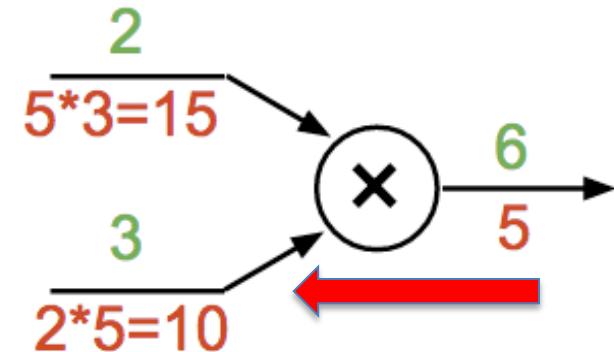
$$w_k \rightarrow w'_k = w_k - \eta \frac{\partial C}{\partial w_k}$$
$$b_l \rightarrow b'_l = b_l - \eta \frac{\partial C}{\partial b_l}.$$

Forward path : backward path

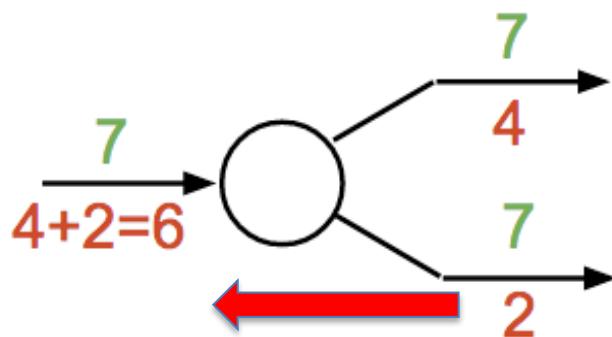
add gate: gradient distributor



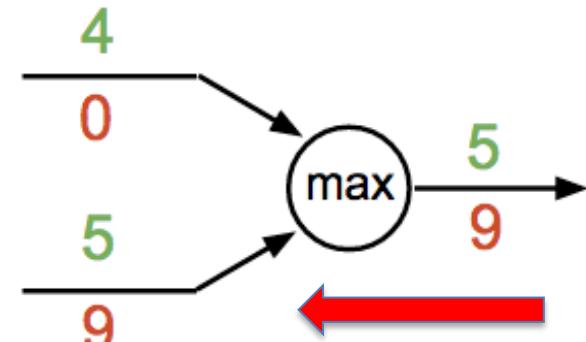
mul gate: “swap multiplier”



copy gate: gradient adder

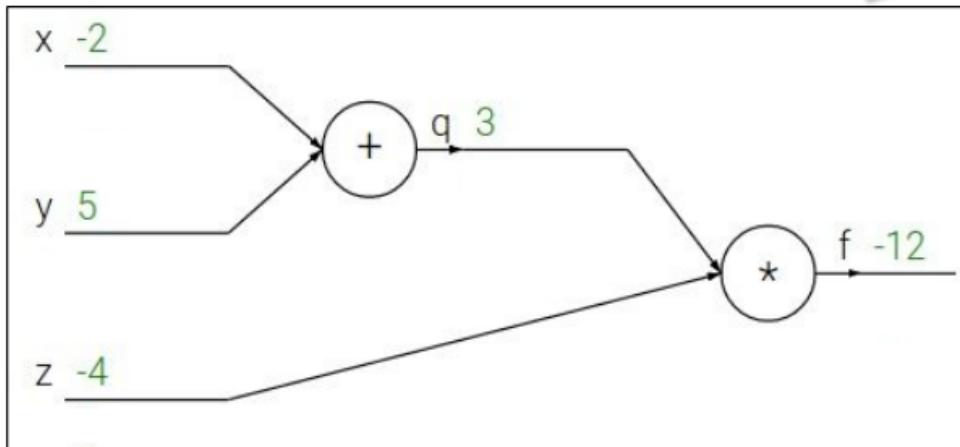


max gate: gradient router



Computational Graph

<http://cs231n.stanford.edu/syllabus.html>



Forward path calculation

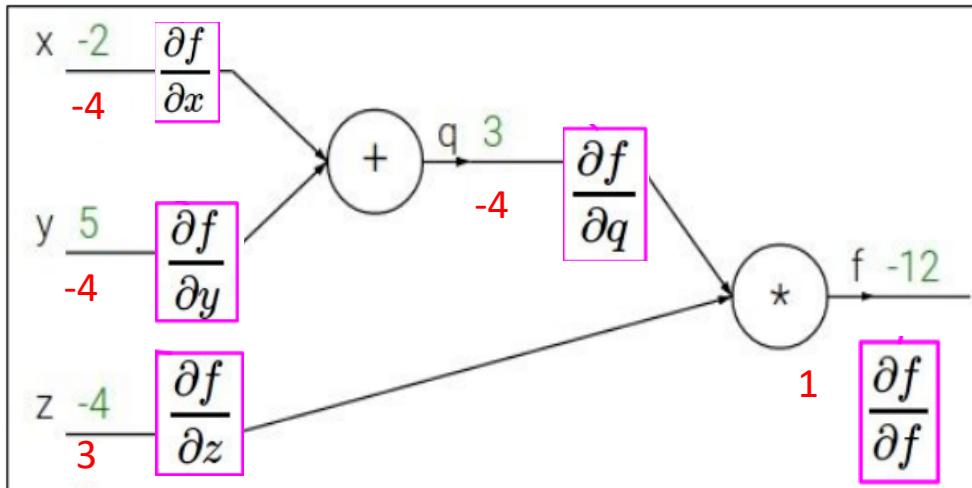
$$f(x, y, z) = (x + y)z$$

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

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

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

Need Gradient of f wrt to inputs



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

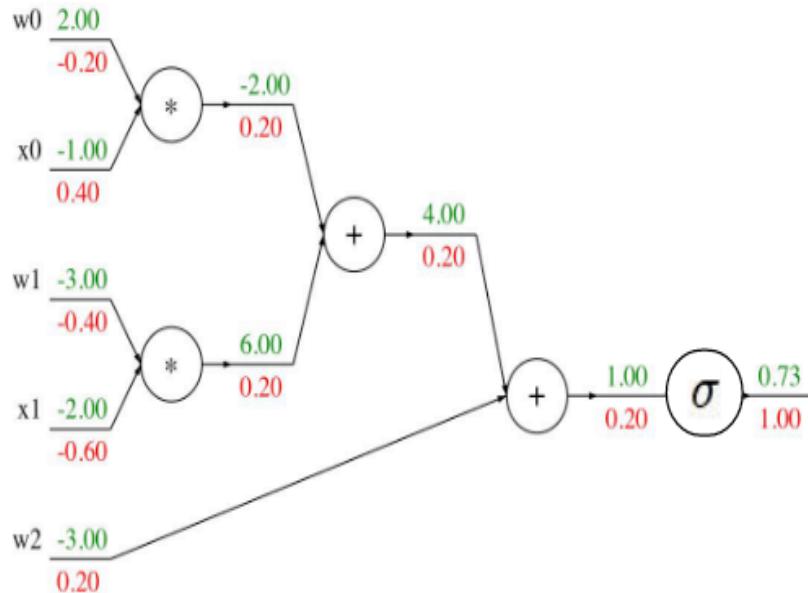
$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x} = (z)(1) = -4$$

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y} = (z)(1) = -4$$

Forward and Backpropagation Code

Forward path : backward path

Backprop Implementation: “Flat” code



Forward pass:
Compute output

```
def f(w0, x0, w1, x1, w2):  
    s0 = w0 * x0  
    s1 = w1 * x1  
    s2 = s0 + s1  
    s3 = s2 + w2  
    L = sigmoid(s3)
```

Backward pass:
Compute grads

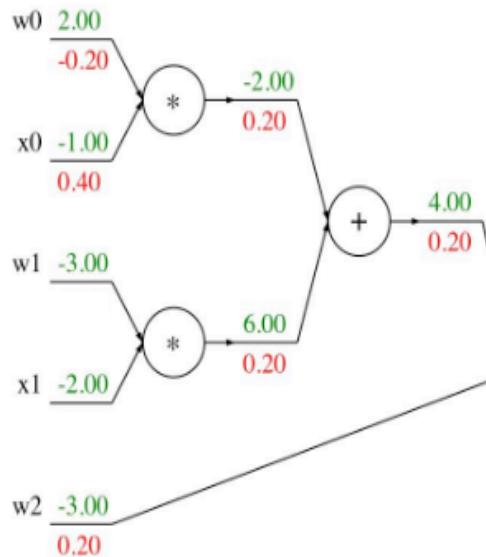
```
grad_L = 1.0  
grad_s3 = grad_L * (1 - L) * L  
grad_w2 = grad_s3  
grad_s2 = grad_s3  
grad_s0 = grad_s2  
grad_s1 = grad_s2  
grad_w1 = grad_s1 * x1  
grad_x1 = grad_s1 * w1  
grad_w0 = grad_s0 * x0  
grad_x0 = grad_s0 * w0
```

Computational Graph

<http://cs231n.stanford.edu/syllabus.html>

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$



Sigmoid
function

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

Computational graph representation may not be unique. Choose one where local gradients at each node can be easily expressed!

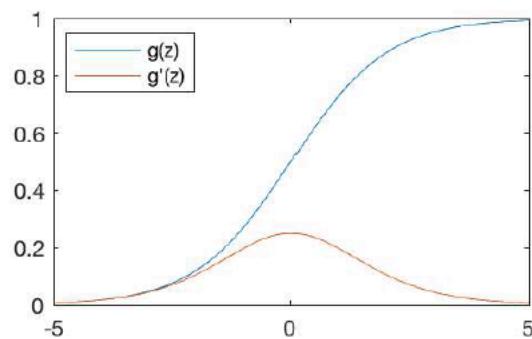
Sigmoid

$$\begin{aligned} & [\text{upstream gradient}] \times [\text{local gradient}] \\ & [1.00] \times [(1 - 0.73)(0.73)] = 0.2 \end{aligned}$$

Sigmoid local
gradient:

$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left(\frac{1 + e^{-x} - 1}{1 + e^{-x}} \right) \left(\frac{1}{1 + e^{-x}} \right) = (1 - \sigma(x))\sigma(x)$$

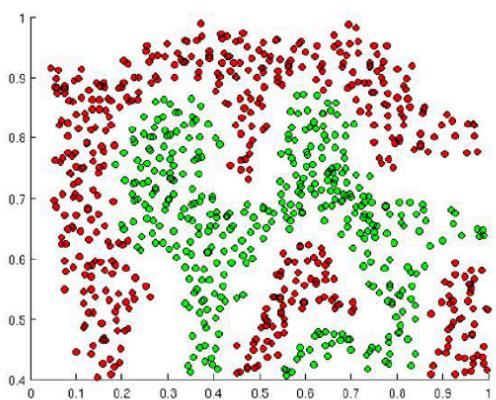
Sigmoid Function



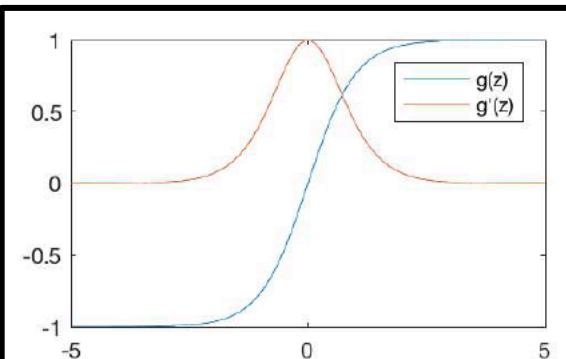
$$g(z) = \frac{1}{1 + e^{-z}}$$

$$g'(z) = g(z)(1 - g(z))$$

`tf.nn.sigmoid(z)`



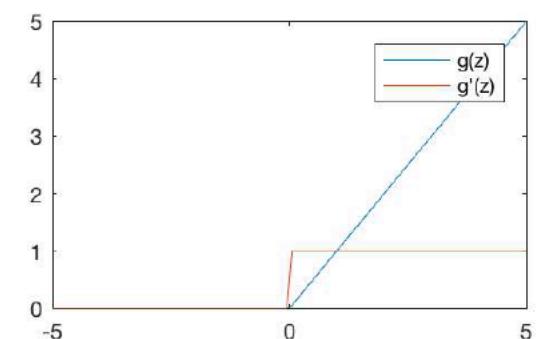
Nonlinear Activation Function



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

$$g'(z) = 1 - g(z)^2$$

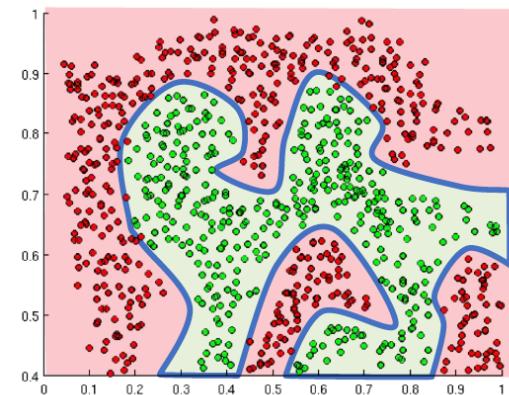
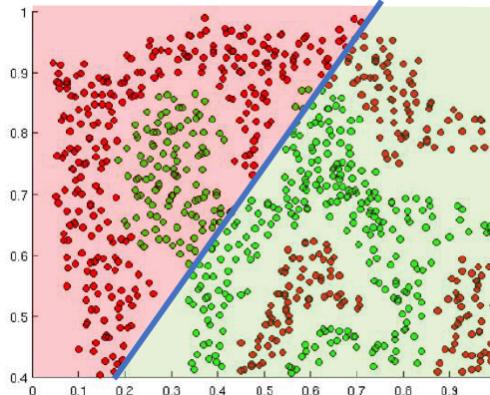
`tf.nn.tanh(z)`



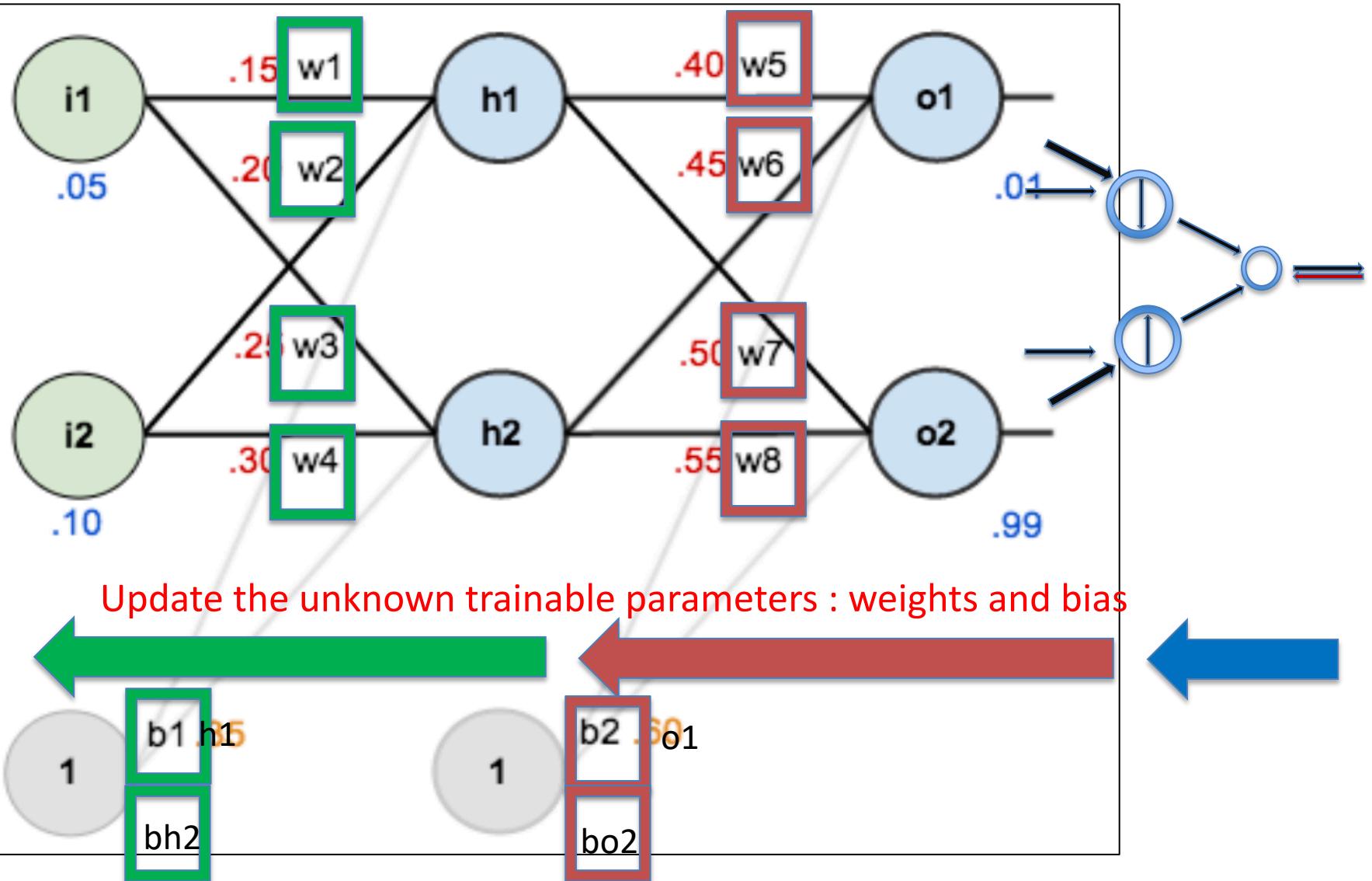
$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

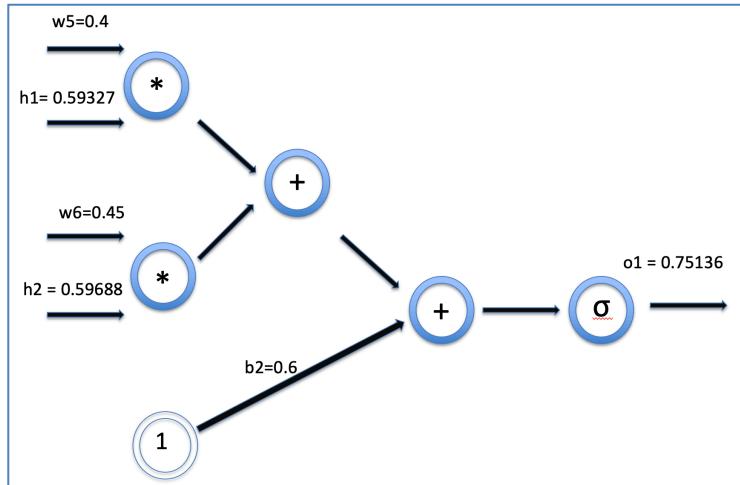
`tf.nn.relu(z)`



Multilayer Perceptron : Backward Path Calculation

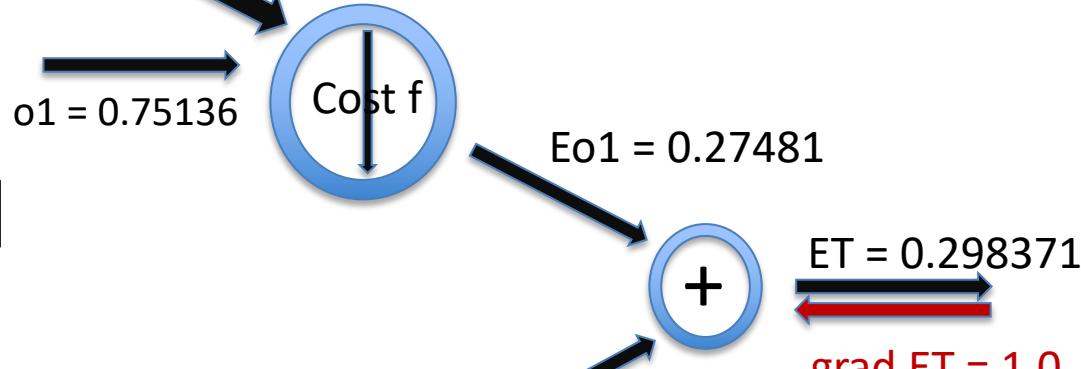
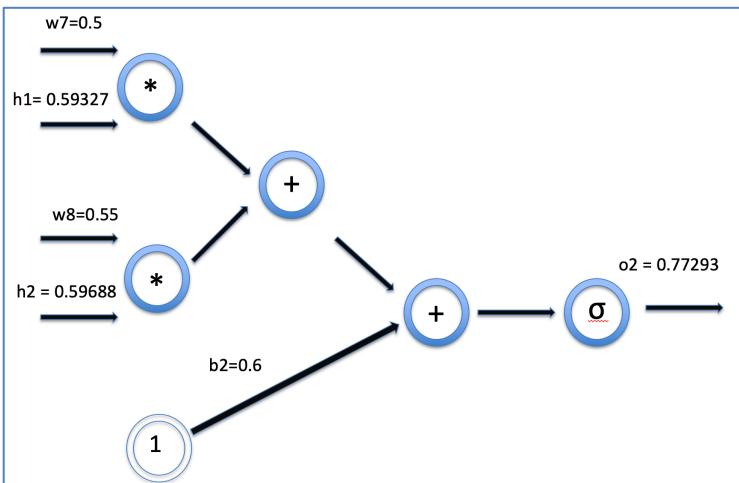
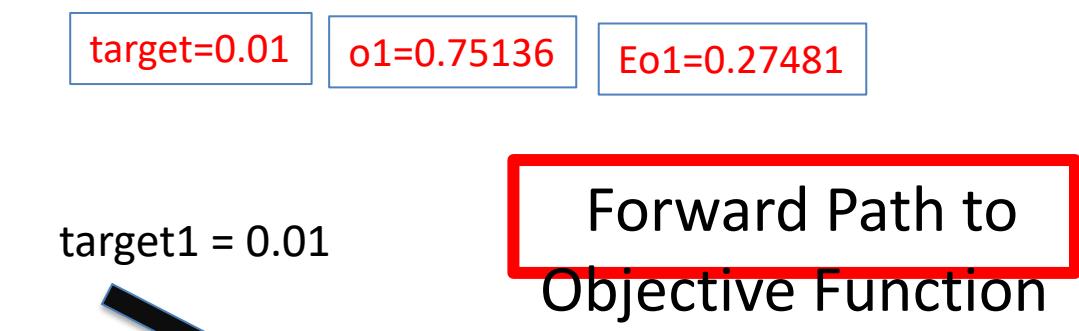


$$E_{o1} = \frac{1}{2}(target_{o1} - out_{o1})^2 = \frac{1}{2}(0.01 - 0.75136507)^2 = 0.274811083$$



$$E_{total} = \sum \frac{1}{2}(target - output)^2$$

$$E_{total} = E_{o1} + E_{o2} = 0.298371109$$



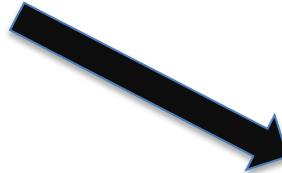
$$target=0.99$$

$$o2=0.77293$$

$$Eo2=0.02356$$

Backward Path from Objective Function

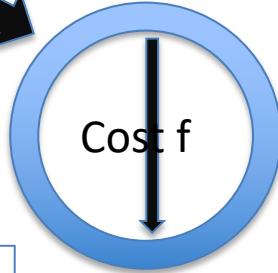
target1 = 0.01



$o_1 = 0.75136$



$\text{grad } o_1 = 0.74136$



composite diff : $\text{grad } o_1 = (\text{grad up}) * (\text{local grad})$

local grad = $(dE_1/do_1) = (o_1 - \text{target1})$

$\text{grad } o_1 = (1) (0.75136 - 0.01) = 0.74136$

$E_{o1} = \frac{1}{2} * (\text{target1} - o_1)^2$

$E_{o1} = 0.27481$

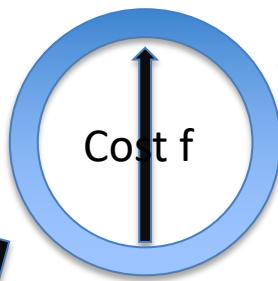
$ET = 0.298371$

$ET = E_1 + E_2$

$o_2 = 0.77293$



$\text{grad } o_2 = -0.21707$



composite diff : $\text{grad } o_2 = (\text{grad up}) * (\text{local grad})$

local grad = $(dE_2/do_2) = (o_2 - \text{target2})$

$E_{o2} = \frac{1}{2} * (\text{target2} - o_2)^2$

$ET = 0.298371$

$\text{grad } ET = 1.0$

target2 = 0.99

composite diff : $\text{grad } o_2 = (\text{grad up}) * (\text{local grad})$

$E_{o2} = 0.02356$

$\text{grad } E_{o2} = 1.0$

$\text{grad } E_{o1} = 1.0$

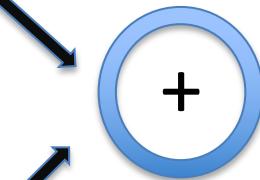
+

Backward Path from o1

$$w5=0.4$$



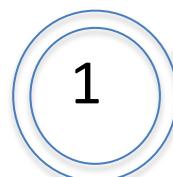
$$h1 = 0.59327$$



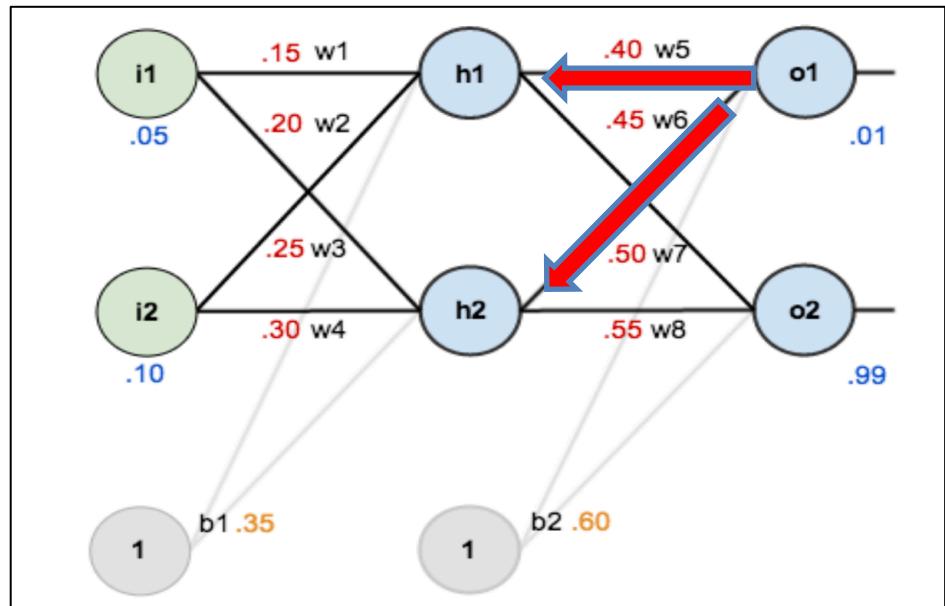
$$w6=0.45$$



$$h2 = 0.59688$$

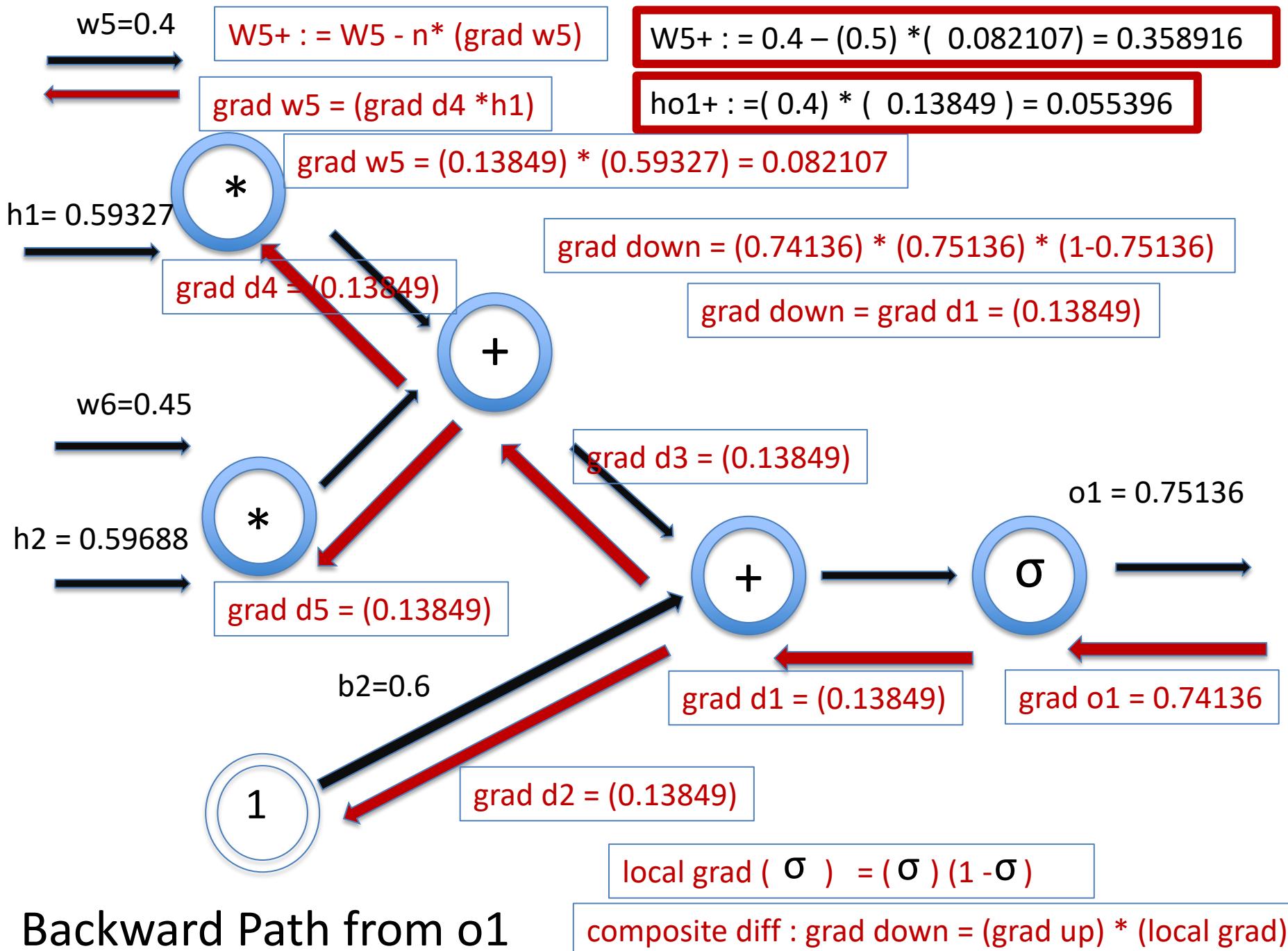


$$b2=0.6$$



$$o1 = 0.75136$$

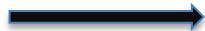




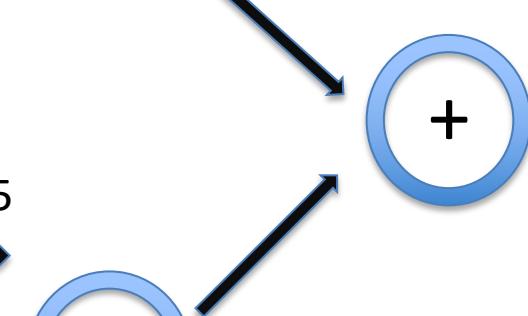
Backward Path from

o_2

$w_7=0.5$



$h_1 = 0.59327$



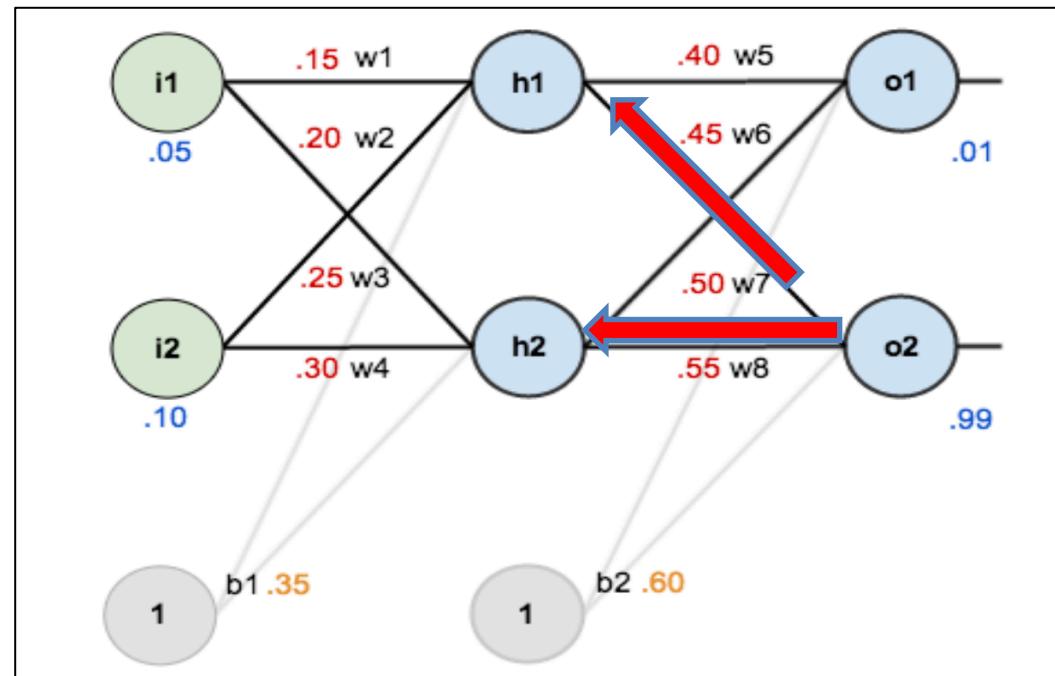
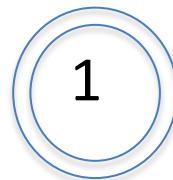
$w_8=0.55$



$h_2 = 0.59688$



$B_2=0.6$



$o_2 = 0.77293$



$w7=0.5$

$W7+ := W7 - n * (\text{grad } w7)$

$W7+ := 0.5 - (0.5) * (-0.022601) = 0.51130$

$\text{grad } w7 = (\text{grad } d9 * h1)$

$h02+ := (0.5) * (-0.038097) = -0.019485$

 $*$

$\text{grad } w7 = (-0.038097) * (0.59327) = -0.022601$

$n=0.5$

$h1 = 0.59327$

$\text{grad } d9 = (-0.038097)$

$\text{grad down} = (-0.21707) * (0.77293) * (1-0.77293)$

$\text{grad down} = \text{grad } d6 = (-0.038097)$

$w8=0.55$

 $+$

$\text{grad } d8 = (-0.038097)$

$o2 = 0.77293$

$h2 = 0.59688$

$\text{grad } d10 = (-0.038097)$

$b2=0.6$

 $+$

$\text{grad } d6 = (-0.038097)$

 σ

$\text{grad } o2 = -0.21707$

 1

$\text{grad } d7 = (-0.038097)$

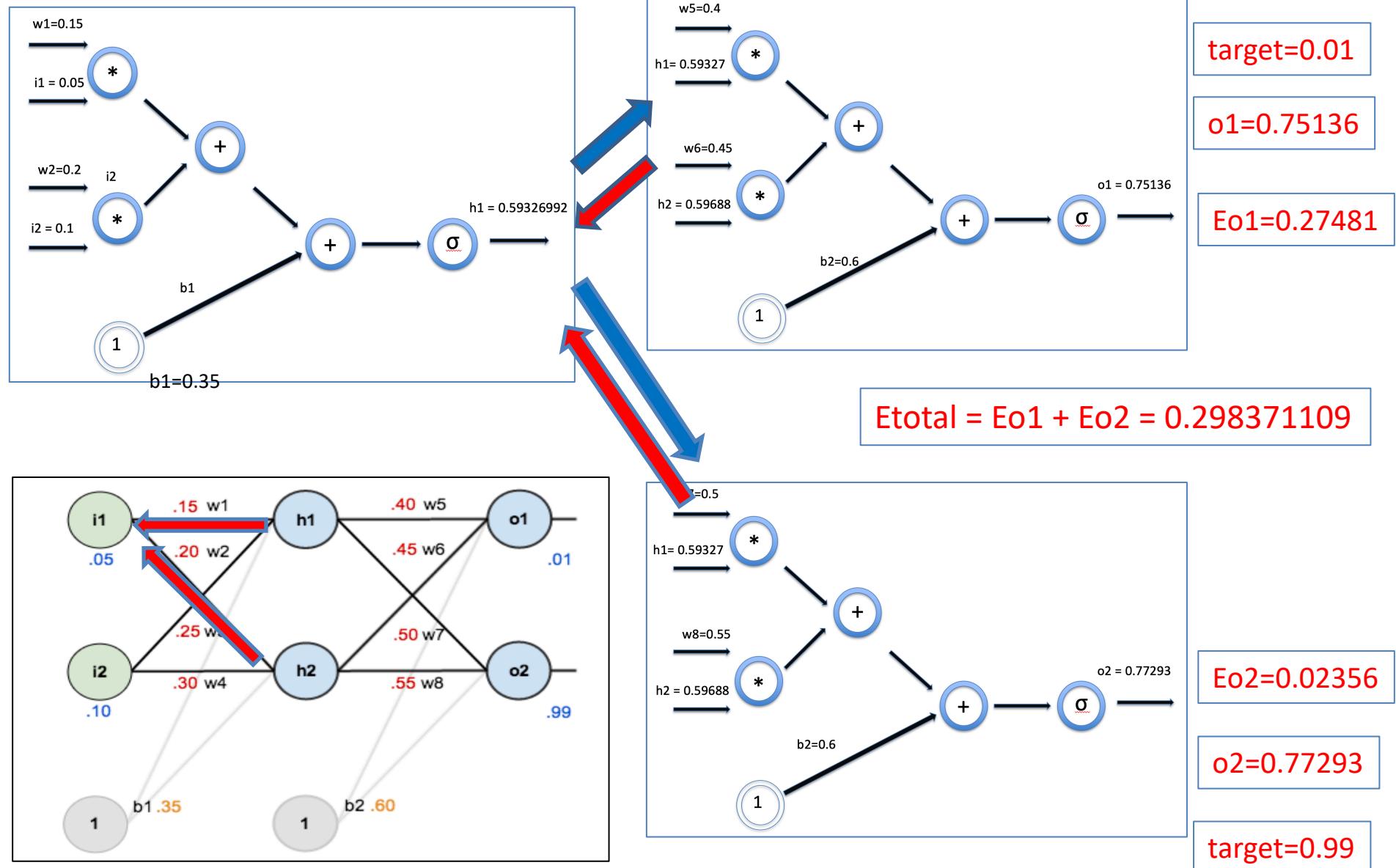
$\text{local grad } (\sigma) = (\sigma)(1-\sigma)$

Backward Path from

$\text{composite diff : grad down} = (\text{grad up}) * (\text{local grad})$

 σ^2

Backward Path from o1 and o2 to h1, from h1 to w1



$w1=0.15$

$W1+ := W1 - n * (\text{grad } w1)$

$\text{grad } w1 = (\text{grad } d14 * i1)$

$W1+ := 0.15 - (0.5) * (0.0004332) = 0.149783$

$i1 = 0.05$

$\text{grad } w1 = (0.008665) * (0.05) = 0.00043326$

$\text{grad } d14 = (0.008665)$

$\text{grad down} = (0.035911) * (0.59327) * (1 - 0.59327)$

$\text{grad down} = \text{grad } d1 = (0.008665)$

$w2=0.20$

$i2 = 0.1$

Backward Path from $h1$ to $w1$

$h1 = 0.59327$

$\text{grad } d15 = (0.008665)$

$\text{grad } d13 = (0.008665)$

$+$

σ

$b1=0.35$

$\text{grad } d11 = (0.008665)$

$\text{grad } h1 = 0.035911$

1

$\text{grad } d12 = (0.008665)$

$\text{grad } h1-o1 := (0.4) * (0.13849) = 0.055396$

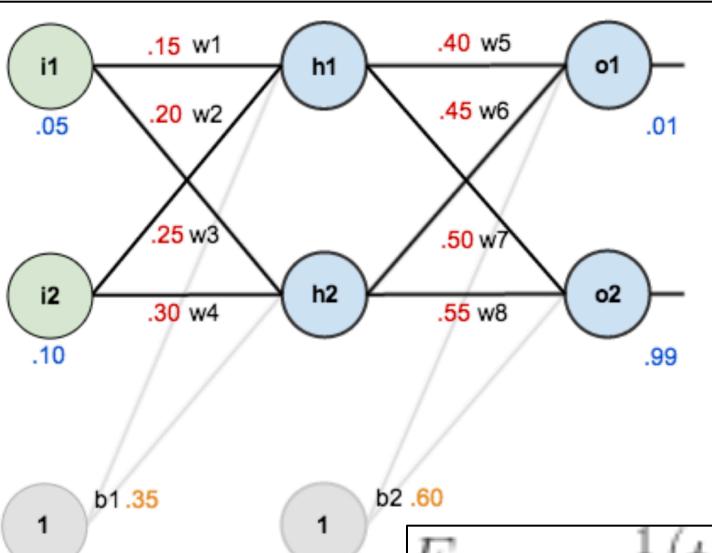
$\text{grad } h1-o2 := (0.5) * (-0.0038097) = -0.019485$

$\text{local grad } (\sigma) = (\sigma)(1-\sigma)$

$\text{grad } h1 := (\text{grad } h1-o1) + (\text{grad } h1-o2) = (0.055396 - 0.019485)$

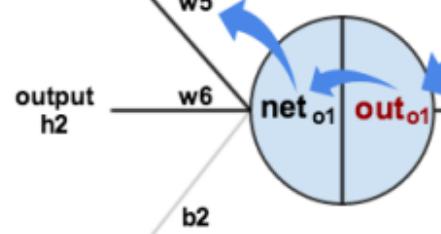
$\text{composite diff : grad down} = (\text{grad up}) * (\text{local grad})$

Backpropagation Example



$$\frac{\partial E_{total}}{\partial w_5} = \frac{\partial E_{total}}{\partial out_{o1}} * \frac{\partial out_{o1}}{\partial net_{o1}} * \frac{\partial net_{o1}}{\partial w_5}$$

$$\frac{\partial net_{o1}}{\partial w_5} * \frac{\partial out_{o1}}{\partial net_{o1}} * \frac{\partial E_{total}}{\partial out_{o1}} = \frac{\partial E_{total}}{\partial w_5}$$



$$E_{o1} = \frac{1}{2}(target_{o1} - out_{o1})^2$$

$$E_{total} = E_{o1} + E_{o2}$$

$$E_{total} = \frac{1}{2}(target_{o1} - out_{o1})^2 + \frac{1}{2}(target_{o2} - out_{o2})^2$$

$$\frac{\partial E_{total}}{\partial out_{o1}} = 2 * \frac{1}{2}(target_{o1} - out_{o1})^{2-1} * -1 + 0 \quad \frac{\partial E_{total}}{\partial out_{o1}} = -(target_{o1} - out_{o1}) = -(0.01 - 0.75136507) = 0.74136507$$

$$out_{o1} = \frac{1}{1+e^{-net_{o1}}} \quad \frac{\partial out_{o1}}{\partial net_{o1}} = out_{o1}(1 - out_{o1}) = 0.75136507(1 - 0.75136507) = 0.186815602$$

$$net_{o1} = w_5 * out_{h1} + w_6 * out_{h2} + b_2 * 1 \quad \frac{\partial net_{o1}}{\partial w_5} = 1 * out_{h1} * w_5^{(1-1)} + 0 + 0 = out_{h1} = 0.593269992$$

$$\frac{\partial E_{total}}{\partial w_5} = \frac{\partial E_{total}}{\partial out_{o1}} * \frac{\partial out_{o1}}{\partial net_{o1}} * \frac{\partial net_{o1}}{\partial w_5} \quad \frac{\partial E_{total}}{\partial w_5} = 0.74136507 * 0.186815602 * 0.593269992 = 0.082167041$$

$$\frac{\partial E_{total}}{\partial w_5} = -(target_{o1} - out_{o1}) * out_{o1}(1 - out_{o1}) * out_{h1} \quad \delta_{o1} = \frac{\partial E_{total}}{\partial out_{o1}} * \frac{\partial out_{o1}}{\partial net_{o1}} = \frac{\partial E_{total}}{\partial net_{o1}}$$

$$\delta_{o1} = -(target_{o1} - out_{o1}) * out_{o1}(1 - out_{o1}) \quad \frac{\partial E_{total}}{\partial w_5} = \delta_{o1} out_{h1} \quad \frac{\partial E_{total}}{\partial w_5} = -\delta_{o1} out_{h1}$$

$$w_5^+ = w_5 - \eta * \frac{\partial E_{total}}{\partial w_5} = 0.4 - 0.5 * 0.082167041 = 0.35891648$$

$$w_6^+ = 0.408666186$$

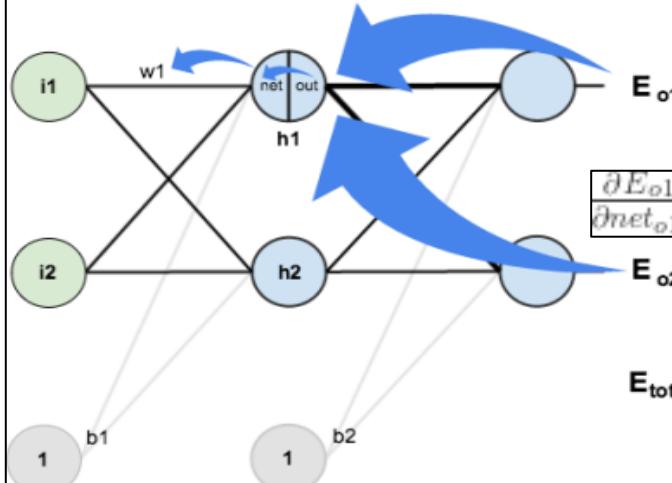
$$w_7^+ = 0.511301270$$

$$w_8^+ = 0.561370121$$

Backpropagation Example

$$\frac{\partial E_{total}}{\partial w_1} = \frac{\partial E_{total}}{\partial out_{h1}} * \frac{\partial out_{h1}}{\partial net_{h1}} * \frac{\partial net_{h1}}{\partial w_1}$$

$$\frac{\partial E_{total}}{\partial out_{h1}} = \frac{\partial E_{o1}}{\partial out_{h1}} + \frac{\partial E_{o2}}{\partial out_{h1}}$$



$$\frac{\partial E_{total}}{\partial w_1} = \frac{\partial E_{total}}{\partial out_{h1}} * \frac{\partial out_{h1}}{\partial net_{h1}} * \frac{\partial net_{h1}}{\partial w_1}$$

$$\frac{\partial E_{total}}{\partial out_{h1}} = \frac{\partial E_{o1}}{\partial out_{h1}} + \frac{\partial E_{o2}}{\partial out_{h1}}$$

$$\frac{\partial E_{o1}}{\partial out_{h1}} = \frac{\partial E_{o1}}{\partial net_{o1}} * \frac{\partial net_{o1}}{\partial out_{h1}}$$

$$\frac{\partial E_{o1}}{\partial net_{o1}} = \frac{\partial E_{o1}}{\partial out_{o1}} * \frac{\partial out_{o1}}{\partial net_{o1}} = 0.74136507 * 0.186815602 = 0.138498562$$

$$net_{o1} = w_5 * out_{h1} + w_6 * out_{h2} + b_2 * 1$$

$$E_{total} = E_{o1} + E_{o2}$$

$$\frac{\partial net_{o1}}{\partial out_{h1}} = w_5 = 0.40$$

$$\frac{\partial E_{o1}}{\partial out_{h1}} = \frac{\partial E_{o1}}{\partial net_{o1}} * \frac{\partial net_{o1}}{\partial out_{h1}} = 0.138498562 * 0.40 = 0.055399425$$

$$\frac{\partial E_{total}}{\partial out_{h1}} = \frac{\partial E_{o1}}{\partial out_{h1}} + \frac{\partial E_{o2}}{\partial out_{h1}} = 0.055399425 + -0.019049119 = 0.036350306$$

$$\frac{\partial E_{o2}}{\partial out_{h1}} = -0.019049119$$

$$\frac{\partial out_{h1}}{\partial net_{h1}} = out_{h1}(1 - out_{h1}) = 0.59326999(1 - 0.59326999) = 0.241300709$$

$$net_{h1} = w_1 * i_1 + w_3 * i_2 + b_1 * 1$$

$$\frac{\partial net_{h1}}{\partial w_1} = i_1 = 0.05$$

$$out_{h1} = \frac{1}{1+e^{-net_{h1}}}$$

$$\frac{\partial E_{total}}{\partial w_1} = \frac{\partial E_{total}}{\partial out_{h1}} * \frac{\partial out_{h1}}{\partial net_{h1}} * \frac{\partial net_{h1}}{\partial w_1} \quad \frac{\partial E_{total}}{\partial w_1} = 0.036350306 * 0.241300709 * 0.05 = 0.000438568$$

$$\frac{\partial E_{total}}{\partial w_1} = \left(\sum_o \frac{\partial E_{total}}{\partial out_o} * \frac{\partial out_o}{\partial net_o} * \frac{\partial net_o}{\partial out_{h1}} \right) * \frac{\partial out_{h1}}{\partial net_{h1}} * \frac{\partial net_{h1}}{\partial w_1} \quad \frac{\partial E_{total}}{\partial w_1} = \left(\sum_o \delta_o * w_{ho} \right) * out_{h1}(1 - out_{h1}) * i_1$$

$$w_1^+ = w_1 - \eta * \frac{\partial E_{total}}{\partial w_1} = 0.15 - 0.5 * 0.000438568 = 0.149780716$$

$$\frac{\partial E_{total}}{\partial w_1} = \delta_{h1} i_1$$

$$w_2^+ = 0.19950143$$

$$w_3^+ = 0.24975114$$

$$w_4^+ = 0.29950229$$

```

# Backpropagation
n = 0.5

grad_o1 = o1 - To1
grad_o2 = o2 - To2

print("grad_o1: "+str(grad_o1))
print("grad_o2: "+str(grad_o2))

grad_d1 = (grad_o1)*o1*(1-o1)
grad_d2 = (grad_o2)*o2*(1-o2)

print("grad_d1: "+str(grad_d1))
print("grad_d2: "+str(grad_d2))

grad_w5 = grad_d1*h1
print("grad_w5: "+str(grad_w5))

w5_final = w5 - n*grad_w5

grad_w6 = grad_d1*h2
print("grad_w6: "+str(grad_w6))
w6_final = w6 - n*grad_w6

grad_w7 = grad_d2*h1
print("grad_w7: "+str(grad_w7))
w7_final = w7 - n*grad_w7

grad_w8 = grad_d2*h2
print("grad_w8: "+str(grad_w8))
w8_final = w8 - n*grad_w8

```

```

grad_h1 = w5*grad_d1 + w7*grad_d2
print("grad_h1: "+str(grad_h1))
grad_d11 = grad_h1*h1*(1-h1)
print("grad_d11: "+str(grad_d11))

grad_w1 = grad_d11*i1
print("grad_w1: "+str(grad_w1))
w1_final = w1 - n*grad_w1

grad_w2 = grad_d11*i1
print("grad_w2: "+str(grad_w2))
w2_final = w2 - n*grad_w2

grad_h2 = w6*grad_d1 + w8*grad_d2
print("grad_h2: "+str(grad_h2))
grad_d22 = grad_h2*h2*(1-h2)
print("grad_d22: "+str(grad_d22))

grad_w3 = grad_d22*i1
print("grad_w3: "+str(grad_w3))
w3_final = w3 - n*grad_w3

grad_w4 = grad_d22*i2
print("grad_w4: "+str(grad_w4))
w4_final = w4 - n*grad_w4

```

```
print("w1+: "+str(w1_final))
print("w2+: "+str(w2_final))
print("w3+: "+str(w3_final))
print("w4+: "+str(w4_final))
print("w5+: "+str(w5_final))
print("w6+: "+str(w6_final))
print("w7+: "+str(w7_final))
print("w8+: "+str(w8_final))
```

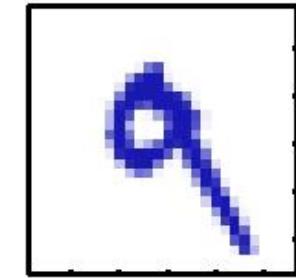
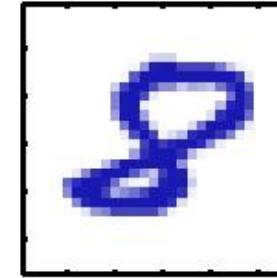
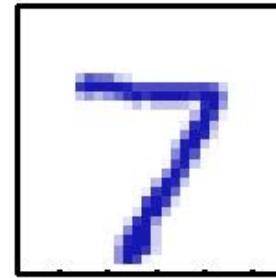
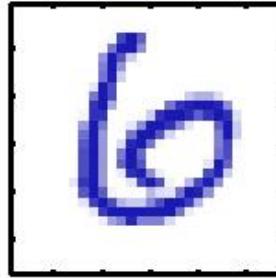
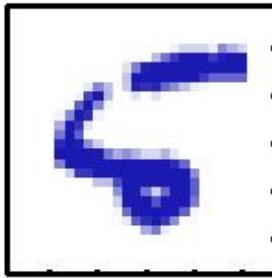
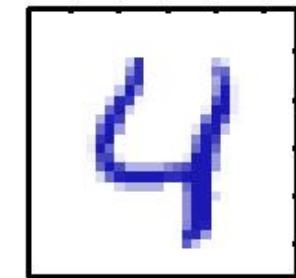
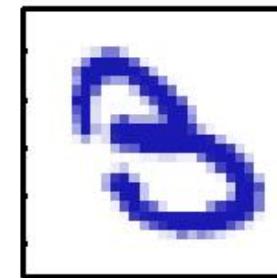
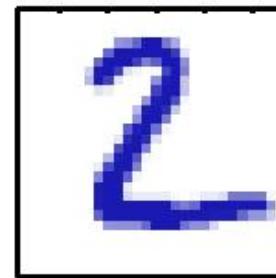
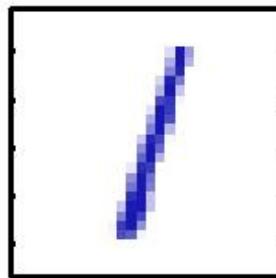
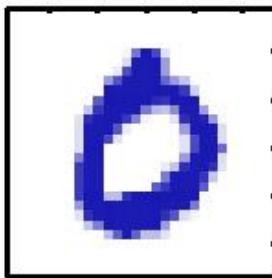
w1+: 0.1497807161327648
w2+: 0.19978071613276482
w3+: 0.24975114363237164
w4+: 0.29950228726474326
w5+: 0.35891647971775653
w6+: 0.4086661860761087
w7+: 0.5113012702391395
w8+: 0.5613701211083925

grad_o1: 0.7413650695475076
grad_o2: -0.21707153468357898
grad_d1: 0.13849856162945076
grad_d2: -0.03809823651803844
grad_w5: 0.08216704056448701
grad_w6: 0.08266762784778263
grad_w7: -0.02260254047827904
grad_w8: -0.02274024221678477
grad_h1: 0.036350306392761086
grad_d11: 0.00877135468940779
grad_w1: 0.0004385677344703895
grad_w2: 0.0004385677344703895
grad_h2: 0.04137032264833171
grad_d22: 0.009954254705134271
grad_w3: 0.0004977127352567136
grad_w4: 0.0009954254705134271

Example: how can computer see images?

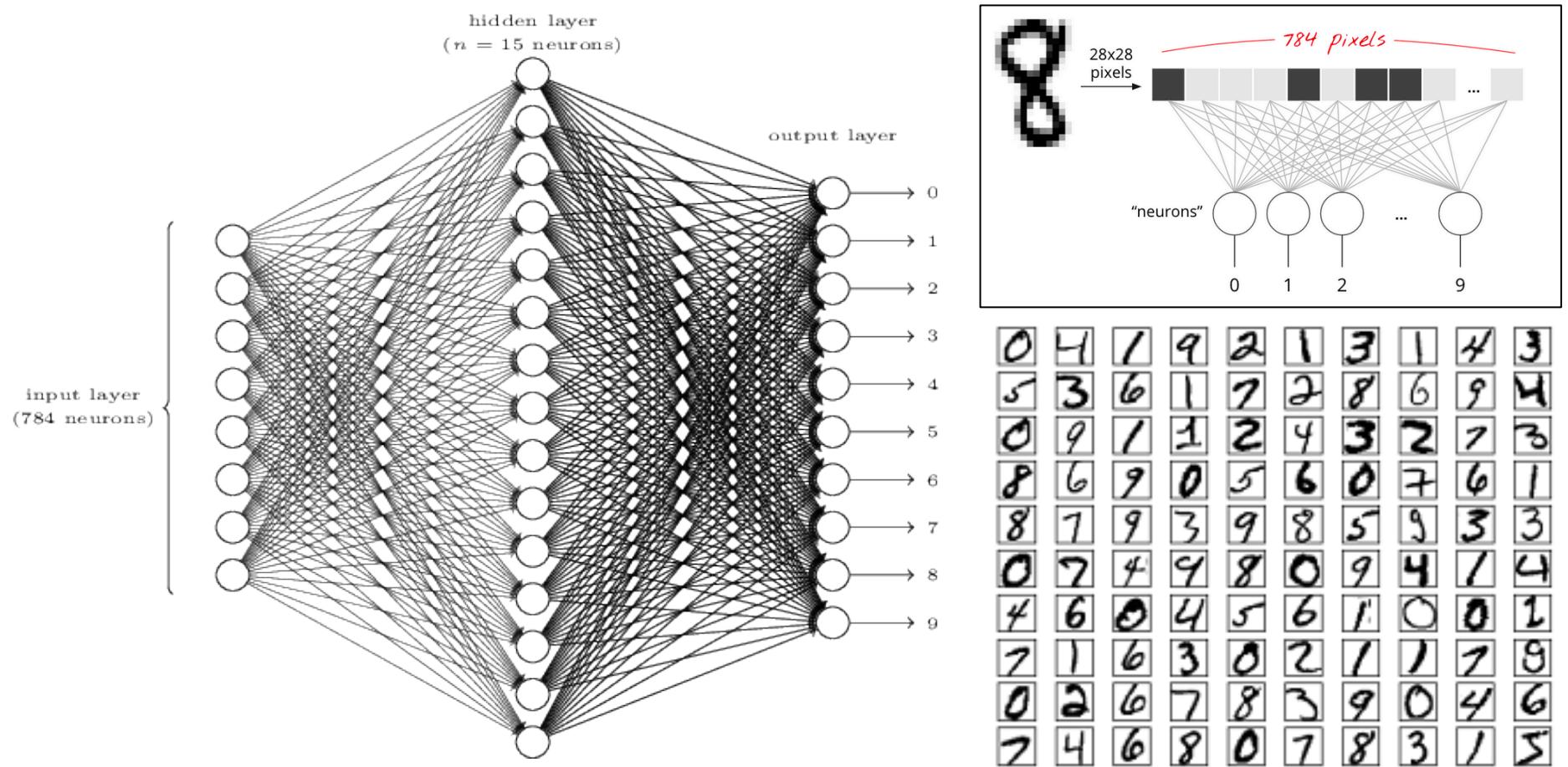
Handwritten Digit Recognition (MNIST data set)

The **MNIST database** (*Modified National Institute of Standards and Technology database*) is a large [database](#) of handwritten digits that is commonly used for [training](#) various [image processing](#) systems. The database is also widely used for training and testing in the field of [machine learning](#).^{[4][5]} It was created by "re-mixing" the samples from NIST's original datasets. The creators felt that since NIST's training dataset was taken from American [Census Bureau](#) employees, while the testing dataset was taken from [American high school](#) students, it was not well-suited for machine learning experiments. Furthermore, the black and white images from NIST were [normalized](#) to fit into a 28x28 pixel bounding box and [anti-aliased](#), which introduced grayscale levels. The MNIST database contains 60,000 training images and 10,000 testing images. – from Wikipedia



MNIST Example (28x28 pixels)

training digits and their labels															
5	9	0	1	5	0	4	2	4	5	7	5	0	6	7	5
validation digits and their labels															
7	2	1	0	4	1	4	9	5	9	0	6	9	0	1	5



0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
13	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
14	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
16	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
17	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
18	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
19	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
20	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
21	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
22	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
23	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
24	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
25	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
26	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
27	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

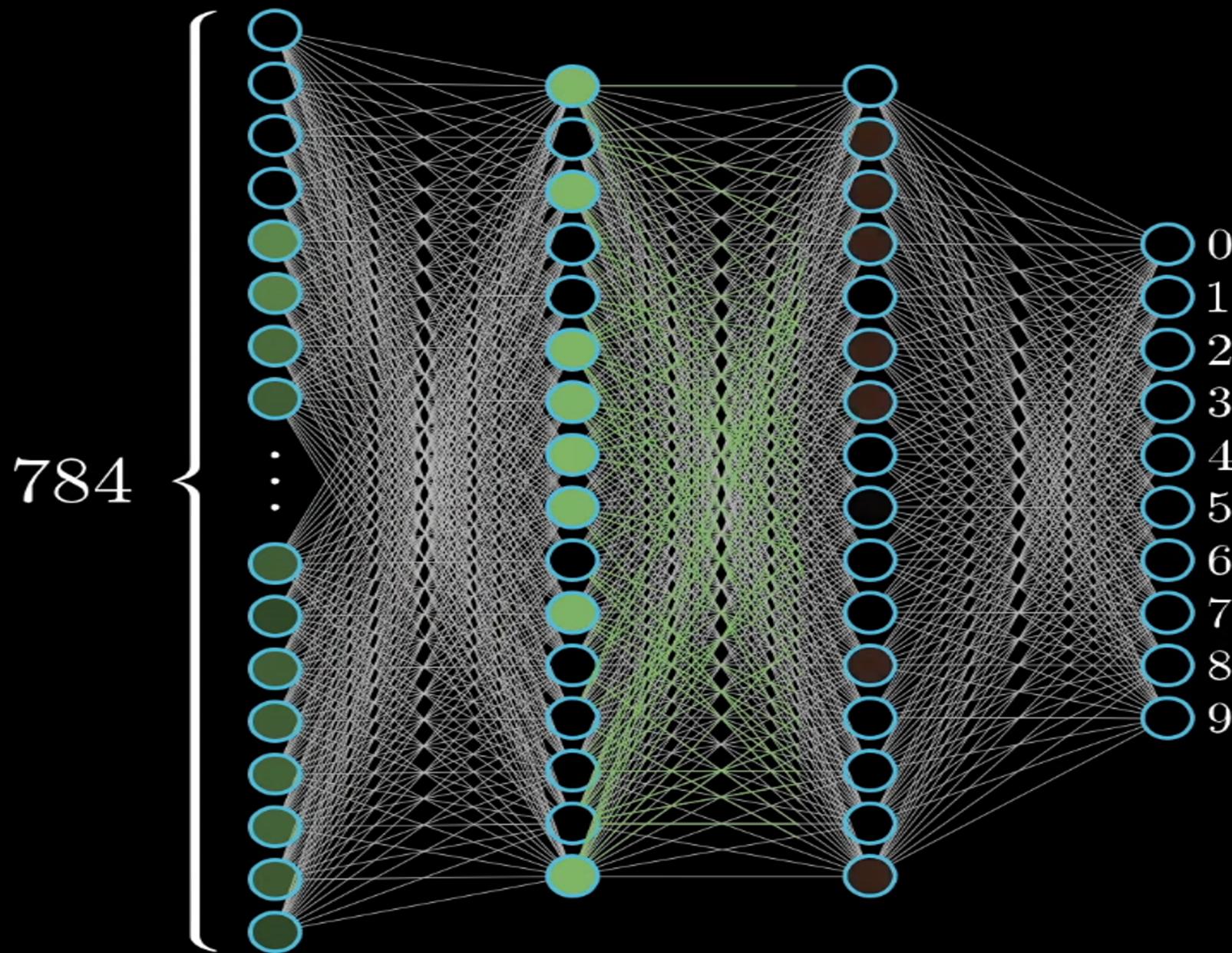
Grayscale image
of 28 x 28 pixels

same as a
28 x 28 matrix

values of
0 - 255

Simple MLP Network

Image:<https://www.3blue1brown.com/>



Tensorflow Example

```
import tensorflow as tf

mnist = tf.keras.datasets.mnist

(x_train, y_train), (x_test, y_test) = mnist.load_data()
x_train, x_test = x_train / 255.0, x_test / 255.0

model = tf.keras.models.Sequential([
    tf.keras.layers.Flatten(input_shape=(28, 28)),
    tf.keras.layers.Dense(16, activation='relu'),
    tf.keras.layers.Dense(16, activation='relu'),
    tf.keras.layers.Dense(10, activation='softmax')
])
```

Tensorflow is imported

Sets the mnist dataset to variable “mnist”

Loads the mnist dataset

Builds the layers of the model
4 layers in this model

```
model.compile(optimizer='sgd', loss='SparseCategoricalCrossentropy',
metrics=['accuracy'])
```

Loss Function

```
model.summary()
```

Compiles the model with the SGD optimizer

```
logdir = os.path.join("logs", datetime.datetime.now().strftime("%Y%m%d-%H%M%S"))
tensorboard_callback = tf.keras.callbacks.TensorBoard(logdir, histogram_freq=1)
```

Print summary

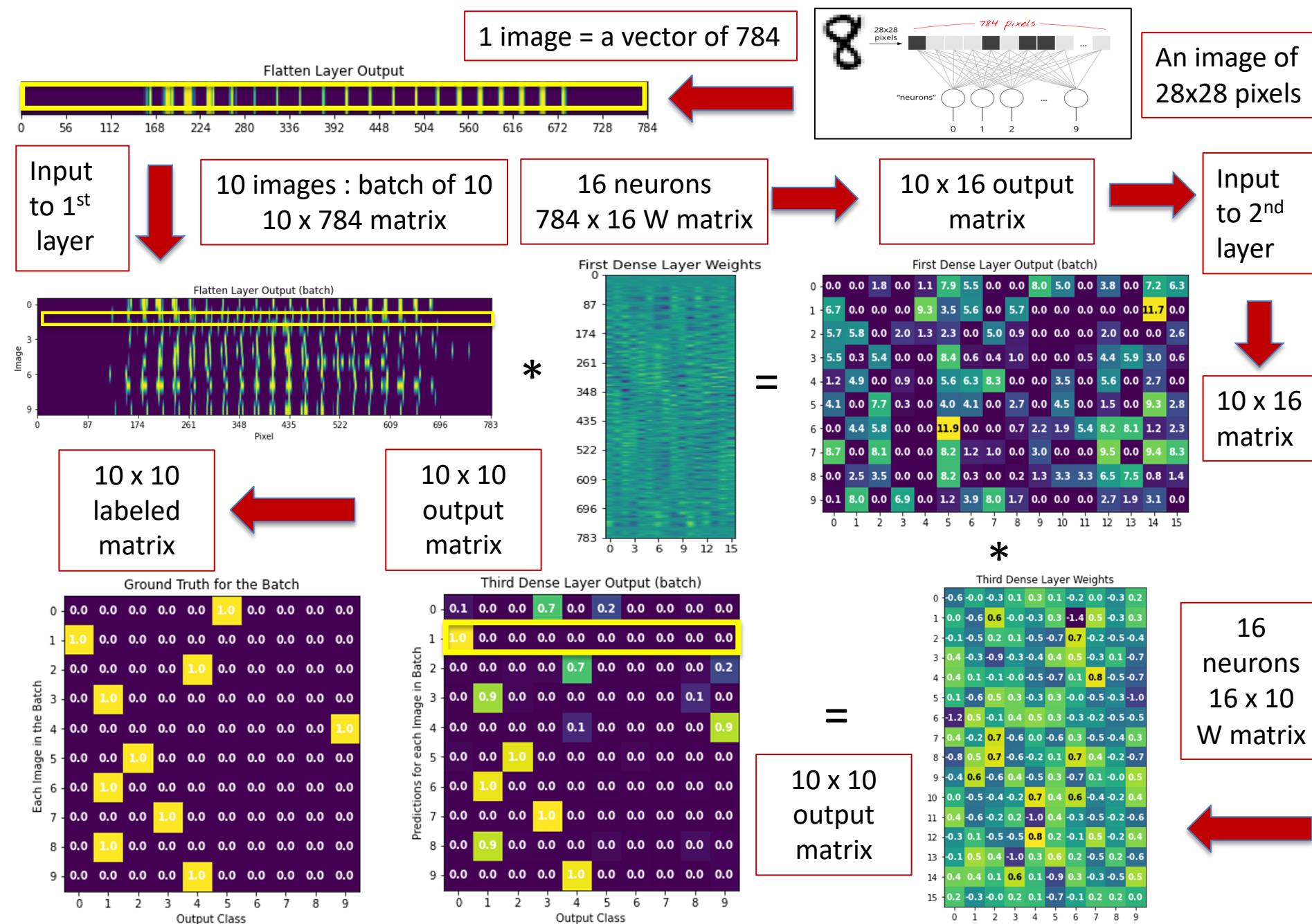
Use tensorborad

```
model.fit(x_train, y_train, epochs=5, batch_size=10,
validation_data=(x_test, y_test), callbacks=[tensorboard_callback])
```

Adjusts model parameters to minimize the loss
Tests the model performance on a test set

```
model.evaluate(x_test, y_test, verbose=2)
%tensorboard --logdir logs
```

Summary : Flow of MLP Dense layer NN



Summary : Flow of NN

```
model.fit(x_train, y_train, epochs=3, batch_size=10)
```

Define Network

Forward Pass

Calculate the analytical gradients

Update weights and bias

backpropagation

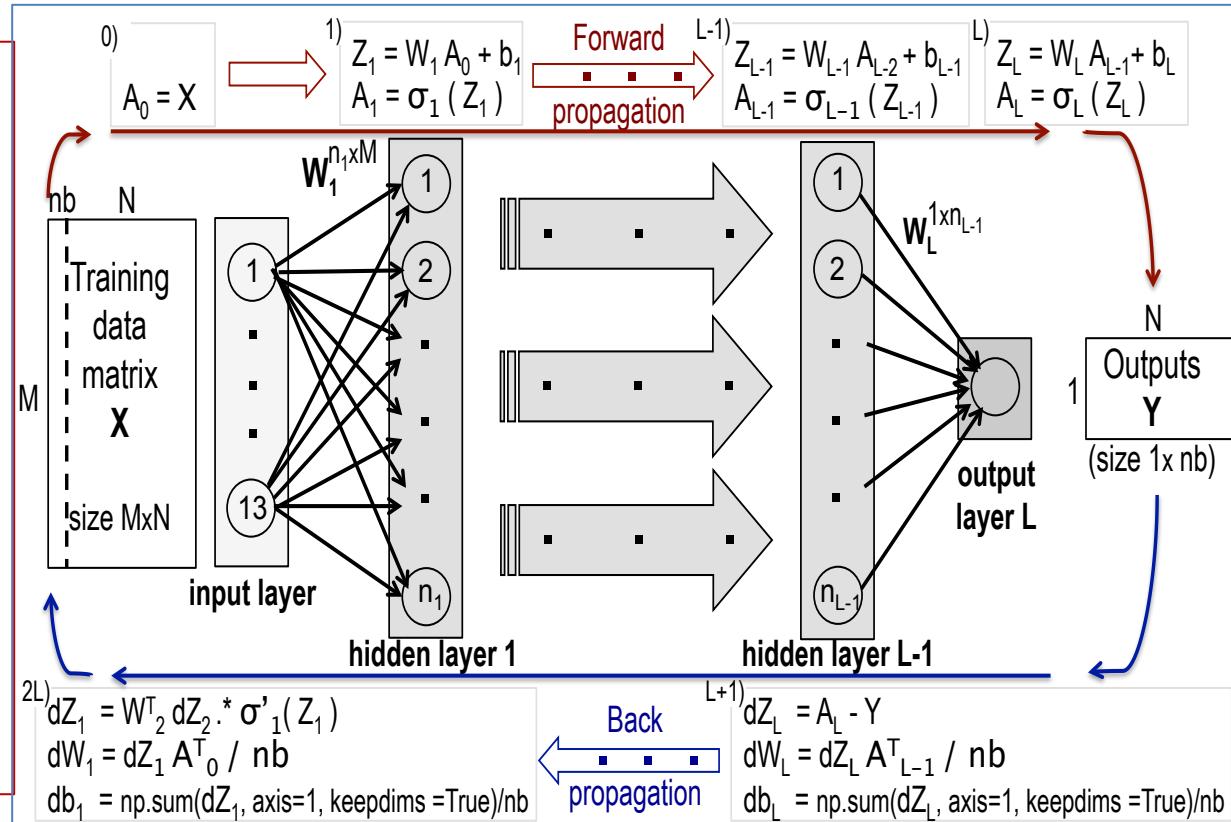
Gradient decent

Quantify loss

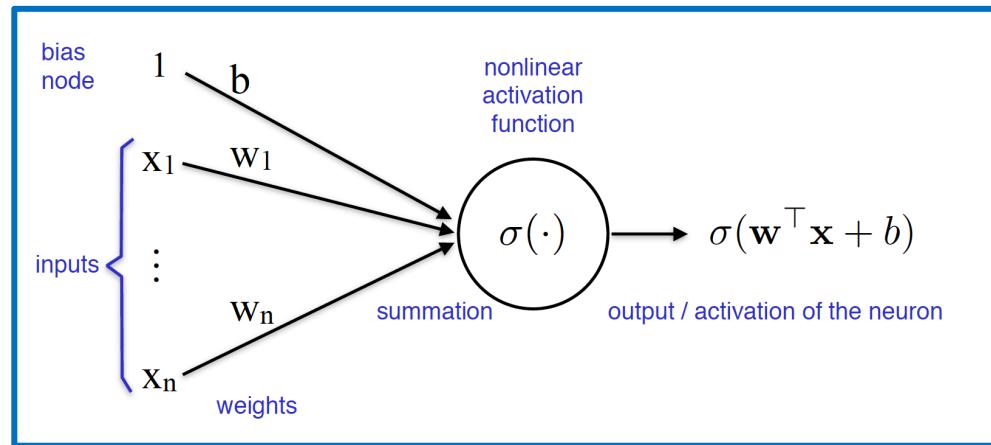
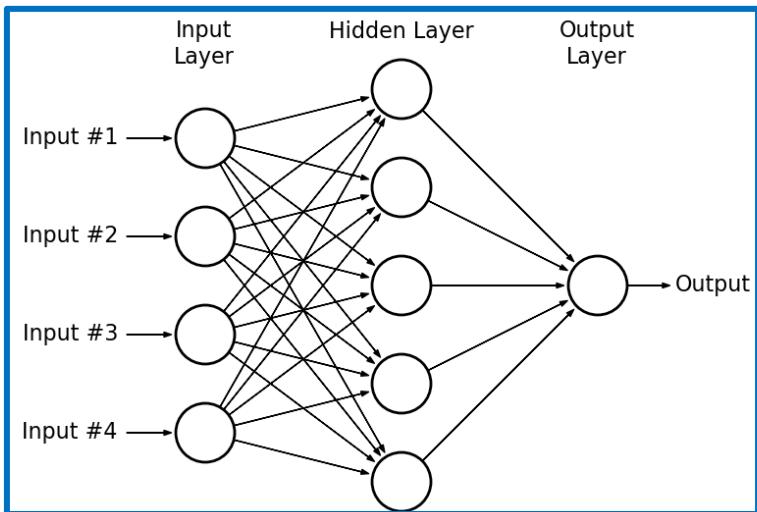
```

1 import numpy as np
2 from numpy.random import randn
3
4 N, D_in, H, D_out = 64, 1000, 100, 10
5 x, y = randn(N, D_in), randn(N, D_out)
6 w1, w2 = randn(D_in, H), randn(H, D_out)
7
8 for t in range(2000):
9     h = 1 / (1 + np.exp(-x.dot(w1)))
10    y_pred = h.dot(w2)
11    loss = np.square(y_pred - y).sum()
12    print(t, loss)
13
14 grad_y_pred = 2.0 * (y_pred - y)
15 grad_w2 = h.T.dot(grad_y_pred)
16 grad_h = grad_y_pred.dot(w2.T)
17 grad_w1 = x.T.dot(grad_h * h * (1 - h))
18
19 w1 -= 1e-4 * grad_w1
20 w2 -= 1e-4 * grad_w2

```



NN Modeling



- ✓ A node in the neural network is a mathematical function or activation function which maps input to output values.
- ✓ Inputs represent a set of vectors containing weights (w) and bias (b). They are the sets of parameters to be determined.
- ✓ Many nodes form a **neural layer**, **links** connect layers together, defining a NN model.
- ✓ Activation function (f or σ), is generally a nonlinear data operator which facilitates identification of complex features.

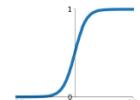
Perceptron

Perceptrons were [developed](#) in the 1950s and 1960s by the scientist [Frank Rosenblatt](#). A perceptron takes several binary inputs, x_1, x_2, \dots , and produces a single binary output. By varying the weights and the threshold, we can get different models of decision-making.

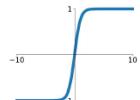
$$\text{output} = \begin{cases} 0 & \text{if } \sum_j w_j x_j \leq \text{threshold} \\ 1 & \text{if } \sum_j w_j x_j > \text{threshold} \end{cases}$$

$$\text{output} = \begin{cases} 0 & \text{if } w \cdot x + b \leq 0 \\ 1 & \text{if } w \cdot x + b > 0 \end{cases}$$

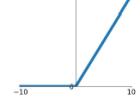
Sigmoid
 $\sigma(x) = \frac{1}{1+e^{-x}}$



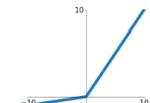
tanh
 $\tanh(x)$



ReLU
 $\max(0, x)$



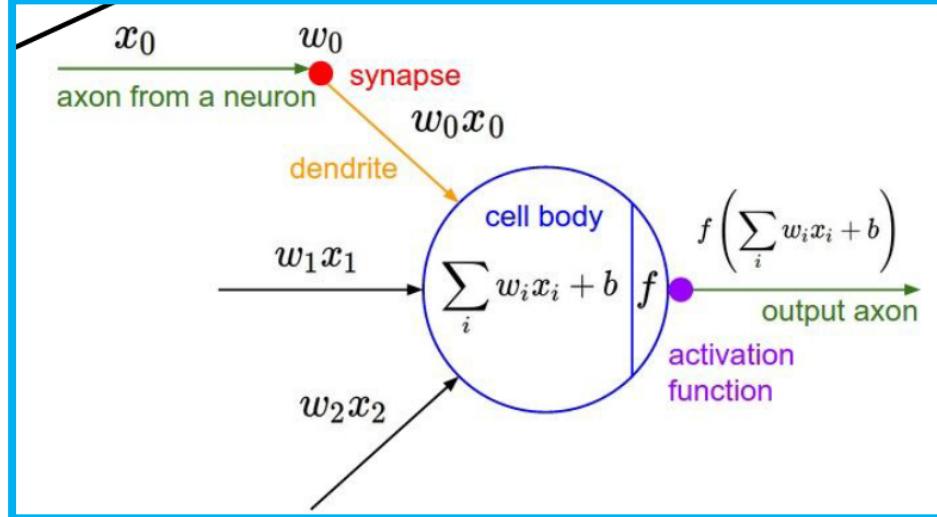
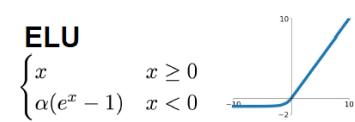
Leaky ReLU
 $\max(0.1x, x)$



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

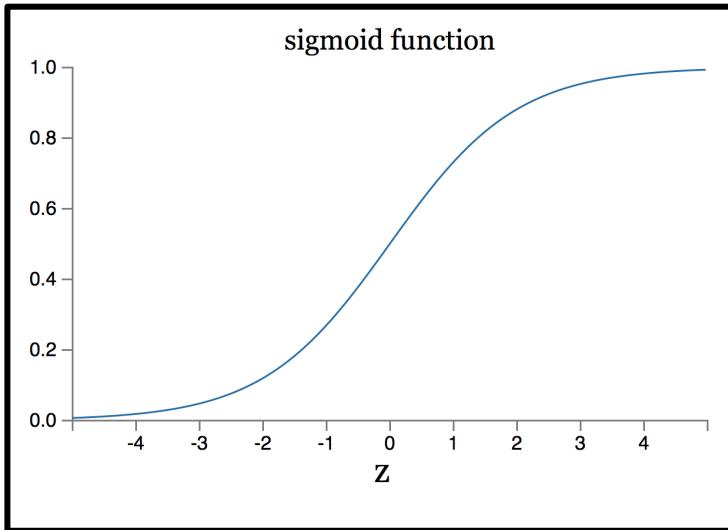
ELU

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

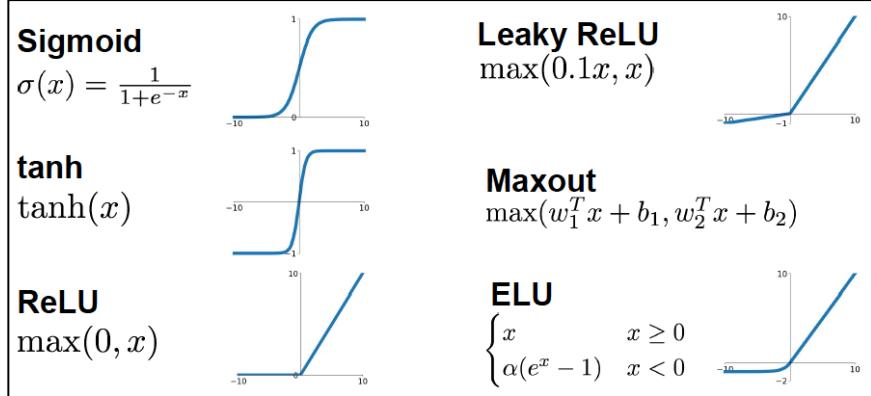


Identity		$f(x) = x$
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Logistic (a.k.a. Soft step)		$f(x) = \frac{1}{1 + e^{-x}}$
TanH		$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$
ArcTan		$f(x) = \tan^{-1}(x)$
Softsign [7][8]		$f(x) = \frac{x}{1 + x }$
Rectified linear unit (ReLU) [9]		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$

Activation Function

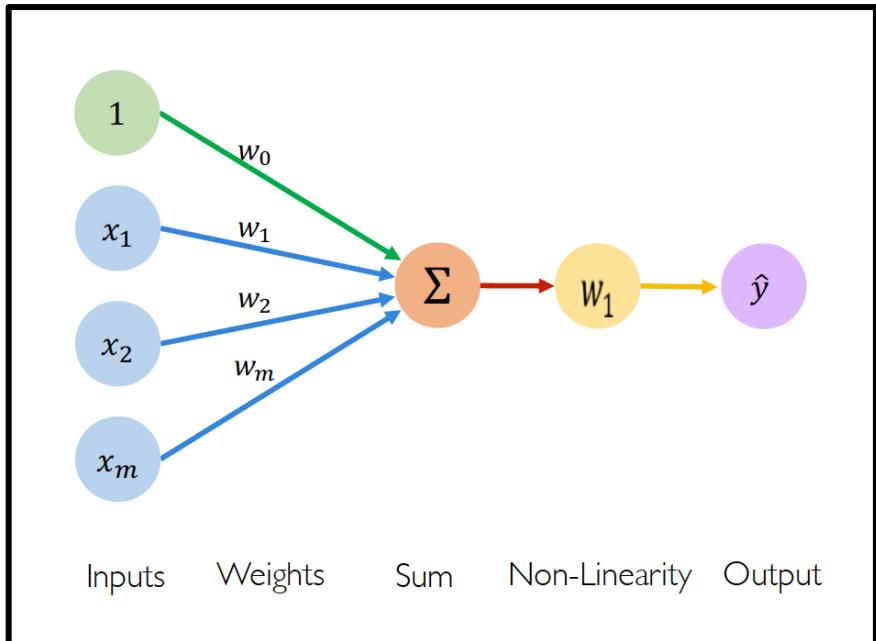


$$\sigma(z) \equiv \frac{1}{1 + e^{-z}}.$$



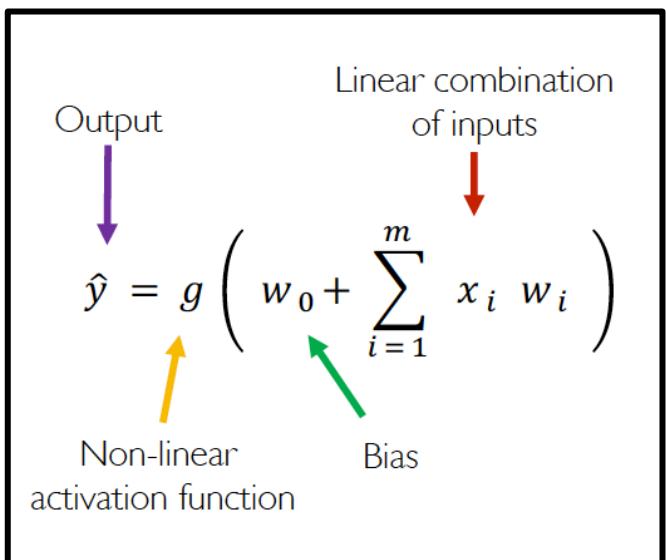
- ✓ In much modern work on neural networks, the main neuron model used is one called the *sigmoid neuron*.
- ✓ Sigmoid (logistic) function, σ , is a smooth out preceptron operating on a set of data, $\sigma(w \cdot x + b)$. The smoothness of σ means that small changes Δw in the weights and Δb in the bias will produce a small change Δoutput in the output from the neuron.
- ✓ If $z=w \cdot x + b$ is a large positive number, then $\exp(-z) \approx 0$ and so $\sigma(z) \approx 1$.
- ✓ If $z=w \cdot x + b$ is very negative, then $\exp(-z) \rightarrow \infty$, and $\sigma(z) \approx 0$. So when $z=w \cdot x + b$ is very negative.
- ✓ It's only when $w \cdot x + b$ is of modest size that there's much deviation from the perceptron model.

DNN MLP Forward Steps



$$\hat{y} = g(w_0 + X^T W)$$

where: $X = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$ and $W = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}$

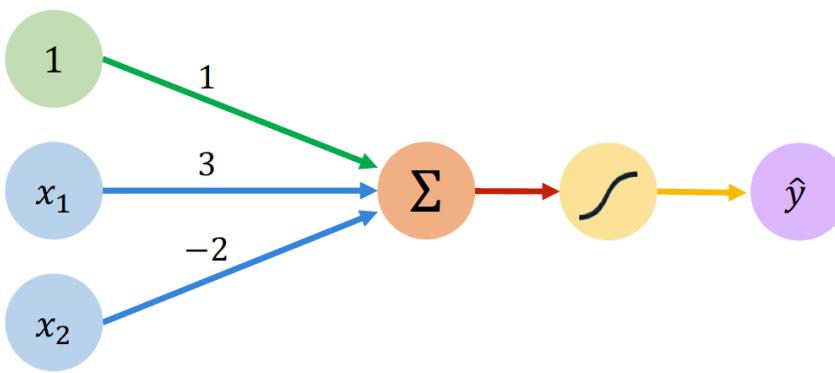


$$\hat{y} = g(w_0 + X^T W)$$

- Example: sigmoid function

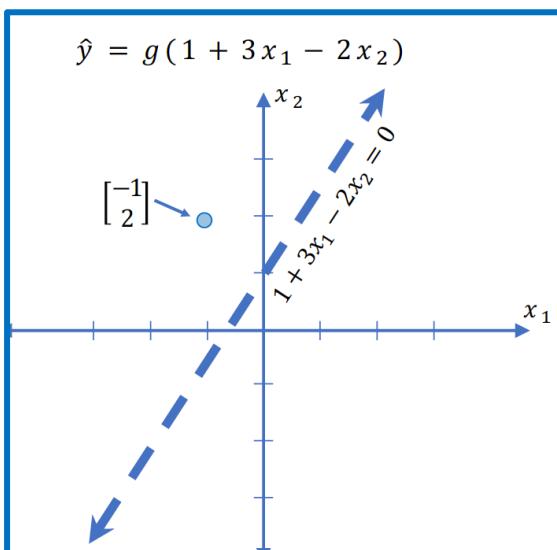
$$g(z) = \sigma(z) = \frac{1}{1 + e^{-z}}$$

Activation Function Example



Assume we have input: $\mathbf{X} = \begin{bmatrix} -1 \\ 2 \end{bmatrix}$

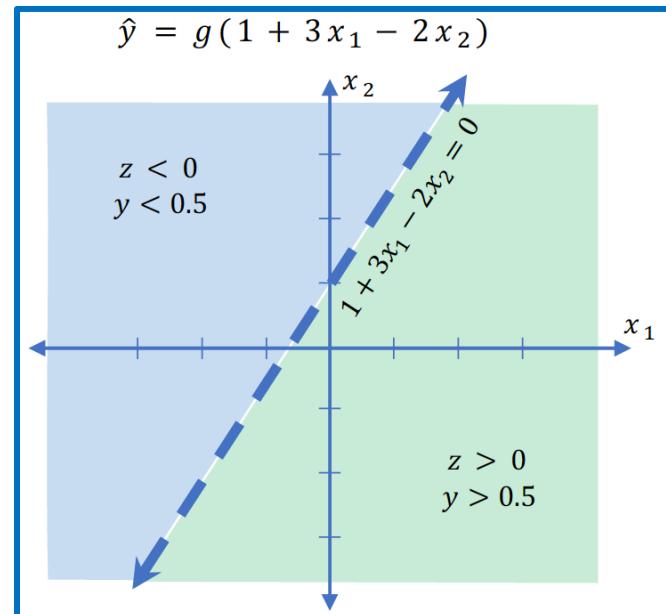
$$\begin{aligned}\hat{y} &= g(1 + (3 * -1) - (2 * 2)) \\ &= g(-6) \approx 0.002\end{aligned}$$



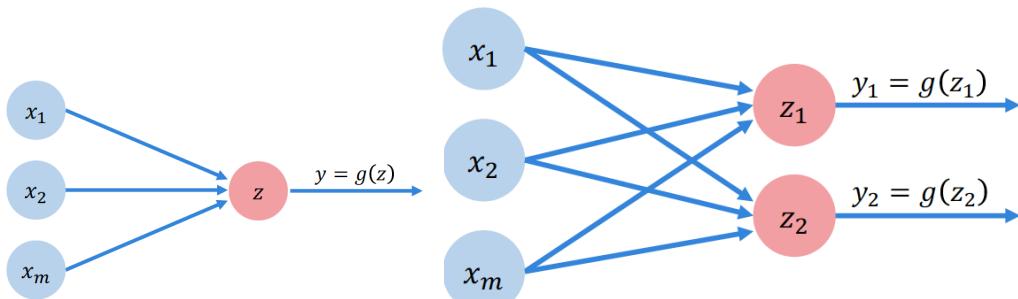
We have: $w_0 = 1$ and $\mathbf{W} = \begin{bmatrix} 3 \\ -2 \end{bmatrix}$

$$\begin{aligned}\hat{y} &= g(w_0 + \mathbf{X}^T \mathbf{W}) \\ &= g\left(1 + \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^T \begin{bmatrix} 3 \\ -2 \end{bmatrix}\right) \\ \hat{y} &= g(1 + 3x_1 - 2x_2)\end{aligned}$$

This is just a line in 2D!

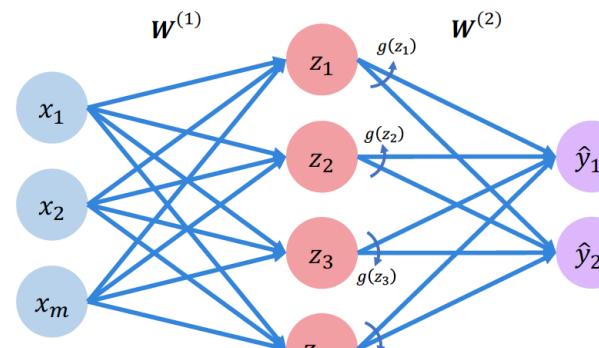


Multilayer Perceptron : Forward Path Calculation



$$z = w_0 + \sum_{j=1}^m x_j w_j$$

$$z_i = w_{0,i} + \sum_{j=1}^m x_j w_{j,i}$$

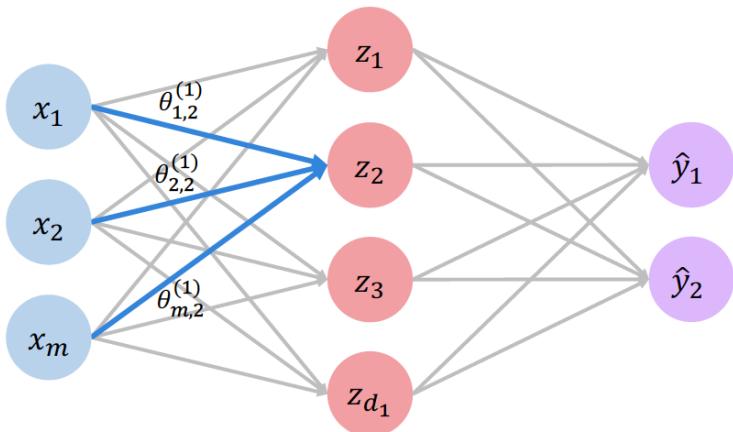


Inputs

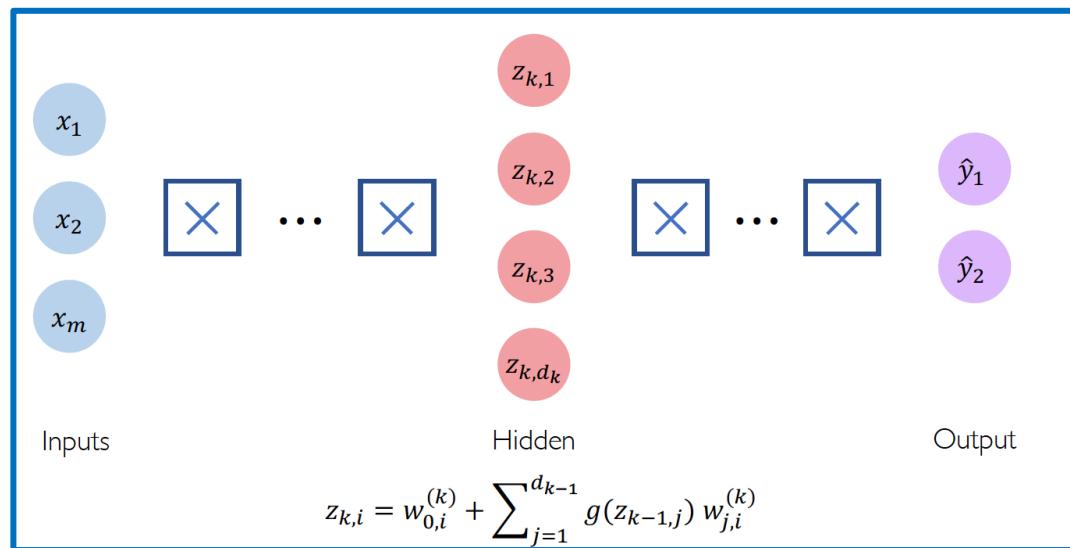
Hidden

Final Output

$$z_i = w_{0,i}^{(1)} + \sum_{j=1}^m x_j w_{j,i}^{(1)} \quad \hat{y}_i = g\left(w_{0,i}^{(2)} + \sum_{j=1}^{d_1} z_j w_{j,i}^{(2)}\right)$$



$$\begin{aligned} z_2 &= w_{0,2}^{(1)} + \sum_{j=1}^m x_j w_{j,2}^{(1)} \\ &= w_{0,2}^{(1)} + x_1 w_{1,2}^{(1)} + x_2 w_{2,2}^{(1)} + x_m w_{m,2}^{(1)} \end{aligned}$$



Inputs

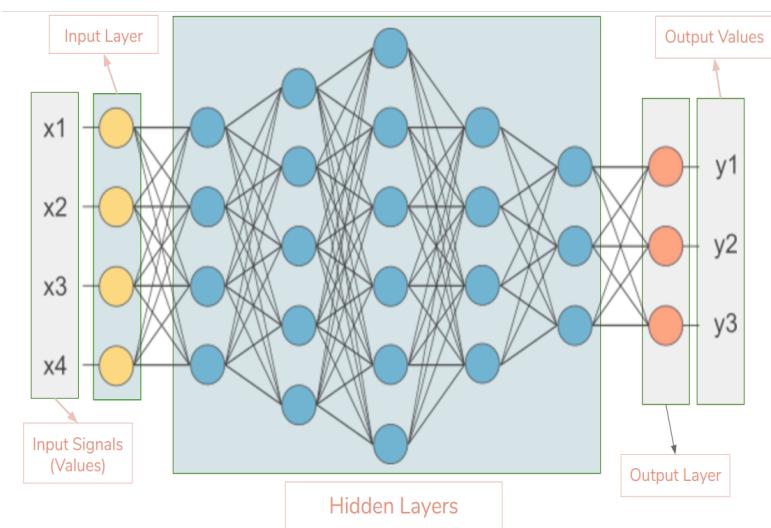
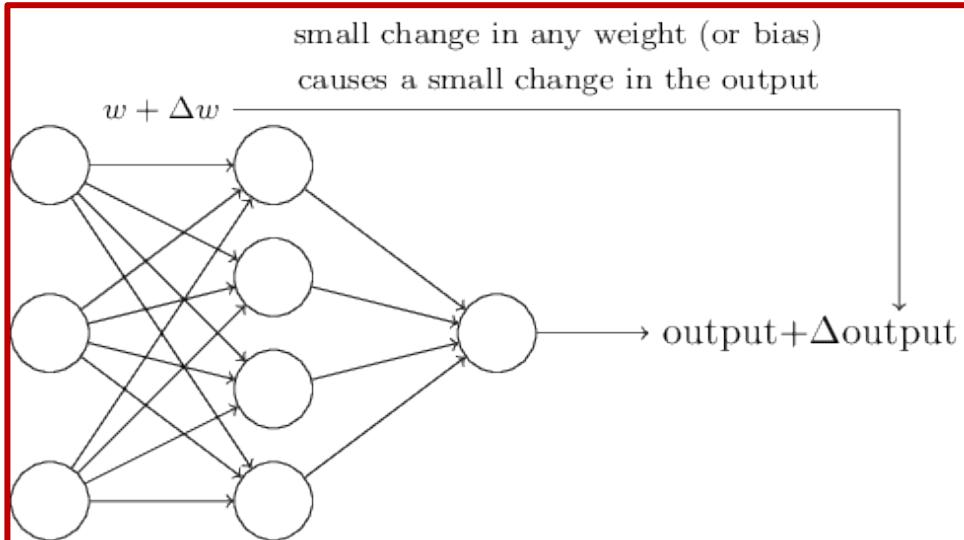
Hidden

Output

$$z_{k,i} = w_{0,i}^{(k)} + \sum_{j=1}^{d_{k-1}} g(z_{k-1,j}) w_{j,i}^{(k)}$$

Representation of NN Model

neuralnetworksanddeeplearning.com



- ✓ Δoutput is approximated using a *linear function* of the changes Δw_j (weights) and Δb (bias). This linearity makes it easy to choose small changes in the weights and biases to achieve any desired small change in the output.

$$\Delta \text{output} \approx \sum_j \frac{\partial \text{output}}{\partial w_j} \Delta w_j + \frac{\partial \text{output}}{\partial b} \Delta b,$$

output= cost function

- ✓ We then cast the model of the NN model as an optimization problem, minimizing the changes of the cost as functions of changes of weight and bias

Minimization : Iteration Scheme

neuralnetworksanddeeplearning.com

$$C(w, b) \equiv \frac{1}{2n} \sum_x \|y(x) - a\|^2.$$

Define a **cost function**, C is the **quadratic cost function** also referred as the *mean squared error(MSE)*.

Find a scheme to minimize the cost $C(w,b)$ as a function of the weights and biases, casting it as an optimization problem using the **gradient descent algorithm**.

Find a way of iterating Δw_j and Δb so as to make ΔC (Δ output) negative,
Pushing MSE (C) smaller and smaller, ideally to zero

Be negative

=

$$\Delta \text{output} \approx \sum_j \frac{\partial \text{output}}{\partial w_j} \Delta w_j + \frac{\partial \text{output}}{\partial b} \Delta b,$$

$$\nabla C \equiv \left(\frac{\partial C}{\partial v_1}, \frac{\partial C}{\partial v_2} \right)^T.$$

For output = C ; $wj = v1$; $b = v2$

$$\Delta C \approx \frac{\partial C}{\partial v_1} \Delta v_1 + \frac{\partial C}{\partial v_2} \Delta v_2.$$

$$\Delta C \approx \nabla C \cdot \Delta v.$$

Be negative

Choose

$$\Delta C \approx \nabla C \cdot \Delta v.$$

$$\Delta v = -\eta \nabla C,$$

$$v \rightarrow v' = v - \eta \nabla C.$$

$$\Delta C \approx -\eta \nabla C \cdot \nabla C = -\eta \|\nabla C\|^2.$$

$$w_k \rightarrow w'_k = w_k - \eta \frac{\partial C}{\partial w_k}$$
$$b_l \rightarrow b'_l = b_l - \eta \frac{\partial C}{\partial b_l}.$$

Mini-Batch Stochastic Gradient Descent

- ✓ Neural network is an optimization process to find a set of parameters (W , weights) that ensure the prediction function $f(x, W)$ to be as close as possible to the true solution (labeled input) y for any input x . We use gradient descent to find the weights w and biases b which minimize the cost function, C .
- ✓ To compute the gradient ∇C we need to compute the gradients ∇C_x separately for each training input, x , and then average them, $\nabla C = 1/n \sum \nabla C_x$. Unfortunately, when the number of training inputs is very large this can take a long time.
- ✓ A way is to use mini-batched **stochastic gradient descent** to speed up learning. The idea is to estimate the gradient ∇C by computing a small sample of randomly chosen training inputs, refer to as a **mini-batch** of input (mini-batched SGD).
- ✓ By averaging over a small sample, we can quickly get a good estimate of the true gradient ∇C , and this helps speed up gradient descent, and thus learning, provided the sample size m is large enough that the average value of the ∇C_{X_j} roughly equals to the average over all ∇C_x .
- ✓ Then, another randomly chosen mini-batch are selected and trained, until all the training inputs are used. It is said to complete an **epoch (iteration)** of training.

$$\frac{\sum_{j=1}^m \nabla C_{X_j}}{m} \approx \frac{\sum_x \nabla C_x}{n} = \nabla C,$$

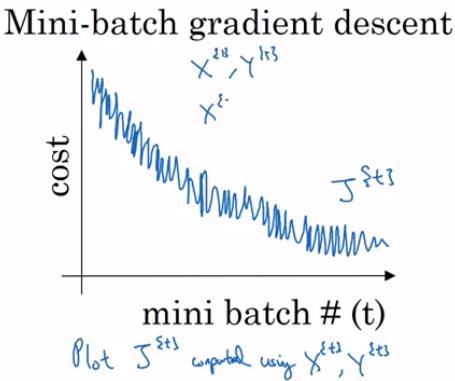
$$w_k \rightarrow w'_k = w_k - \eta \frac{\partial C}{\partial w_k}$$
$$b_l \rightarrow b'_l = b_l - \eta \frac{\partial C}{\partial b_l}.$$

Optimization: Mini-Batch SGD

Instead of computing a gradient across the entire dataset with size N , divides the dataset into a fixed-size subset (“minibatch”) with size m , and computes the gradient for each update on this smaller batch:

In one epoch

1. Pick a mini-batch
2. Feed it to Neural Network
3. Calculate the mean gradient of the mini-batch
4. Use the mean gradient calculated in step 3 to update the weights
5. Repeat steps 1–4 for the mini-batches we created



(N is the dataset size, m is the minibatch size)

Given: training set $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$

Initialize all $\Theta^{(l)}$ randomly (NOT to 0!)

Loop // each iteration is called an epoch

Loop // each iteration is a mini-batch

Set $\Delta_{i,j}^{(l)} = 0 \quad \forall l, i, j$ (Used to accumulate gradient)

Sample m training instances $\mathcal{X} = \{(\mathbf{x}'_1, y'_1), \dots, (\mathbf{x}'_m, y'_m)\}$ without replacement

For each instance in \mathcal{X} , (\mathbf{x}_k, y_k) :

Set $\mathbf{a}^{(1)} = \mathbf{x}_k$

Compute $\{\mathbf{a}^{(2)}, \dots, \mathbf{a}^{(L)}\}$ via forward propagation

Compute $\delta^{(L)} = \mathbf{a}^{(L)} - y_k$

Compute errors $\{\delta^{(L-1)}, \dots, \delta^{(2)}\}$

Compute gradients $\Delta_{ij}^{(l)} = \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$

Compute mini-batch regularized gradient $D_{ij}^{(l)} = \begin{cases} \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)} & \text{if } j \neq 0 \\ \frac{1}{m} \Delta_{ij}^{(l)} & \text{otherwise} \end{cases}$

Update weights via gradient step $\Theta_{ij}^{(l)} = \Theta_{ij}^{(l)} - \alpha D_{ij}^{(l)}$

Until all training instances are seen

Until weights converge or max #epochs is reached

DNN Forward Backward Calculation : Weights and bias

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
4. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
5. Return weights



```
weights = tf.random_normal(shape, stddev=sigma)
```



```
grads = tf.gradients(ys=loss, xs=weights)
```



```
weights_new = weights.assign(weights - lr * grads)
```

In one epoch

1. Pick a mini-batch
2. Feed it to Neural Network
3. Forward path calculation
4. Calculate the mean gradient of the mini-batch
5. Use the mean gradient calculated in step 4 to update the weights
6. Repeat steps 1–5 for the mini-batches we created

$$\frac{\sum_{j=1}^m \nabla C_{X_j}}{m} \approx \frac{\sum_x \nabla C_x}{n} = \nabla C,$$

$$w_k \rightarrow w'_k = w_k - \eta \frac{\partial C}{\partial w_k}$$
$$b_l \rightarrow b'_l = b_l - \eta \frac{\partial C}{\partial b_l}.$$

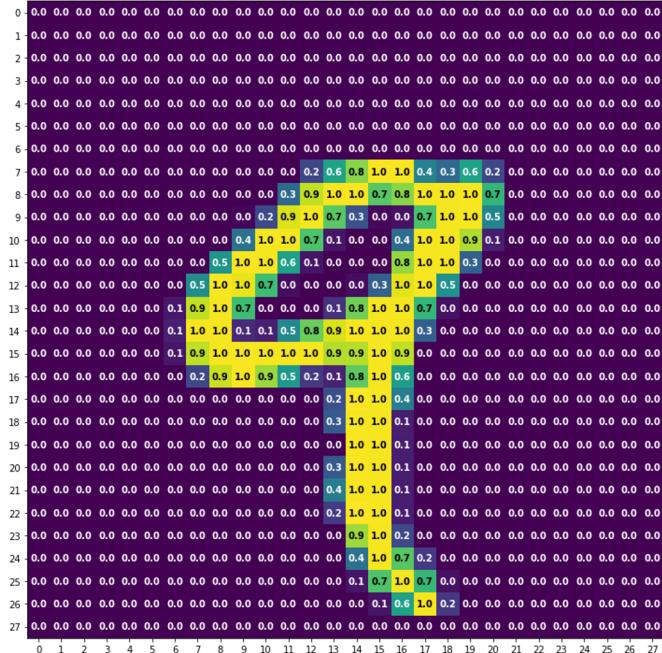
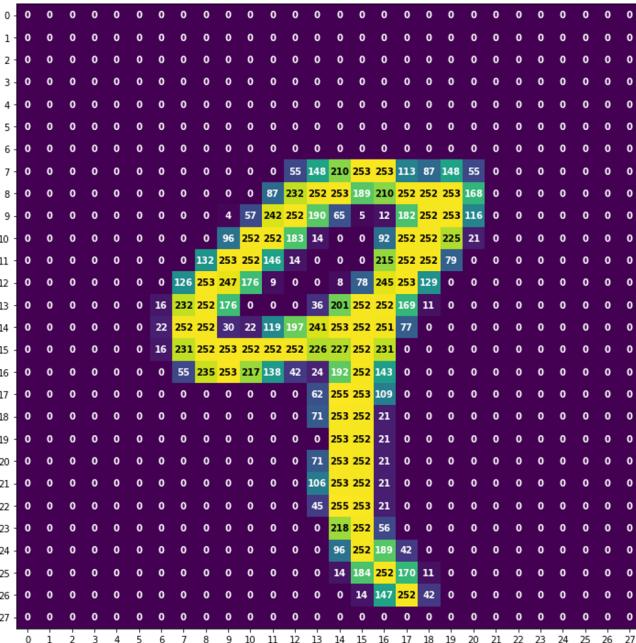
TensorFlow 2.0

Loading and Preparing the MNIST Dataset

```
import tensorflow as tf  
  
mnist = tf.keras.datasets.mnist  
  
#download the images  
(x_train, y_train), (x_test, y_test) =  
mnist.load_data()  
  
x_train, x_test = x_train / 255.0,  
x_test / 255.0 #normalize
```

Each Image is a 28x28 image of a handwritten digit

x 60,000 images

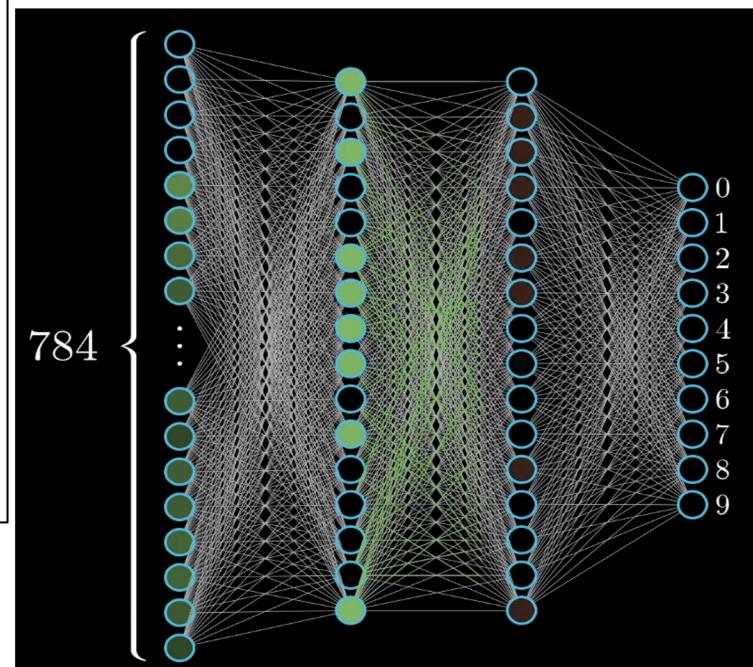


TensorFlow 2.0 : NN Network

Let's use this simple network to demonstrate what TensorFlow is doing behind the scenes

```
model = tf.keras.models.Sequential([
    tf.keras.layers.Flatten(input_shape=(28, 28)),
    tf.keras.layers.Dense(16, activation='relu'),
    tf.keras.layers.Dense(16, activation='relu'),
    tf.keras.layers.Dense(10, activation='softmax')
])
model.compile(optimizer='sgd',
              loss='Categorical_Crossentropy',
              metrics=['accuracy'])
Model.summary
model.fit(x_train, y_train, epochs=3, batch_size=10)
```

Simple MLP Network



- ✓ **Batch Size:** how many images do we pass through the network for each iteration (10)
- ✓ **Epochs:** how many times we pass over the entire dataset (3)
- ✓ **# Iterations per Epoch** = number of images / batch size (60,000/10=6,000)

Image:
<https://www.3blue1brown.com/>

<https://codelabs.developers.google.com/codelabs/cloud-tensorflow-mnist/#0>

TensorFlow 2.0 : Forward Pass

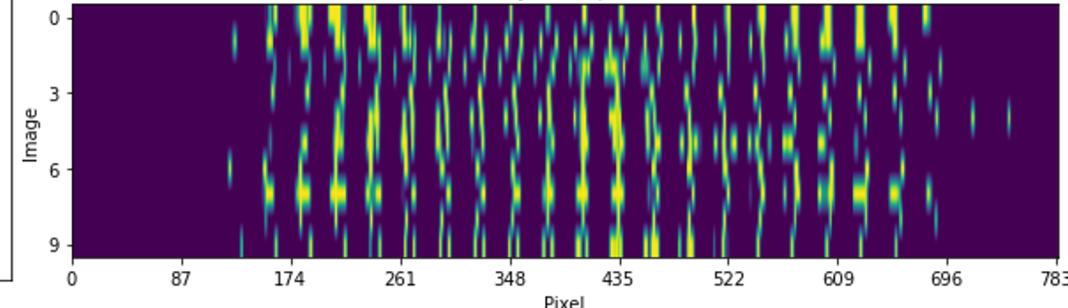
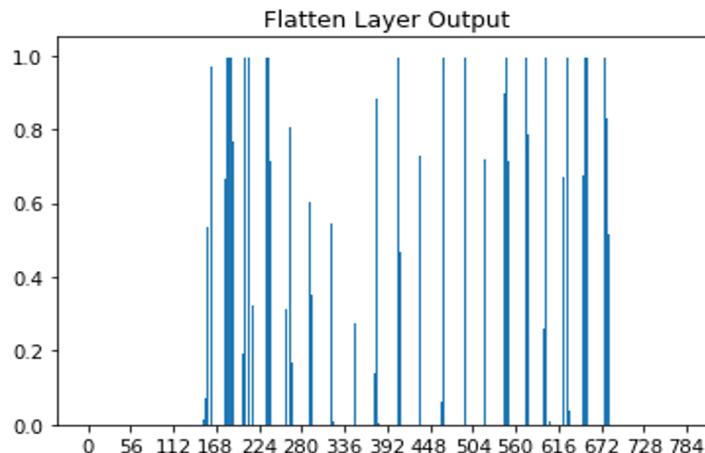
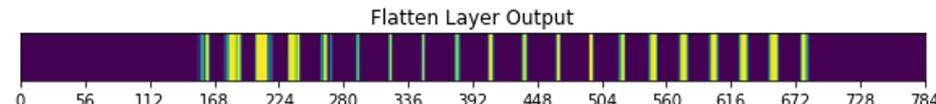
```
tf.reshape(tensor, shape)
#example
tf.reshape(x_train[0], [784])
```

Flatten a single image

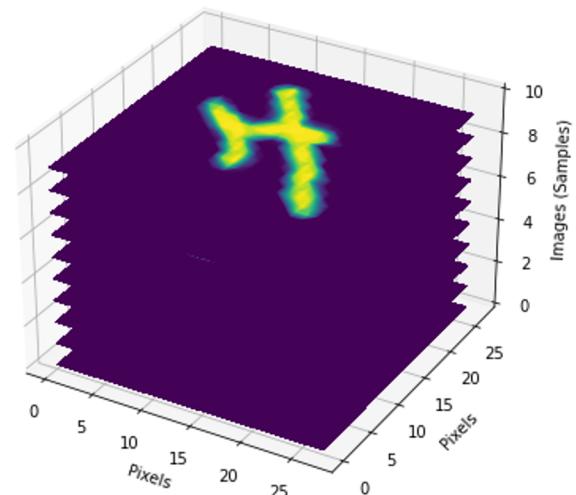
```
tf.reshape(tensor, shape)
#example
tf.reshape(x_train[0:10], [10: 784])
```

Flatten a batch of images

```
model = tf.keras.models.Sequential([
    tf.keras.layers.Flatten(input_shape=(28, 28))
```



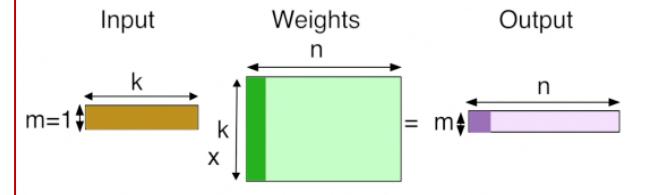
Network Input Per Iteration



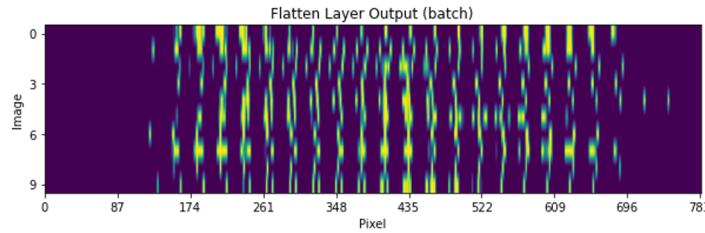
TensorFlow 2.0 : Forward Pass

1ST Fully Connected Layer with 16 neurons Matrix multiplication

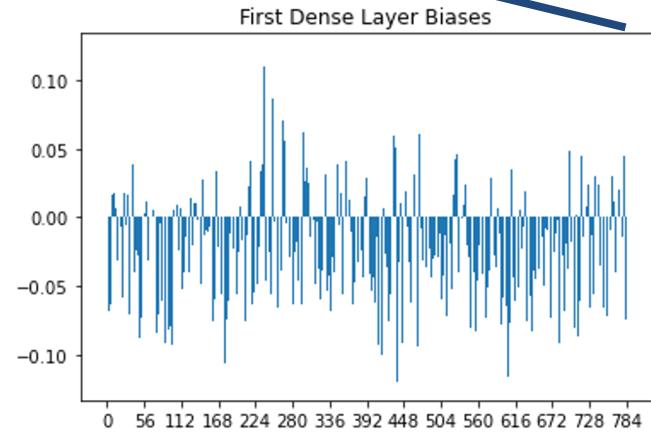
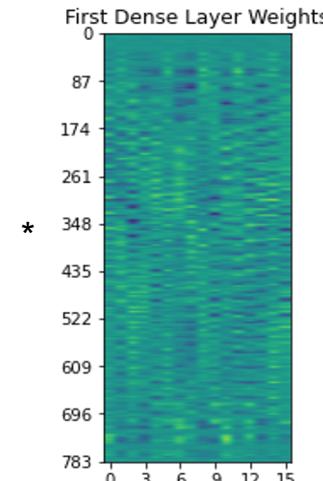
```
...  
tf.keras.layers.Flatten(input_shape=(28, 28)),  
tf.keras.layers.Dense(16, activation='relu'),  
...
```



$$a^1 = \text{ReLU}(x * W^1 + b^1)$$



Weights are initialized with random values

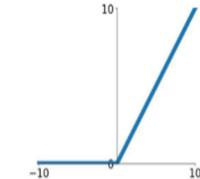


TensorFlow 2.0 : Forward Pass

Second Fully Connected Layer with 16 neurons

```
model = tf.keras.models.Sequential([
    tf.keras.layers.Flatten(input_shape=(28, 28)),
    tf.keras.layers.Dense(16, activation='relu'),
    tf.keras.layers.Dense(16, activation='relu'),
    ...
```

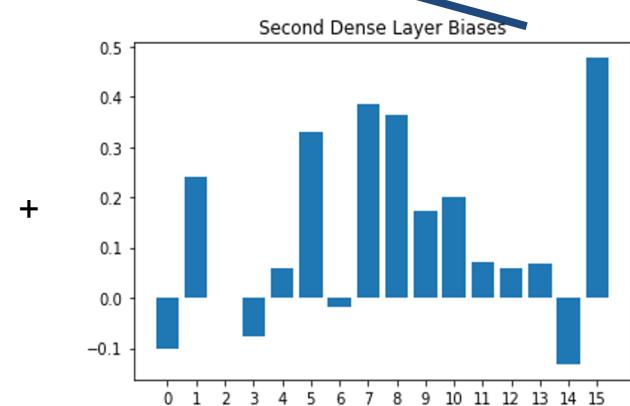
ReLU
 $\max(0, x)$



$$a^2 = \text{ReLU}(a^1 * W^2 + b^2)$$

First Dense Layer Output (batch)																
0	0.0	0.0	1.8	0.0	1.1	7.9	5.5	0.0	0.0	8.0	5.0	0.0	3.8	0.0	7.2	6.3
1	6.7	0.0	0.0	0.0	9.3	3.5	5.6	0.0	5.7	0.0	0.0	0.0	0.0	0.0	11.7	0.0
2	5.7	5.8	0.0	2.0	1.3	2.3	0.0	5.0	0.9	0.0	0.0	0.0	2.0	0.0	0.0	2.6
3	5.5	0.3	5.4	0.0	0.0	8.4	0.6	0.4	1.0	0.0	0.0	0.5	4.4	5.9	3.0	0.6
4	1.2	4.9	0.0	0.9	0.0	5.6	6.3	8.3	0.0	0.0	3.5	0.0	5.6	0.0	2.7	0.0
5	4.1	0.0	7.7	0.3	0.0	4.0	4.1	0.0	2.7	0.0	4.5	0.0	1.5	0.0	9.3	2.8
6	0.0	4.4	5.8	0.0	0.0	11.9	0.0	0.0	0.7	2.2	1.9	5.4	8.2	8.1	1.2	2.3
7	8.7	0.0	8.1	0.0	0.0	8.2	1.2	1.0	0.0	3.0	0.0	0.0	9.5	0.0	9.4	8.3
8	0.0	2.5	3.5	0.0	0.0	8.2	0.3	0.0	0.2	1.3	3.3	3.3	6.5	7.5	0.8	1.4
9	0.1	8.0	0.0	6.9	0.0	1.2	3.9	8.0	1.7	0.0	0.0	0.0	2.7	1.9	3.1	0.0

Second Dense Layer Weights																
15	0.0	0.4	0.1	-0.2	-0.4	0.9	0.6	-0.4	-0.1	0.2	-0.3	0.9	-0.6	-0.3	-0.2	0.0
14	-0.3	0.2	0.0	0.7	-0.1	0.5	-0.1	0.5	-0.1	0.1	0.6	0.4	-0.3	0.3	0.2	-0.1
13	0.0	-0.6	-0.3	0.8	-0.4	-0.8	0.1	-0.3	0.2	0.2	0.3	0.2	0.5	0.9	-0.2	-0.6
12	-0.5	-0.2	-0.3	-0.3	0.1	0.1	0.2	-0.4	-0.5	0.7	0.3	0.0	0.6	-0.4	0.8	0.6
11	0.2	-0.2	0.8	0.2	0.9	-0.2	0.0	-0.1	1.1	-0.1	0.4	-0.7	0.0	0.4	-0.2	0.0
10	0.1	1.2	0.2	-0.6	0.2	-0.1	-0.8	0.1	-0.0	1.0	-0.1	0.2	-0.3	-0.4	0.2	0.2
9	0.1	-0.1	-0.2	0.1	-0.3	0.4	0.9	-0.7	-0.1	0.7	0.3	0.1	-0.6	0.1	-0.2	-0.1
8	0.2	-0.3	-0.3	0.1	0.4	0.2	0.2	0.6	0.6	-0.5	0.4	0.2	-0.4	-0.2	0.4	0.3
7	-0.1	0.4	-0.2	-0.4	-0.4	-0.5	0.1	-0.1	-0.6	0.0	0.8	-0.3	1.1	0.2	-0.2	0.2
6	0.2	1.0	-0.3	0.2	0.1	-0.1	-0.3	-0.2	-0.2	0.4	-0.2	0.6	0.7	-0.0	-0.5	0.2
5	0.0	0.1	-0.0	0.0	0.4	-0.2	-0.5	-0.1	0.5	0.6	0.2	0.2	-0.1	0.4	0.6	-0.4
4	-0.3	0.4	-0.3	0.7	0.8	-0.2	-0.2	-0.1	-0.5	-0.3	-0.4	0.8	0.2	-0.4	-0.1	0.5
3	-0.5	0.1	0.1	-0.4	0.1	0.3	0.8	0.1	0.3	-0.5	0.1	0.2	0.1	0.6	-0.5	-0.2
2	-0.1	0.8	-0.1	-0.4	0.4	0.7	0.1	-0.1	0.6	-0.1	-1.0	-0.2	-0.2	-0.1	0.2	0.3
1	-0.4	-0.1	0.6	0.2	0.1	0.2	0.4	0.2	0.4	-0.5	0.5	-0.8	0.6	-0.2	0.6	0.6
0	0.1	0.3	-0.2	-0.6	0.0	-0.2	-0.3	0.5	-0.1	0.4	-0.5	0.1	0.2	0.1	0.8	0.4

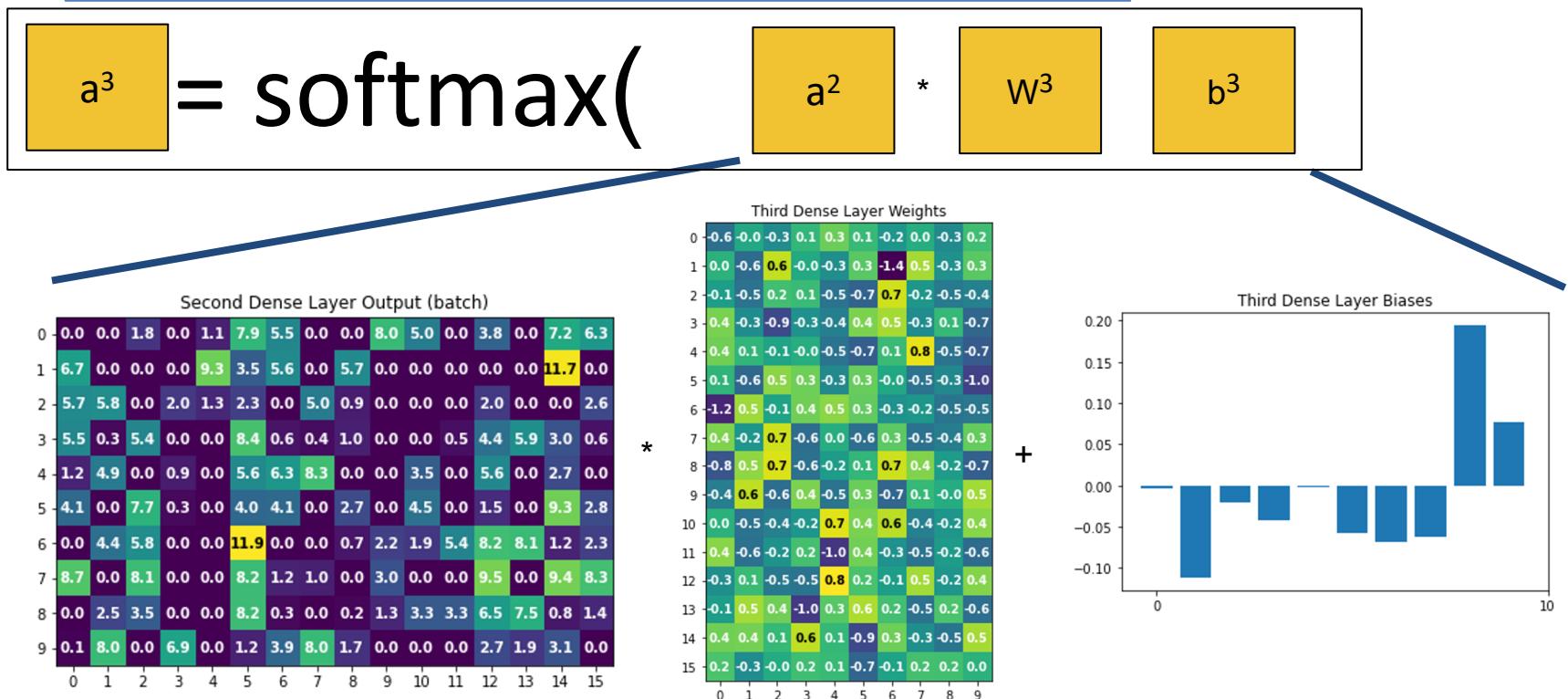


Final Fully Connected Layer with 10 neurons : 10 classes

$$\text{Softmax} = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

- ✓ The softmax function is a function that turns a vector of K real values into values between 0 and 1, and the values sum up to 1, so that they can be interpreted as probabilities.
- ✓ This is because the softmax is a generalization of logistic regression that can be used for multi-class classification.

<https://deeppi.org/machine-learning-glossary-and-terms/softmax-layer>



TensorFlow 2.0 : Forward Pass

`tf.losses.mean_squared_error(y_true_tf, y_pred_tf)`

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

Evaluating the Error

Softmax Classifier (Multinomial Logistic Regression)



Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

Probabilities
must be ≥ 0

cat	3.2
car	5.1
frog	-1.7

Unnormalized log-probabilities / logits

24.5
164.0
0.18

unnormalized probabilities

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

Softmax Function

Probabilities
must sum to 1

0.13
0.87
0.00

probabilities

1.00
0.00
0.00

Cross Entropy

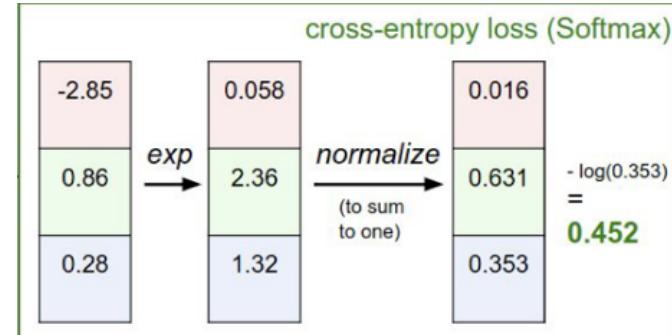
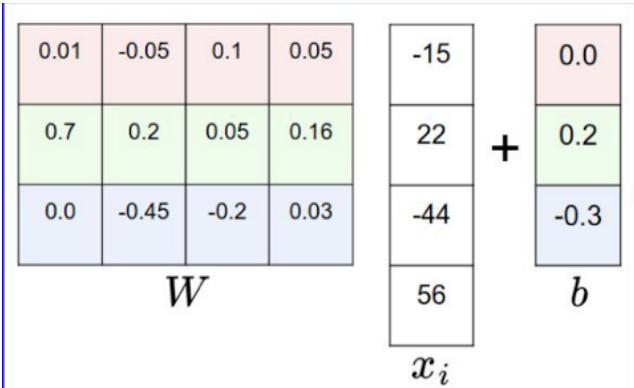
$$H(P, Q) = H(p) + D_{KL}(P||Q)$$

Correct probs

exp

normalize

compare

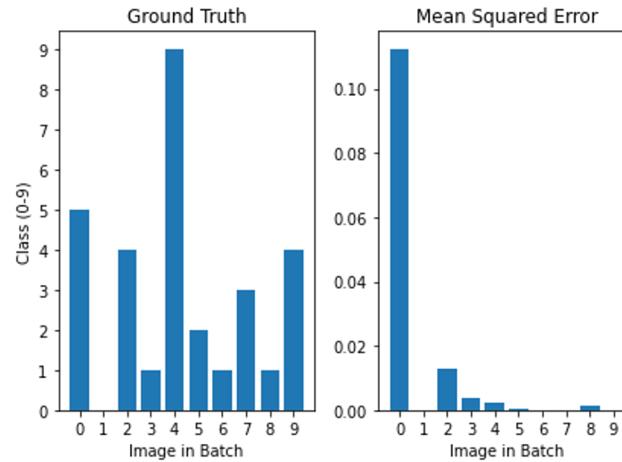
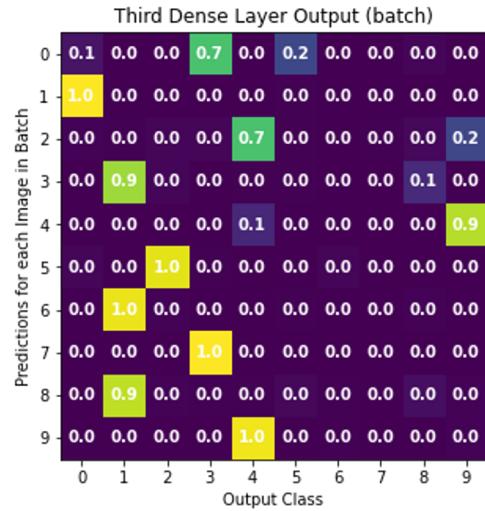


TensorFlow 2.0 : Forward Pass

`tf.losses.mean_squared_error(y_true_tf, y_pred_tf)`

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

Evaluating the Error



Softmax Classifier (Multinomial Logistic Regression)



Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

Softmax Function

cat
car
frog

Unnormalized log-probabilities / logits	exp	normalize	probabilities
3.2	5.1	24.5	0.13
5.1	-1.7	164.0	0.87
-1.7	0.18	0.00	0.00

Probabilities must be ≥ 0

Probabilities must sum to 1

$$L_i = -\log P(Y = y_i | X = x_i)$$

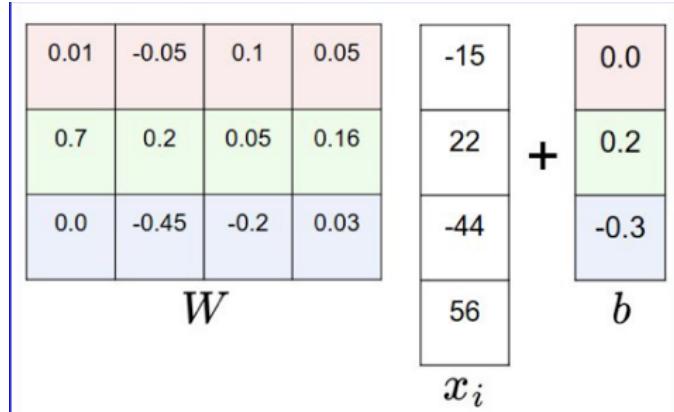
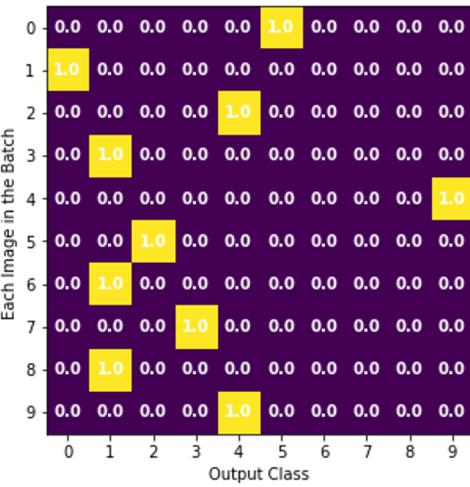
Cross Entropy

$$H(P, Q) = H(p) + D_{KL}(P || Q)$$

Correct probs

$$L_i = -\log\left(\frac{e^{sy_i}}{\sum_j e^{sj}}\right)$$

Ground Truth for the Batch



cross-entropy loss (Softmax)

$$L_i = -\log\left(\frac{e^{sy_i}}{\sum_j e^{sj}}\right)$$

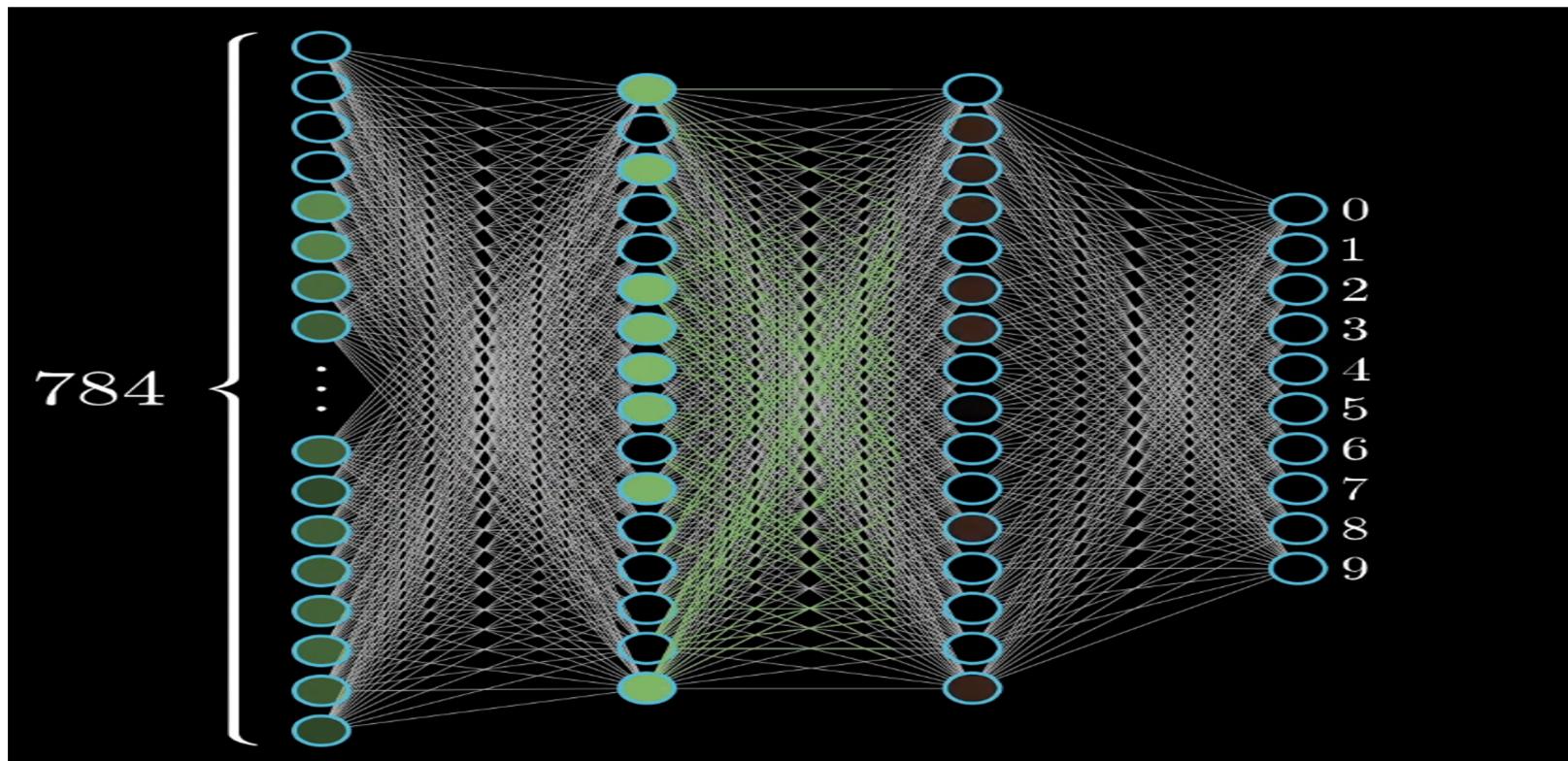
exp → normalize (to sum to one)

-2.85	0.058	0.016
0.86	2.36	0.631
0.28	1.32	0.353

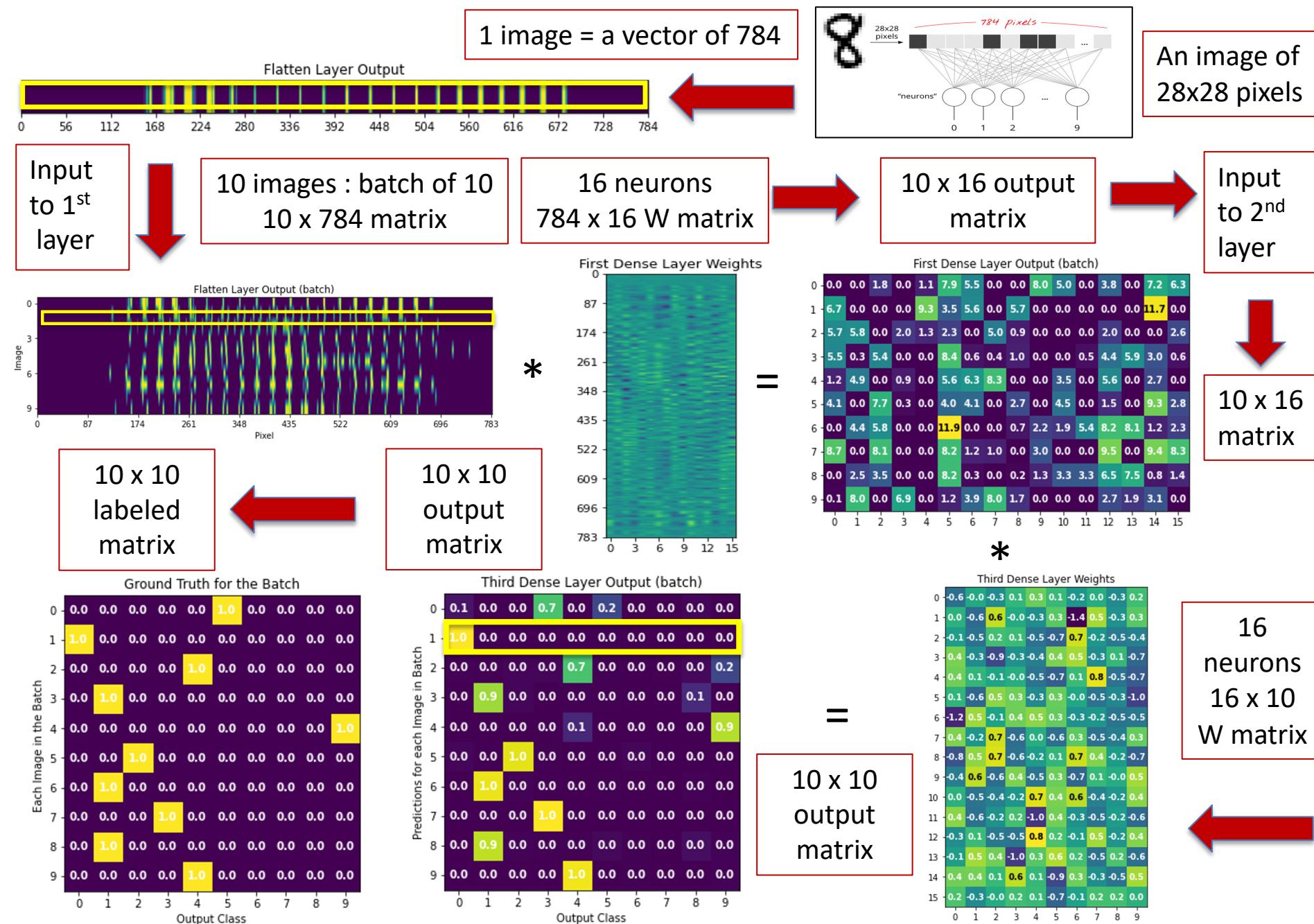
$= -\log(0.353) = 0.452$

Simple MNIST MLP Network

Google Colab Code



Summary : Flow of MLP Dense layer NN

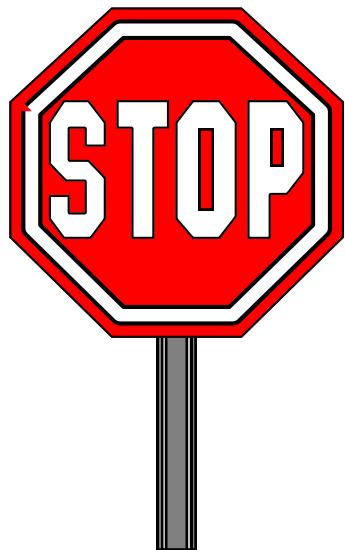


Acknowledgements and References

This portion of the tutorial contains many extracted materials from many online websites and courses. Listed below are the major sites. This portion of the materials is not intended for public distribution. Please visit the sites websites for detail contents. If you are beginner users of DNN, I suggest to read the following list of websites in its order.

- 1) <http://neuralnetworksanddeeplearning.com>
- 2) <https://codelabs.developers.google.com/codelabs/cloud-tensorflow-mnist>
- 3) <https://www.deeplearning.ai/deep-learning-specialization/>
- 4) <http://cs231n.stanford.edu/>
- 5) MIT 6.S191, <https://www.youtube.com/watch?v=njKP3FqW3Sk>
- 6) <https://livebook.manning.com/book/deep-learning-with-python/about-this-book/>
- 7) <https://www.fast.ai/>
- 8) <https://petewarden.com/2015/04/20/why-gemm-is-at-the-heart-of-deep-learning/>
- 9) <https://www.nersc.gov/users/training/gpus-for-science/gpus-for-science-2020/>
- 10) <https://oneapi-src.github.io/oneDNN/>
- 11) <http://www.cs.cornell.edu/courses/cs4787/2020sp/>
- 12) More sites and free books are listed in www.jics.utk.edu/actia → ML
- 13) <https://machinelearningmastery.com>

The End



- The End!

