

# SCRP HPC Cluster

HPC cluster managed by the CUHK Department of Economics

- Run your code on the cloud
- All necessary software installed and maintained
- GPU access
- scrp.econ.cuhk.edu.hk



#### **Step 1: Account Creation**

Accounts have been already been created for all participants.

#### Step 2: Create Password

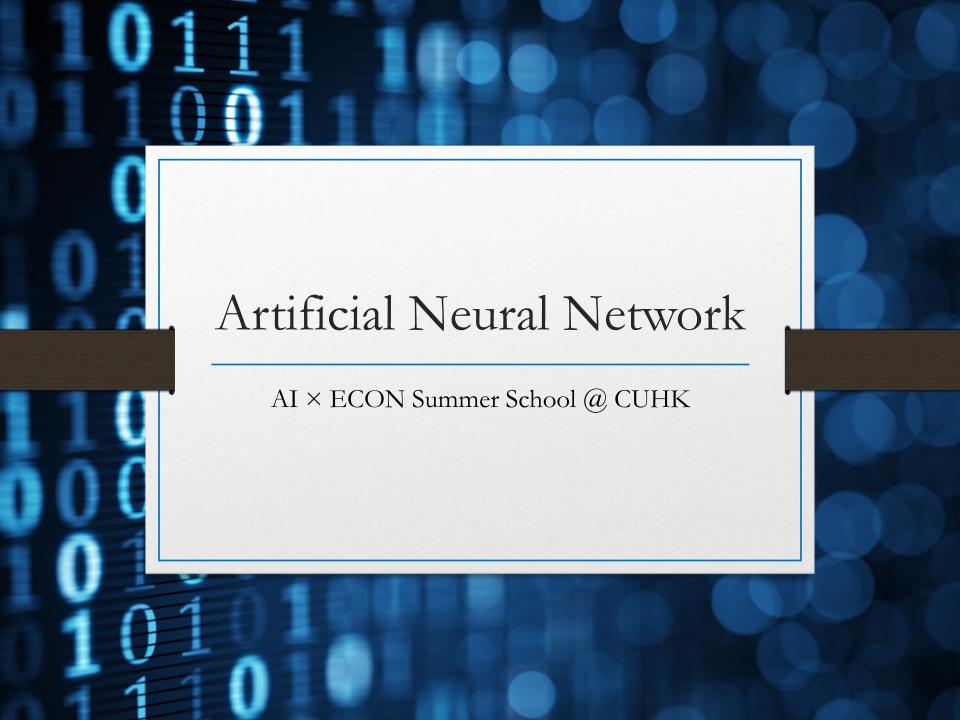
Follow instructions in account creation email.

#### Step 3: Connect to SCRP

#### Web access:

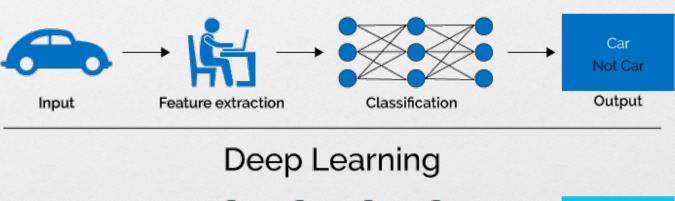
- <a href="https://scrp-login.econ.cuhk.edu.hk">https://scrp-login.econ.cuhk.edu.hk</a>
- <a href="https://scrp-login-2.econ.cuhk.edu.hk">https://scrp-login-2.econ.cuhk.edu.hk</a>

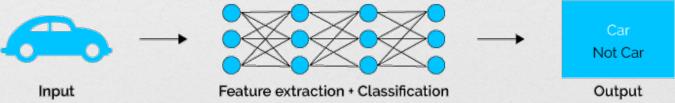




- The most prominent A.I. systems today are powered by artificial neural networks (ANN).
- ANN is the pinnacle of model complexity

#### Machine Learning





#### This tutorial:

- What is ANN?
- What are the common types of ANN?
- How to write an ANN in Python?

# This is a Linear Regression

Output 
$$\dot{y} = \alpha + \dot{x}\dot{\beta}$$

Where  $\vec{x} = [x_1 \dots x_k]$  is one sample of features (a.k.a. one observation of independent variables)

$$\vec{\beta} = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_k \end{bmatrix}$$
 is a vector of weights (a.k.a. coefficients)

# This is a Logistic Regression

$$P(y=1) = G(\alpha + \vec{x}\vec{\beta})$$
Data

Where 
$$G(z) = \frac{e^z}{1+e^z}$$

# This is a Neuron

Output 
$$h = G(\alpha + \vec{x}\vec{\beta})$$
Input

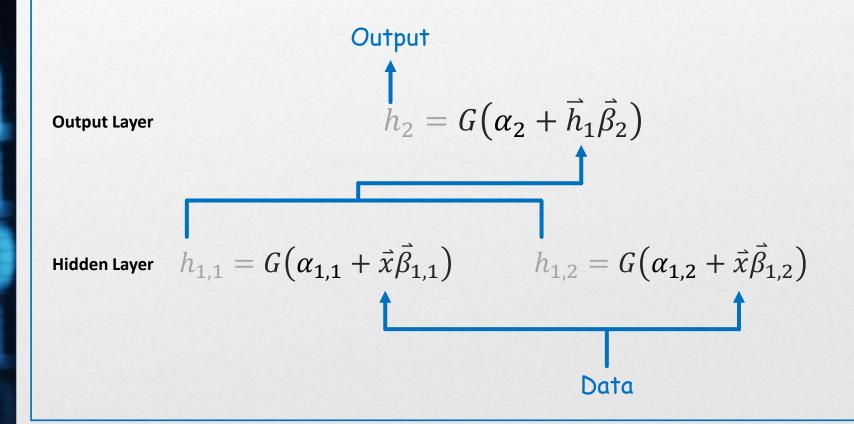
Where G(z) is a non-linear function G(z) is called the activation function and h the neuron's activation

(A Very Simple One)

**Output Layer** 

Output 
$$h_2 = G_2(\alpha_2 + h_1\beta_2)$$
 
$$h_1 = G_1(\alpha_1 + \vec{x}\vec{\beta}_1)$$
 Data

(With Two Hidden Neurons)

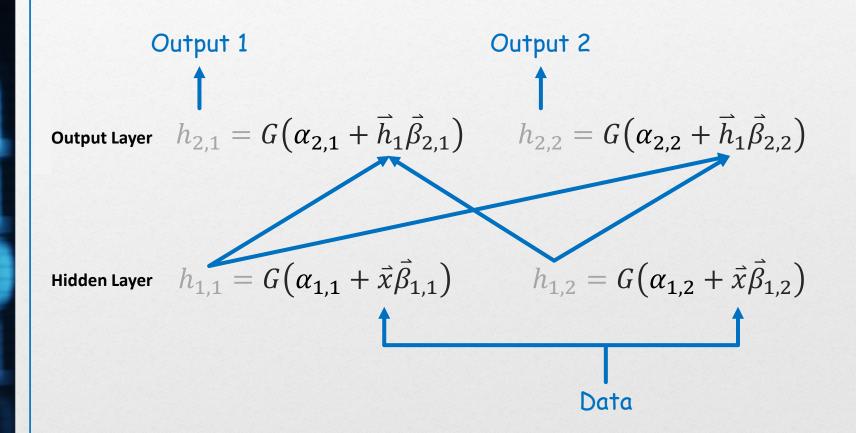


(With Multiple Hidden Neurons)

**Output Layer** 

Output 
$$h_2 = G_2 (\alpha_2 + \vec{h}_1 \vec{\beta}_2)$$
 
$$\vec{h}_1 = G_1 (\vec{\alpha}_1 + \vec{x} \boldsymbol{\beta}_1)$$
 Data

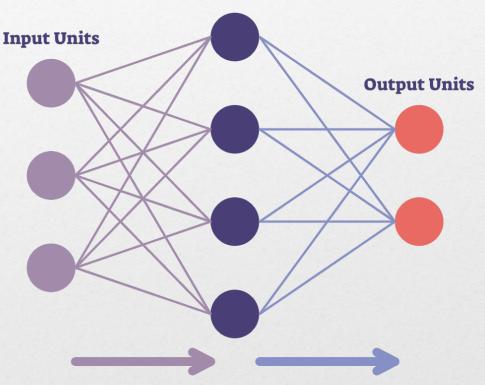
(With Two Hidden Neurons and Two Output Neurons)



# **AI Neural Networks**

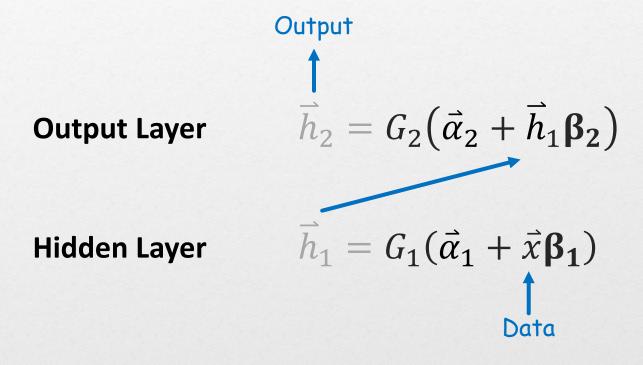


#### **Hidden Units**



**Flow of Activation** 

(With Multiple Hidden Neurons and Outputs)



#### Softmax

(a.k.a. Multinomial Logit)

### Classification Task

Output 1

$$h_{2,1} = F(\alpha_{2,1} + h_1 \beta_{2,1})$$

Output 2

Output Layer 
$$h_{2,1} = Fig(lpha_{2,1} + ec{h}_1ec{eta}_{2,1}ig) \qquad h_{2,2} = Fig(lpha_{2,2} + ec{h}_1ec{eta}_{2,2}ig)$$

$$h_{1,1} = F(\alpha_{1,1} + \vec{x}\bar{\beta}_{1,1})$$

Hidden Layer 
$$h_{1,1} = Fig(lpha_{1,1} + ec{x}ec{eta}_{1,1}ig) \qquad h_{1,2} = Fig(lpha_{1,2} + ec{x}ec{eta}_{1,2}ig)$$

Data

# Why Does Activation Have to be Non-Linear?

$$h_2 = G_2(\alpha_2 + \beta_2 h_1)$$
  
$$h_1 = G_1(\alpha_1 + \vec{\beta}_1 \vec{x})$$

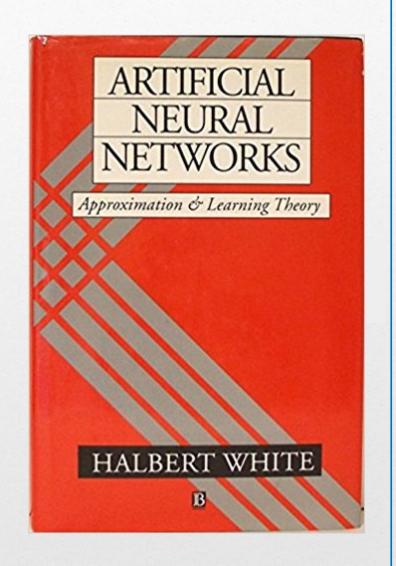
If  $G_1$  is linear, you get

$$h_2 = G_2 \left( \alpha_2 + \beta_2 (\alpha_1 + \vec{\beta}_1 \vec{x}) \right)$$
$$= G_2 \left( \alpha'_2 + \vec{\beta}'_2 \vec{x} \right)$$

So there is no point in having the hidden layer.

An econometrician's view of artificial neural network: a bunch of regressions stacked together

- Halbert White made significant contribution to the theoretical foundation of ANN in the 1980s
- Application was rare because the lack of computational power



# Why is ANN Powerful?

**Definition** A function  $\Psi: R \to [0,1]$  is a squashing function if it is non-decreasing,  $\lim_{\lambda \to \infty} \Psi(\lambda) = 1$  and  $\lim_{\lambda \to -\infty} \Psi(\lambda) = 0$ .

#### Universal approximation theorem

(Hornik, Stinchcombe and White 1989 Theorem 2.4)

Let 
$$\Sigma^r(G) = \{f: \mathbb{R}^r \to \mathbb{R}: f(x) = \sum_{i=1}^q \beta_i G(a_i + b_i x), a_i, b_i, \beta_i \in \mathbb{R}^r, q = 1, 2, \dots\}$$

For every squashing function  $\Psi$ , every continuous function  $g: \mathbb{R}^r \to \mathbb{R}$ , every compact subset  $K \subset \mathbb{R}^r$  and every  $\varepsilon > 0$ , there exists an  $f \in \Sigma^r(\Psi)$  such that

$$\sup_{x \in K} |f(x) - g(x)| < \varepsilon$$

# Why is ANN Powerful?

The universal approximation theorem says ANN can approximate any continuous function arbitrarily well.

The theorem can be extended to multiple outputs and non-continuous functions under additional assumptions.

#### Universal approximation theorem

(Hornik, Stinchcombe and White 1989 Theorem 2.4)

Let 
$$\Sigma^{r}(G) = \{f : \mathbb{R}^{r} \to \mathbb{R} : f(x) = \sum_{i=1}^{q} \beta_{i} G(a_{i} + b_{i}x), a_{i}, b_{i}, \beta_{i} \in \mathbb{R}^{r}, q = 1, 2, ...\}$$

For every squashing function  $\Psi$ , every continuous function  $g: \mathbb{R}^r \to \mathbb{R}$ , every compact subset  $K \subset \mathbb{R}^r$  and every  $\varepsilon > 0$ , there exists an  $f \in \Sigma^r(\Psi)$  such that

$$\sup_{x \in K} |f(x) - g(x)| < \varepsilon$$

# Why is ANN Powerful?

Note that this is an existence theorem—it does not say how many neurons q is needed, which squashing/activation function  $\Psi$  works best, or how to get the biases and weights.

#### Universal approximation theorem

(Hornik, Stinchcombe and White 1989 Theorem 2.4)

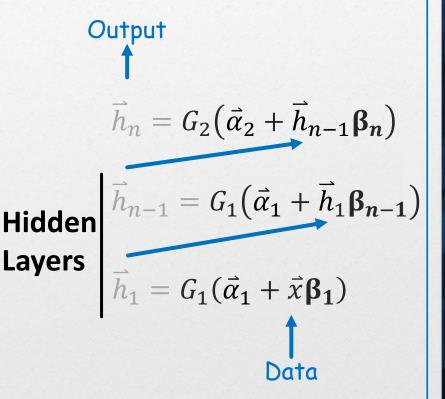
Let 
$$\Sigma^r(G) = \{f: \mathbb{R}^r \to \mathbb{R}: f(x) = \sum_{i=1}^q \beta_i G(a_i + b_i x), a_i, b_i, \beta_i \in \mathbb{R}^r, q = 1, 2, \dots\}$$

For every squashing function  $\Psi$ , every continuous function  $g: \mathbb{R}^r \to \mathbb{R}$ , every compact subset  $K \subset \mathbb{R}^r$  and every  $\varepsilon > 0$ , there exists an  $f \in \Sigma^r(\Psi)$  such that

$$\sup_{x \in K} |f(x) - g(x)| < \varepsilon$$

# Deep Learning

- Deep Learning refers to the stacking of multiple hidden layers
- Typical layer count is in the single digit but can go as high as a hundred.



# How Does an ANN Learn?

Gradient Descent

- **Gradient descent** is a very general optimization algorithm, usable with pretty much all models.
- It uses the <u>first derivative</u> of the loss function with respect to a parameter to adjust the parameter.

- Gradient descent uses the <u>first derivative</u> of the loss function with respect to a parameter to adjust the parameter.
- E.g. suppose our model is

$$\hat{y} = \alpha + x$$

This is a (simplified) regression task, which we can solve by minimizing the squared error:

$$c = (y - \hat{y})^2$$

• For illustration purpose, let us assume the data is generated by y = 5 + x, so that if x = 1, y = 6.

Error 
$$\epsilon = y - \hat{y}$$
  
Loss  $= \epsilon^2$ 

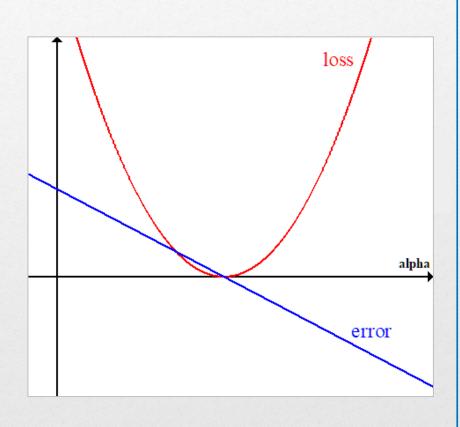
First derivative of loss function:

$$\frac{d}{d\alpha} \epsilon^2$$

$$= 2\epsilon \frac{d}{d\alpha} [y - (\alpha + x)]$$

$$= -2\epsilon$$

This is the marginal effect, or gradient, of  $\alpha$  on loss.



We have an initial guess of what  $\alpha$  is (usually just a random number.)

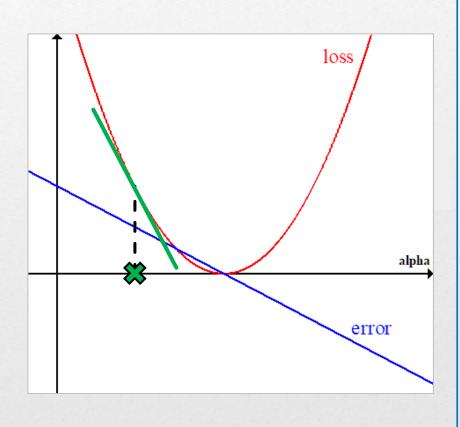
This gives us an initial prediction  $\hat{y}$  and corresponding error and loss.

$$\hat{\alpha}_0 = 2$$

$$\hat{y}_0 = 2 + 1 = 3$$

$$\epsilon_0 = y - \hat{y} = 6 - 3 = 3$$

$$\frac{d}{d\alpha}\epsilon^2 = -2\epsilon = -6$$

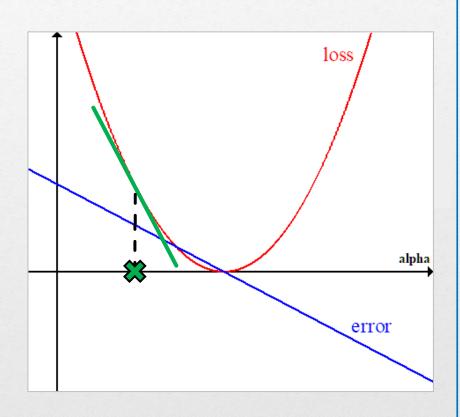


$$\epsilon_0 = 3$$

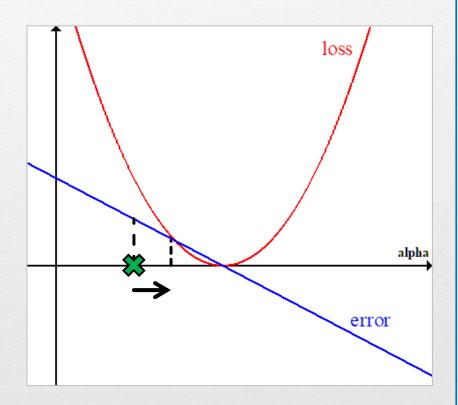
$$\frac{d}{d\alpha}\epsilon^2 = -6 > 0$$

In our example,  $\epsilon$  is positive, so the gradient is negative.

This means loss is decreasing in alpha.

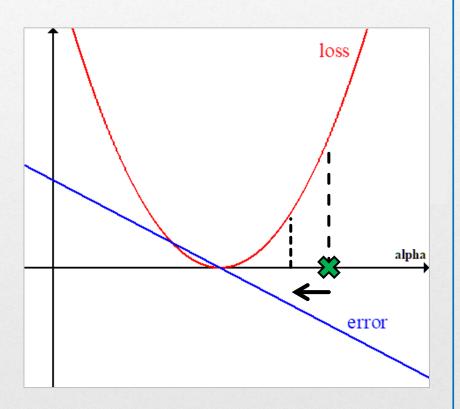


This makes sense, because if  $y > \hat{y} = \alpha + x$ increasing  $\alpha$  will bring  $\hat{y}$  closer to y.



Conversely, if  $\epsilon < 0$ , then

$$y < \hat{y} = \alpha + x$$
  
decreasing  $\alpha$  will bring  $\hat{y}$  closer to  $y$ .



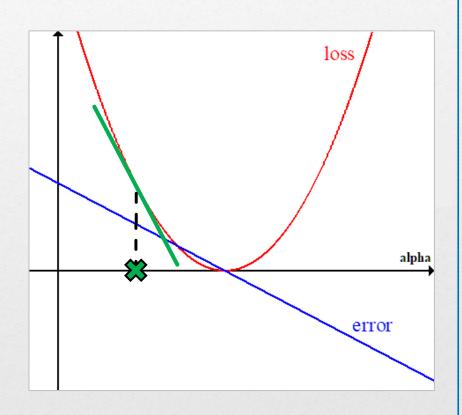
The gradient tells us the direction we need to adjustment our parameter.

The amount we need to adjust is, to a first-order approximation, inversely proportional to the gradient.

E.g.

$$\frac{d}{d\alpha}\epsilon^2 = -6$$

 $\alpha$  needs to be bigger.



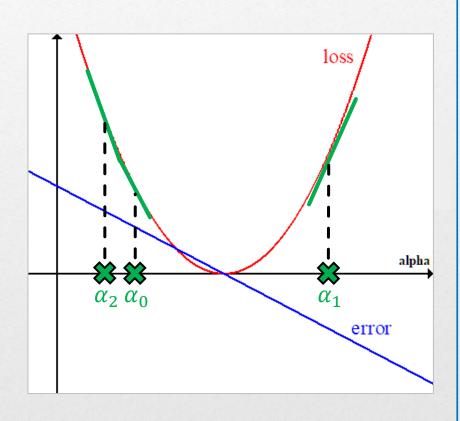
If we adjust the parameter by the exact amount of the gradient, we might overshoot:

$$\hat{\alpha}_1 = \hat{\alpha}_0 - \frac{d}{d\alpha} \epsilon^2$$

$$= 2 - (-6)$$

$$= 8 > 5 = \text{true } \alpha$$

Overshooting could result in our parameters bouncing around the true value, never to converge.



To prevent overshooting, we moderate the gradient by a learning rate  $\gamma$ :

$$\hat{\alpha}_1 = \hat{\alpha}_0 - \gamma \frac{d}{d\alpha} \epsilon^2$$

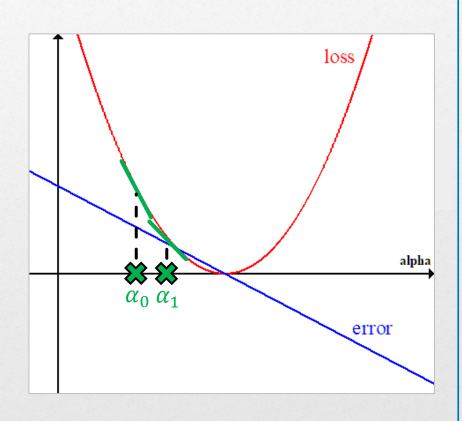
E.g. if  $\gamma = 0.1$ :

$$\hat{\alpha}_1 = \hat{\alpha}_0 - \gamma \frac{d}{d\alpha} \epsilon^2$$

$$= 2 - 0.1(-6)$$

$$= 2.6 < 5 = \text{true } \alpha$$

Now we are not overshooting, but it will take us more iterations to get to the true  $\alpha$  value.



- The gradient tells us the direction we need to adjustment our parameter. In the above example, the amount we need to adjust is, to a first-order approximation, proportional to  $-dc/d\alpha$ .
- Moderating the adjustment by a learning rate  $\gamma$ , we have the following update rule:

$$\alpha_t = \alpha_{t-1} - \gamma \frac{dc}{d\alpha} \bigg|_{\alpha_{t-1}}$$

Or more typical in computer science:

$$\alpha \leftarrow \alpha - \gamma \frac{dc}{d\alpha}$$

• With more than one sample, we take the average.

# Stochastic Gradient Descent

- Learning is quite slow if we only update model weights after we go through all data.
- We could instead update every time after we gone through a given number of samples. This is **Stochastic Gradient Descent (SGD).**
- Because we are not using all data, we could be updating towards the wrong direction sometimes, but on average the updates will be correct. Hence, *stochastic*.
- The stochastic nature is actually a good property because it helps the model escape from local optima.

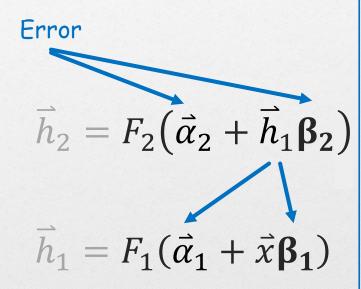
# How Does an ANN Learn?

- Gradient Descent
- Back Propagation

• Because neural network have multiple layers, the gradient of lower layers needs be computed using the chain rule. This process is called back propagation.

Output 
$$\vec{h}_2 = F_2(\vec{lpha}_2 + \vec{h}_1 oldsymbol{eta}_2)$$
  $\vec{h}_1 = F_1(\vec{lpha}_1 + \vec{x} oldsymbol{eta}_1)$  Data

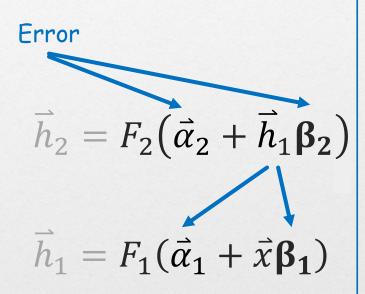
• Because neural network have multiple layers, the gradient of lower layers needs be computed using the chain rule. This process is called **back propagation**.



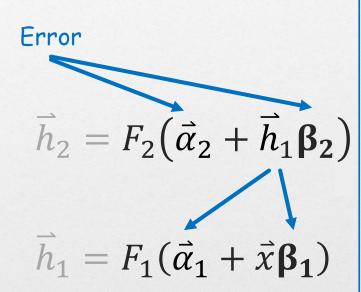
E.g. single target regression task, the gradient of the second weight of the first neuron is:

$$\begin{split} &\frac{\partial}{\partial \beta_{1}^{(1,2)}} (y - h_{2})^{2} \\ &= -2(y - h_{2}) F_{2}' (\vec{\alpha}_{2} + \vec{h}_{1} \vec{\beta}_{2}) \frac{\partial}{\partial \beta_{1}^{(1,2)}} (\vec{\alpha}_{2} + \vec{h}_{1} \vec{\beta}_{2}) \\ &= -2(y - h_{2}) F_{2}' (\vec{\alpha}_{2} + \vec{h}_{1} \vec{\beta}_{2}) \beta_{2}^{(1)} \\ &\cdot \frac{\partial}{\partial \beta_{1}^{(1,2)}} F_{1} \left( \alpha_{1} + \vec{x} \vec{\beta}_{1}^{(1)} \right) \\ &= -2(y - h_{2}) F_{2}' (\vec{\alpha}_{2} + \vec{h}_{1} \vec{\beta}_{2}) \beta_{2}^{(1)} \\ &\cdot F_{1}' \left( \vec{\alpha}_{1} + \vec{x} \vec{\beta}_{1}^{(1)} \right) x_{2} \end{split}$$

An important feature of any neural network library is to automate the computation of gradient.



- As the number of layers go up, the process becomes quite brittle.
- Remember the chain rule involves a lot of multiplication.
- Small times small means very small, resulting in vanishing gradient.
- Big times big means very big, resulting in exploding gradient.
- Neither is good for learning.
- A lot of research goes into finding ways to combat the issue: different activation functions, different randomization distributions...



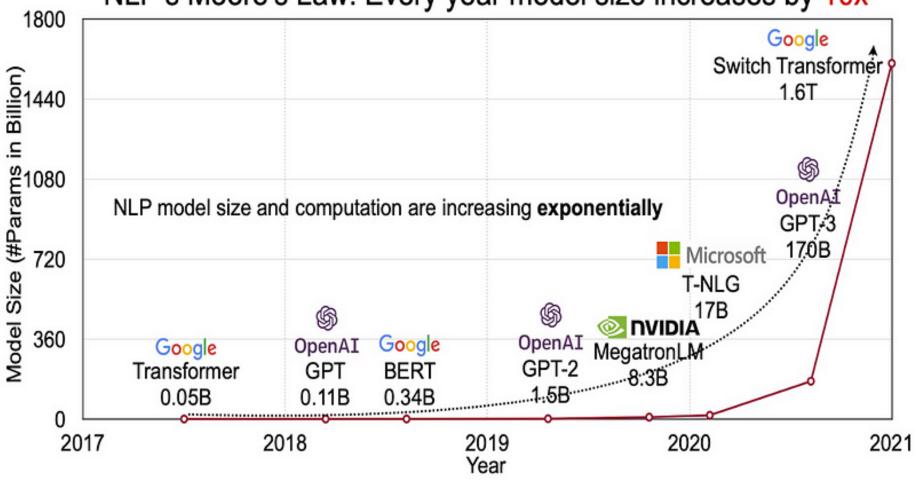
Hands-on Demo

## Computation

- The idea of artificial neural network can be traced back to as far back as the 1940s.
- Due to the large number of parameters and large data size involved, effective use of ANN was prohibitive until 2010s.

# 11 Model Size

NLP's Moore's Law: Every year model size increases by 10x



# It Takes a Lot to Train These Models!

Model size	Hidden size	Number of layers	Number of parameters (billion)	Model-parallel size	Number of GPUs	Batch size	Achieved teraFIOPs per GPU	Percentage of theoretical peak FLOPs	Achieved aggregate petaFLOPs
1.7B	2304	24	1.7	1	32	512	137	44%	4.4
3.6B	3072	30	3.6	2	64	512	138	44%	8.8
7.5B	4096	36	7.5	4	128	512	142	46%	18.2
18B	6144	40	18.4	8	256	1024	135	43%	34.6
39B	8192	48	39.1	16	512	1536	138	44%	70.8
76B	10240	60	76.1	32	1024	1792	140	45%	143.8
145B	12288	80	145.6	64	1536	2304	148	47%	227.1
310B	16384	96	310.1	128	1920	2160	155	50%	297.4
530B	20480	105	529.6	280	2520	2520	163	52%	410.2
1T	25600	128	1008.0	512	3072	3072	163	52%	502.0

Source: Nvidia



HOME COMPUTE STORE CONNECT CONTROL CODE AI HPC ENTERPRISE HYPERSCALE CLOUD

EDGE

LATEST > Ampere Arm Server CPUs To Get 512 Cores, AI Accelerator > COMPUTE

Search ...

**HOME** > **AI** > So Who Is Building That 100,000 GPU Cluster For xAI?

### SO WHO IS BUILDING THAT 100,000 GPU CLUSTER FOR XAI?

July 30, 2024 Timothy Prickett Morgan



ANN took off due to massive increase in computational capabilities, particularly in the use of **graphic processing unit** (GPU) for computation.



# NVIDIA® DGX STATION™ YOUR PERSONAL AI SUPERCOMPUTER

### Announcing NVIDIA H200 Tensor Core GPU

Supercharging the highest-performing generative AI and HPC platforms



141GE

Memory Bandwidth
4.8 TB/S

GPT-3 175B Inference

1.6X
Performance vs H100

Llama2 70B Inference

1.9X
Performance vs H100

**HPC Simulation** 

2.0X

HGX H200



1.1 TB HBM3e | 32 PF FP8 fore Memory Capacity | 1.4X More HBM Bandwidth

### Announcing GB200 NVL72 Delivers New Unit of Compute



36 GRACE CPUs GB200 NVL72 -

72 BLACKWELL GPUs

Fully Connected NVLink Switch Rack

720 PFLOPs Training Inference 1,440 PFLOPs **NVL Model Size** 27T params Multi-Node Bandwidth 130 TB/s Multi-Node All-Reduce 260 TB/s



**NVIDIA** Corp @ NVIDIA Financials Overview Compare NASDAQ: NVDA : Market Summary > NVIDIA Corp 134.91 USD + Follow +130.72 (3,119.81%) ↑ past 5 years Closed: 10 Jul, 7:59 pm GMT-4 • Disclaimer After hours 135.60 +0.69 (0.51%) 5D 1M 6M YTD 1Y <u>5Y</u> 1D 150 100 50 2021 2023 2022 2024

