

# History of Neural Networks

Deep learning

→ Perception

→ 1957

[Rosenblatt]

logistic regression

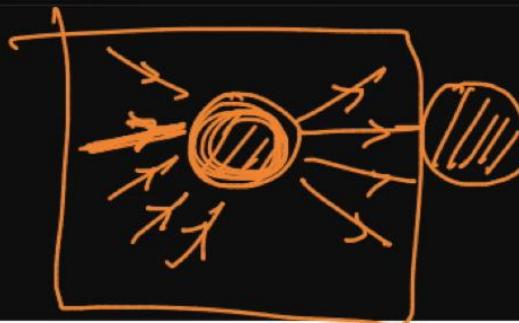
Russian - English

→ Alan Turing

[Computer]

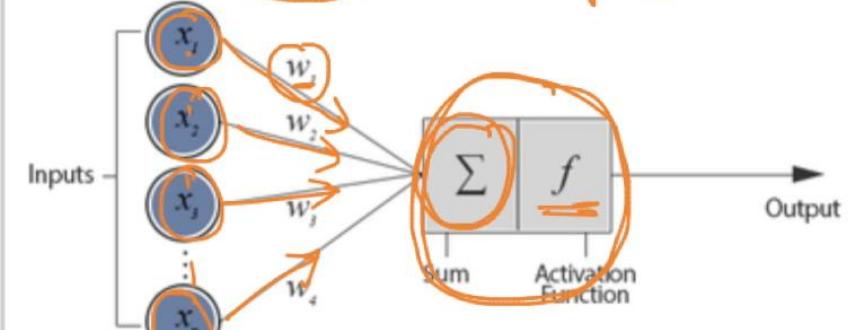
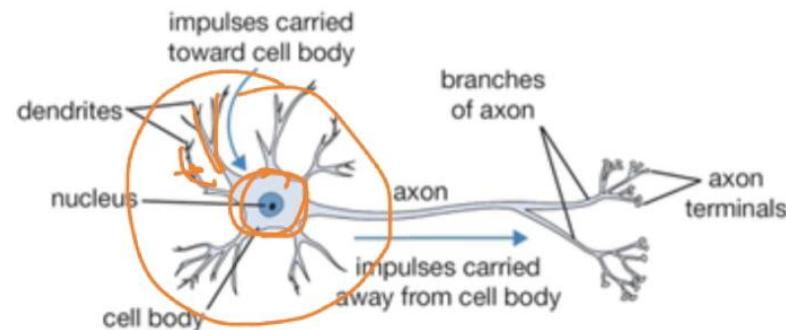
→ "loosely" inspired from biology



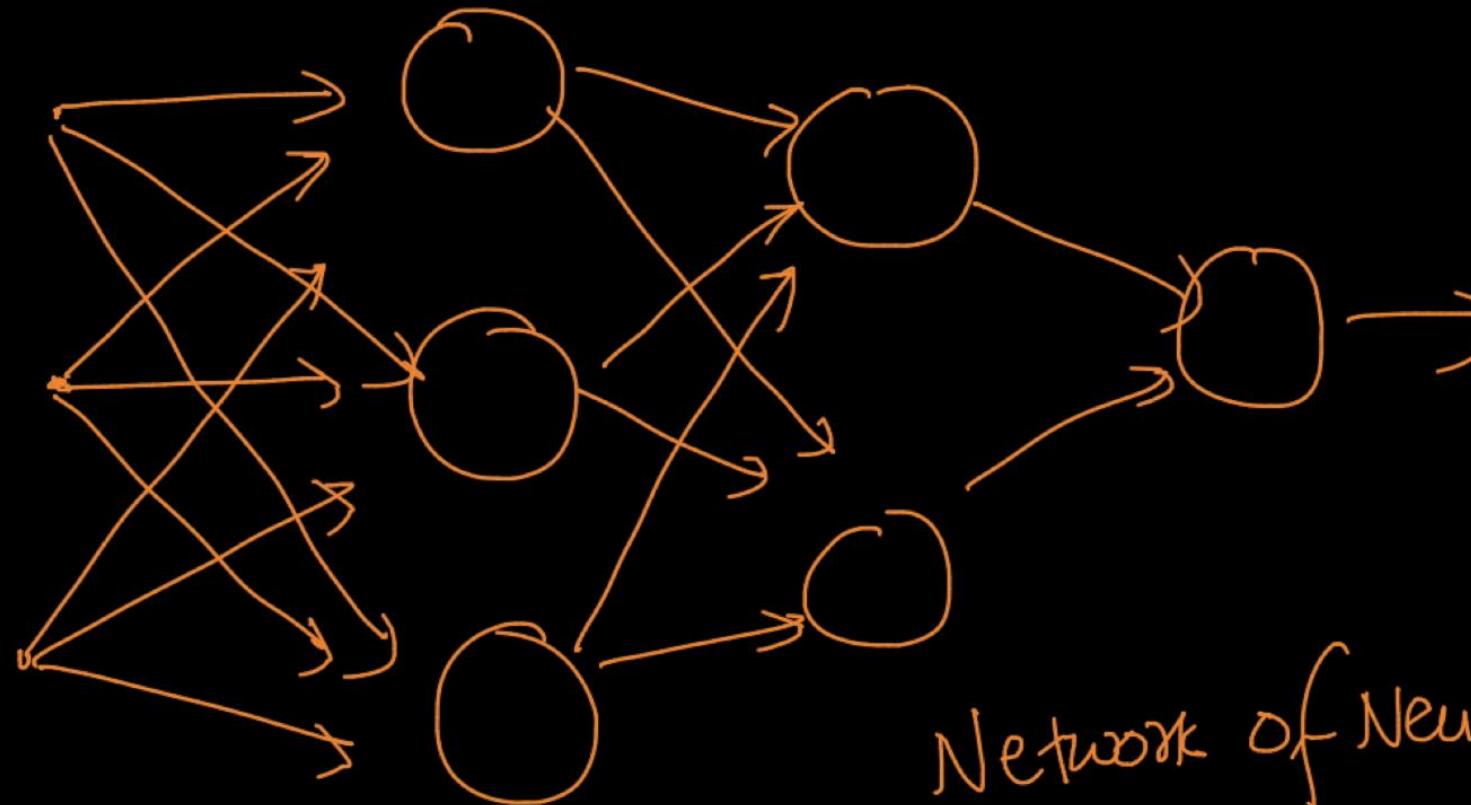


most simple

## Biological Neuron versus Artificial Neural Network $f(w_1x_1 + w_2x_2 + \dots + w_nx_n)$



Brain



Network of Neurons



1960s →

Artificial NN

1986

Hinton & others

late  
1990's

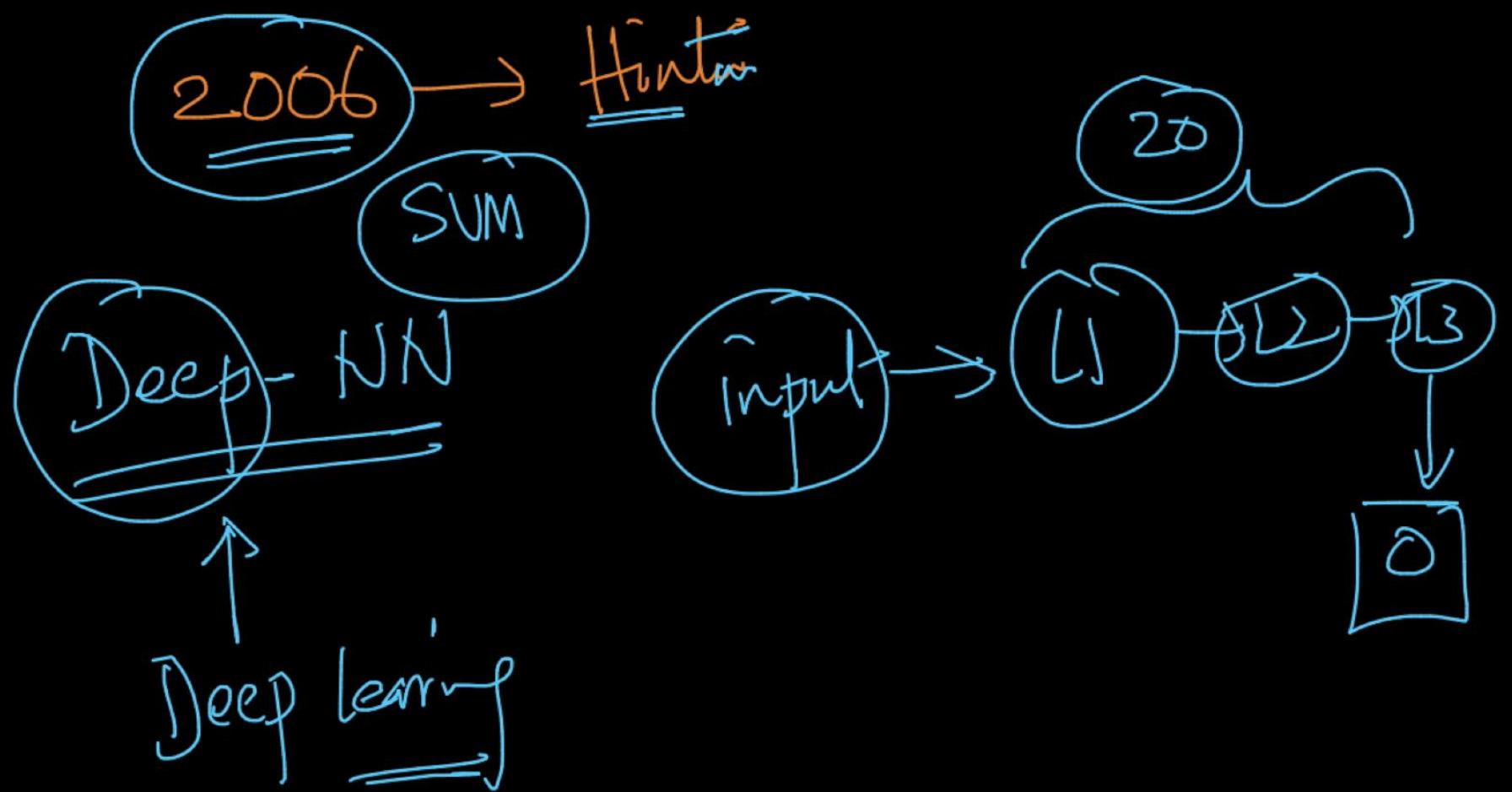
Back propagation algo

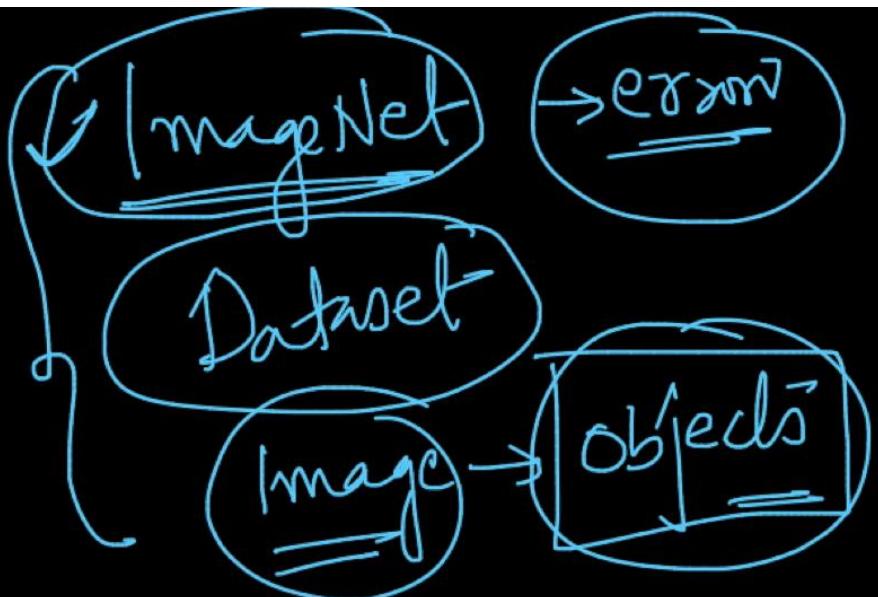
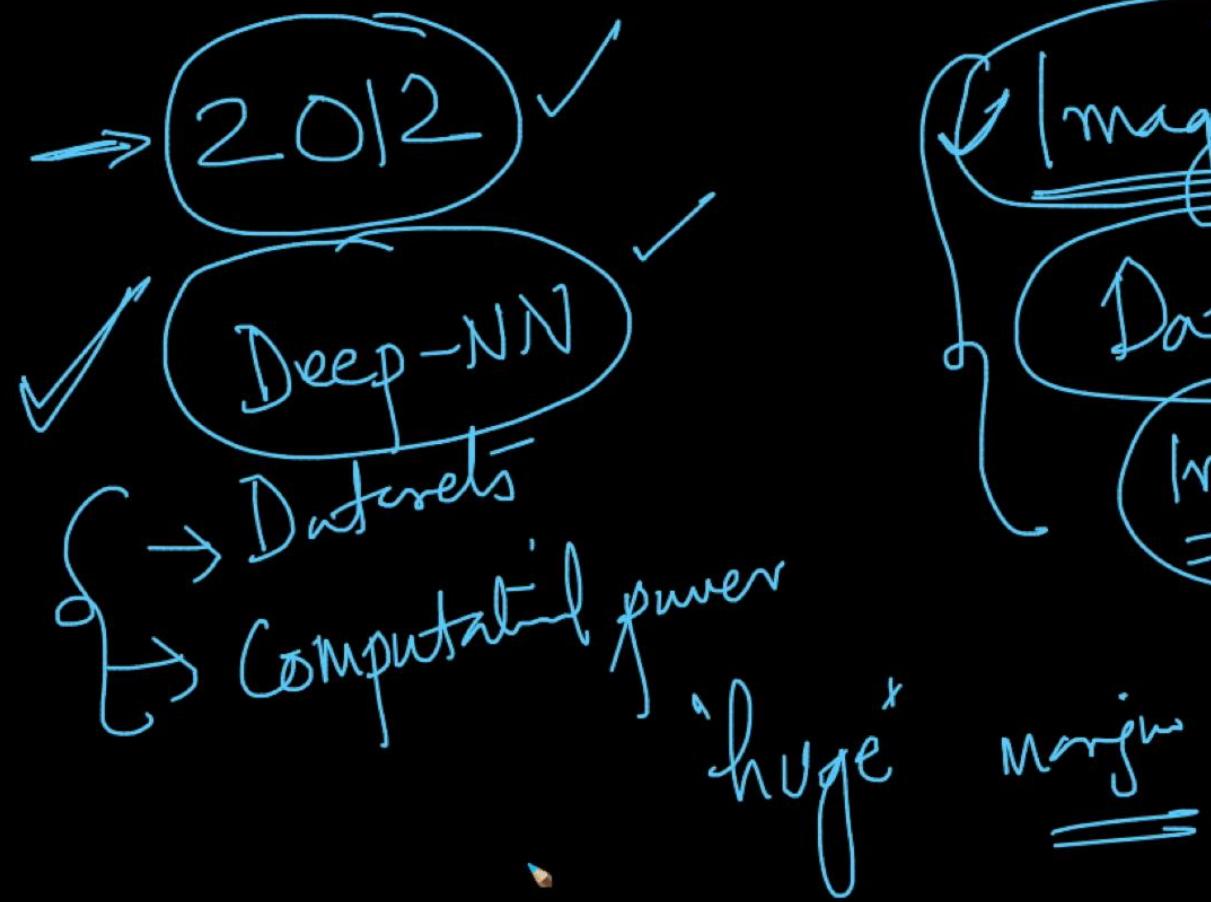
Chain rule  
of diff

AI Hype

loops

{ computing power }  
Data





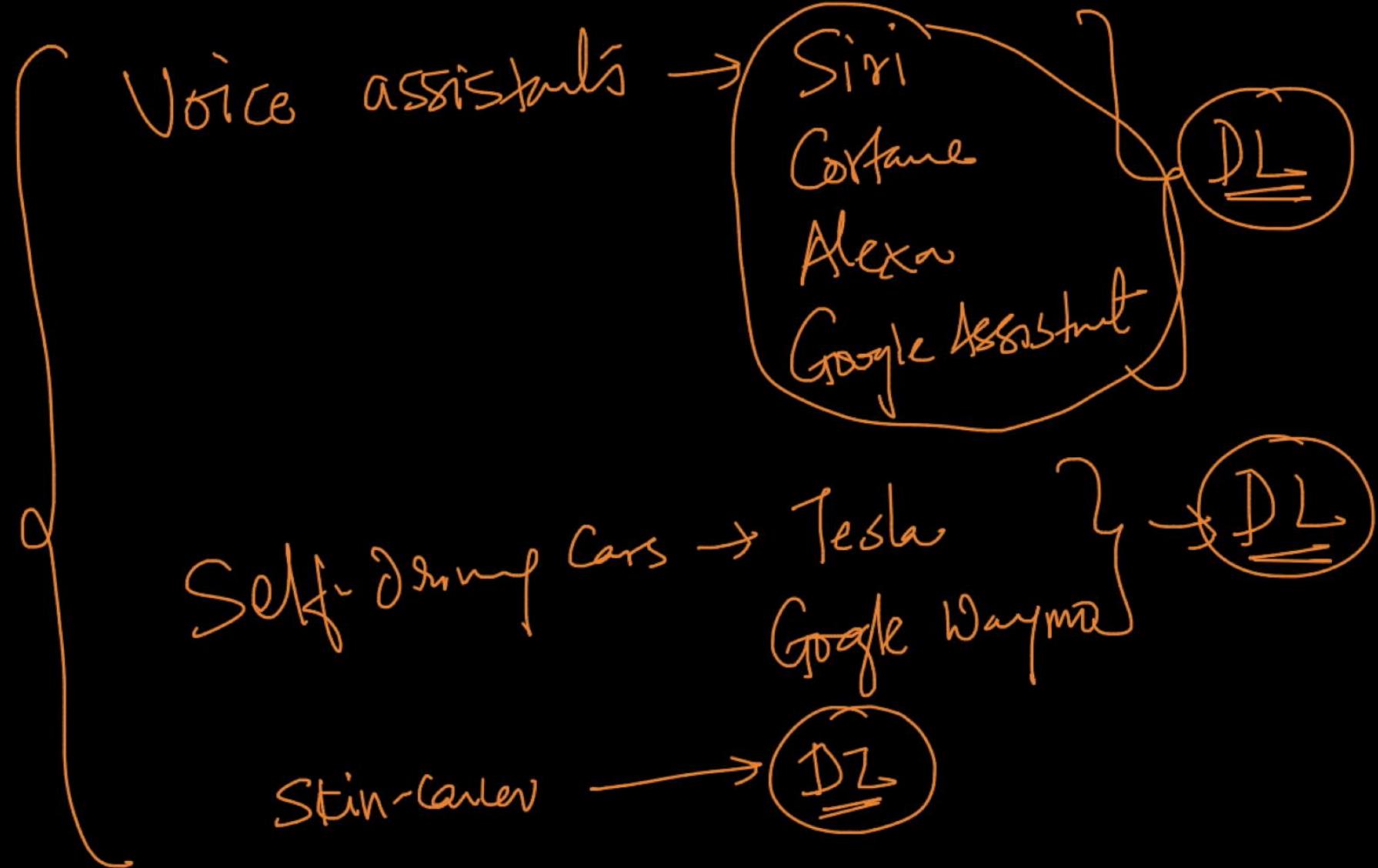
Google, MS, FB, Amazon, Baidu

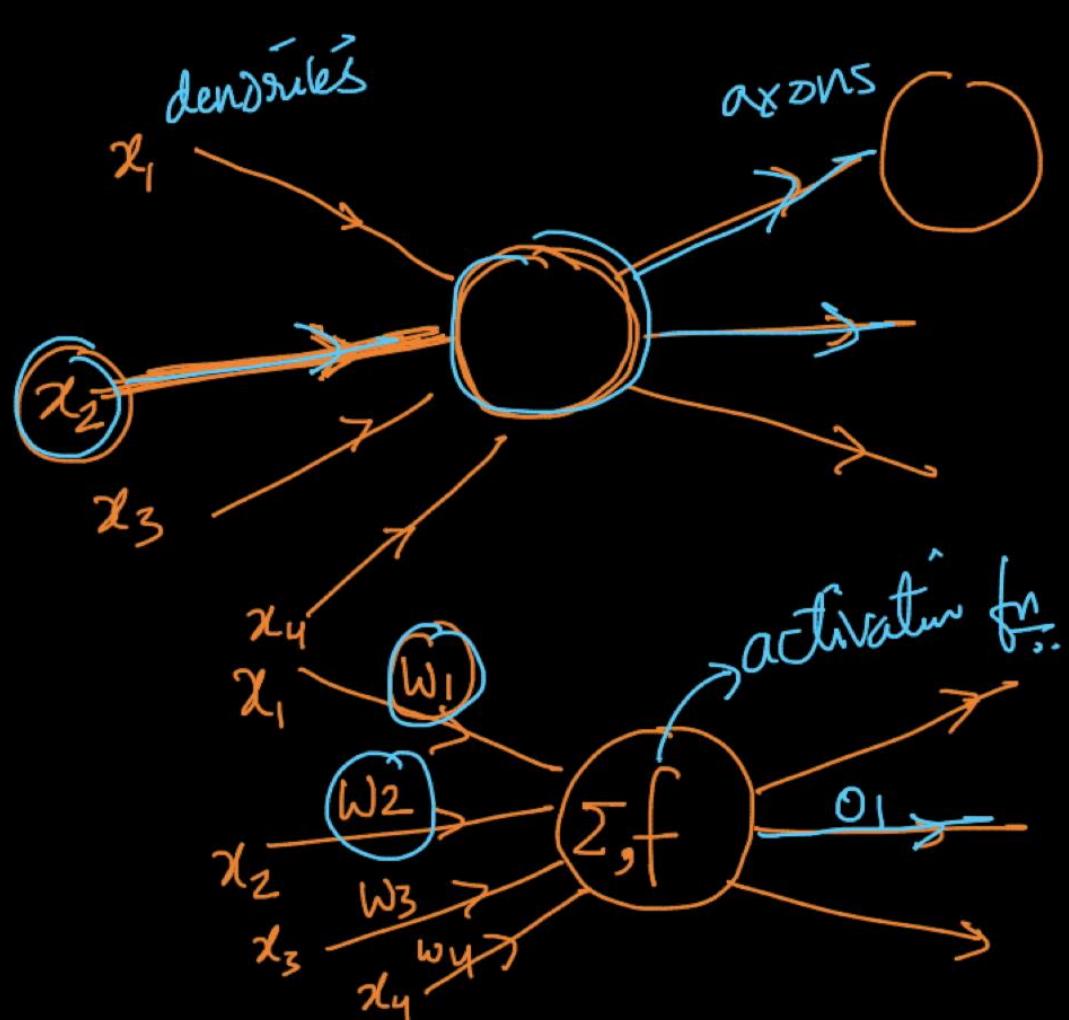
{ → Lots of Data  
→ lots of Computational resources

2012

→ Golden Days of AI





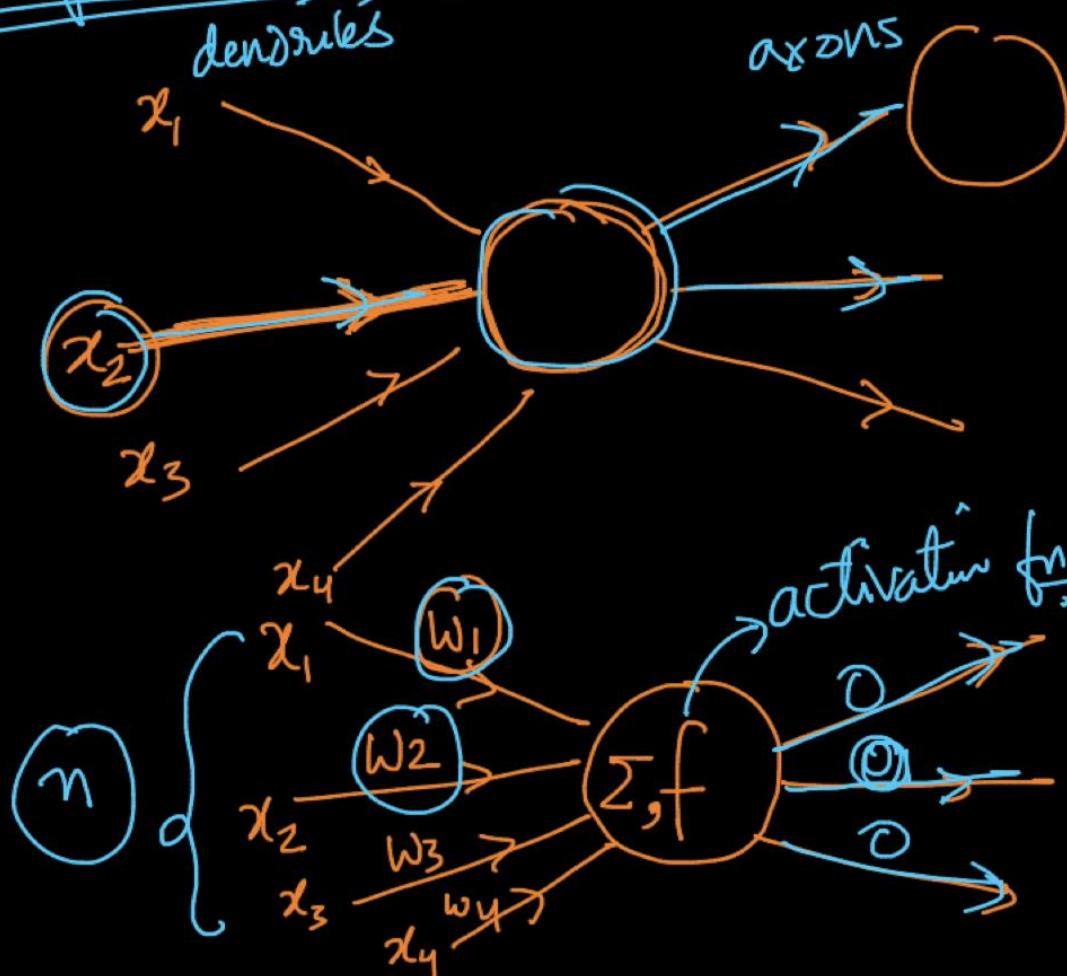


{ simplified view  
of a Neuron

{ activated  
fired

$$o_1 = f \left( w_1 x_1 + w_2 x_2 + \dots + w_n x_n \right)$$

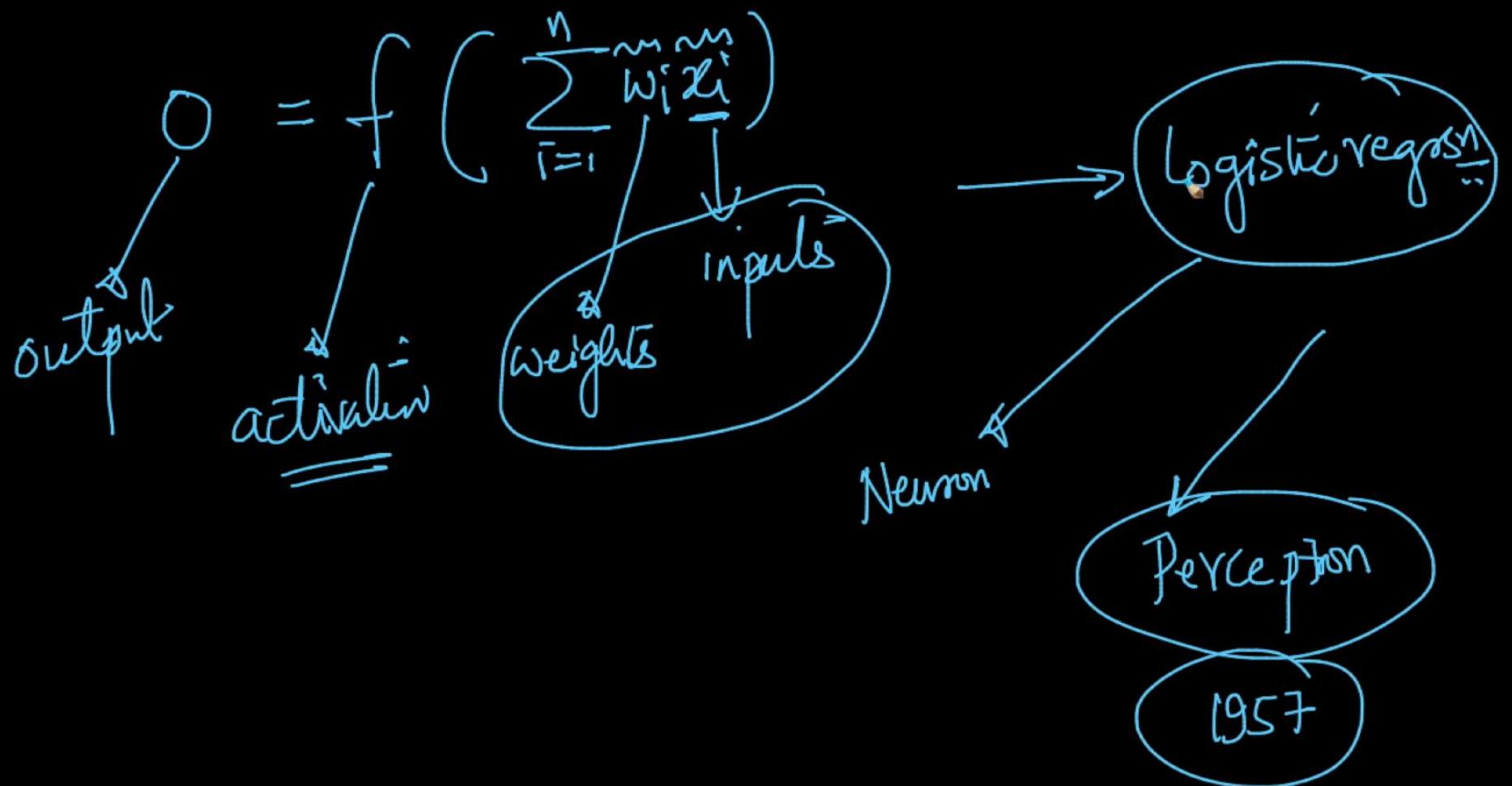
outgoing vertices / edges / lines → axons



{ simplified view  
of a Neuron

{ activated  
fired

$$o = f(w_1x_1 + w_2x_2 + \dots + w_nx_n)$$



Secure | https://developingchild.harvard.edu/resources/five-numbers-to-remember-about-early-childhood-development/

**language →**

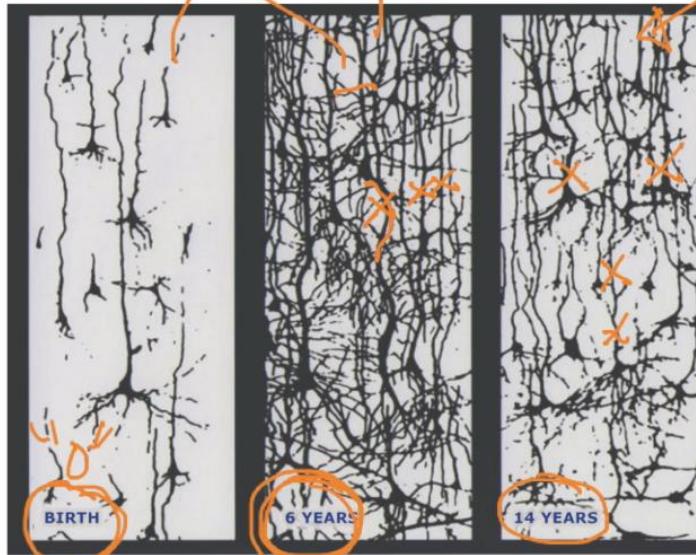
**Speech objects**

**more connections**

↓

**Calories**

**1 More Than 1 Million New Neural Connections Per Second\***



BIRTH      6 YEARS      14 YEARS

Image source: Coneil, JL. The postnatal development of the human cerebral cortex. Cambridge, Mass: Harvard University Press, 1959.

million new neural connections are formed every second.\* Neural connections are formed through the interaction of genes and a baby's environment and experiences, especially "serve and return" interaction with adults, or what developmental researchers call contingent reciprocity. These are the connections that build brain architecture – the

**Teen-age**

1 **More Than 1 Million New Neural Connections Per Second\***

2 **18 Months: Age At Which Disparities in Vocabulary Begin to Appear**

3 **90 - 100% Chance of Developmental Delays When Children Experience 6 - 7 Risk Factors**

4 **3:1 Odds of Adult Heart Disease After 7 - 8 Adverse Childhood Experiences**

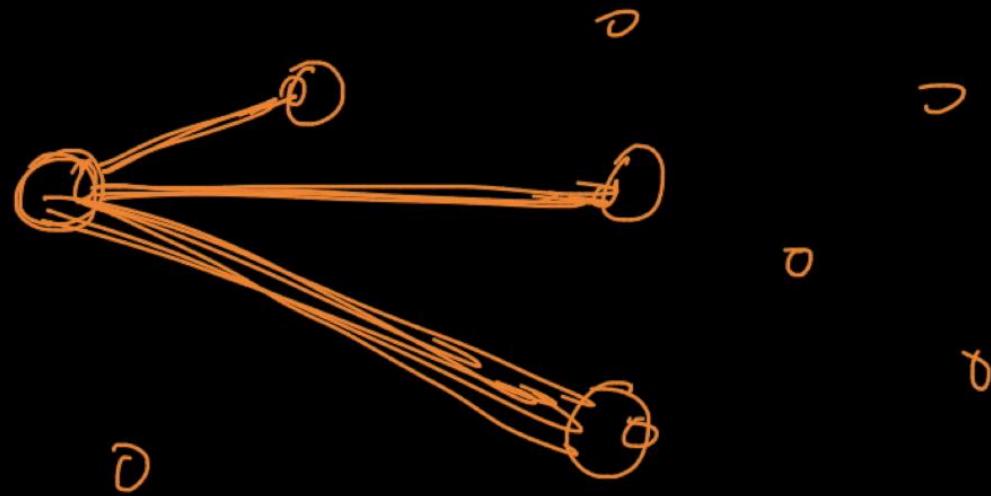
5 **\$4 - \$9 in Returns For Every Dollar Invested in Early Childhood Programs**

Join our mailing list to receive news and other updates from the Center.

Your email address

**Subscribe**

all biological brain  $\Rightarrow$  Weights on neural connections



LR:

$x_i \rightarrow \hat{y}_i \rightarrow$  predicted value of  $y_i$

$$\hat{y}_i = \text{Sigmoid}((w^T x_i) + b)$$

$$\mathcal{D} = \{x_i, y_i\}$$

Train LR

$x_i \in \mathbb{R}^d$

$w, b$   
 $w \in \mathbb{R}^d$   
 $b \in \mathbb{R}$

$$x_i = [x_{i1}, x_{i2}, x_{i3}, x_{i4}, \dots, x_{id}]$$



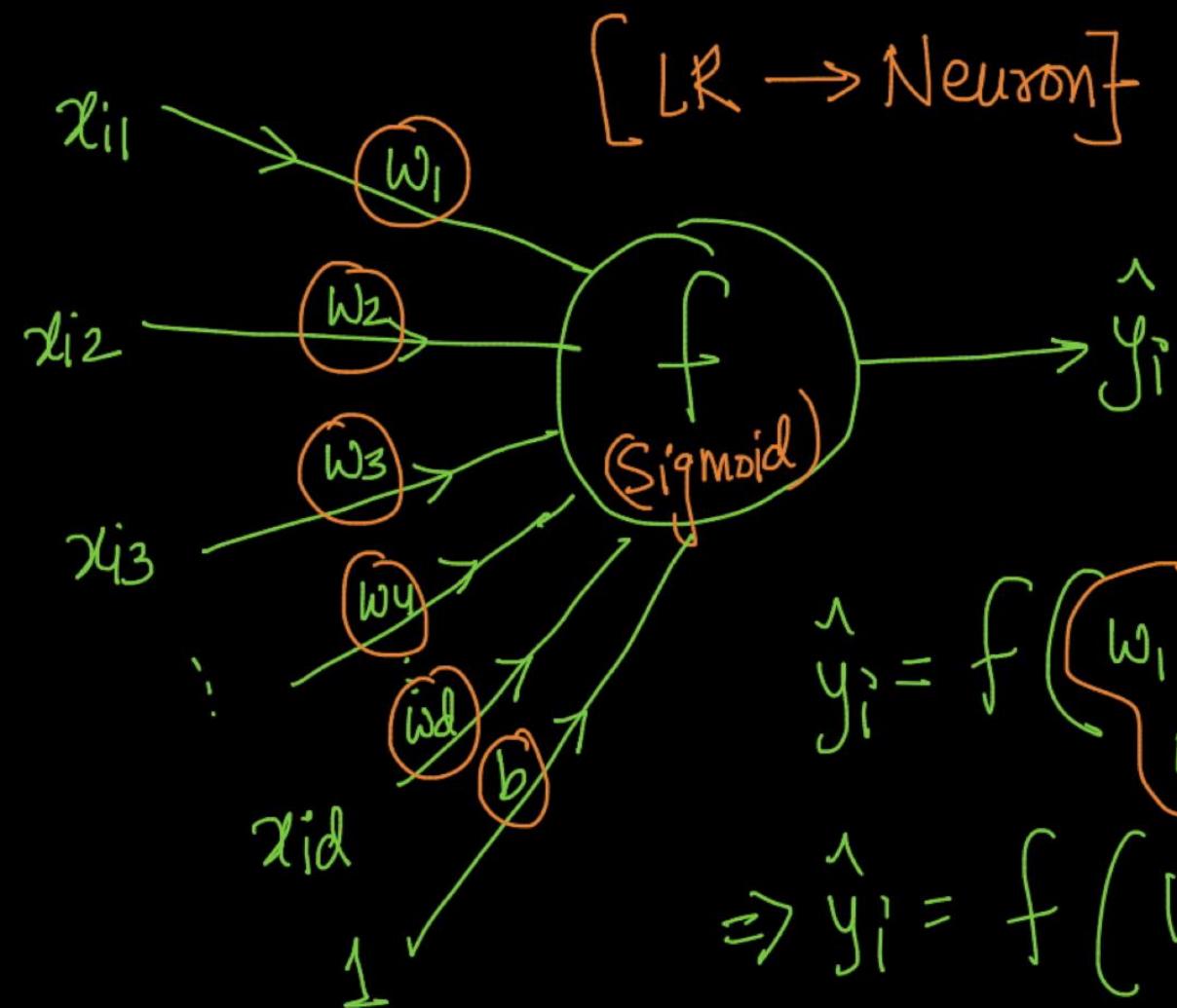
$\hat{y}_i = \text{Sigmoid} \left( \sum_{j=1}^d w_j x_{ij} + b \right)$

$x_i = [x_{i1}, x_{i2}, \dots, x_{id}]$

$w = [w_1, w_2, \dots, w_d]$

$o = f \left( \sum_{j=1}^d w_j x_{ij} \right)$





$$\mathcal{D} = \{(x_i, y_i)\}$$

Task :-  $w_i, b$        $i = 1 \rightarrow d$

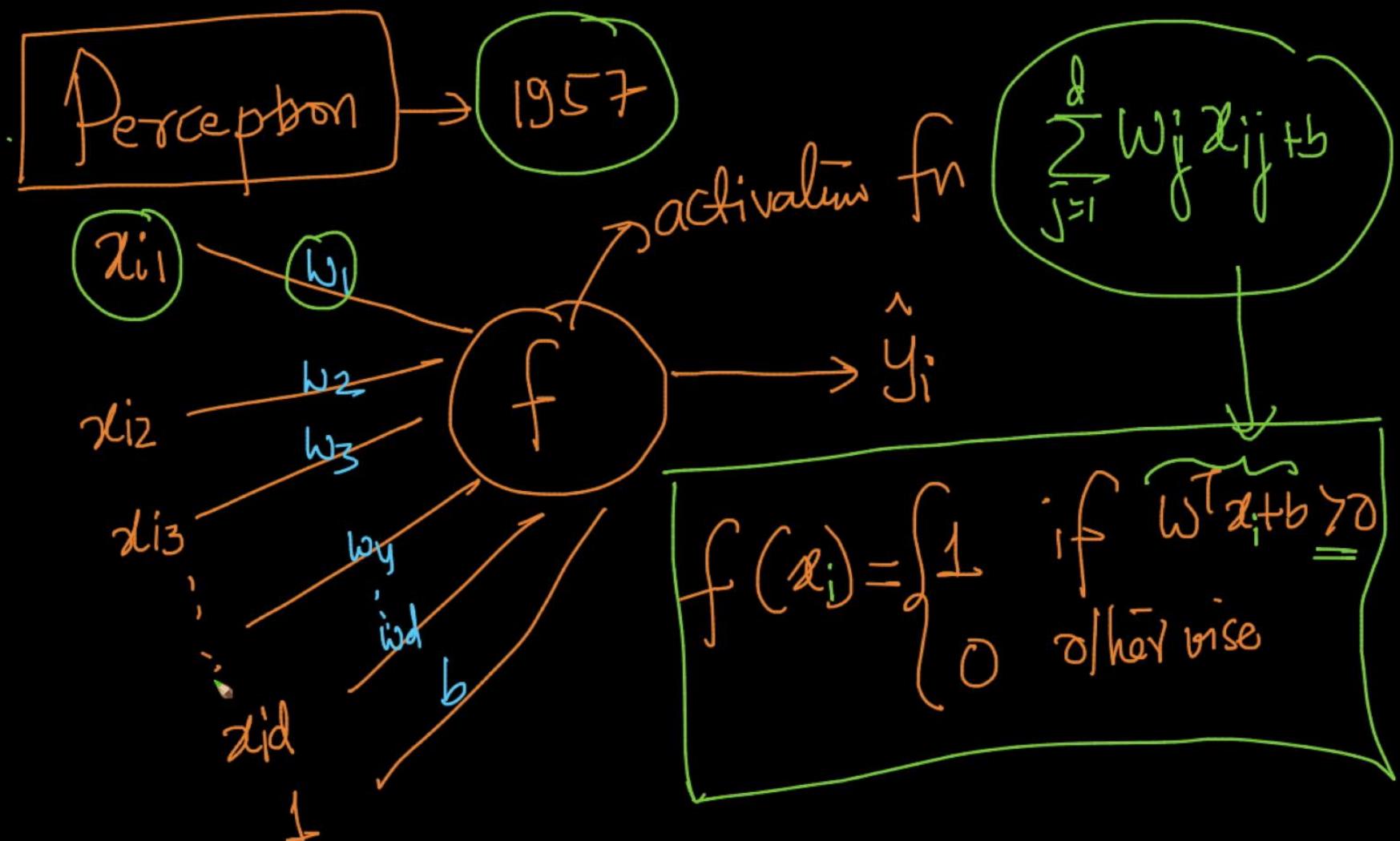
SGD  $\rightarrow$  lr. regression & optimiza<sup>n</sup>

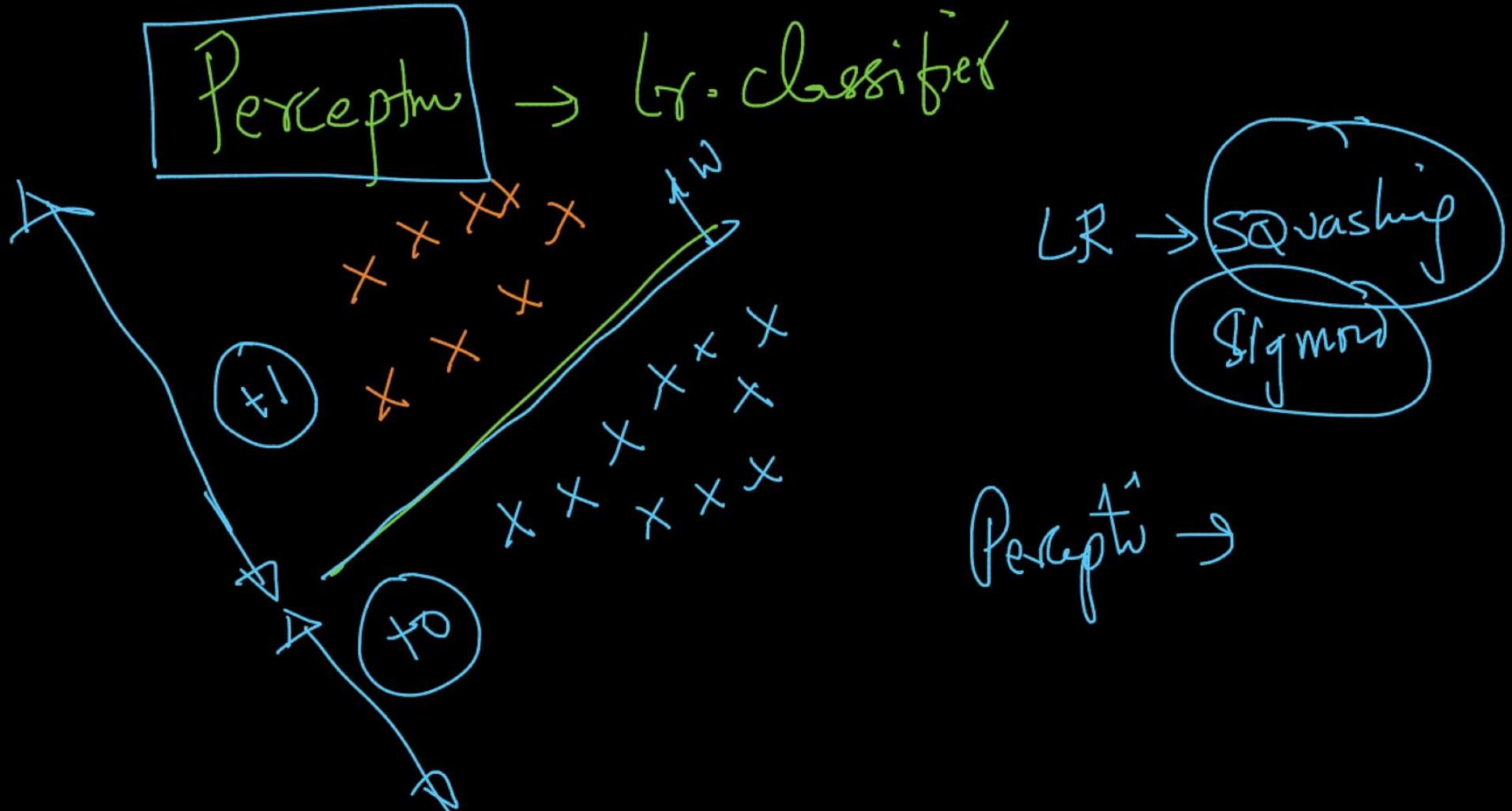


Train a NN  $\Rightarrow$  Weights on edges/vertices

LR  $\rightarrow$  Neuron  
activation fn  $\rightarrow$  Sigmoid







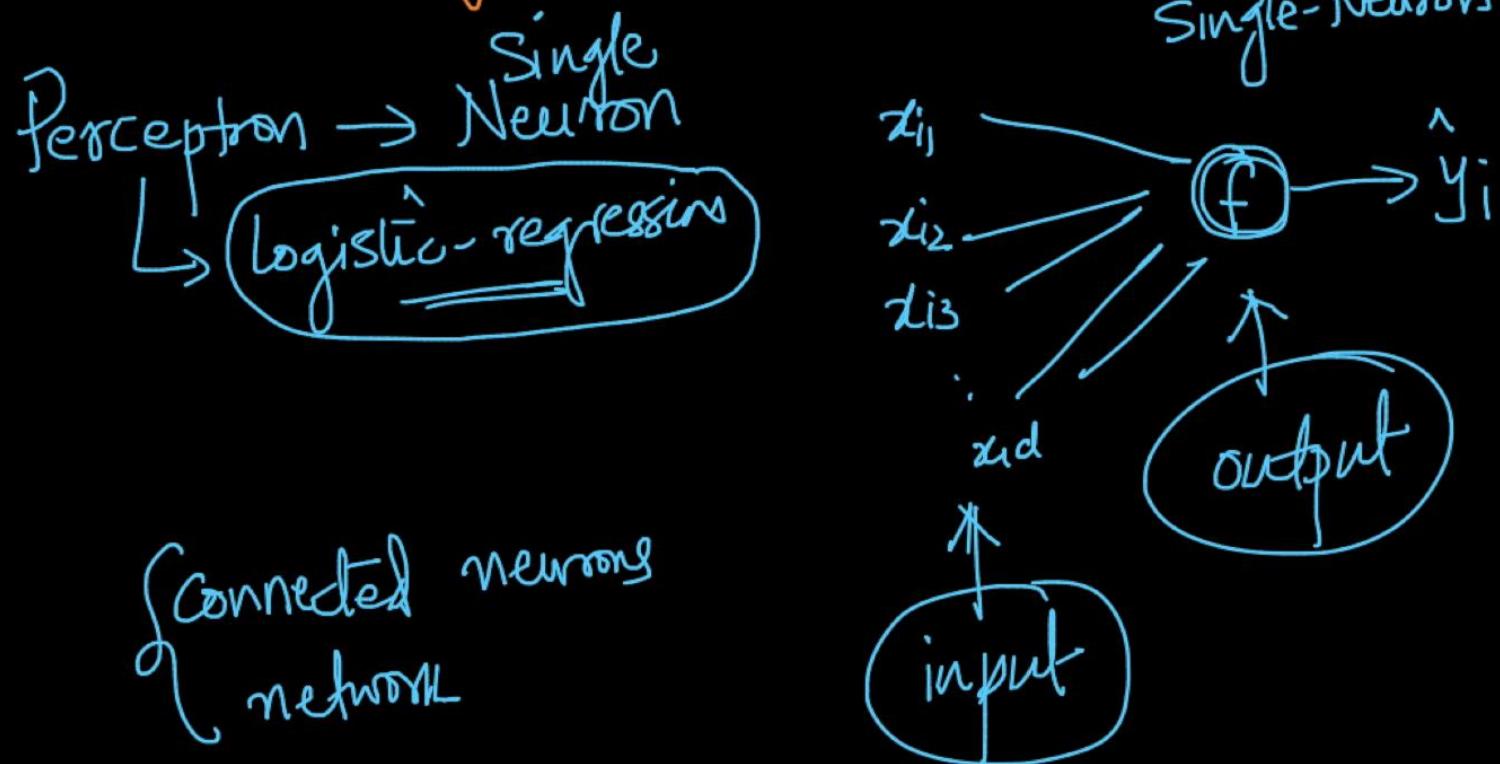
LR, perceptron



Simple single Neuron models



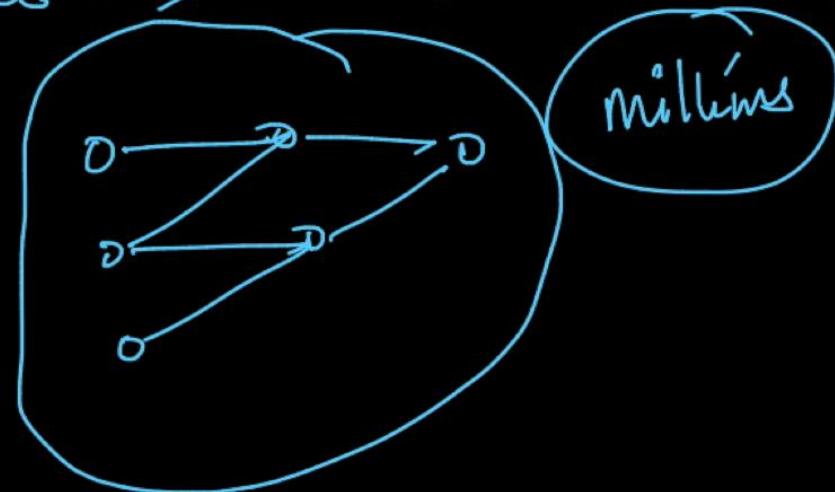
# Multi-layered Perception (MLP)



(Q) Why should we care about MLP?

(a) biological inspiration  
Neuroscience → humans, rats, monkeys

Neuron → Perception



c)

Mathematical

regression

$$\{x_i, y_i\} = \mathcal{D}$$

$$y_i = F(\underline{x_i})$$

$y_i \in \mathbb{R}$

Contents - Google Docs

plot(2\*sin(x^2)+sqrt(x^5)) - Google Search

Secure | https://www.google.co.in/search?ei=EzqvWpC6HrcvATRkJG4Aw&q=plot%282\*sin%28x%5E2%29%2Bsqr... Chekuri Srikan...

Google

plot(2\*sin(x^2)+sqrt(x^5))

All Maps Images News Videos More Settings Tools

About 1,25,00,000 results (0.75 seconds)

Graph for  $2\sin(x^2) + \sqrt{x^5}$

$f(x_i) = y_i$

$F(x) = 2\sin(x^2) + \sqrt{x^5}$

GraphSketch

<https://graphsketch.com/>

GraphSketch.com. Graph Plot Save Click here to download this graph. Permalink Permanent link to this graph page. ... graph of a function, you can just go to [http://graphsketch.com/\[function\]](http://graphsketch.com/[function]), like [http://graphsketch.com/sin\(x\)](http://graphsketch.com/sin(x)). You can even separate multiple equations with commas, like [http://graphsketch.com/sin\(x\),cos\(x\)](http://graphsketch.com/sin(x),cos(x))

APPLIED COURSE

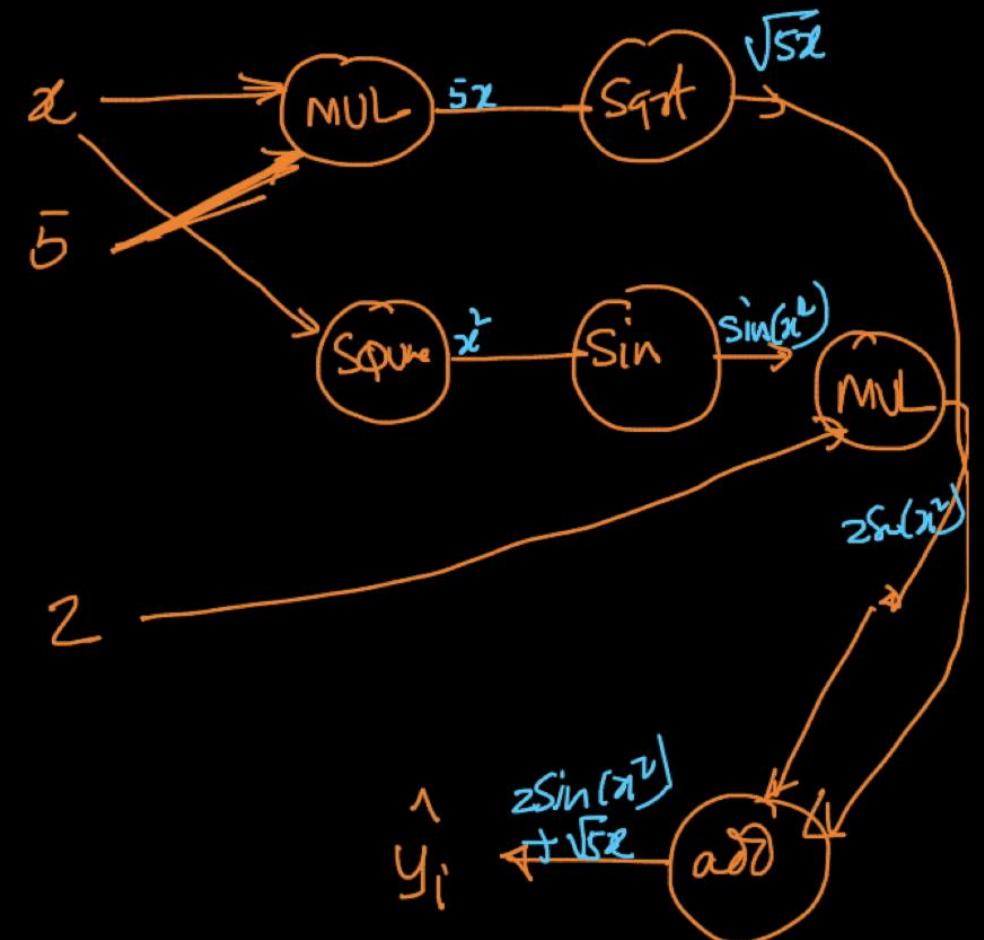
$f_1 \rightarrow \text{add}()$

$f_2 \rightarrow \text{square}()$

$f_3 \rightarrow \text{sqrt}()$

$f_4 \rightarrow \text{Sin}()$

$f_5 \rightarrow \text{MUL}()$

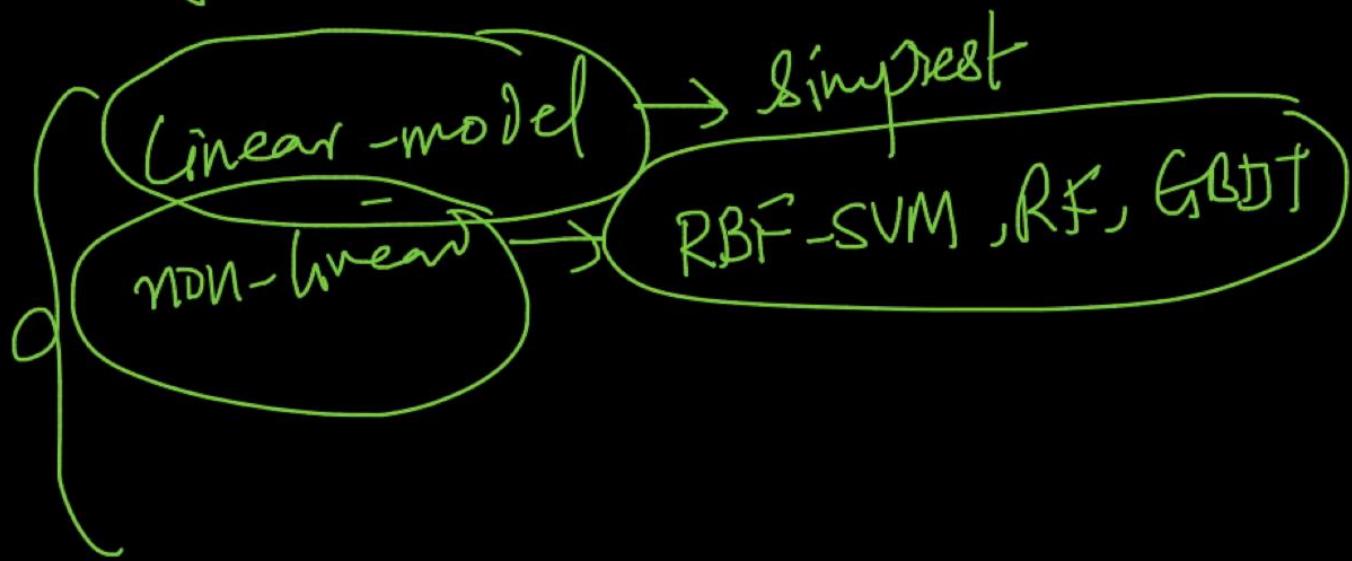


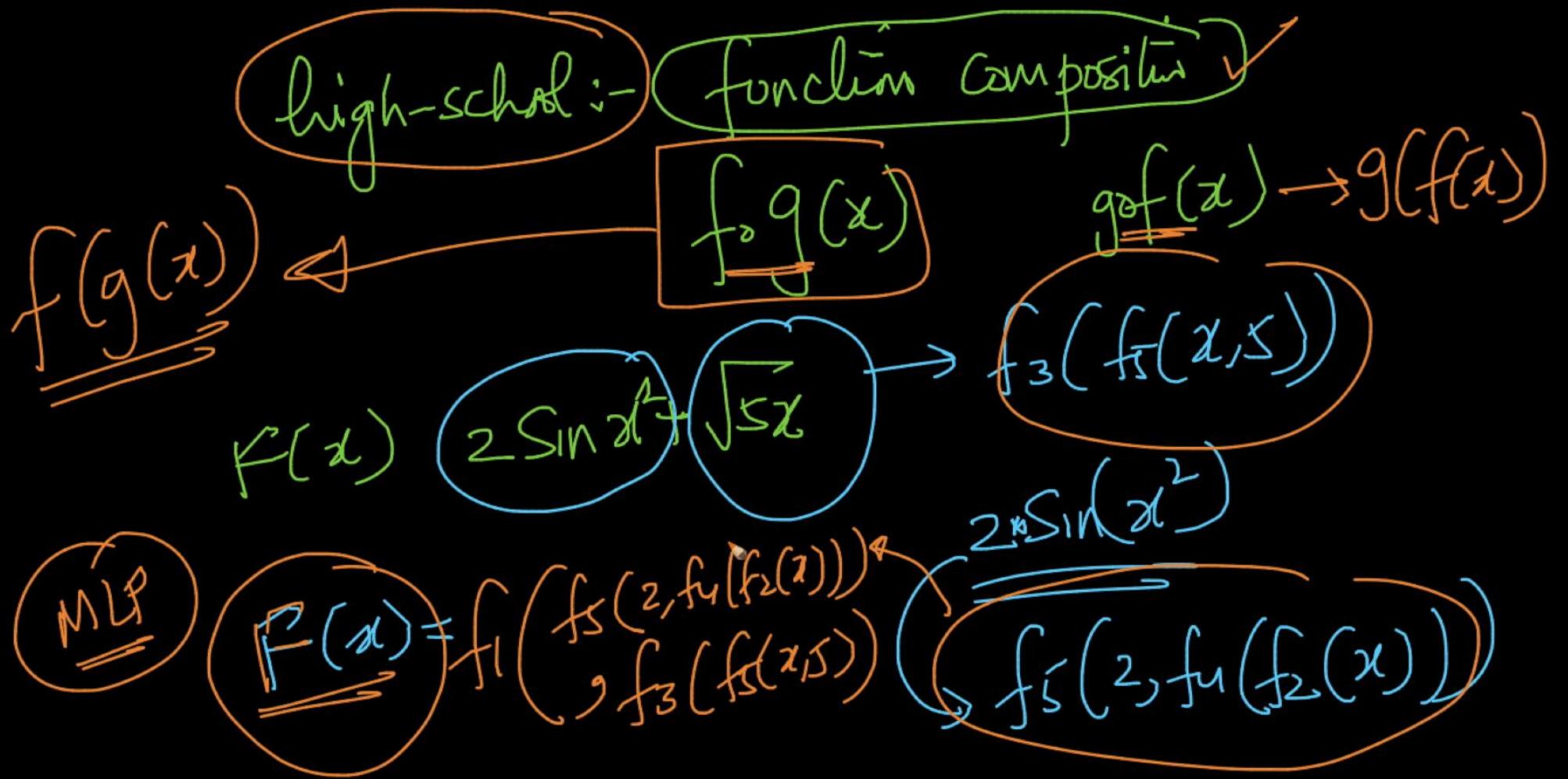
$$F(x) = 2 \sin(x^2) + \sqrt{5}x$$

{ By using multi-layered stretching  
we can come up with complex  
math-fns to solve your reg'n  
problems }

Multi-layered

→ enormous power

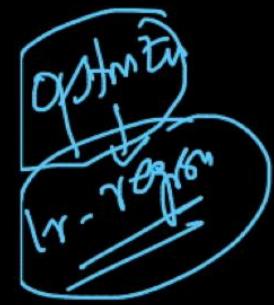




MLP :- graphical way of  
representing



ML - Struct → Powerful models  
(Overfit - very easily)

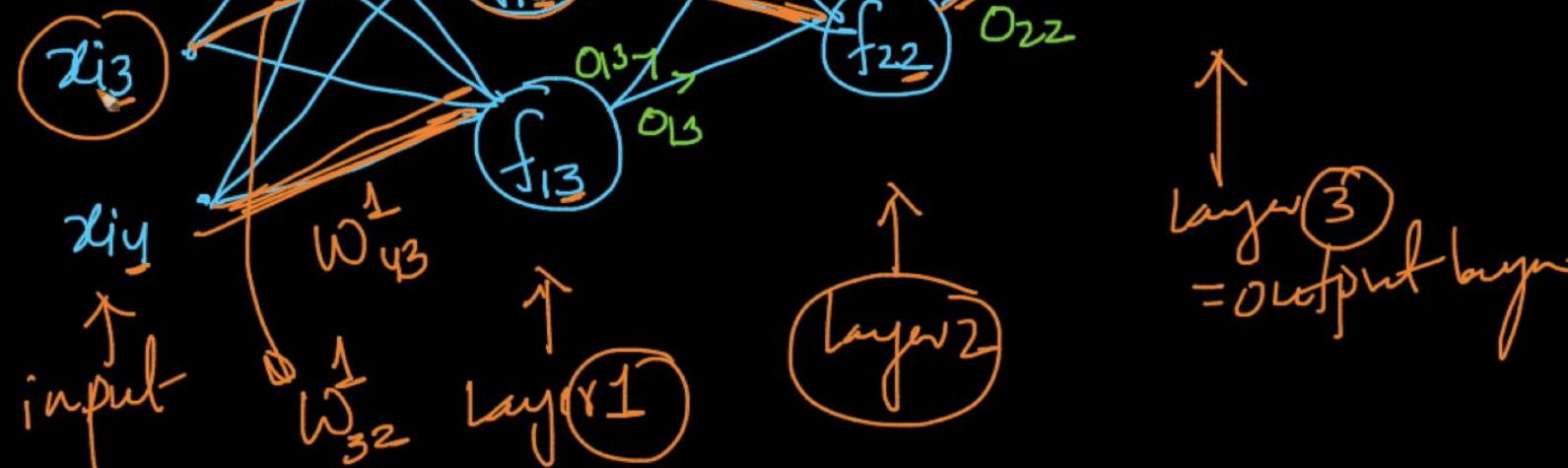


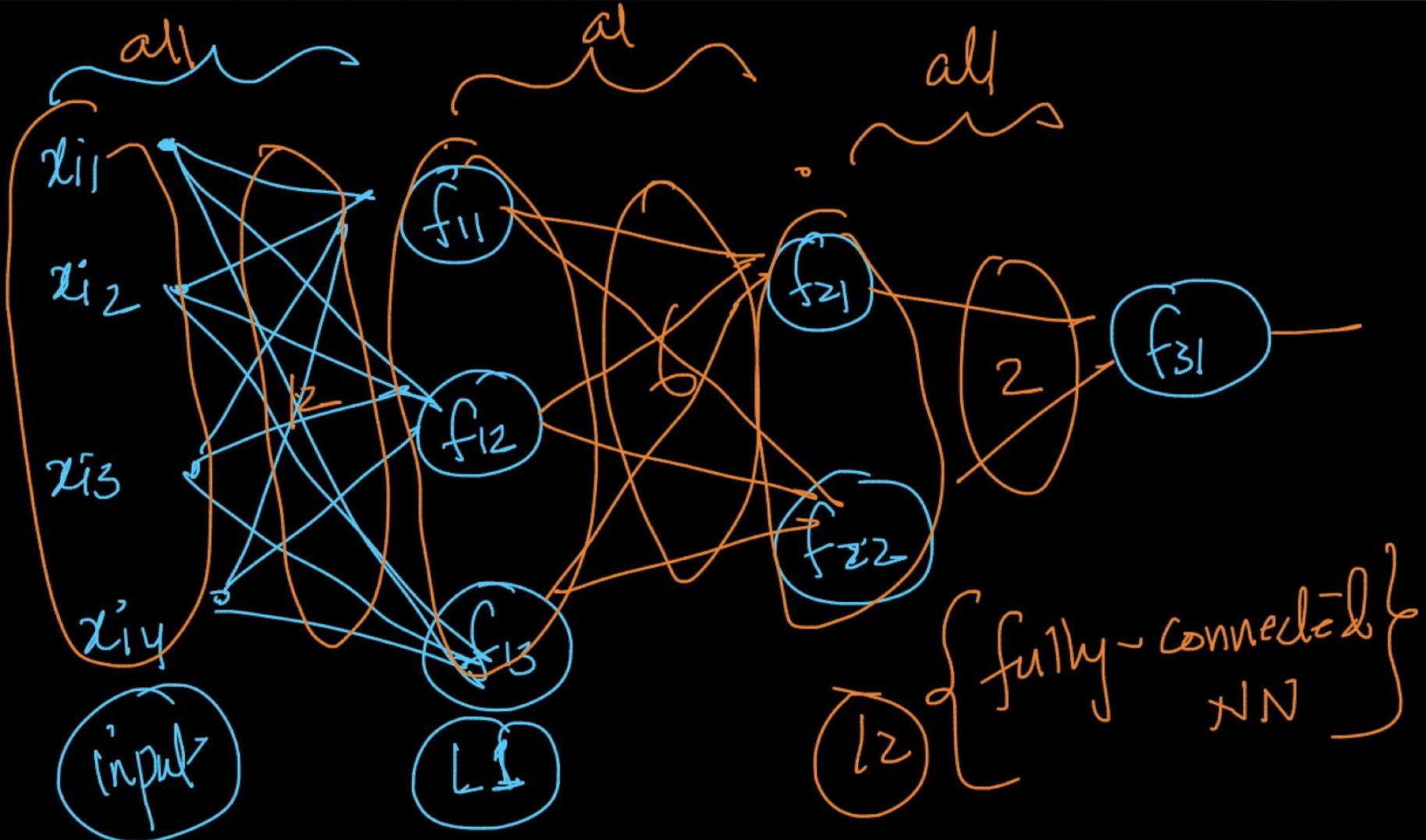
$x_{ij}$   
pt feature

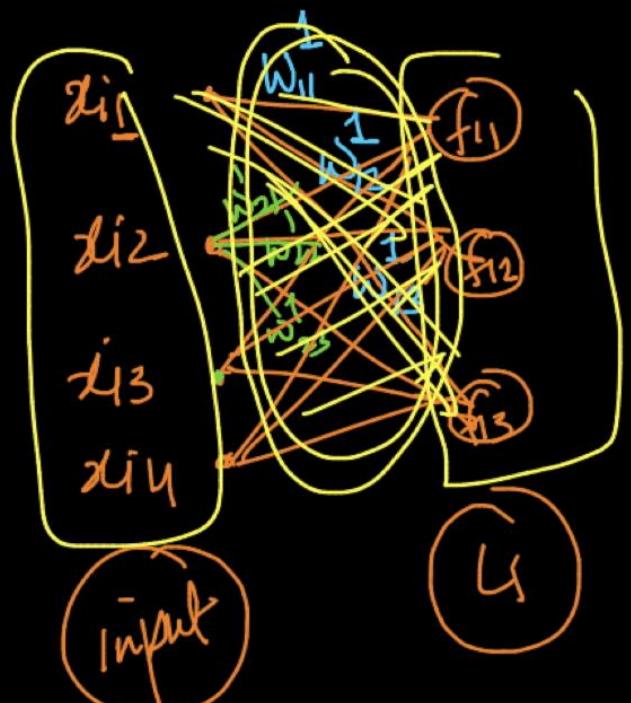


## Notation

$\mathcal{D} = \{(x_i, y_i)\}; x_i \in \mathbb{R}^M; y_i \in \mathbb{R}$  (Regn)





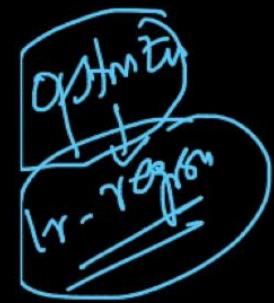


$$W = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \\ w_{41} & w_{42} & w_{43} \end{bmatrix}$$

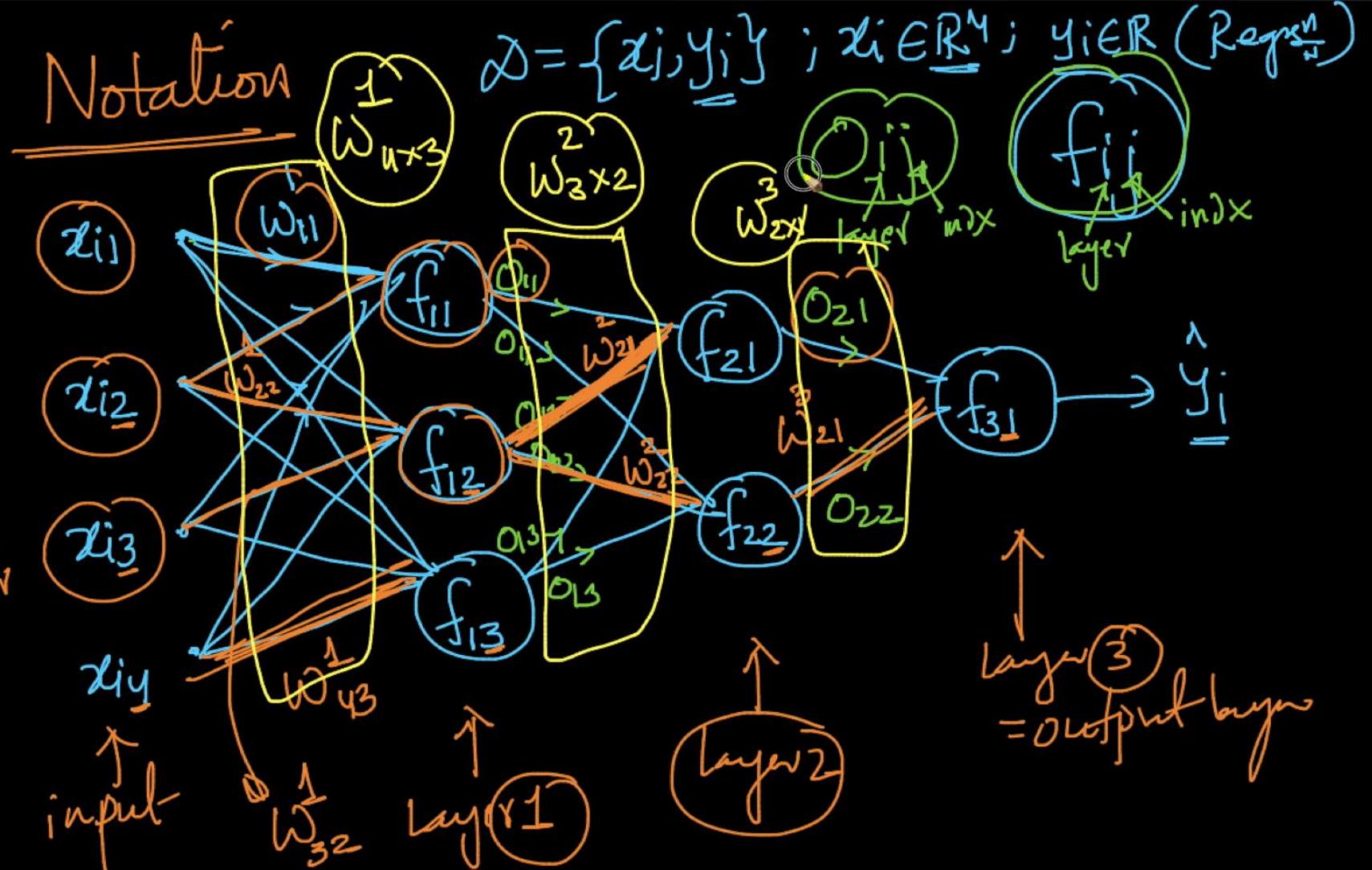
$4 \times 3$

$$(4) \times (3) \neq 12$$





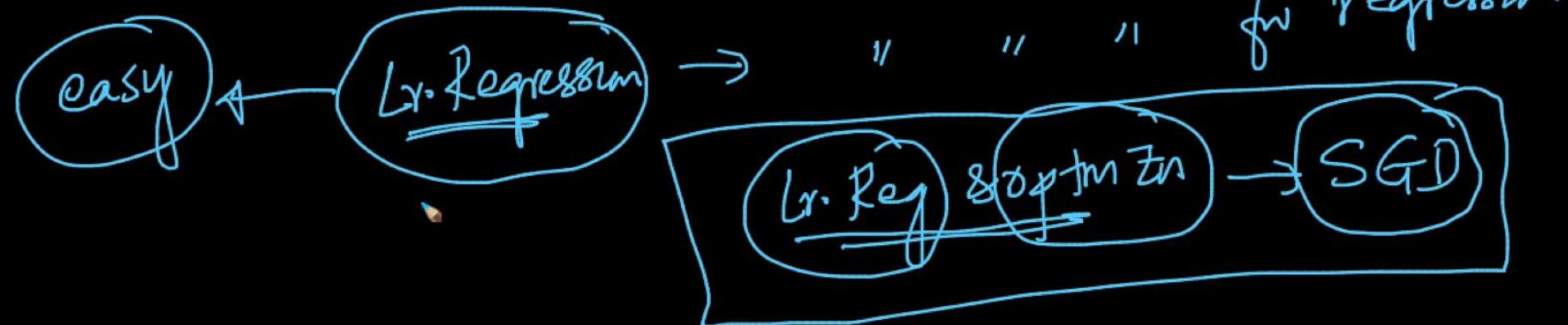
$x_{ij}$   
pt feature  
next layer



Train a single Neuron model

↓  
find the best edge-weights using  $D_{TR}$

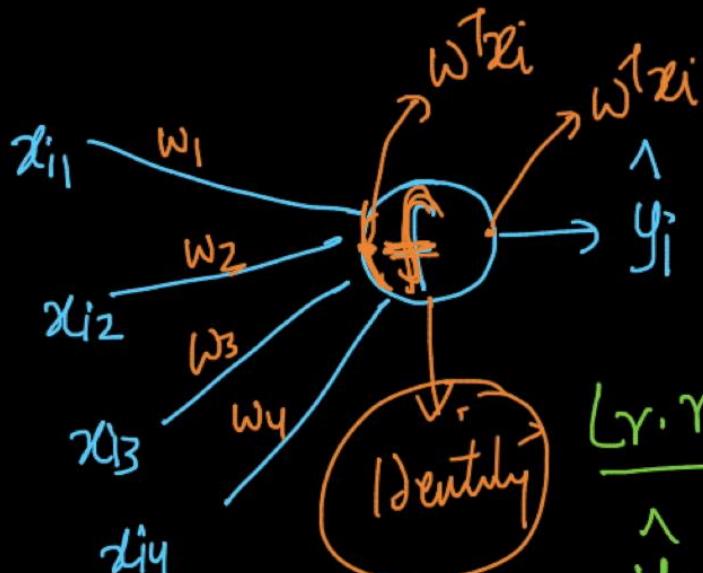
✓ Perception & LR → Single Neuron models  
for classif.



$$\mathcal{D} = \{(x_i, y_i)\}$$

$\log\text{-logit}$   
 $f = \text{logit}^{-1}$

$$y = \frac{\sum_{j=1}^d w_j x_{ij}}{f(z) = z}$$



lr. regression

$$\hat{y}_i = \sum_{j=1}^d w_j x_{ij}$$

$$\hat{y}_i = \boxed{w^T x_i}$$

$$x_i \in \mathbb{R}^d$$

$$y_i \in \mathbb{R}$$

lr. optimz

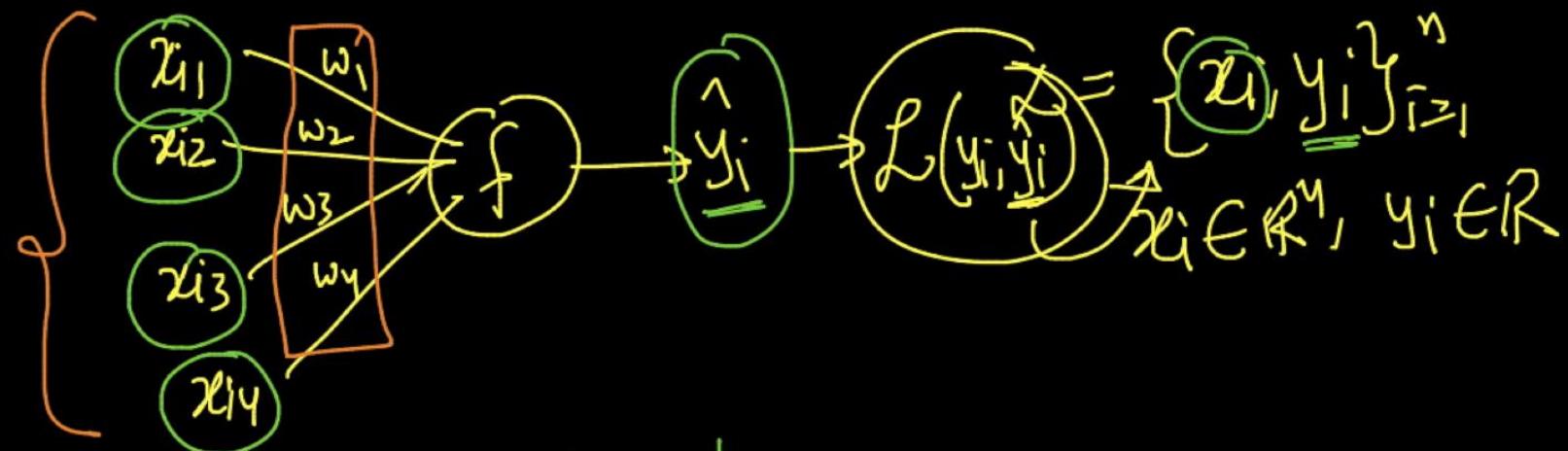
Ly. regres<sup>n</sup>:

$$\min_{w_i} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \text{reg}$$

$\hat{y}_i = w^T x_i$

OPTIMISATION

$$\min_{w_i} \sum_{i=1}^n (y_i - w^T x_i)^2 + \|w\|_2^2$$



① Define loss-fn

$$L = \sum_{i=1}^n L_i$$

$$\begin{aligned} L &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \text{reg} \\ L_i &= (y_i - \hat{y}_i)^2 \end{aligned}$$

$$\hat{y}_i = \mathbf{w}^T \mathbf{x}_i$$

②

optimization

$$\min_{w_i} \sum_{i=1}^n (y_i - \underbrace{w_i^\top x_i}_{\hat{y}_i = f(w^\top x_i)})^2 + \text{reg}$$

$\hookrightarrow$  Lr. reg  $f$  is I

$$\min_{w_i} \sum_{i=1}^n (y_i - \underbrace{f(w^\top x_i)}_{\text{percept}})^2 + \text{reg}$$

$\hookrightarrow$  I : l.y. reg  
 $\hookrightarrow$  Sigmoid : - log-reg



$$\hat{\omega} = \arg \min_{\omega} \sum_{i=1}^n (y_i - f(\omega^T x_i))^2 + \text{reg}$$

③ Solve the optimization problem

(a) Initialization of  $\underline{\omega}$

(b)



$$L = \begin{bmatrix} \frac{\partial L}{\partial \omega_1} \\ \frac{\partial L}{\partial \omega_2} \\ \vdots \\ \frac{\partial L}{\partial \omega_d} \end{bmatrix}$$



$$\omega \in \mathbb{R}^Y$$

$$x_i \in \mathbb{R}^Y$$

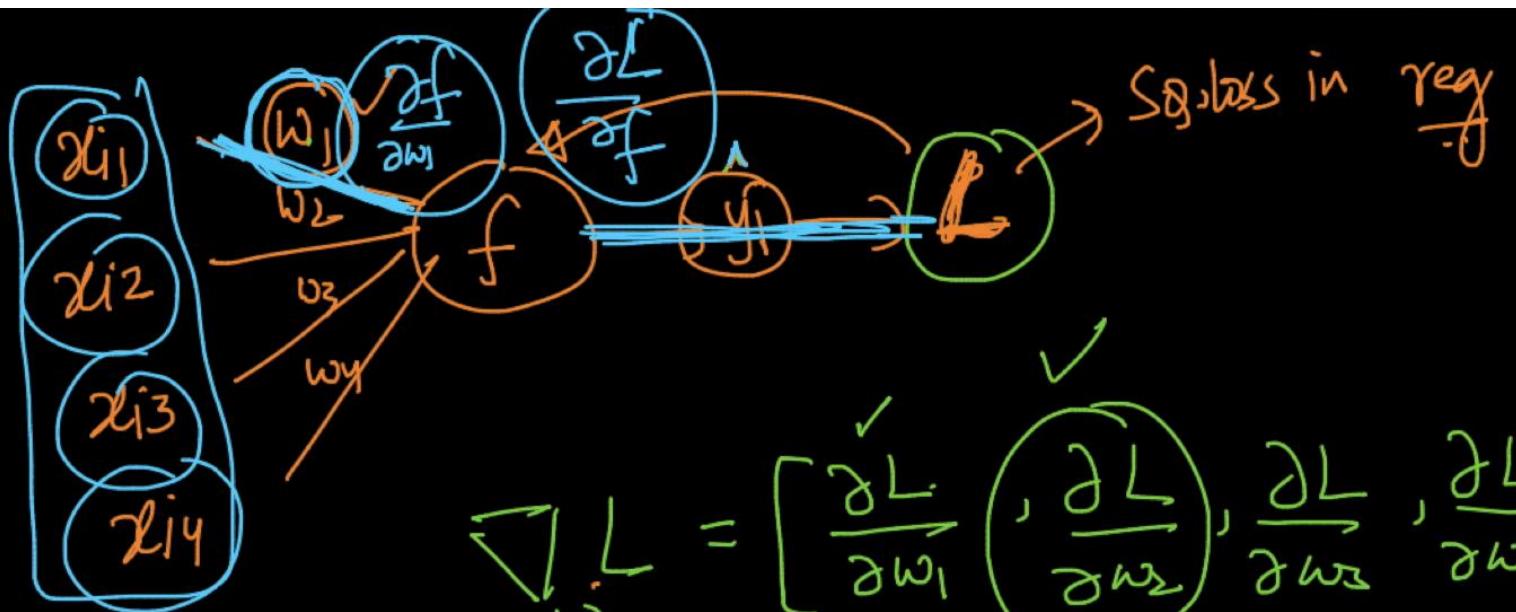
⑥  $w_{\text{new}} = w_{\text{old}} - \eta [\nabla_L]_{\text{old}}$

$$\boxed{(w_i)_{\text{new}} = (w_i)_{\text{old}} - \eta \left[ \frac{\partial L}{\partial w_i} \right]_{(w_i)_{\text{old}}}}$$

$f_w - \text{iter} = 1 \text{ to } K$

GD:   $\rightarrow x_i's \& y_i's$

SGD:   $\approx$  One pt  $\{x_i, y_i\} \leftarrow$   
Small batch of pts  $\leftarrow$  batch SGD

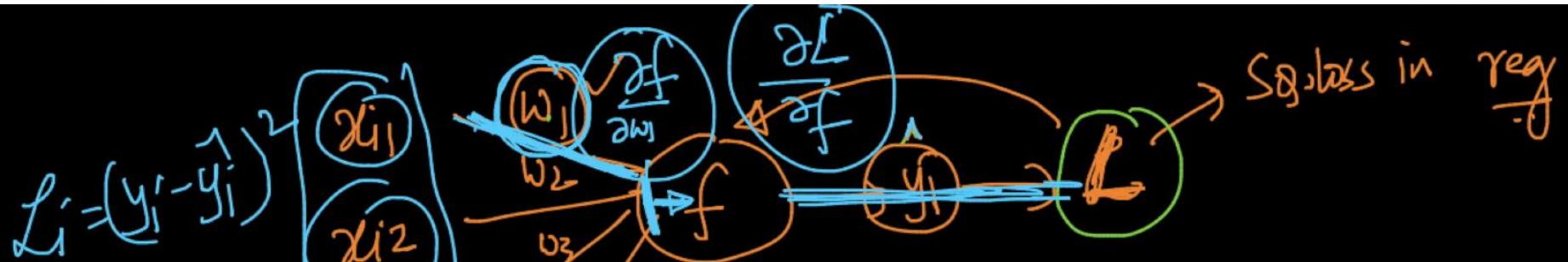


$$\nabla_{\omega} L = \left[ \frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}, \frac{\partial L}{\partial w_3}, \frac{\partial L}{\partial w_4} \right]^T$$

$$\frac{\partial L}{\partial w_1} =$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial w_1}$$

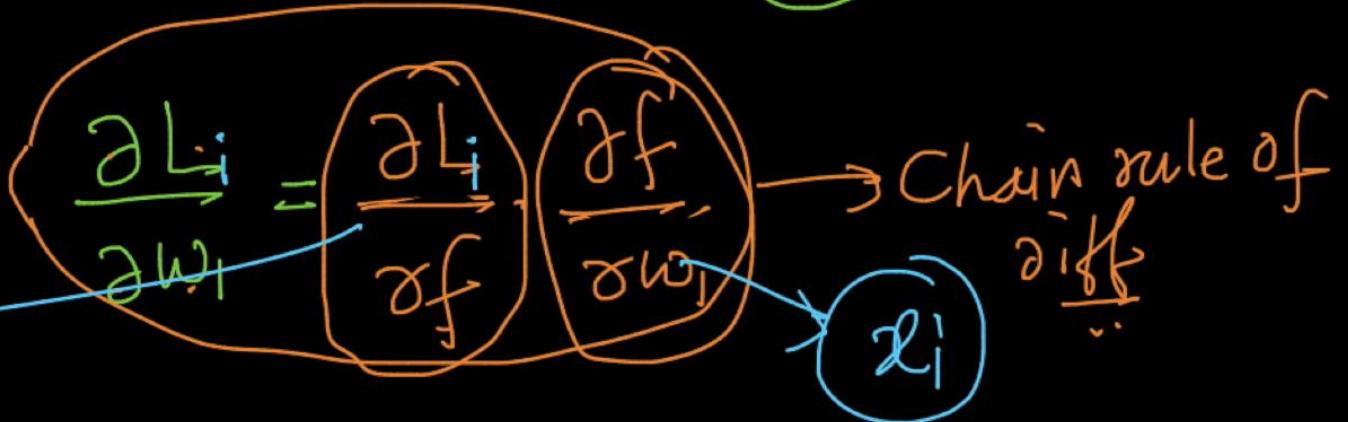
Chain rule of diff.



$$\nabla_w L = \left[ \frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}, \frac{\partial L}{\partial w_3}, \frac{\partial L}{\partial w_4} \right]^T$$

$$\frac{\partial L}{\partial w_1} =$$

$$-2(y_i - \hat{y}_i)$$



$$w_{\text{new}} = w_{\text{old}} - \eta \left( \nabla L_w \right)_{\text{old}}$$

gradien $\leftarrow$

$$(w_i)_{\text{new}} = (w_i)_{\text{old}} - \eta \left( \frac{\partial L}{\partial w_i} \right)_{(w_i)_{\text{old}}}$$

fw-iter = 1 to K

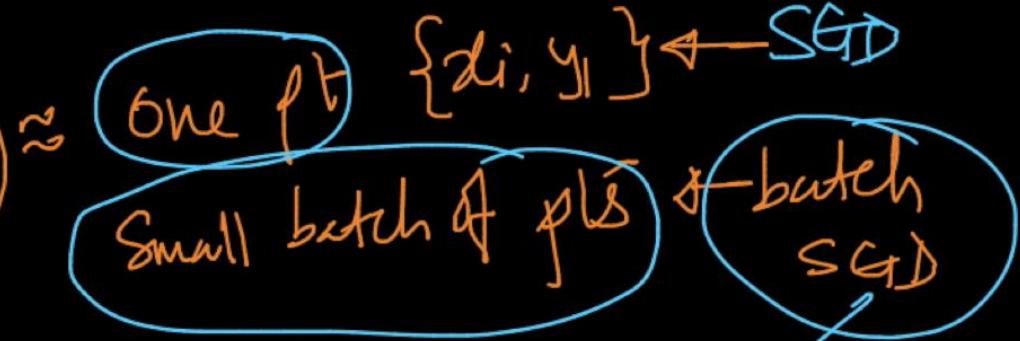
learning rate

GD:



$\rightarrow$  all  $x_i$ 's &  $y_i$ 's

SGD:





Training an MLP → Back propagation algo

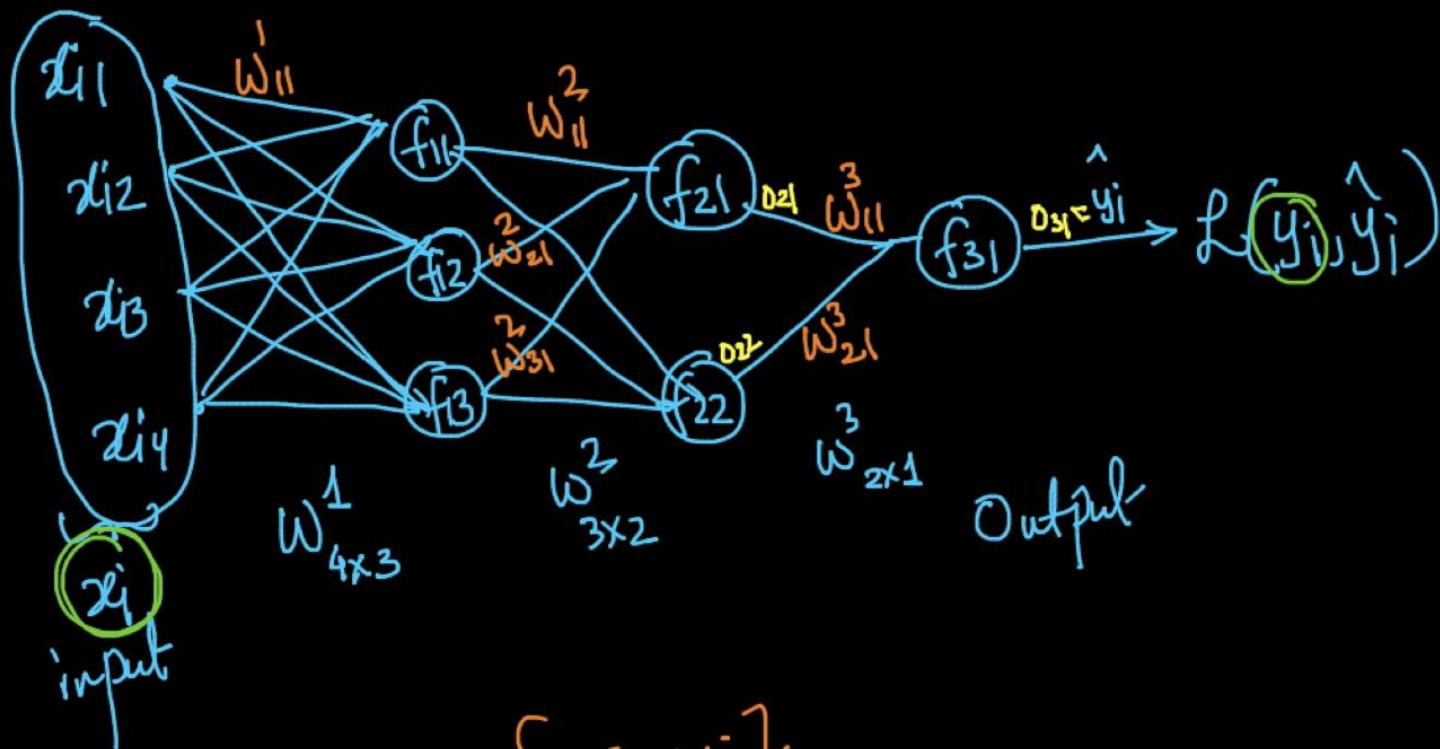
✓ MLP

Train single Neuron - model  
SGD ✓  
chain-rule

$$\mathcal{D} = \{x_i, y_i\}$$

$x_i \in \mathbb{R}^4$  } regression prob → squared loss  
 $y_i \in \mathbb{R}$  }





$$\mathcal{D} = \{x_i, y_i\}$$

determining  $\omega^1_{4 \times 3}, \omega^2_{3 \times 2}, \omega^3_{2 \times 1}$  = 20 weights

(12) (6) (2)

①

$$L = \sum_{i=1}^n \underbrace{(y_i - \hat{y}_i)^2}_{\text{Sq. loss}} + \gamma \text{reg}$$

$$\sqrt{L_i} = (y_i - \hat{y}_i)^2$$

$$L = \sum_{i=1}^n L_i + \gamma \text{reg}$$

~~Optimization:~~

$$\min_{w^1, w^2, w^3} L$$

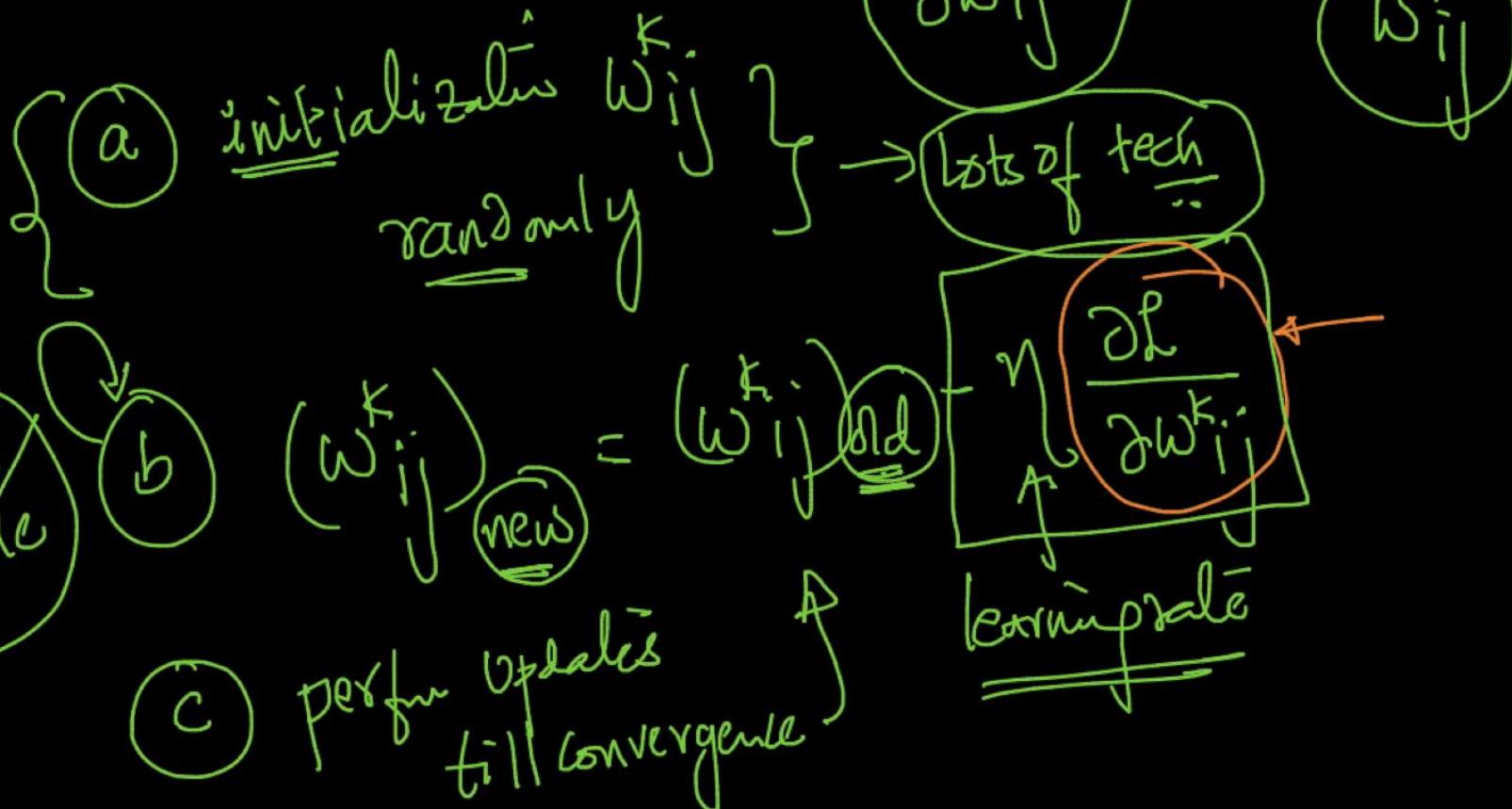
$$L_2 = \sum_{i,j,k} (w_{ijk}^k)^2$$

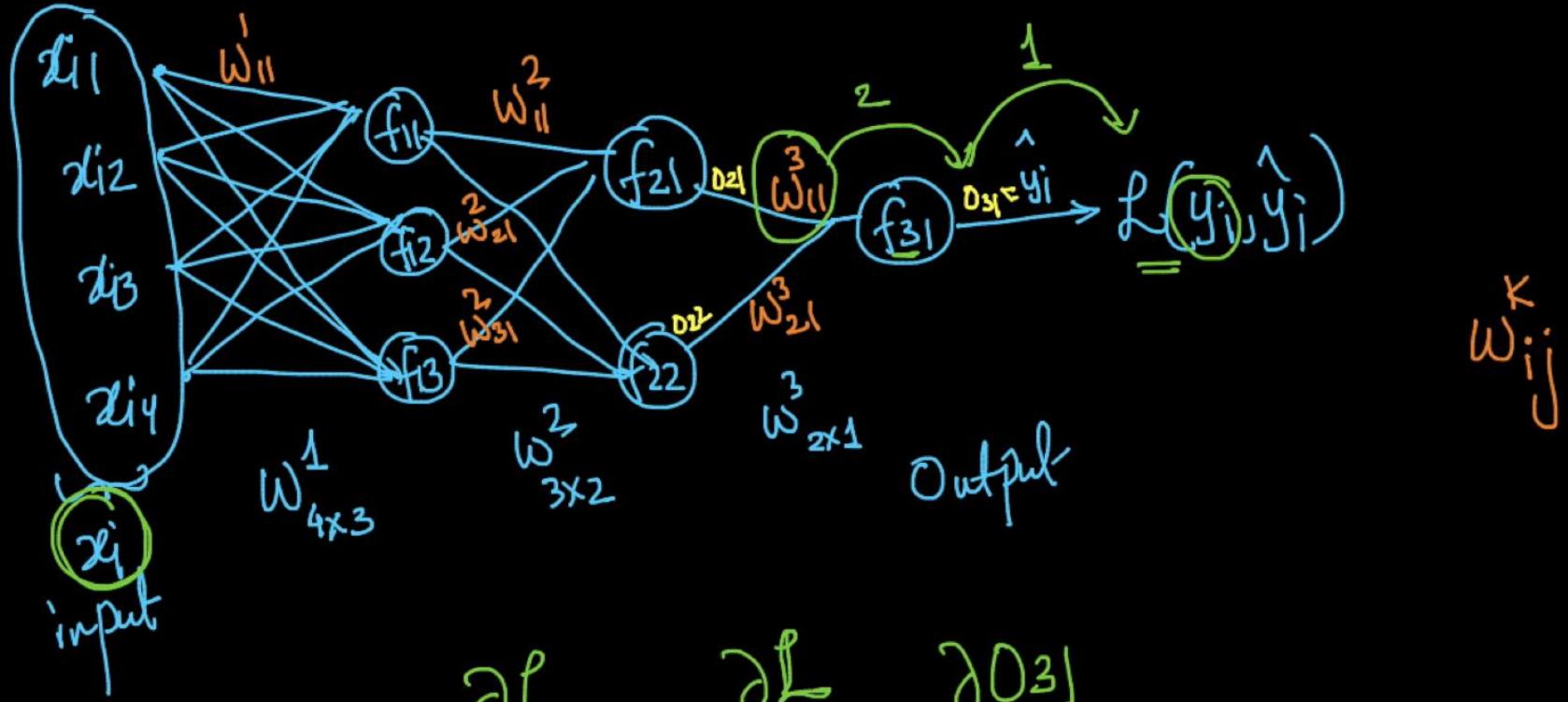
$$\sum_{i,j,k} |w_{ijk}| \quad \uparrow L_1$$

$$\boxed{\min_{w_{ijk}} L}$$

②

SGD & GD





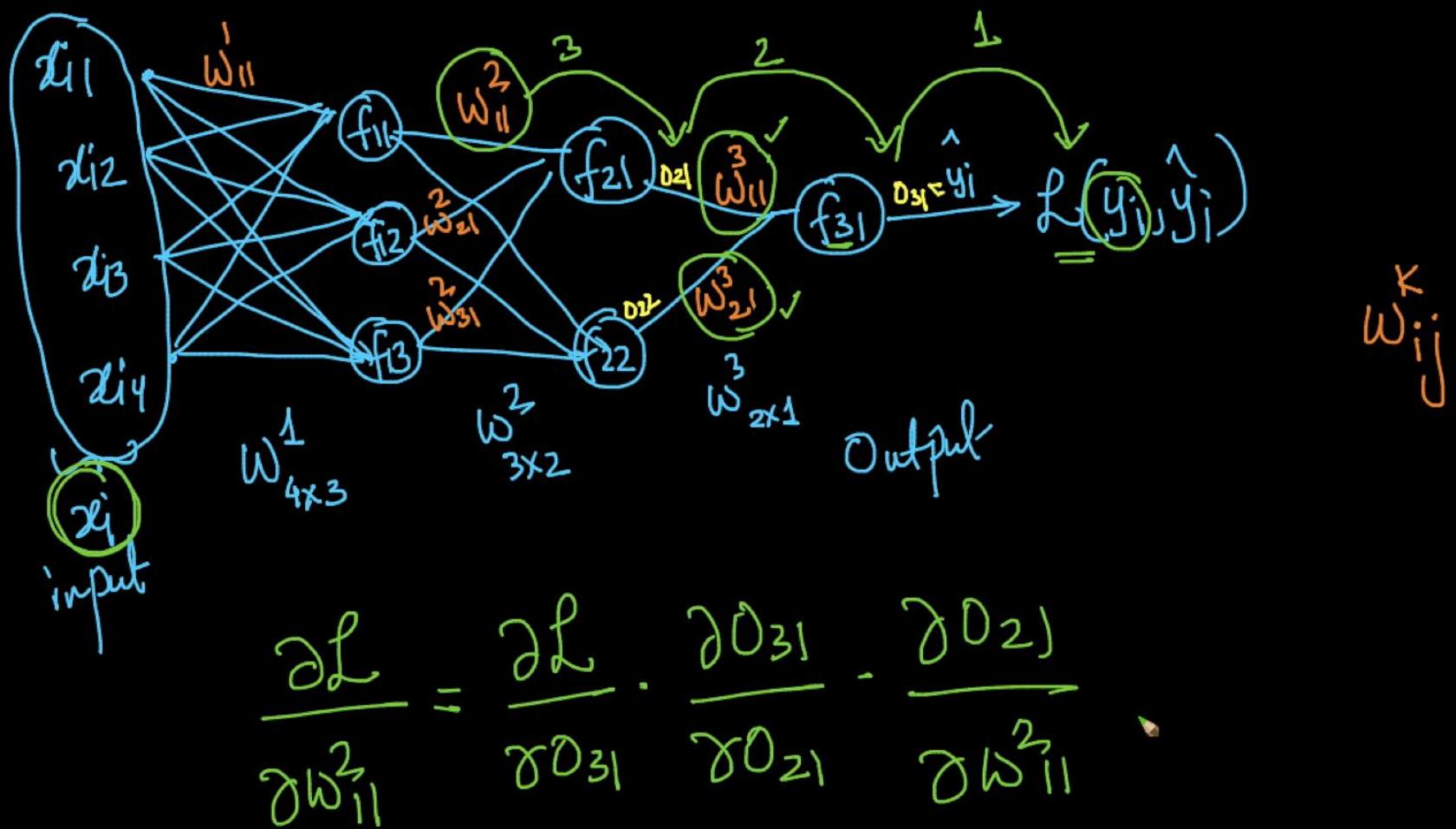
$$\frac{\partial L}{\partial w_{11}^3} = \underbrace{\frac{\partial L}{\partial o_{31}}}_{1} \cdot \underbrace{\frac{\partial o_{31}}{\partial w_{11}^3}}_{2}$$

$w^3$

$$\frac{\partial L}{\partial w_{11}^3} = \frac{\partial L}{\partial \theta_{31}} \cdot \frac{\partial \theta_{31}}{\partial w_{11}^3} \quad \leftarrow \text{chain rule}$$

$$\frac{\partial L}{\partial w_{21}^3} = \frac{\partial L}{\partial \theta_{31}} \cdot \frac{\partial \theta_{31}}{\partial w_{21}^3} \quad \leftarrow \text{chain rule}$$





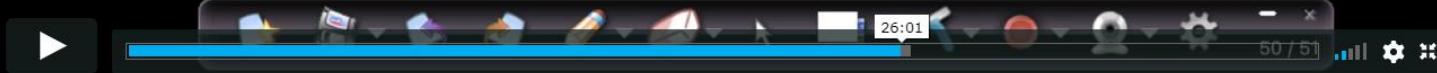
$$\frac{\partial L}{\partial w_{11}^3} = \frac{\partial L}{\partial \delta_{31}} \cdot \frac{\partial \delta_{31}}{\partial \delta_{21}} \cdot \frac{\partial \delta_{21}}{\partial w_{11}^3}$$

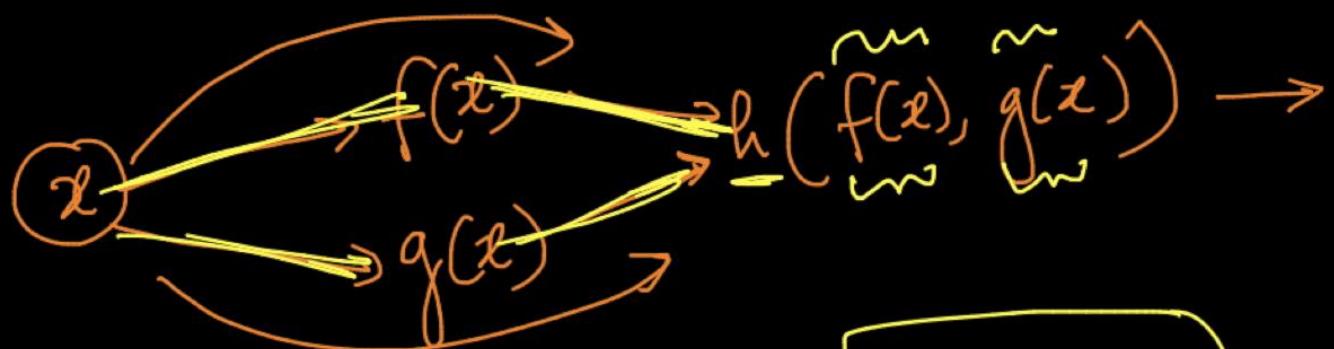
$\tilde{w}^2$

$$\frac{\partial \hat{L}}{\partial w_{11}^2} = \frac{\partial L}{\partial \theta_{31}} \cdot \frac{\partial \theta_{31}}{\partial \theta_{21}} \cdot \frac{\partial \theta_{21}}{\partial w_{11}^2}$$

$$\frac{\partial \hat{L}}{\partial w_{21}^2} = \frac{\partial L}{\partial \theta_{31}} \cdot \frac{\partial \theta_{31}}{\partial \theta_{21}} \cdot \frac{\partial \theta_{21}}{\partial w_{21}^2}$$

$$\frac{\partial \hat{L}}{\partial w_{31}^2} = \frac{\partial L}{\partial \theta_{31}} \cdot \frac{\partial \theta_{31}}{\partial \theta_{21}} \cdot \frac{\partial \theta_{21}}{\partial w_{31}^2}$$





$$\frac{\partial h}{\partial x} = \left[ \frac{\partial h}{\partial f} \cdot \frac{\partial f}{\partial x} \right] + \left[ \frac{\partial h}{\partial g} \cdot \frac{\partial g}{\partial x} \right]$$

+  
↓  
SUM

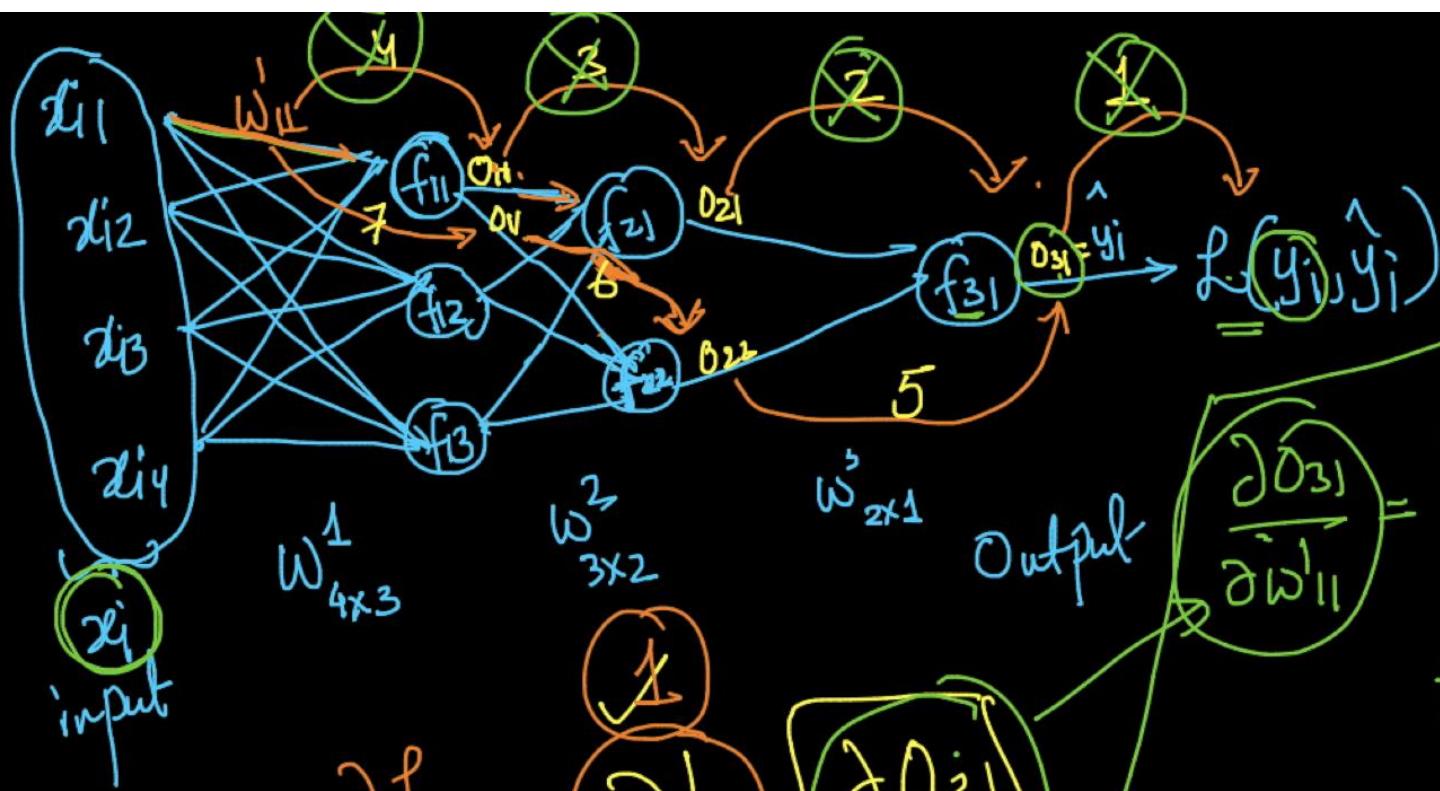
← chain  
 rule  
 =

$$\begin{array}{ccccc}
 & f & & h & \\
 x & \nearrow & \searrow & \swarrow & k \\
 & g & & &
 \end{array}$$

$$\frac{\partial k}{\partial x} = \frac{\partial k}{\partial h} \cdot \frac{\partial h}{\partial x} + \frac{\partial h}{\partial g} \cdot \frac{\partial g}{\partial x}$$

$$\frac{\partial h}{\partial x} = \frac{\partial h}{\partial f} \cdot \frac{\partial f}{\partial x} + \frac{\partial h}{\partial g} \cdot \frac{\partial g}{\partial x}$$

$$\frac{\partial k}{\partial x} = \frac{\partial k}{\partial h} \cdot \left\{ \frac{\partial h}{\partial f} \cdot \frac{\partial f}{\partial x} + \frac{\partial h}{\partial g} \cdot \frac{\partial g}{\partial x} \right\}$$



$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial D_{31}} \cdot \frac{\partial D_{31}}{\partial w_{11}}$$

Output

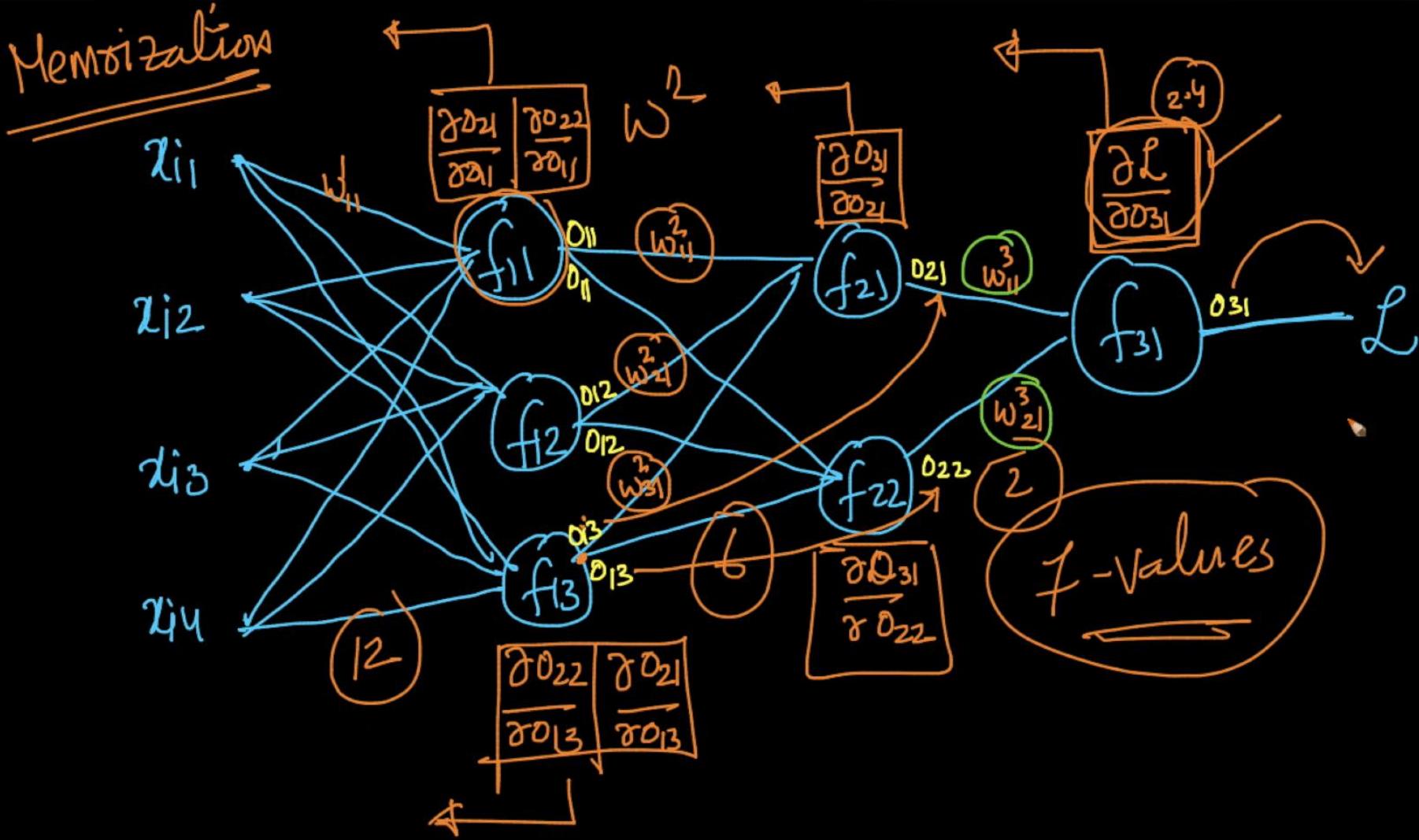
$$\begin{aligned} \frac{\partial D_{31}}{\partial w_{11}} &= \frac{\partial D_{31}}{\partial D_{21}} \cdot \frac{\partial D_{21}}{\partial w_{11}} \\ &+ \frac{\partial D_{31}}{\partial D_{22}} \cdot \frac{\partial D_{22}}{\partial w_{11}} \end{aligned}$$

$\hat{\omega}^1$

$$\partial \omega_{11}^1 = \frac{\partial L}{\partial \dot{\theta}_{31}} \cdot \left\{ \frac{\partial \theta_{31}}{\partial \theta_{21}} \cdot \frac{\partial \theta_{21}}{\partial \theta_{11}} \frac{\partial \theta_{11}}{\partial \omega_{11}^1} \right\} +$$

$$\frac{\partial \theta_{31}}{\partial \theta_{22}} \cdot \frac{\partial \theta_{22}}{\partial \theta_{11}} \cdot \frac{\partial \theta_{11}}{\partial \omega_{11}^1} \}$$



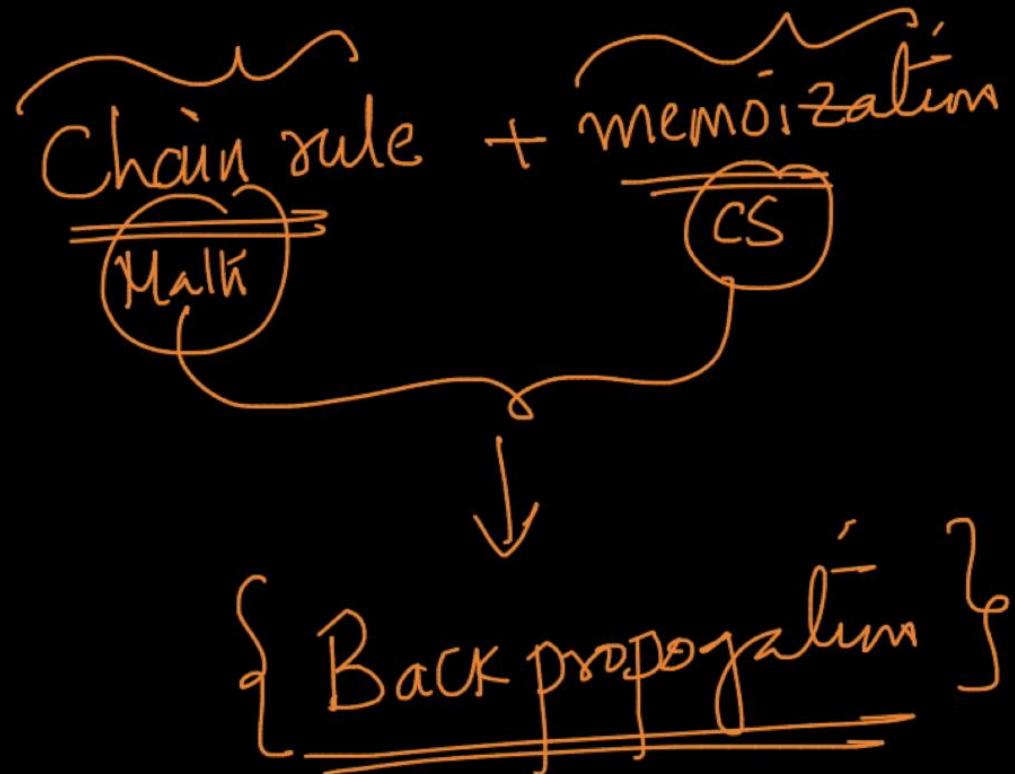


## Memoization

{ if there is any operation that is used  
many times repeatedly,  
{ - it's a good idea to compute it  
once → save it → send it

huge Speed-up ; slightly more memory

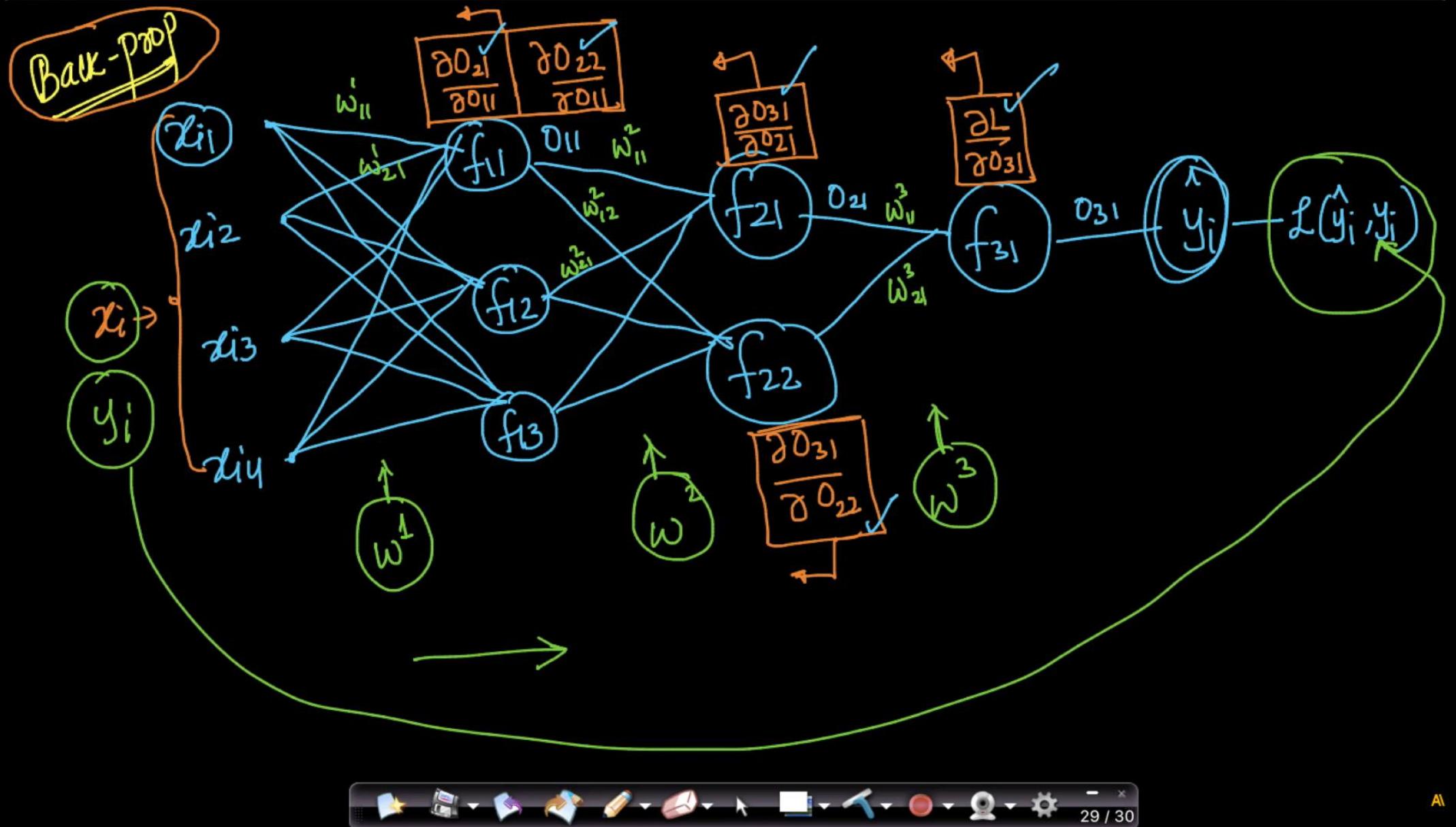


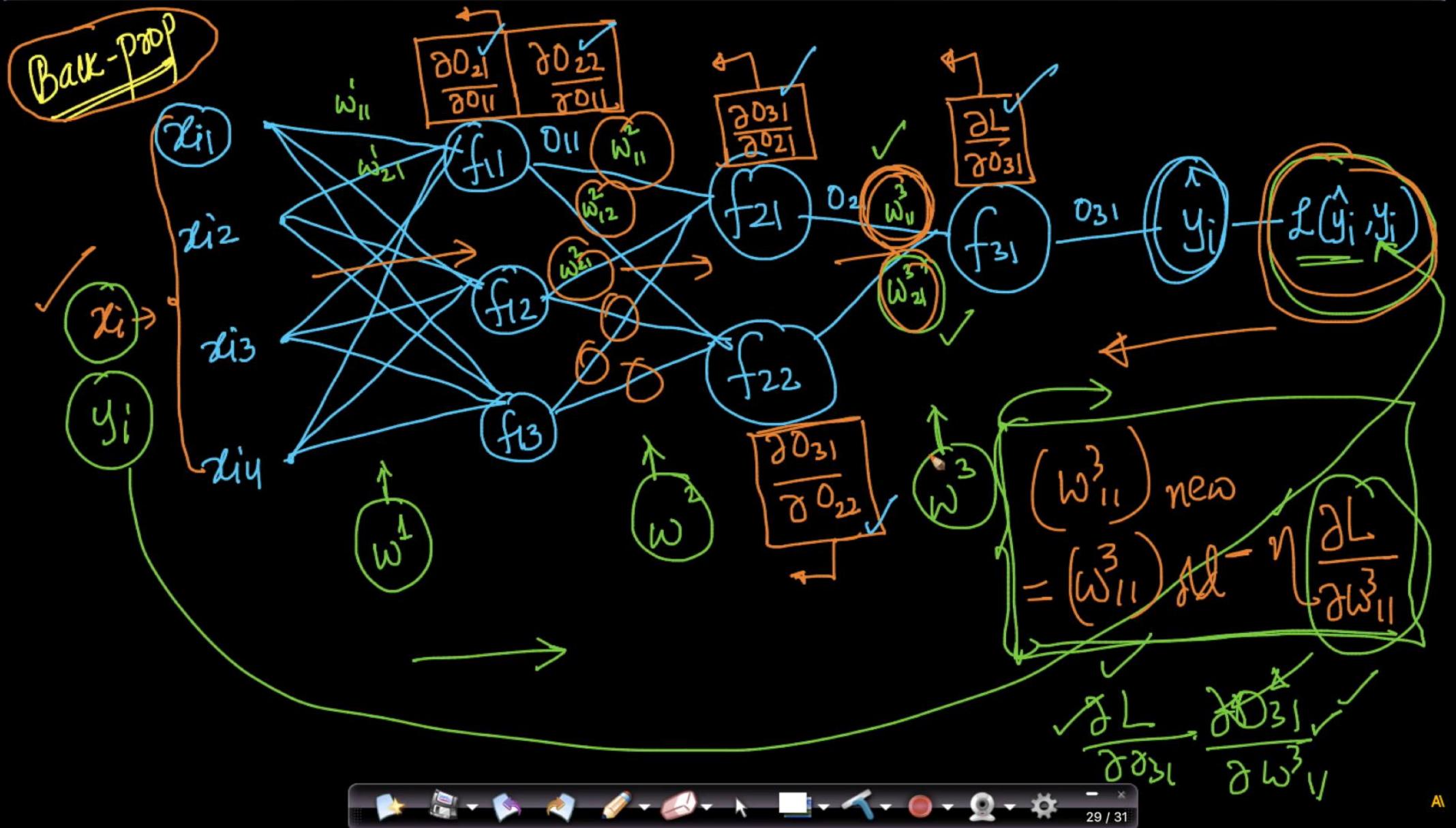


$$\mathcal{D} = \{\underline{x_i}, \underline{y_i}\}$$

- ① initialize  $w_{ij}$ 's
- ② for each  $x_i$  in  $\mathcal{D}$ 
  - a pass  $x_i$  forward through the network  $\rightarrow$  forward prop
  - b Compute  $L(\hat{y}_i, y_i)$
  - c Compute all the derivatives chain-rules,  
& memorization







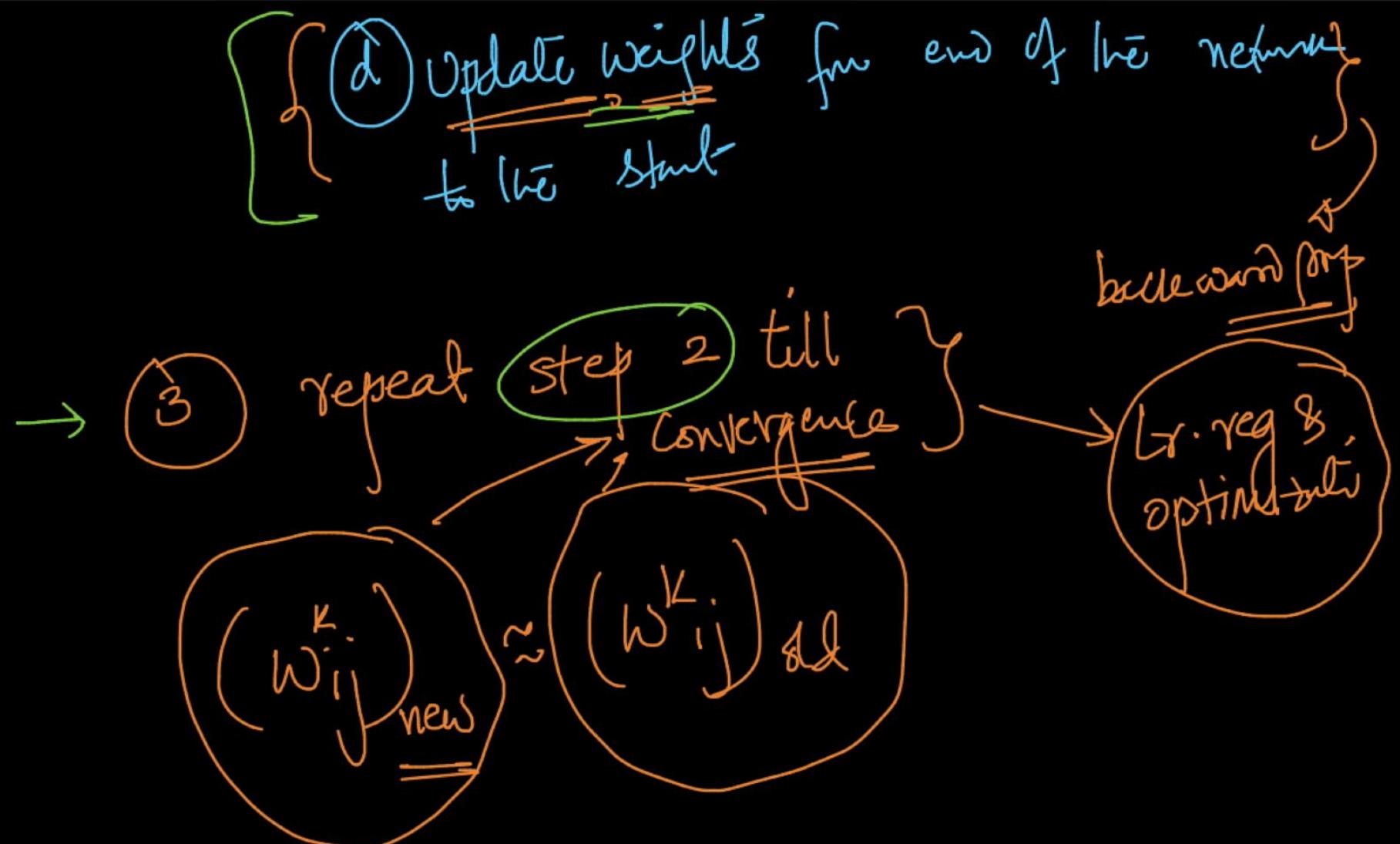
$\mathcal{D} \Rightarrow \{\underline{x_i}, \underline{y_i}\}$

✓ ① initialize  $w_{ij}$ 's

epoch

② one through NN

- ✓ ② for each  $\underline{x_i}$  in  $\mathcal{D}$
- a pass  $\underline{x_i}$  forward through the network  $\rightarrow$  forward prop
  - b Compute  $L(\hat{y}_i, y_i)$
  - c Compute all the derivatives
- chain-rule,  
& memorization

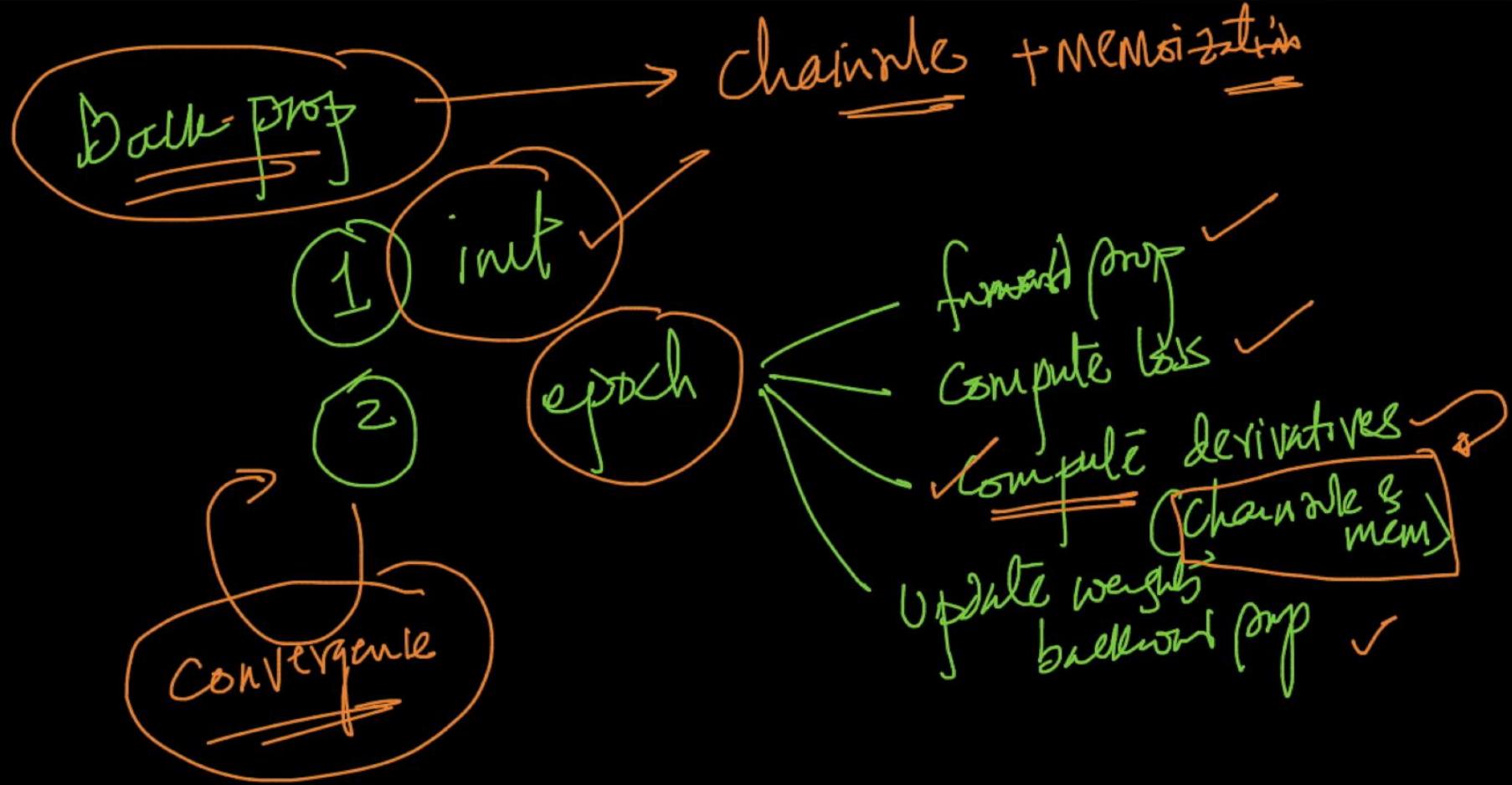


$\mathcal{D} \rightarrow \{(x_i, y_i)\} \rightarrow \text{epochs}$

5 times  $\rightarrow$  5 epochs

multiple epochs





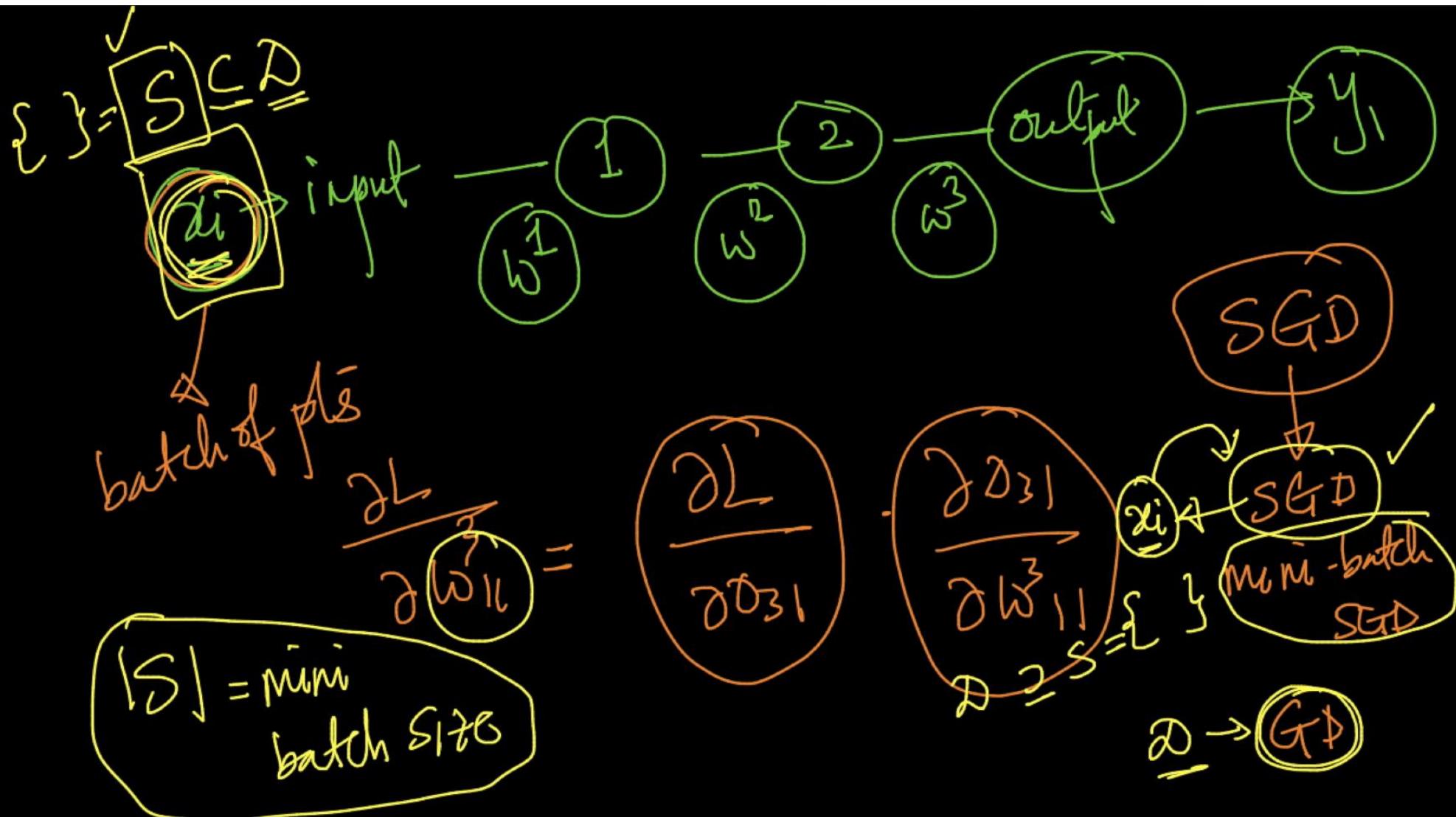
V.V. · IMP

Back-prop  
is training NN

\* Activations functions are  
differentiable

If it is easily / fast  
differentiable

{ Speed up the training of the  
NN



{ all of the data points in RAM  
is computing "J" using  $\mathcal{D}$

mini-batch based back-prop  $\rightarrow$  most popular approach

$$S = \{ \gamma \} \subseteq \mathcal{D}$$


$10^k$  points in  $\mathcal{D}$   
mini-batch =  $\underline{\underline{100}}$  →  $[64, 128, 256, 32]$  → RAM

epoch →  $10,000 - 100$   
 $100$   
for each site  
mini-batches of size  $\underline{\underline{100}}$   
forward pass  $L$ ,  $\frac{\partial L}{\partial w_{ij}}$ , update, backprop

(2) for each  $10^k \rightarrow 1$

## Activation functions

$$f_{ij}$$

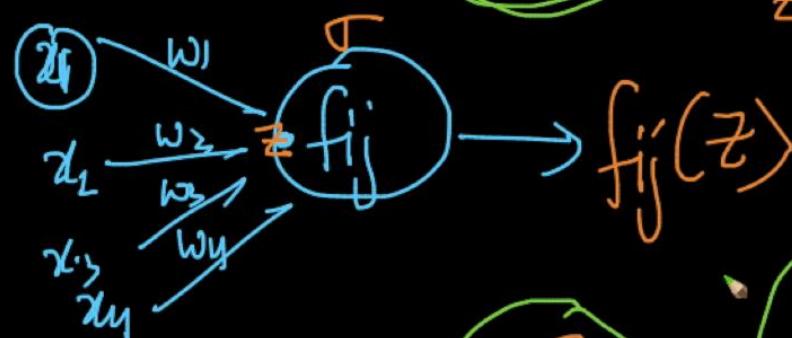
classical NN  $\xrightarrow{\text{2006}}$   $\xrightarrow{\text{2012}}$  DL

1980 & 90s

Sigmoid & tanh

$$z = \sum_i w_i x_i = \omega^T x$$

- ✓ differentiable
- ✓ easy to diff



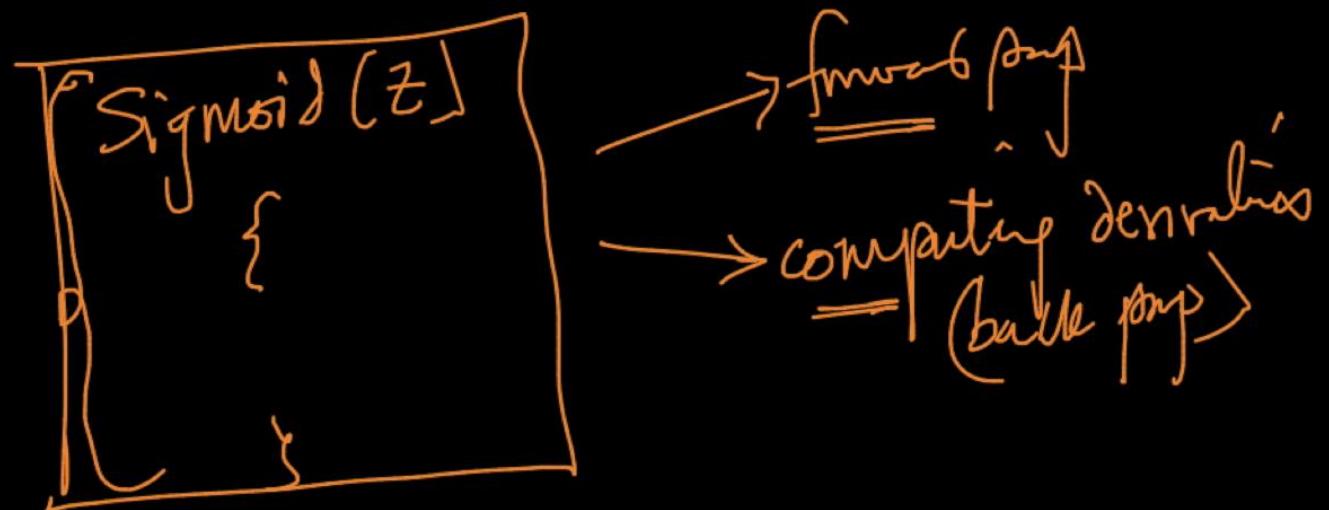
$$\underline{\sigma(z)} = \frac{e^z}{1+e^z} \text{ or } \frac{1}{1+e^{-z}}$$

$$\omega = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \end{bmatrix}$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}$$

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

✓  $\frac{\partial \sigma}{\partial z} = \sigma(z)(1-\underline{\sigma(z)})$



Contents - Google Docs calculus - Derivative of sigmoid Ronny Restrepo Ronny Restrepo Chekuri Srikan...

ronny.rest/blog/post\_2017\_08\_10\_sigmoid/

# Ronny Restrepo

Portfolio Blog Tutorials Contact

Comments

Note you can comment without any login by:

1. Typing your comment
2. Selecting "sign up with Disqus"
3. Then checking "I'd rather post as a guest"

3 Comments ronny.rest 1 Login

Recommend 11 Share Sort by Best

Join the discussion...

LOG IN WITH OR SIGN UP WITH DISQUS

D f G Name 556 / 551

APPLIED COURSE A

Contents - Google Docs calculus - Derivative of sigmoid Ronny Restrepo Ronny Restrepo Chekuri Srikan...

ronny.rest/blog/post\_2017\_08\_10\_sigmoid/

# Ronny Restrepo

Portfolio Blog Tutorials Contact

$$0 \leq \frac{dF}{dz} < 1$$

Comments

Note you can comment without any login by:

1. Typing your comment
2. Selecting "sign up with Disqus"
3. Then checking "I'd rather post as a guest"

3 Comments ronny.rest 1 Login

Recommend 11 Share Sort by Best

Join the discussion...

LOG IN WITH OR SIGN UP WITH DISQUS

D f G Name 556 / 551

APPLIED COURSE A

$$\tanh =$$
$$\frac{d \tanh}{dz} = 1 - \tanh^2(z)$$



$$\tanh =$$
$$\frac{d \tanh}{dz} = 1 - \tanh^2(z)$$

$\text{80's, 90's}$

more popular  $\rightarrow$  ReLU

Contents - Google Docs    calculus - Derivative of sigmoid    Ronny Restrepo    Ronny Restrepo    Chekuri Srikan...

ronny.rest/blog/post\_2017\_08\_16\_tanh/

# Ronny Restrepo

Portfolio - Blog Tutorials Contact

```
ax.xaxis.set_ticks_position('bottom')
ax.yaxis.set_ticks_position('left')

# Create and show plot
ax.plot(z,a, color="#307EC7", linewidth=3, label="tanh")
ax.plot(z,dz, color="#9621E2", linewidth=3, label="derivative")
ax.legend(loc="upper right", frameon=False)
fig.show()
```

optimization

① -1 to 1

②  $\frac{dz}{d\lambda} \leq 1$

## Comments

Note you can comment without any login by:

1. Typing your comment
2. Selecting "sign up with Disqus"
3. Then checking "I'd rather post as a guest"

0 Comments ronny.rest 13:43 558 / 551

1 Loading APPLIED COURSE A

{ Vanishing gradients }

$\rightarrow \{ \delta_{D^L}, \delta_D^L, \delta_D^L \rightarrow NN \}$

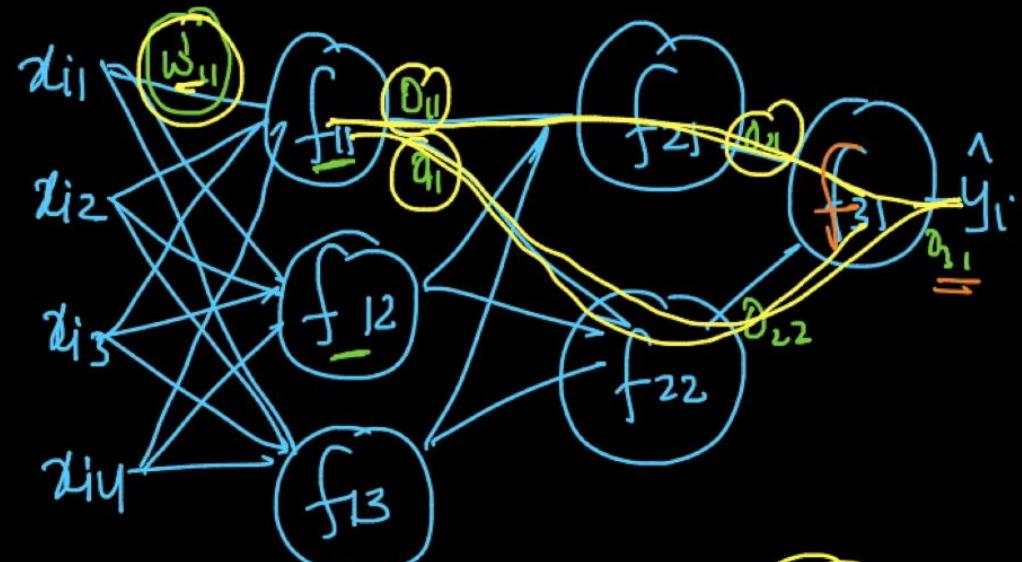
$f_{ij} \rightarrow \underline{\text{Sigmoid}}$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial \delta_{31}} \left[ \frac{\partial \delta_{31}}{\partial \delta_{21}} \cdot \frac{\partial \delta_{21}}{\partial \delta_{11}} \cdot \frac{\partial \delta_{11}}{\partial w_{11}} \right.$$

$$\left. + \frac{\partial \delta_{31}}{\partial \delta_{22}} \cdot \frac{\partial \delta_{22}}{\partial \delta_{11}} \cdot \frac{\partial \delta_{11}}{\partial w_{11}} \right]$$

$$\frac{\partial \delta_{31}}{\partial \delta_{21}} = \boxed{\frac{\partial f_{31}}{\partial \delta_{21}}}$$

$$\delta_{31} = f_{31}()$$



$$(\omega_{11})_{\text{new}} = (\omega_{11})_{\text{old}} - \eta \left( \frac{\partial L}{\partial \omega_{11}} \right)$$

03:35

41 / 42

A1

$$2.5 - 1 + 0.001 \\ = 2.5 - 0.001 \\ = 2.499$$

$0 \leq \frac{\partial f_3}{\partial w_{21}} < 1$

$$(w_{11})_{\text{new}} = (w_{11})_{\text{old}} - \eta \frac{\partial L}{\partial w_{11}}$$

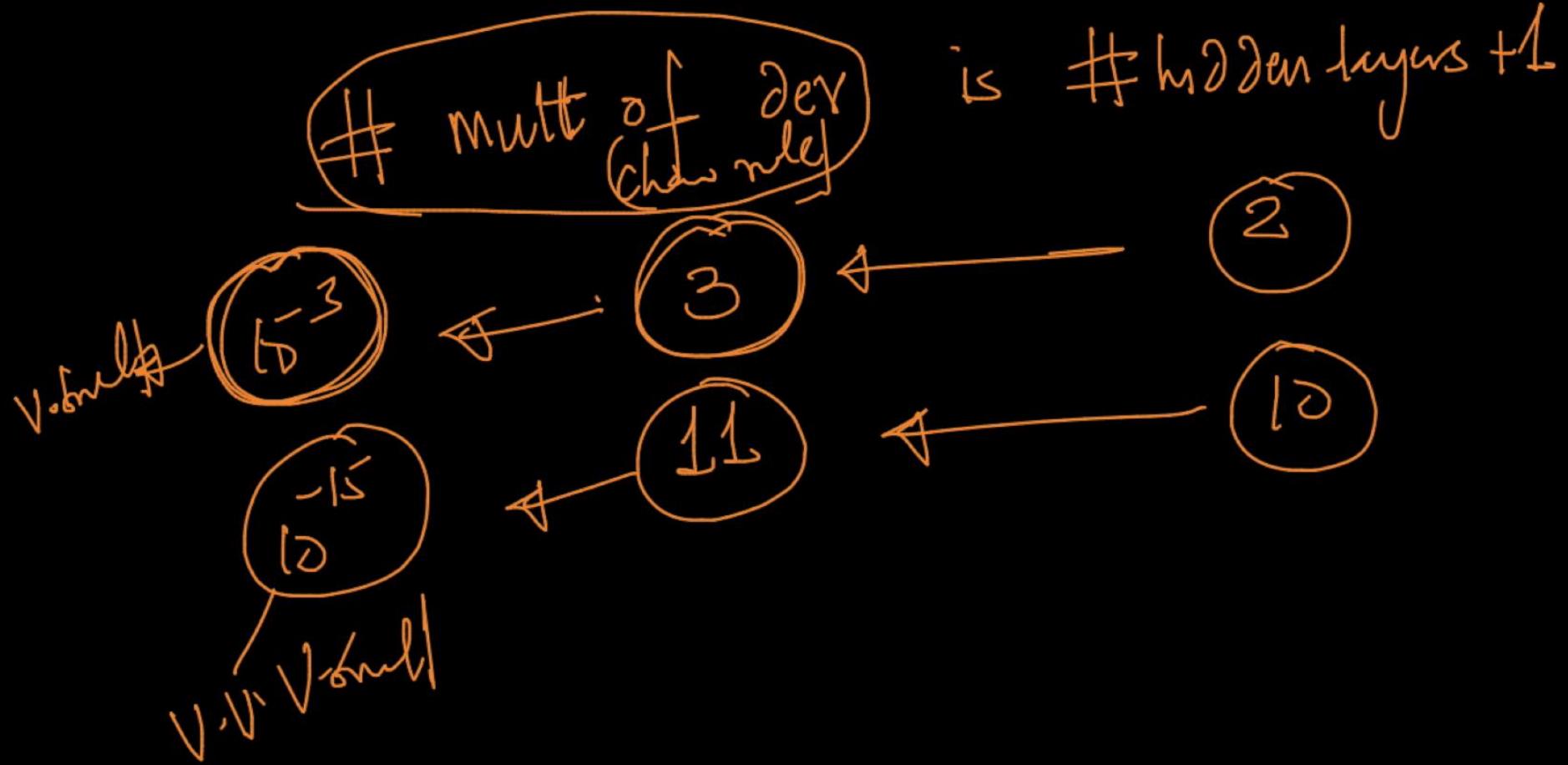
$\eta = 0.001$

✓

$\frac{\partial f_3}{\partial w_{21}} \cdot \frac{\partial f_2}{\partial w_{11}} \cdot \frac{\partial f_1}{\partial w_{11}}$

$0.2 * 0.1 * 0.05 = 0.0010$   $\approx 1D^{-3}$

$\ll 1$  V-level



$$2.5 - 1 + 0.001$$

$$= 2.5 - 0.001$$

$$0 \leq$$

$$\frac{\partial f_3}{\partial w_2} < 1$$

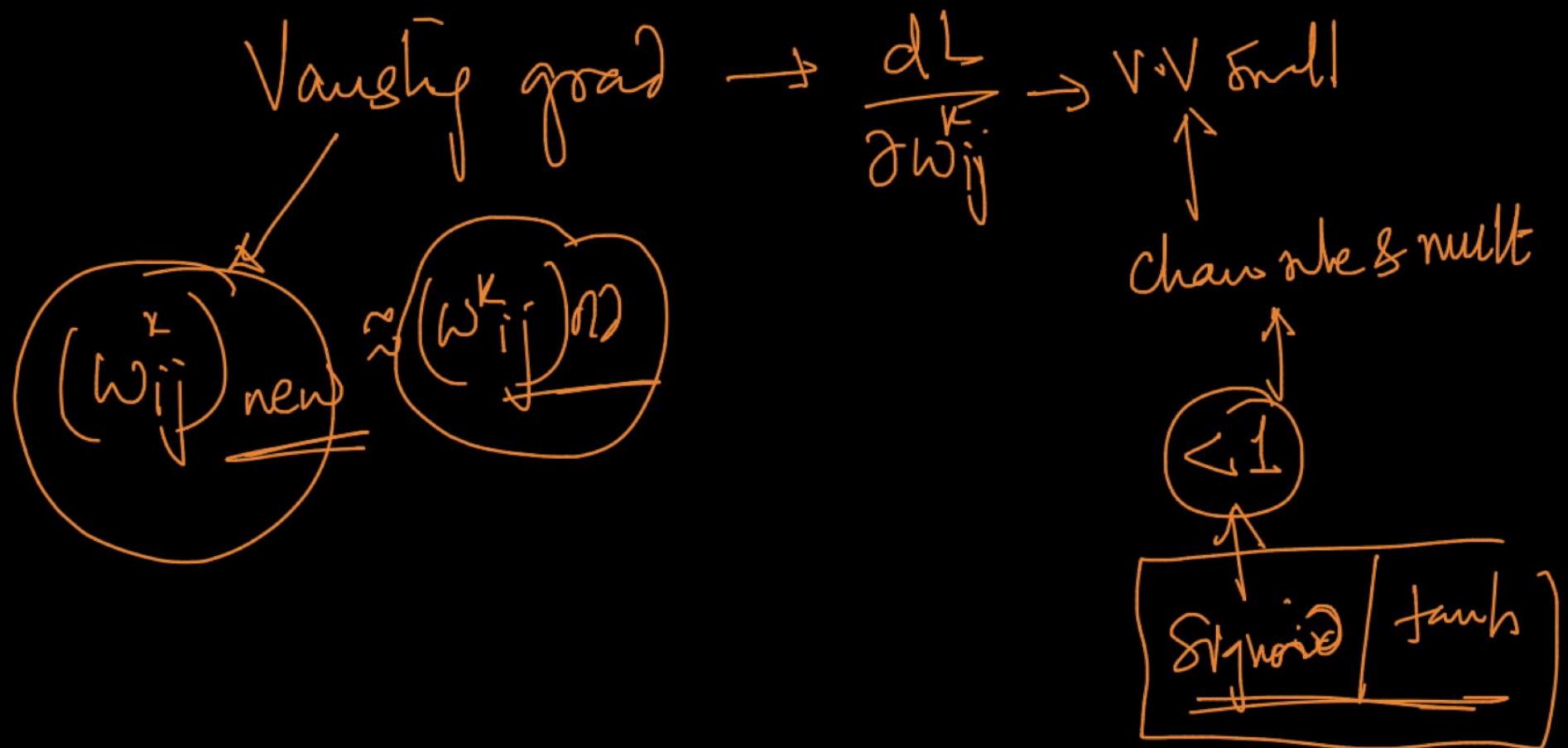
$$(w_{11})_{\text{new}} = (w_{11})_{\text{old}}$$

$$- \eta \frac{\partial L}{\partial w_{11}} = 0.001$$

$$V_{\text{final}} = 0.6 \text{ mV}$$

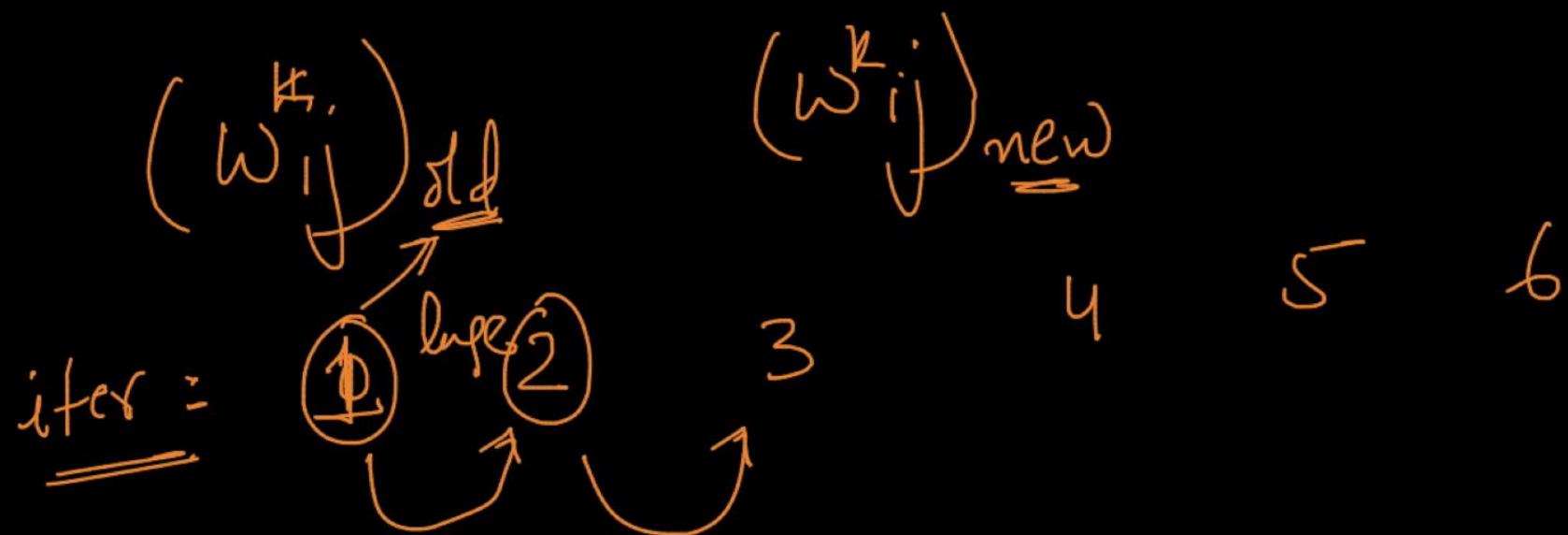
$$\begin{aligned} & \frac{\partial f_3}{\partial w_{21}}, \frac{\partial f_2}{\partial w_{11}}, \frac{\partial f_1}{\partial w_{11}} \\ & 0.2 * 0.1 + 0.05 = 0.0010 \end{aligned}$$

$$= 10^{-3} \text{ V-fall}$$



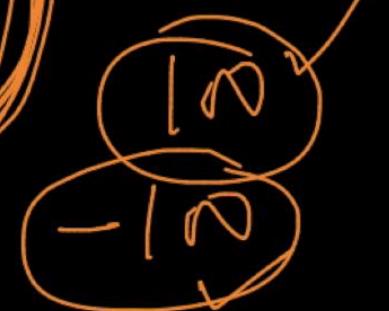
$$(\underline{w}_{ij}^k)_{\text{new}} = (\underline{w}_{ij}^k)_{\text{old}} - \gamma \frac{\partial L}{\partial w_{ij}^k}$$

W.V. tape



$$(\underline{w}_{ij}^k)_{\text{new}} = \underline{w}_{ij}^k - \gamma \frac{\partial L}{\partial w_{ij}^k}$$

V.V. tape



$w_{ij}^k$   $\delta L$

iter = 1 step 2

3

$$(\underline{w}_{ij}^k)_{\text{new}}$$

4

5 6

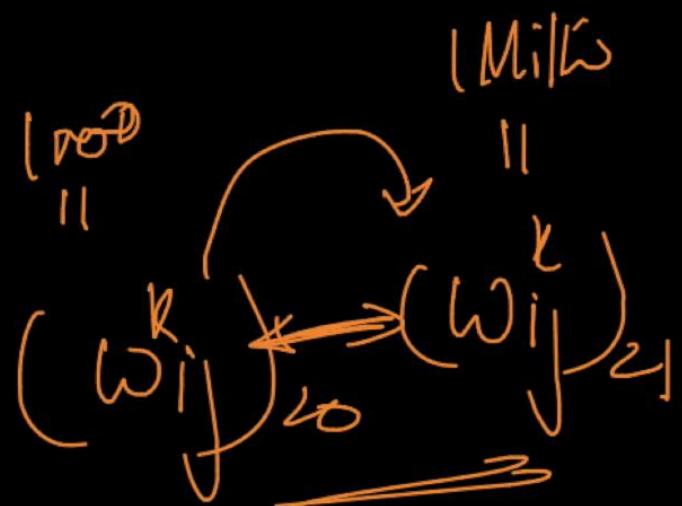


$$(\omega_{ij}^k)_1 = 2.2$$

$$(\omega_{ij}^k)_2 = 10.8$$

$$(\omega_{ij}^k)_3 = 100.6$$

⋮



Exploding grad!

chain rule: mult lots of derivatives

$$f_{ij} : \text{activation} =$$

$$\frac{\partial L}{\partial w_{11}} =$$

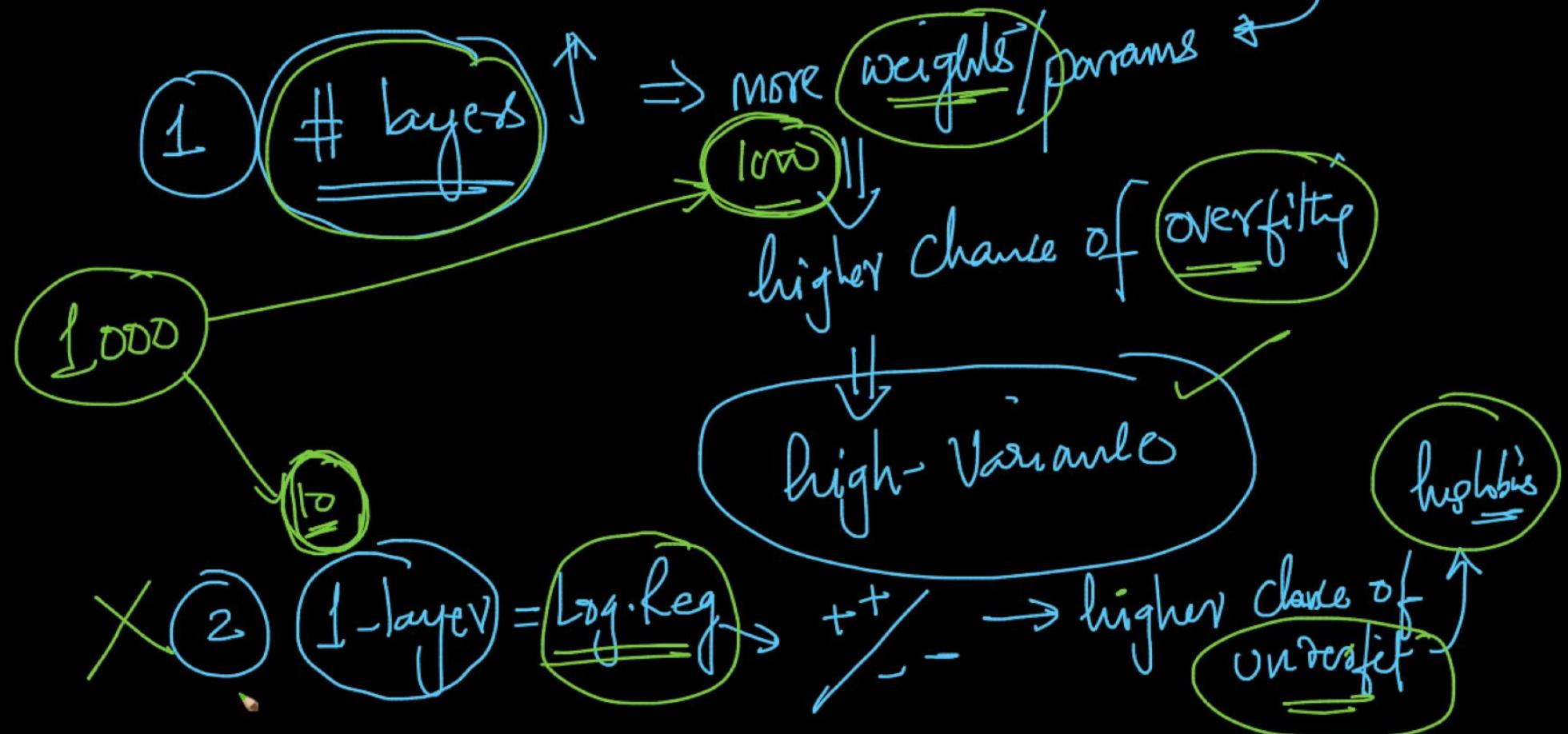
$$\frac{\partial \delta_{31}}{\partial o_{21}} \cdot \frac{\partial \delta_{21}}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial w_{11}}$$
$$\gamma_1 \quad \gamma_1 \quad \gamma_1$$
$$\delta = 2 \quad 2 \quad 2$$



## Bias - Variance Tradeoff:

NN

MLP



MLP

multiple layers

→ overfitting / high-Var

optimization

regularization

$$L(\hat{y}_i, y_i)$$

$$L = \sum_{i=1}^n \text{loss}(y_i, \hat{y}_i) + \sum_{i,j,k} (w_{ij}^k)^2$$

~~log-reg & L1-reg~~

$L_2\text{-reg}$

$L_1\text{-reg}$



$$L = \sum_{i=1}^n \text{loss}_i$$

reg in weights

larger  $\lambda \Rightarrow$  lesser overfit

$L_1$  reg  $\Rightarrow$  Sparsity  $\downarrow$  some  $w_{ij} = 0$

MLP Sparse

$$\mathcal{L} = \sum_i \text{loss}_i + \lambda \mathcal{L}_{\text{reg}}$$

MLP

1)  $\lambda$ :

hyperparam

$\lambda \uparrow \Rightarrow \text{Var} \downarrow$

2)

# layers: ↑

$\Rightarrow \text{Var} \uparrow$

Sigmoid