

Career outcomes



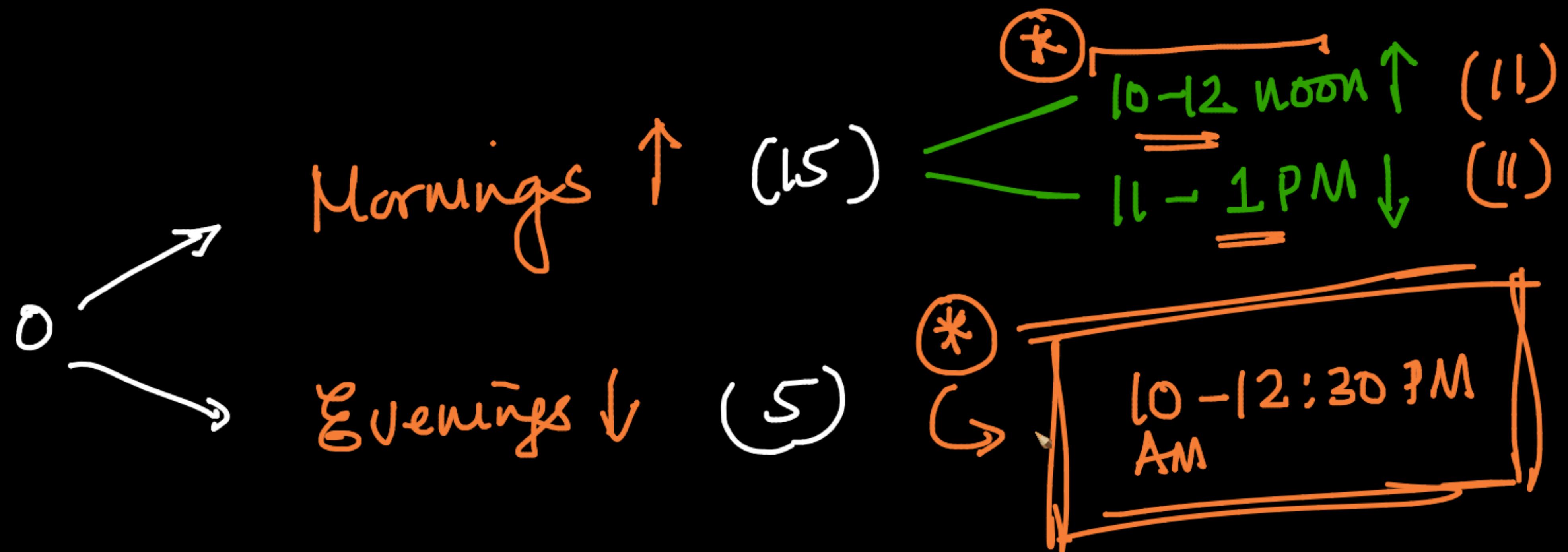
{ Sat & Sun
(recap + Interview Prep)
Sessions

MWF → DL

- { 1. Programming ✓
- 2. Prob + Stats ✓
- 3. Math for ML
- 4. Classification & Regression ML-1

{ 5. ML-2 (classical ML)
rest of
6. SQL

Suggestion : Hold off Case-study sessions





Please complete

mock-interviews



Careers

Practice

Tabular Data → some DL-NN

{ Image / Speech / Complex Time Series
RecSys / Text }

{ Neural Networks (NN) }

(2006-...)
(2012-...)

DL (NN)

[lots of research papers]

any new
model

Obj: equip you with strong foundations →
+ cover most widely used models
+ research papers

Assumption: classical ML + Math (recap)



Perceptron → 1957 (Rosenblatt)

✓ { NN loosely
inspired from
brains)

Back-propagation → 1980's (Geoff Hinton & others)

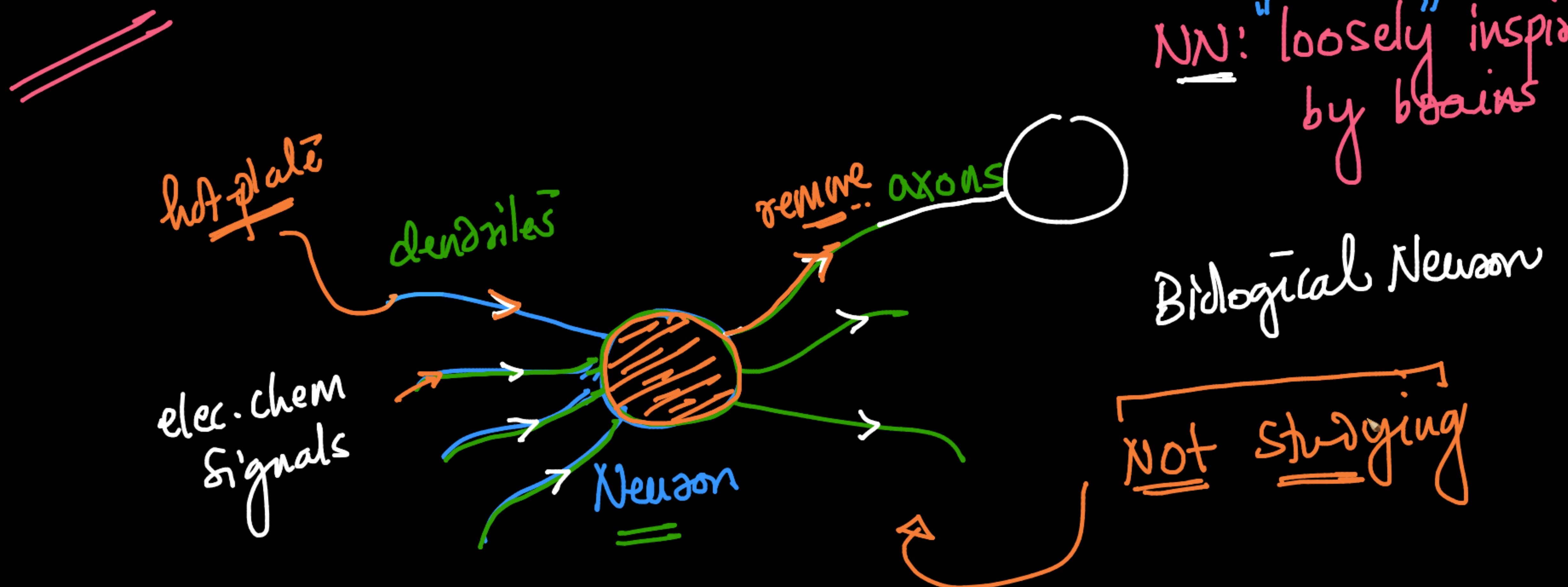
modern DL → 2006 ✓
AlexNet (CV) → 2012 - - - - ~2022

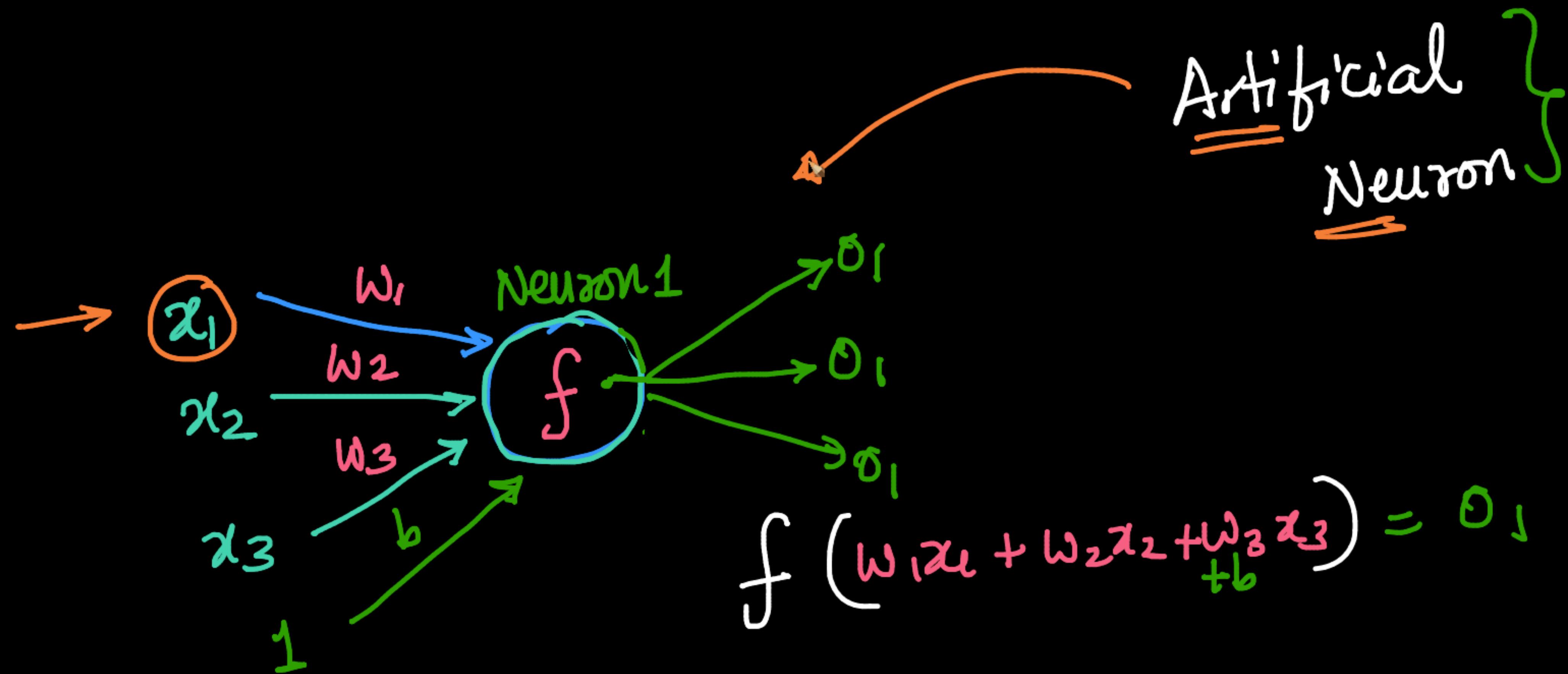
✓ TensorFlow 2 & Keras
=====

Google - - -
;

Py-Torch
=====

FB
openAI
Tesla





f : activation fn

Terminology:

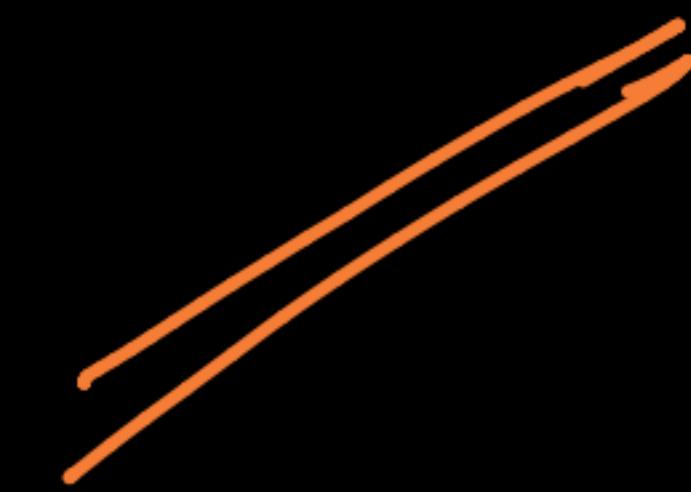
$$o_i = f \left(\sum_{j=1}^d x_{ij} w_j \right)$$

Diagram illustrating the components of the equation:

- $x_i \in \mathbb{R}^d$ (Input vector)
- w_i (Weights)
- f (activation function)
- x_{ij} (Inputs)

Annotations:

- x_i is labeled as $\underline{x_i} \in \mathbb{R}^d$
- w_i is labeled as $\underline{w_i}$
- f is labeled as \underline{f}
- x_{ij} is labeled as $\underline{x_{ij}}$
- w_j is labeled as $\underline{w_j}$

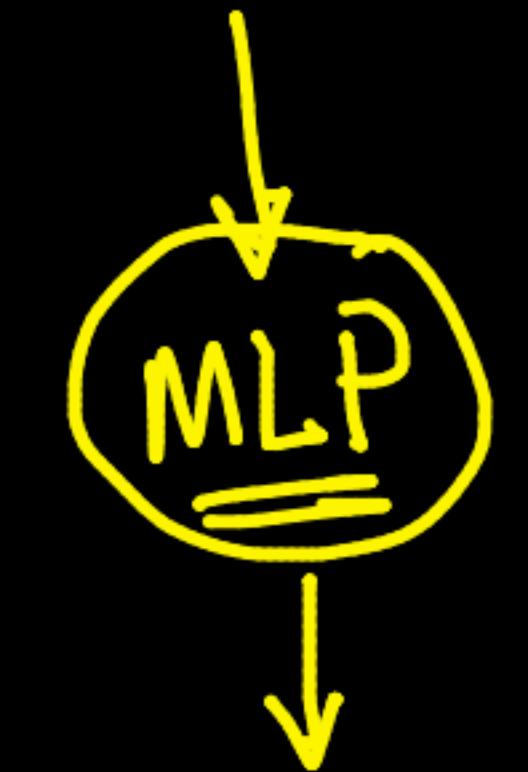
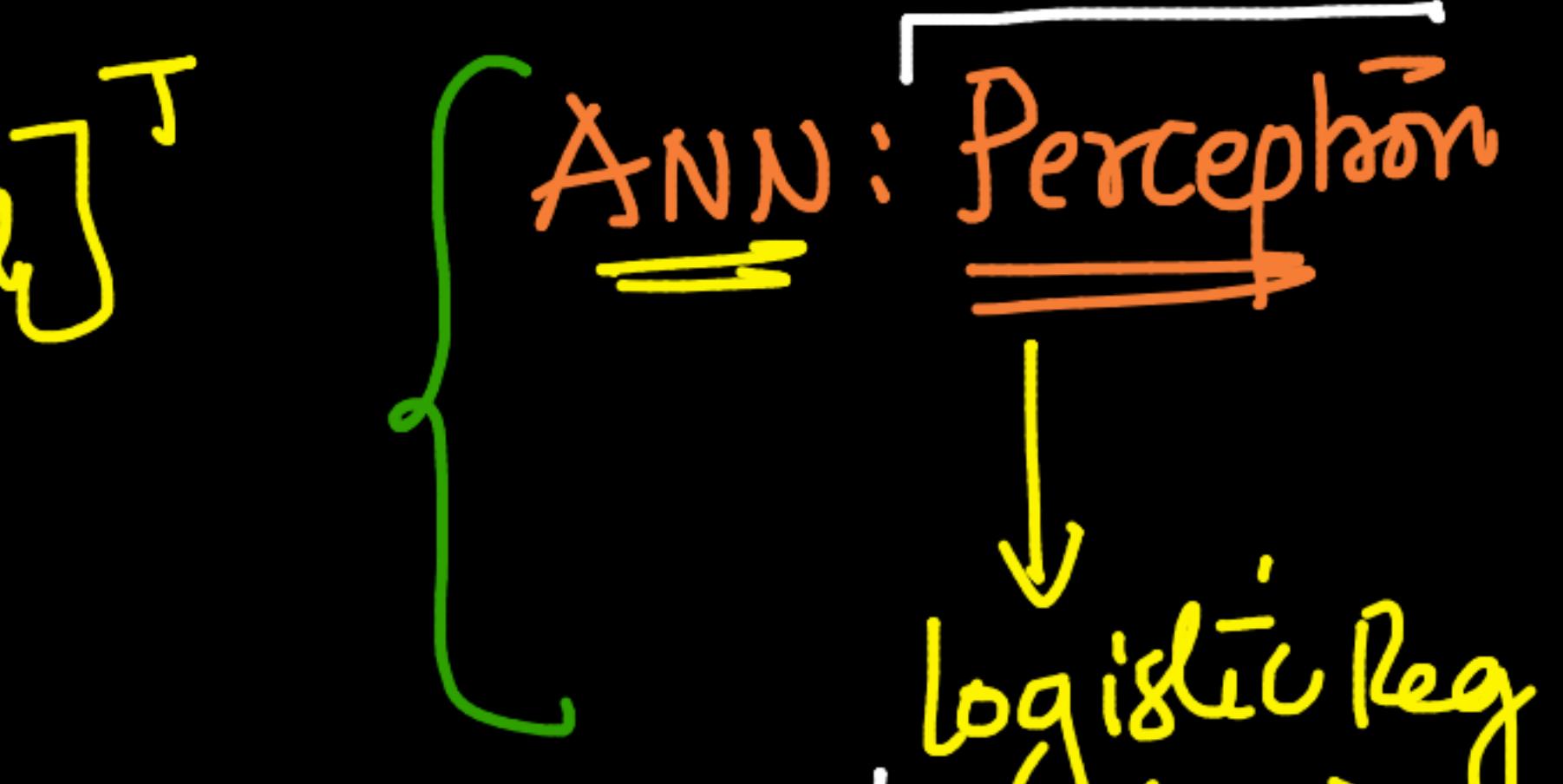


Logistic Regr!

$$\mathcal{D}_{\text{Tr}} = \left\{ \left(\underline{x}_i \in \mathbb{R}^d, y_i \right) \right\}_{i=1}^n$$

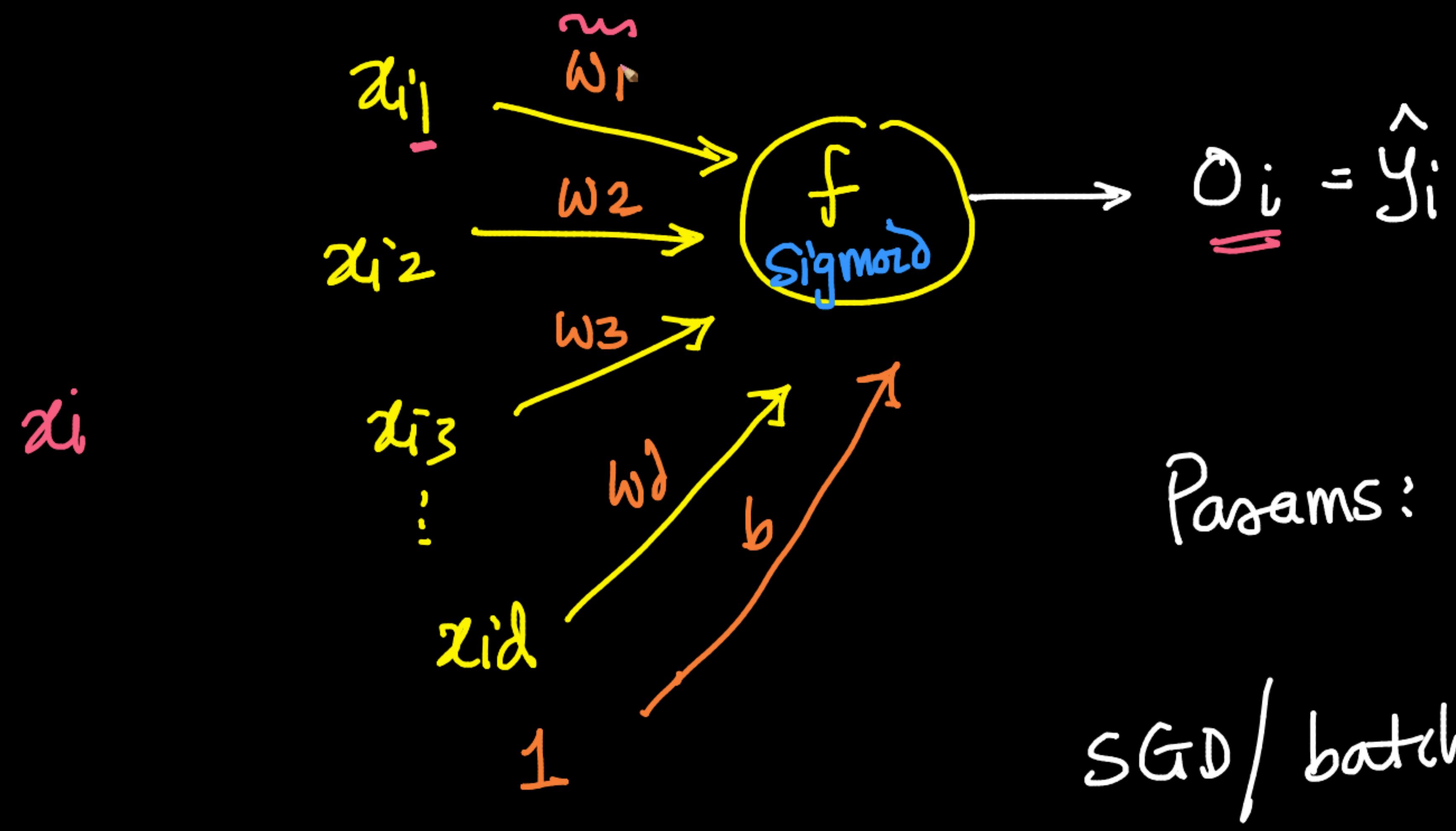
$$\hat{y}_i = \text{Sigmoid} \left(\omega^T \underline{x}_i + b \right)$$

↓ ↓
d-dim scalar

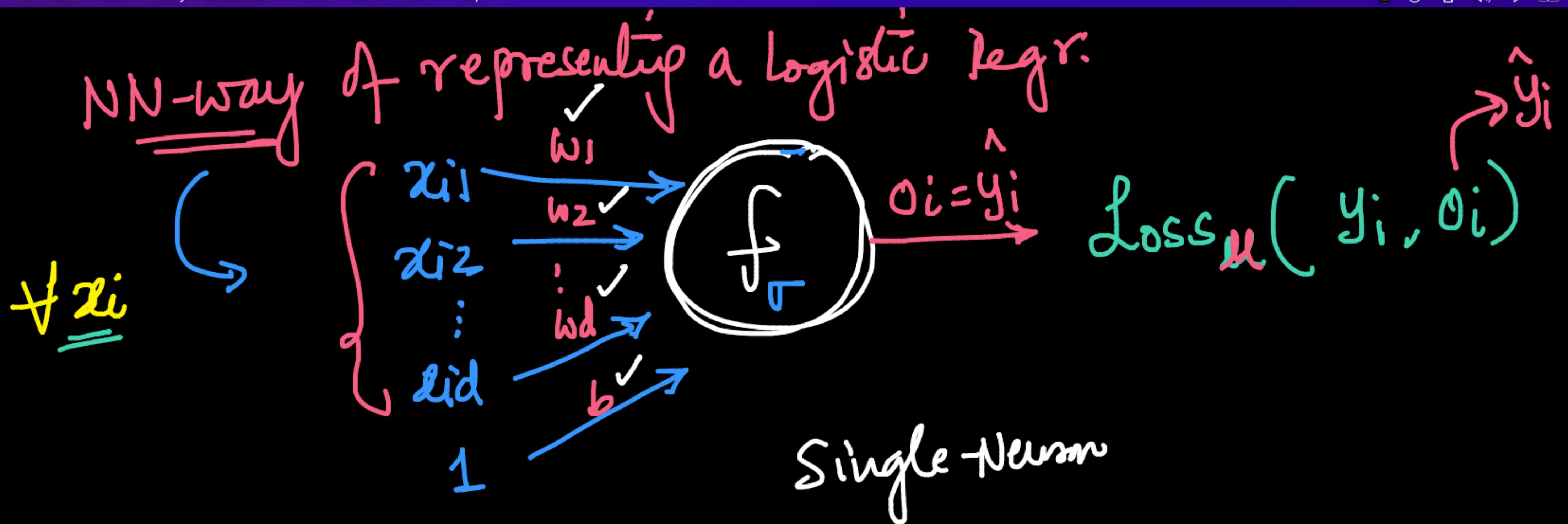


back prop

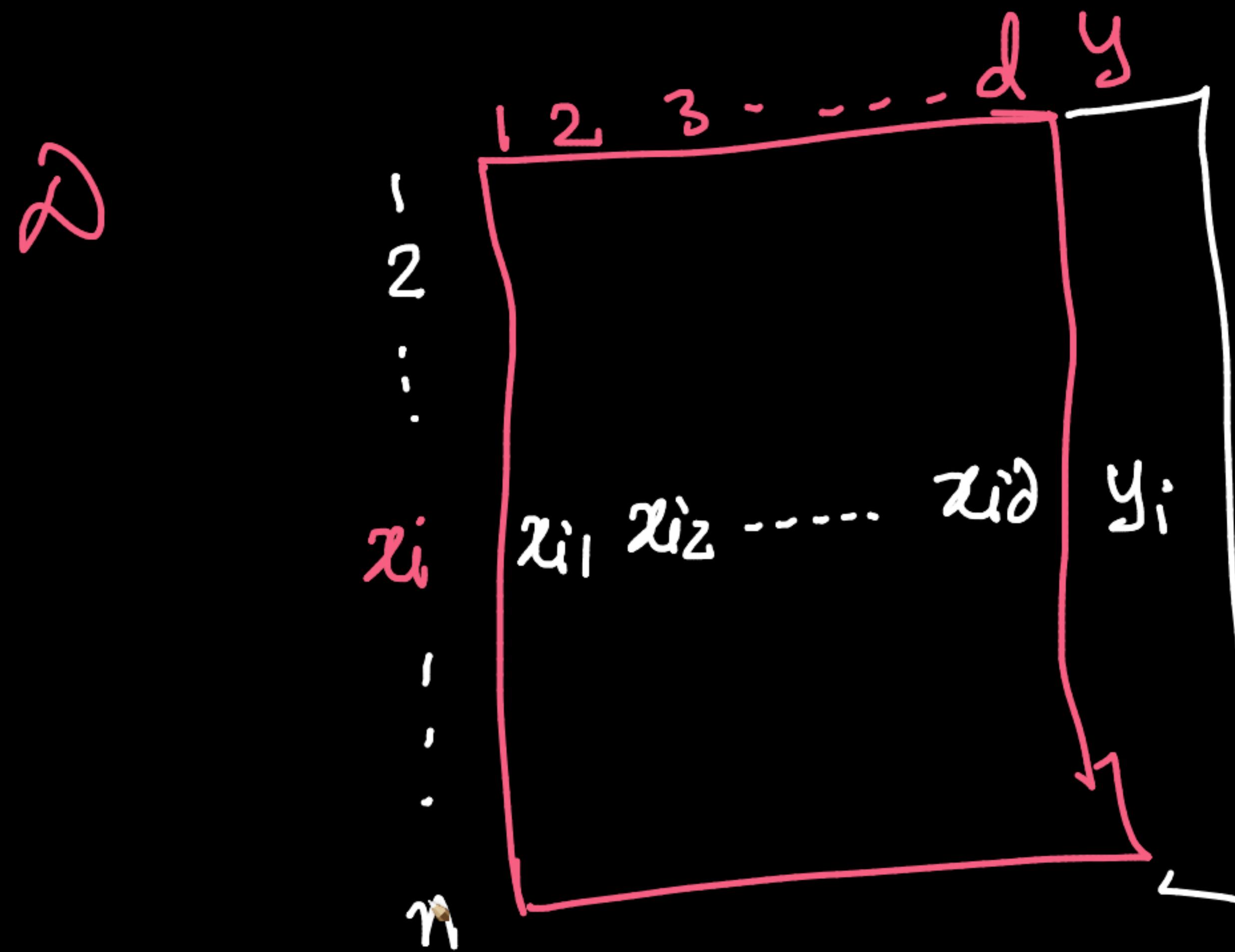
Code

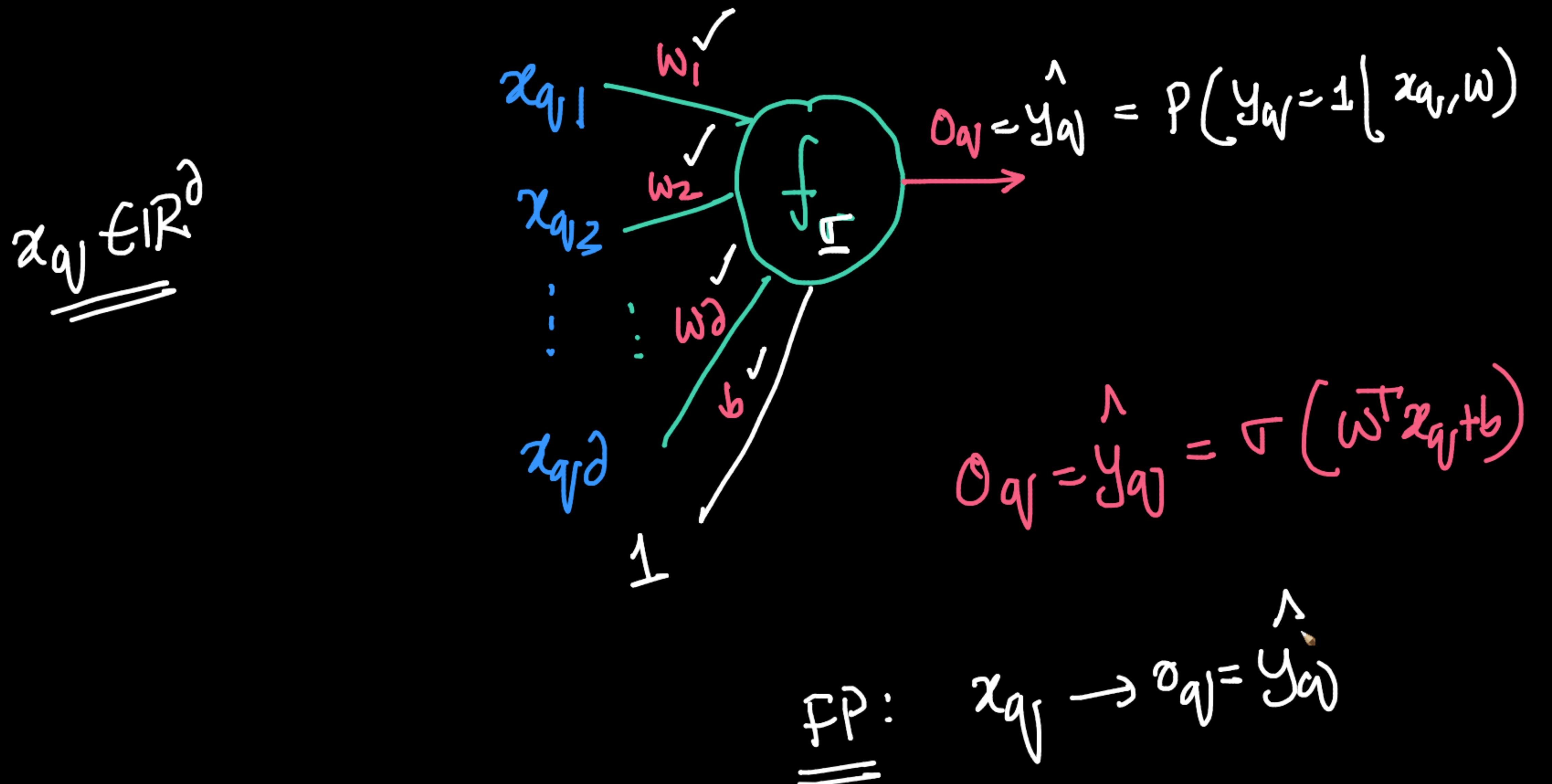


SGD / batch GD



forward propagation : $f_\Gamma\left(\sum_{j=1}^d w_j x_{ij} + b\right) = f_\sigma(\underline{w^T x + b})$

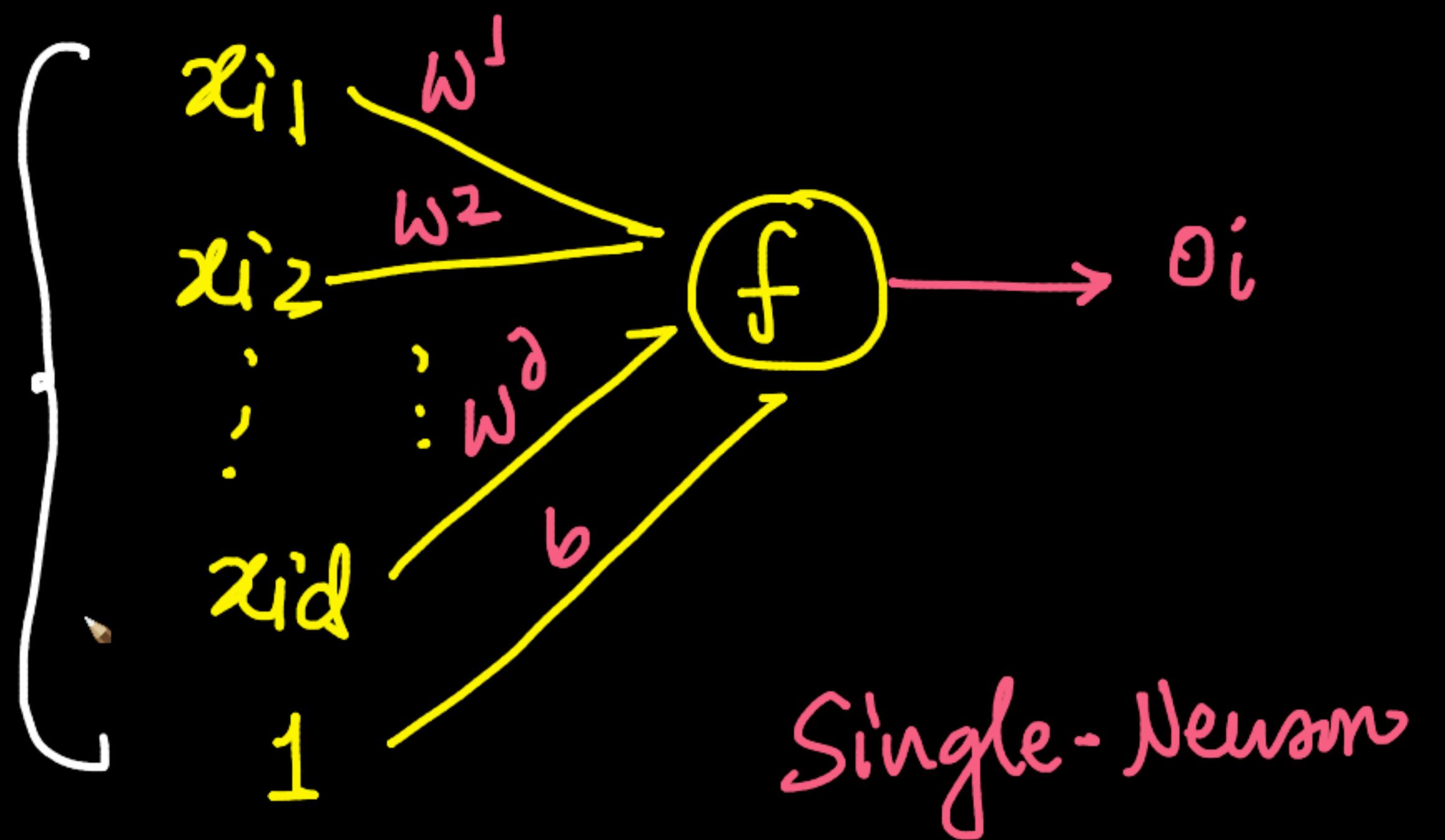




Perceptron

$$x_i \in \mathbb{R}^D, y_i \in \{0, 1\}$$

$$\begin{aligned} o_i &= f_{\text{perception}}(x_i, w, b) \\ &= \begin{cases} 1 & \text{if } w^T x_i + b > 0 \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

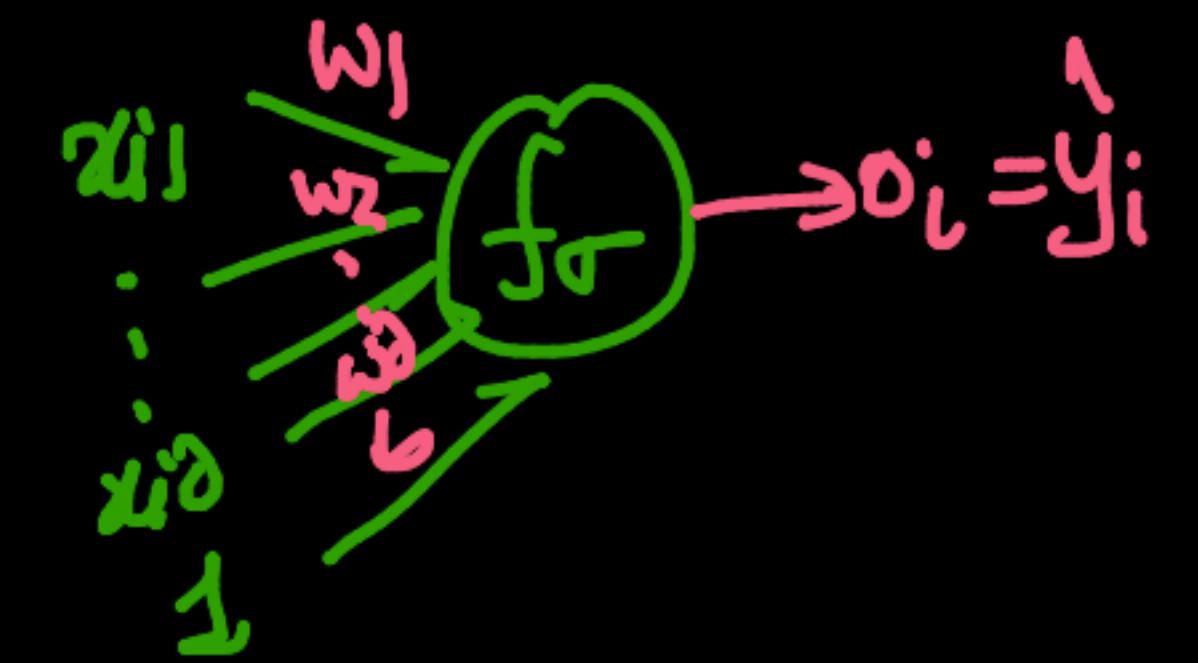


Eqns \rightarrow Geom \rightarrow NN-diagram

logistic reg:-

$$\sigma(\omega^T x_i + b) = \hat{y}_i$$

π^θ ; squashing

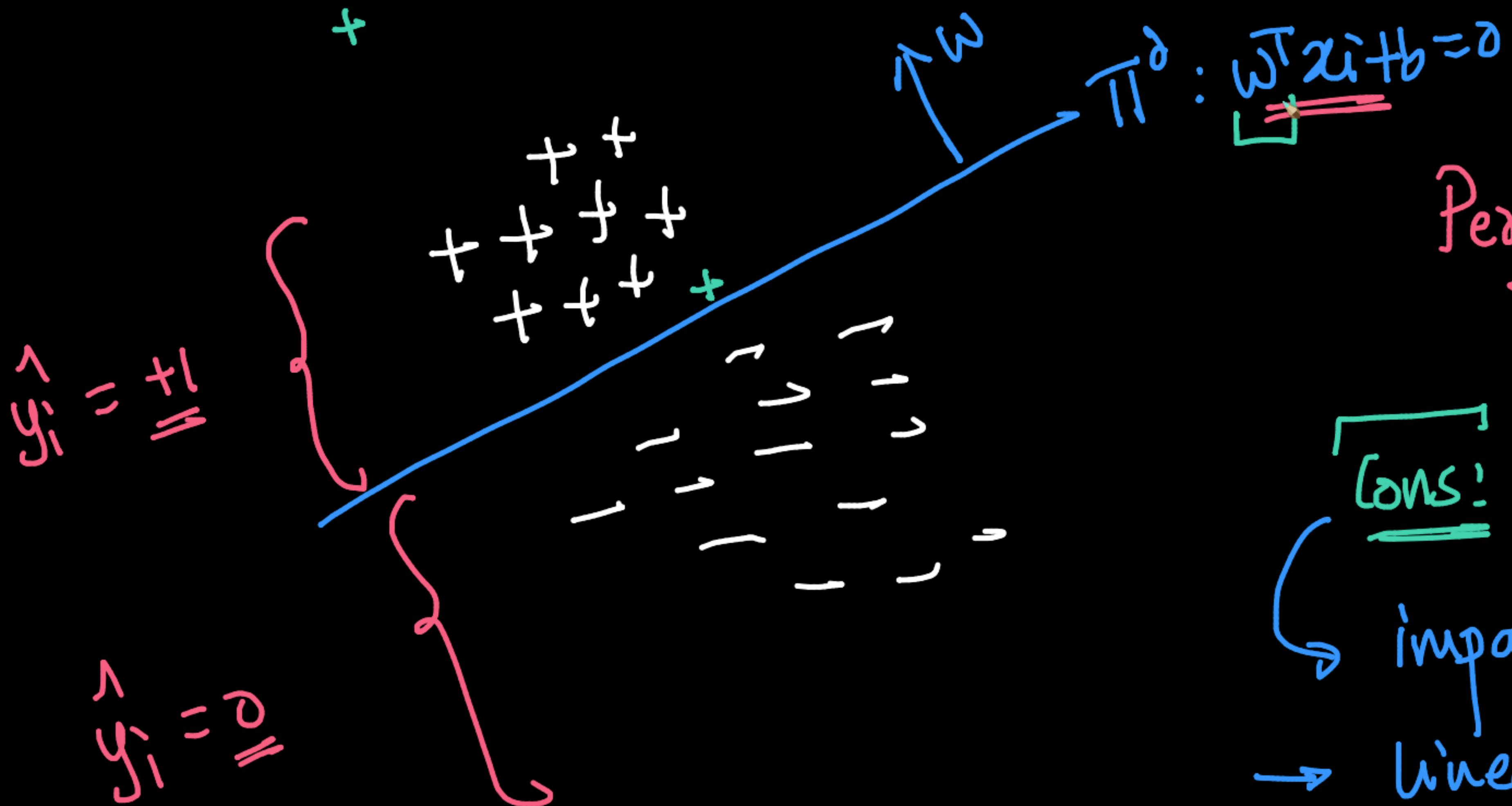


Perceptron :-

$$f_{\text{perceptron}}(\omega^T x_i + b) = \hat{y}_i$$

$\underline{\pi^\theta}$; no-squashing

$$f(x, \underline{\omega}, b) = \begin{cases} 1; & \omega^T x_i + b \geq 0 \\ 0 & \text{otherwise} \end{cases}$$



Perception 1957

cons:

impact of outliers

→ linear model

→ no probabilities

→

TrainPerception

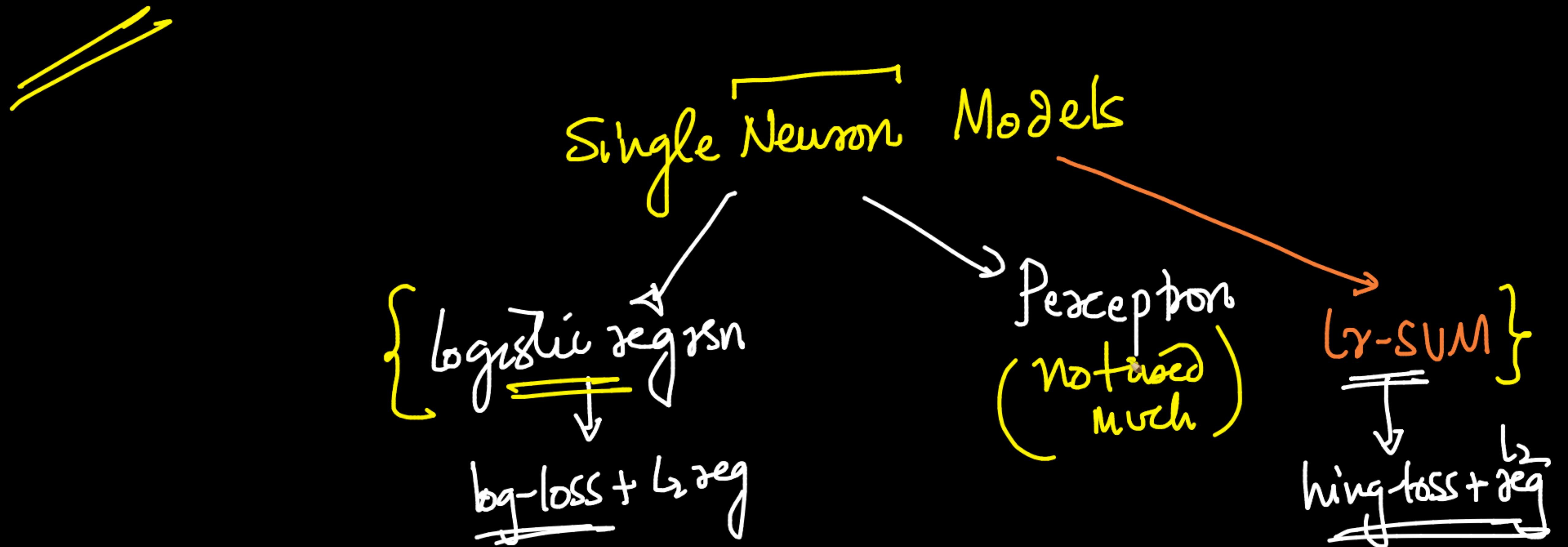
$$x_i \rightarrow y_i \xrightarrow{f_p} 0 \rightarrow 1$$

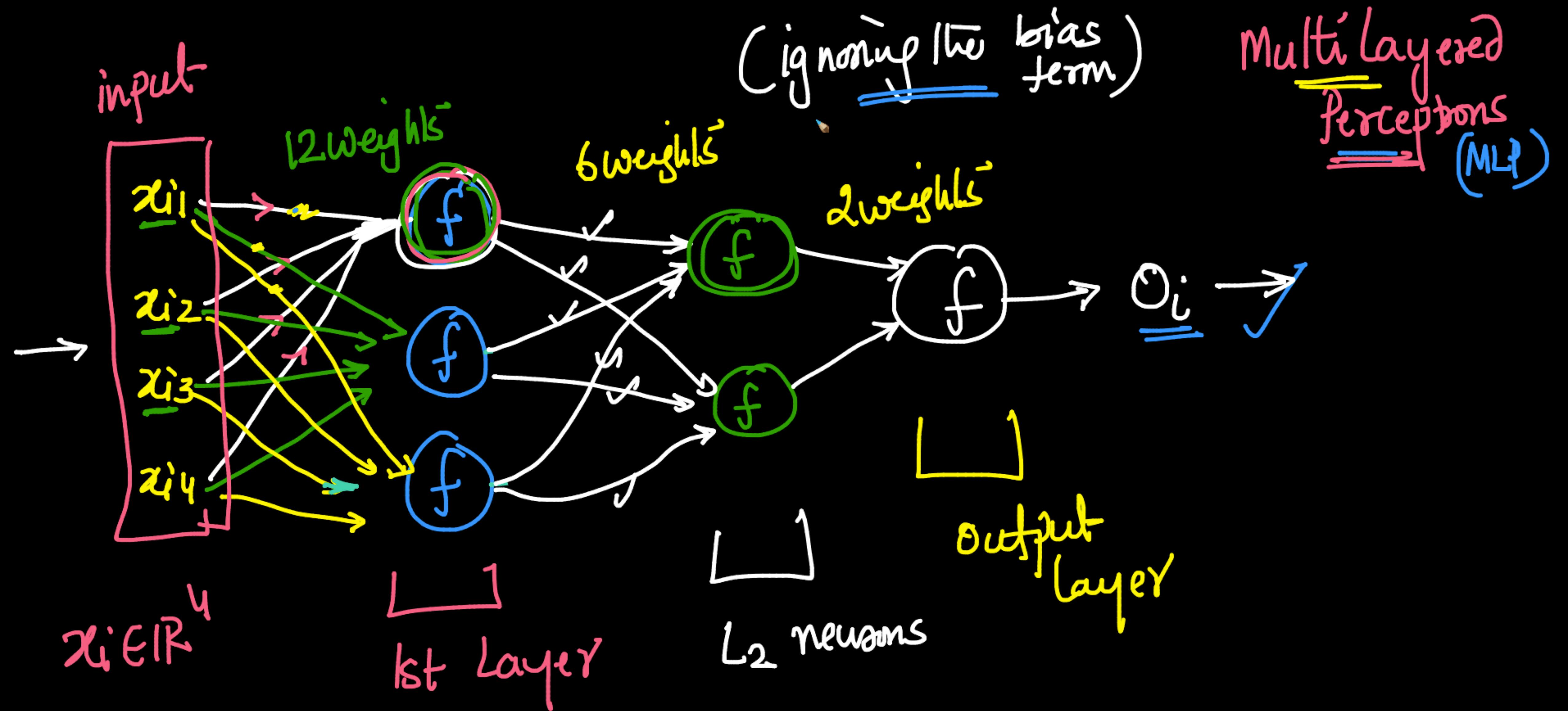
$$\hat{y}_i = f_p(w^T x_i + b) \xrightarrow{0} 1$$

GD

$$\min_{w, b} \sum_{i=1}^n \text{Loss}(y_i, \hat{y}_i) + \lambda \underbrace{\|w\|_2}_{L_2-\text{reg}}(w)$$

even use MSE







Why do MLP?

====

(Algebra)

fog got

$f_1: +$

$f_3: \sin$

$f_2: -$

$f_4: \cos$

$f_5: *$

$f_1(x_{i1}, x_{i2})$

$f_6: /$

$= x_{i1} + x_{i2}$

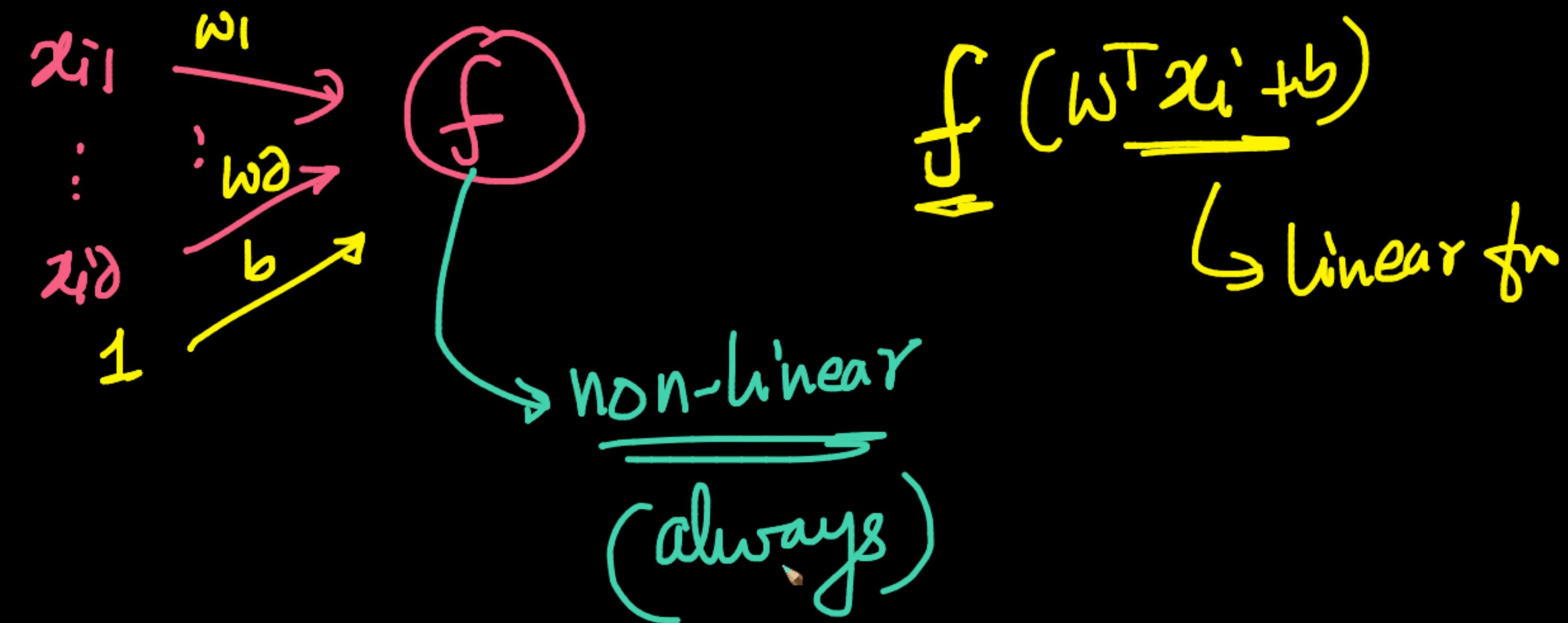
$f_7: \text{sqrt}$

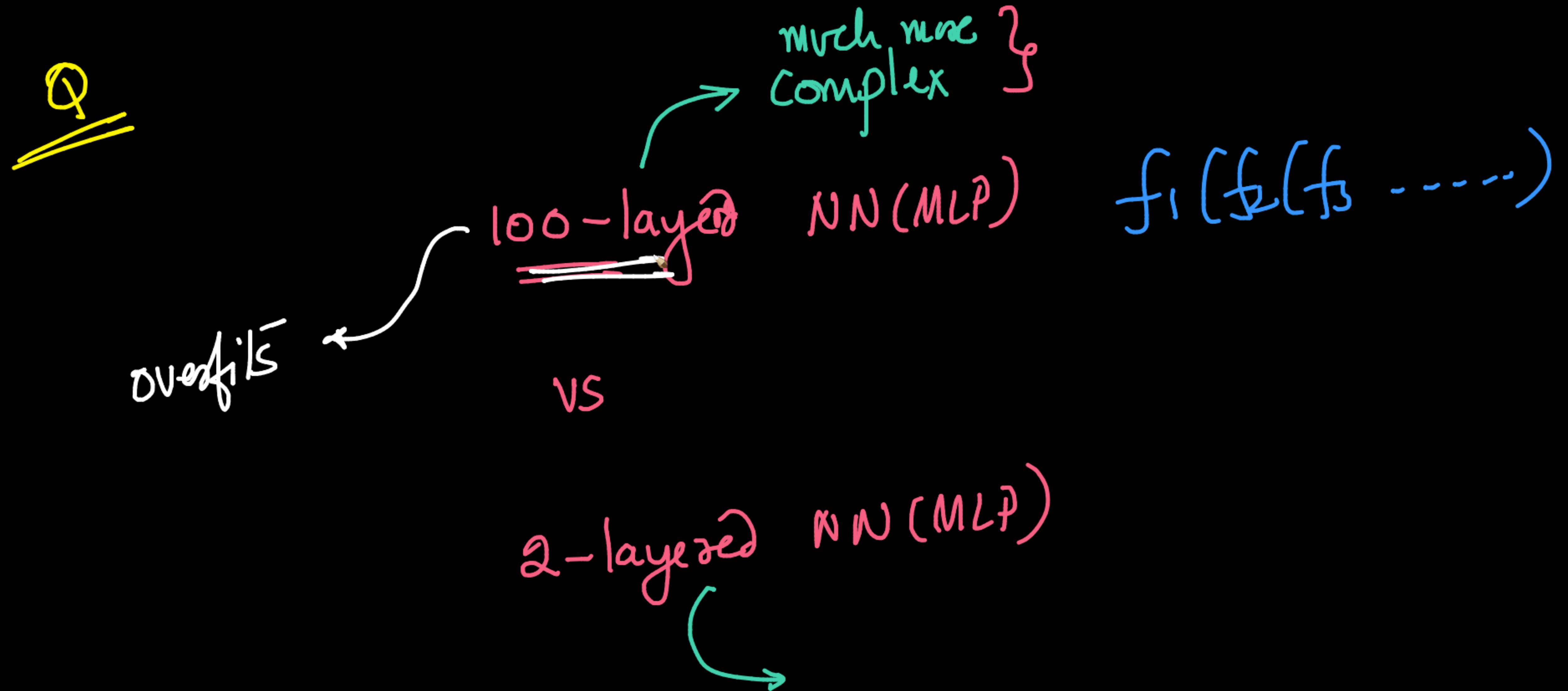
$f_6(f_1(x_{i1}, x_{i2})) = \sin(x_{i1} + x_{i2})$

Composition of fn

Construct complex fns as gof or fog over & over again.

$$fog(x) = f(g(x))$$
$$gaf(x) = g(f(x))$$

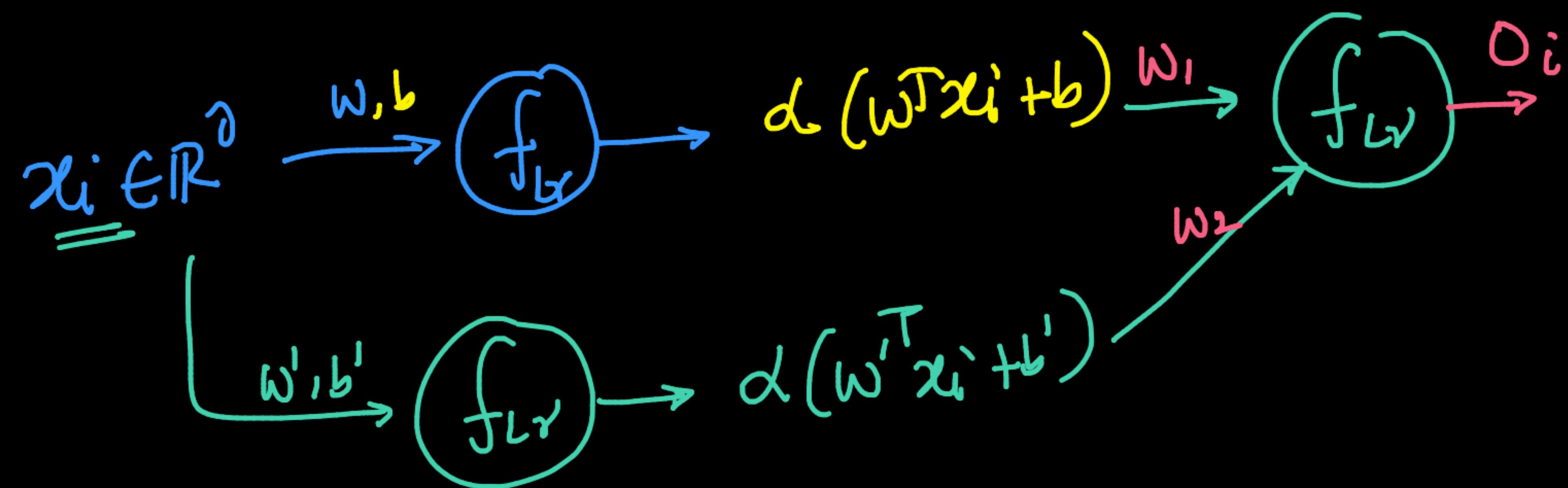






what happens if each actn. fn in a
MLP is a linear-fn-

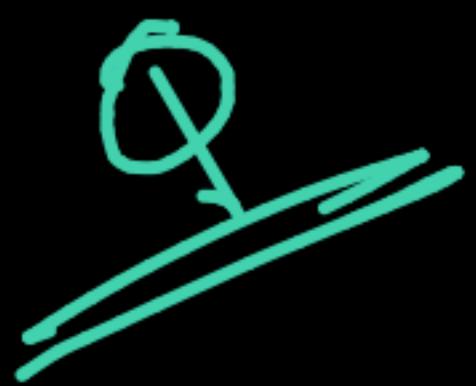
$$f(x_i, w, b) = f_{Lr}(w^T x_i + b) = d(w^T x_i + b) + \beta$$



$$\begin{aligned}
 o_i &= \alpha(\omega_1(w^T x_i + b) + \omega_2(w'^T x_i + b')) \\
 &= \underline{\alpha \omega_1 (w^T x_i + b)} + \underline{\alpha \omega_2 (w'^T x_i + b')}
 \end{aligned}$$

if all of our action-fns are linear
then $\text{MLP} = \text{linear-model}$

Hence always use non-linear action
fns ...
in DL & NN



different

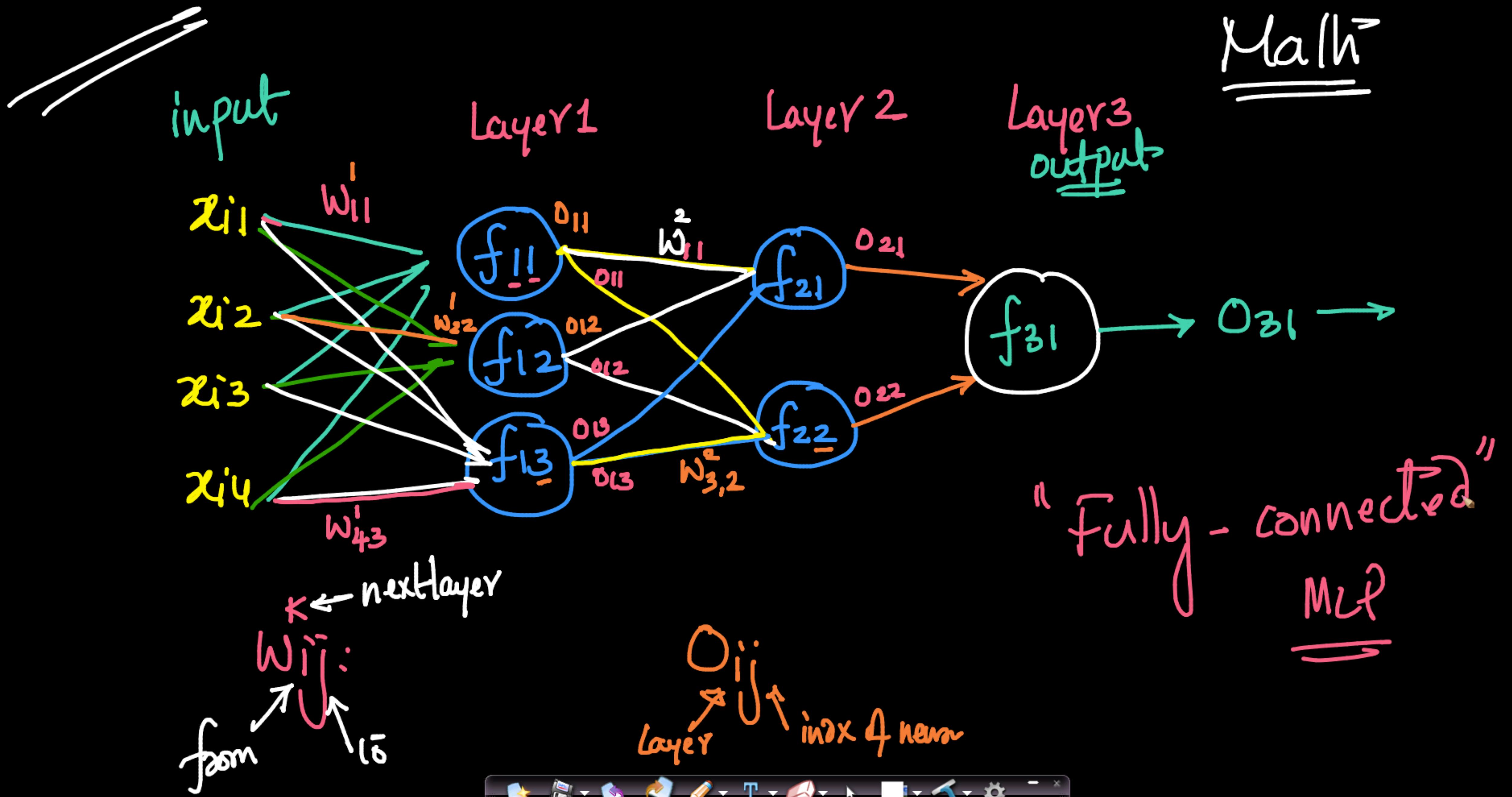
activn fns in various layers

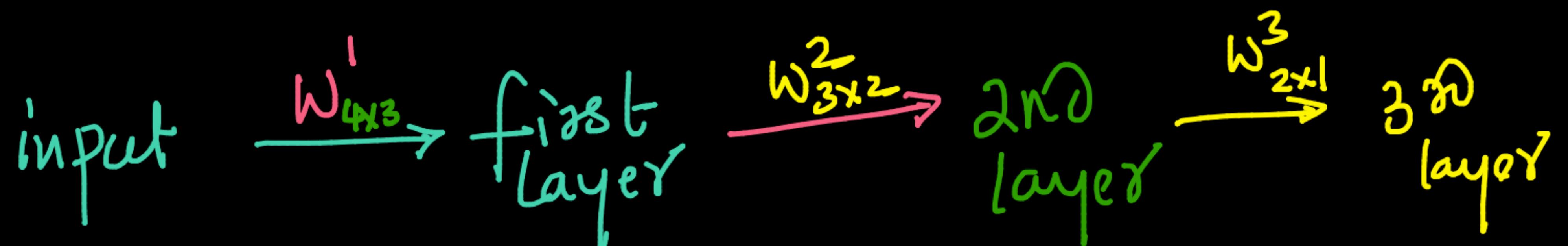


You-can (theoretically)

Practically: not much incremental
value

e.g.: ~~Sigmoid~~ → 1980's & 90's



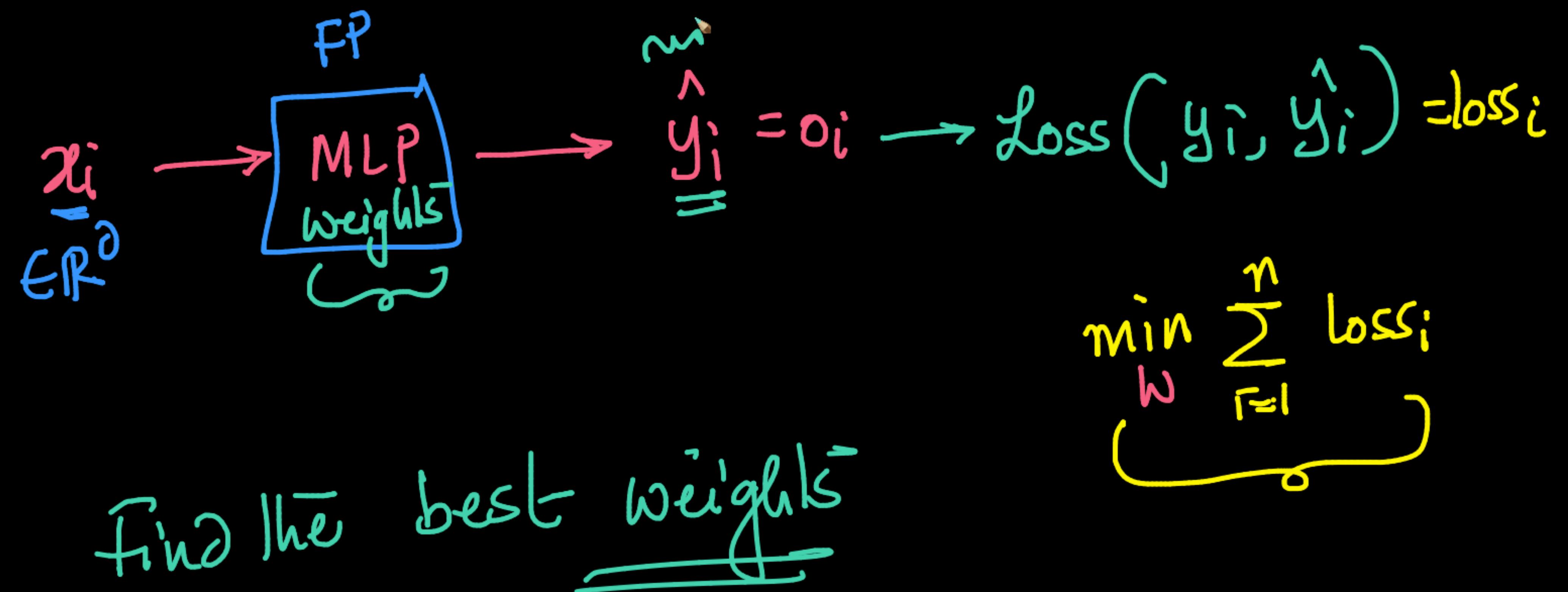


$$W^1 = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ \vdots & & \vdots \\ w_{41} & w_{42} & w_{43} \end{bmatrix}_{4 \times 3} \rightarrow \text{Matrices}$$

Classfn:

$$\text{Data} \rightsquigarrow \{(x_i, y_i)\}_{i=1}^n$$

Train an MLP

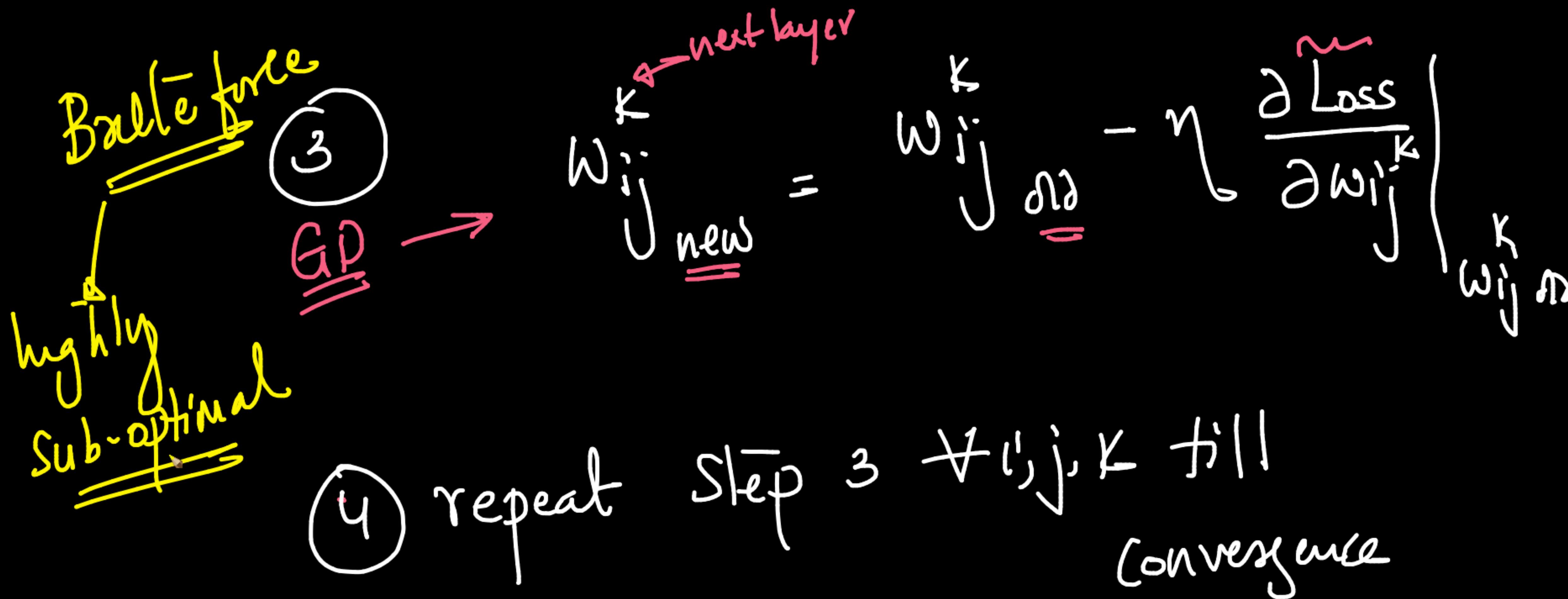




$$\min_w \sum_{i=1}^n \underline{\underline{\text{loss}_i}}$$

Loss ✓

{ all activation &
loss fn have
to be
differentiable



Back-prop → Math
Back-prop → CS
Back-prop → code [from scratch]

Q
/ \

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

$n = 1 \text{Million}$

GD:

$$\frac{\partial \text{Loss}}{\partial w}$$

all the n datapts

SGD:

$$\frac{\partial \text{Loss}}{\partial w}$$

using

one random pt

Most widely used

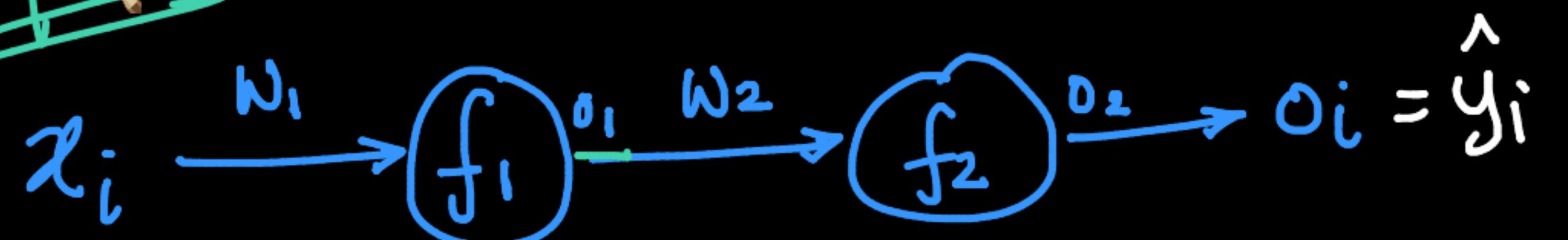
{ minibatch-SGD:
(k)

$$\frac{\partial L}{\partial w}$$

: random batch of
k-pts

Q

Toy-ex:

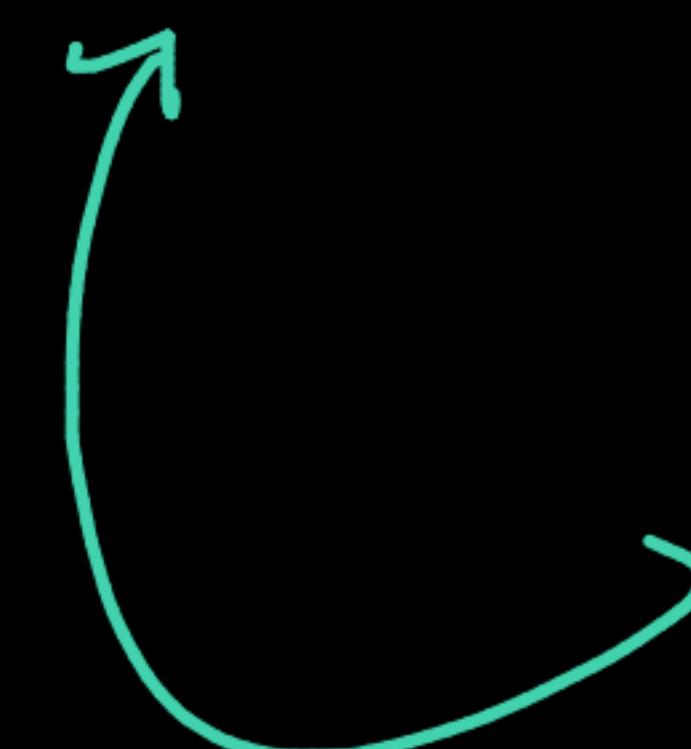


$$\text{loss} = \sum_{i=1}^n \text{loss}_i$$

$$\text{loss}_i = (y_i - \hat{y}_i)^2$$

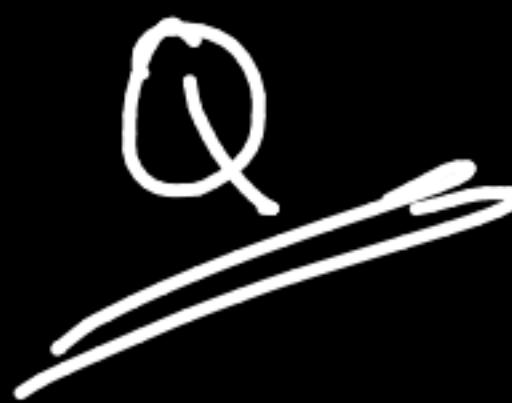
$$\frac{\partial \text{Loss}}{\partial w_1} = \sum_{i=1}^n \frac{\partial \text{loss}_i}{\partial w_1}$$

$$\frac{\partial \text{Loss}_i}{\partial w_1} = \frac{\partial}{\partial w_1} \left(y_i - \hat{y}_i - \left\{ f_1(x_i w_1) \times w_2 \right\} \right)^2$$

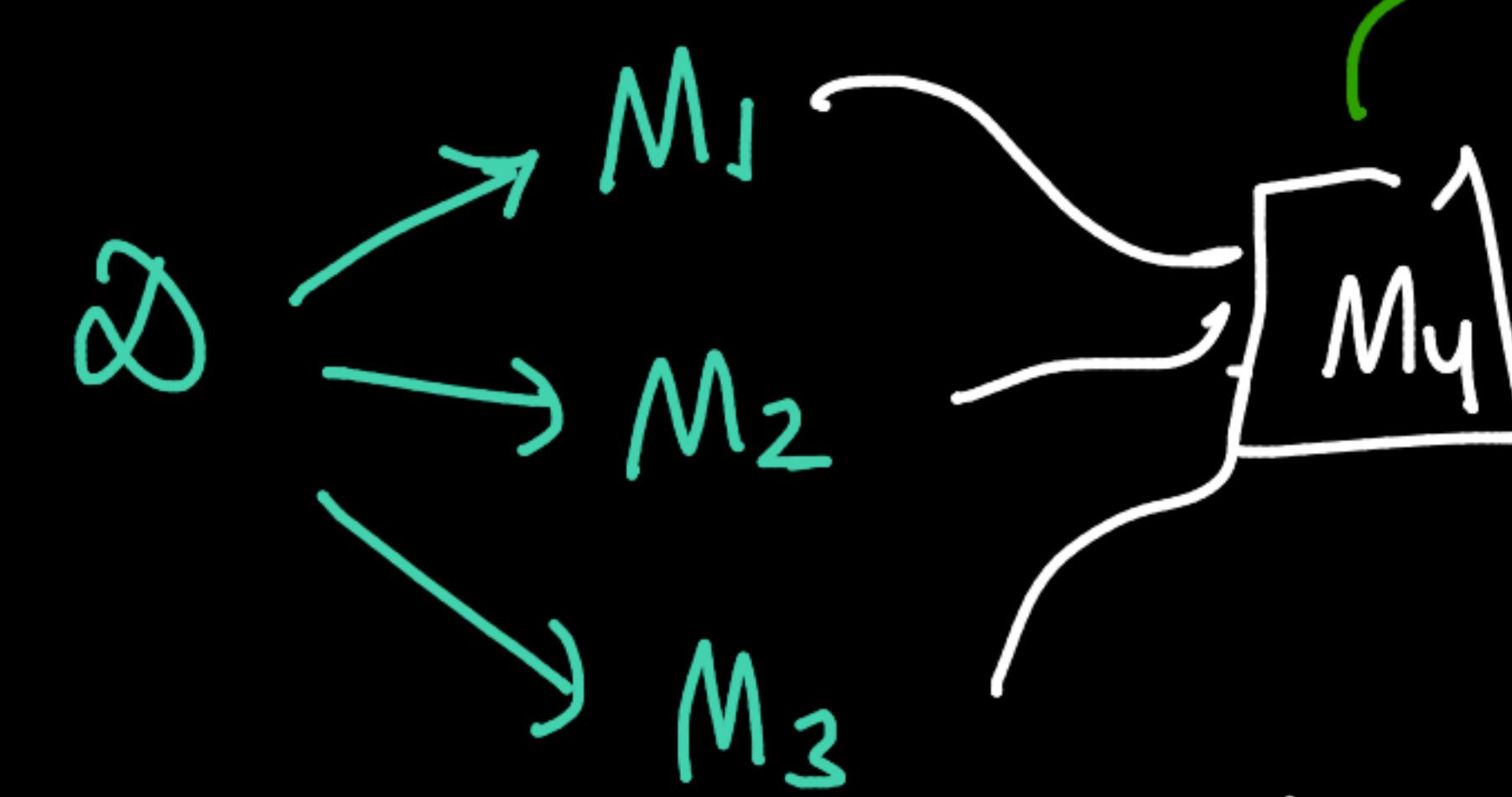


f_1, f_2 are differentiable

Chain-rule

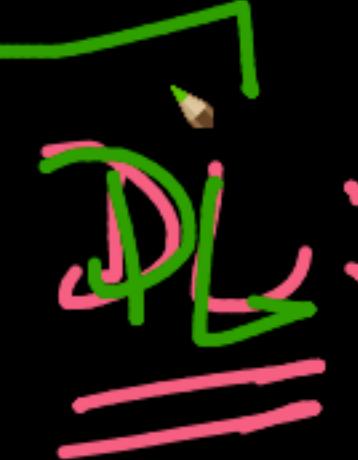
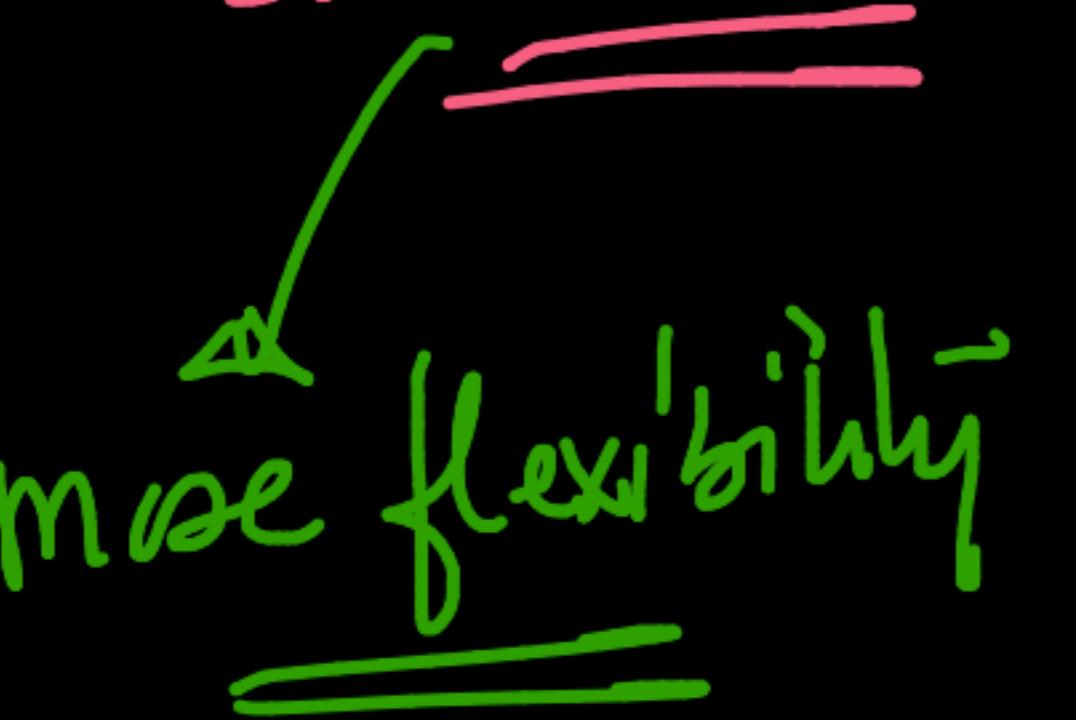


X { Stacking:



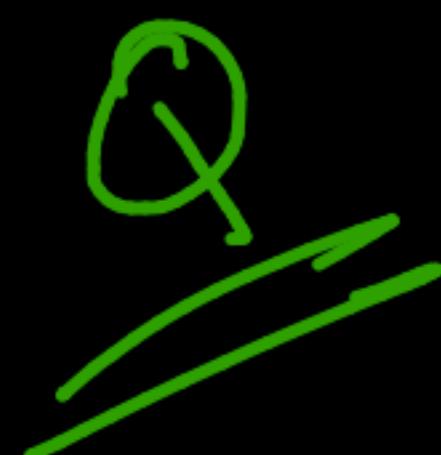
was built after
 $M_1, M_2 \& M_3$

all models were built
indep

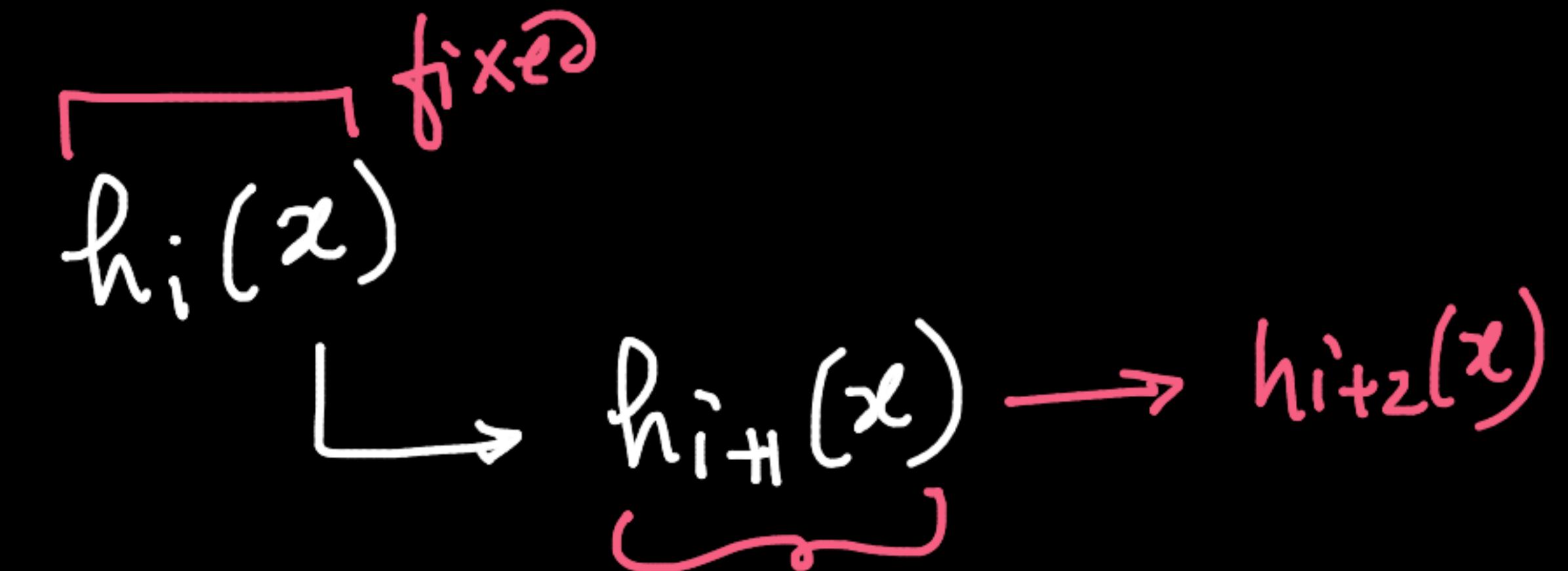
 DL: all params across multiple actions are learnt simultaneously 

→ models

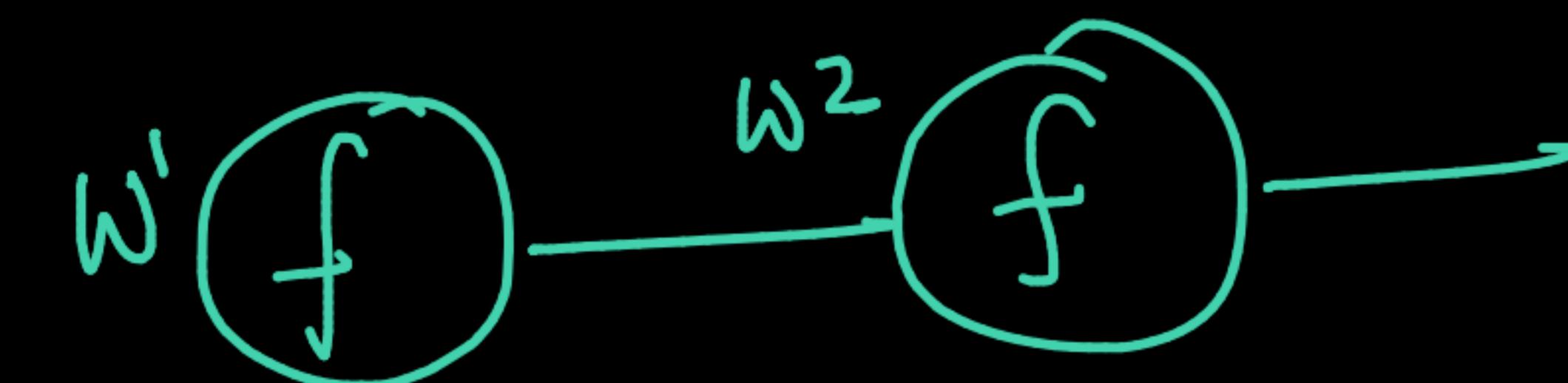
more flexibility



Boosting:



DL:



simultaneously trained

DL

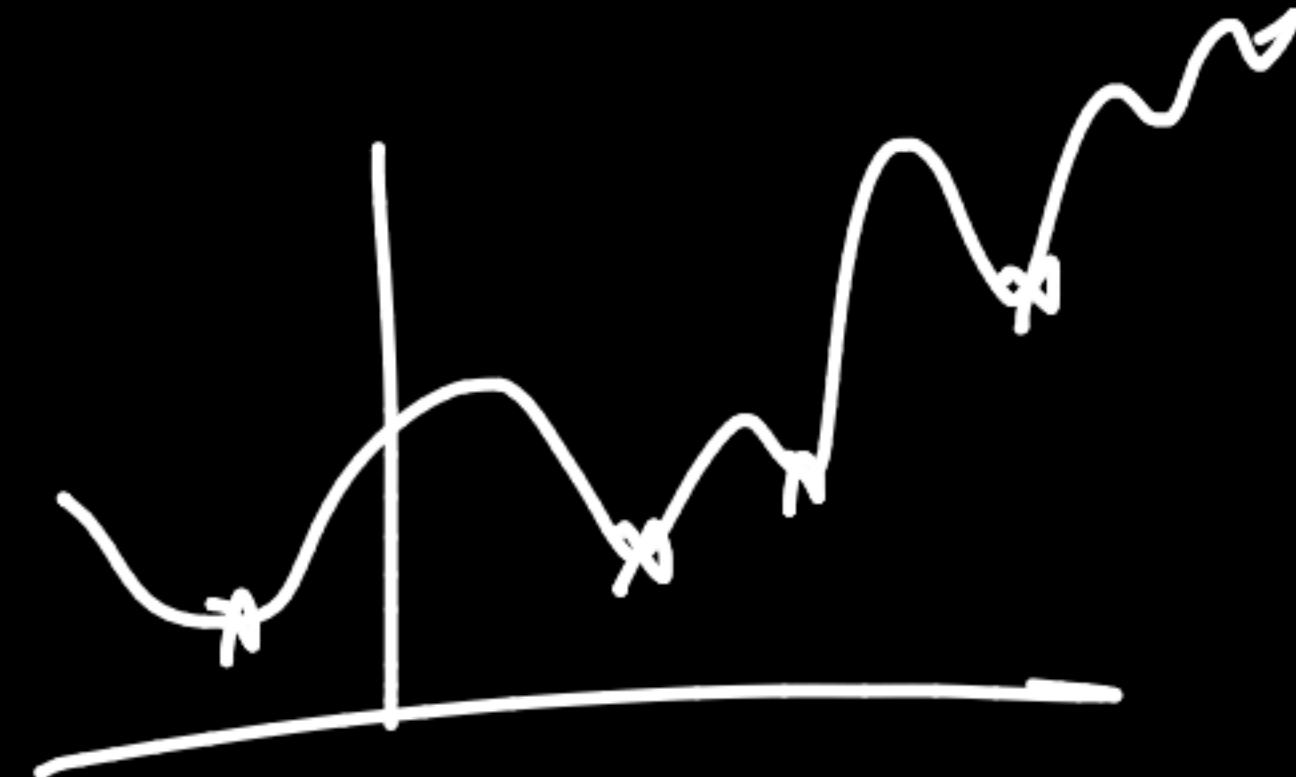
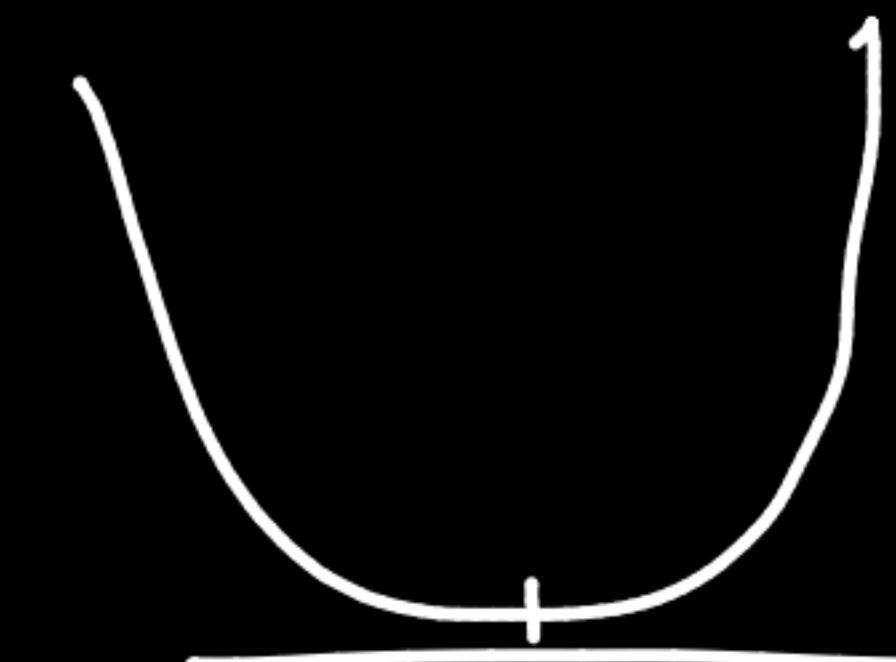
$$f(w^3 f(w^2 f(w^1 x)))$$

Very complex

non-convex

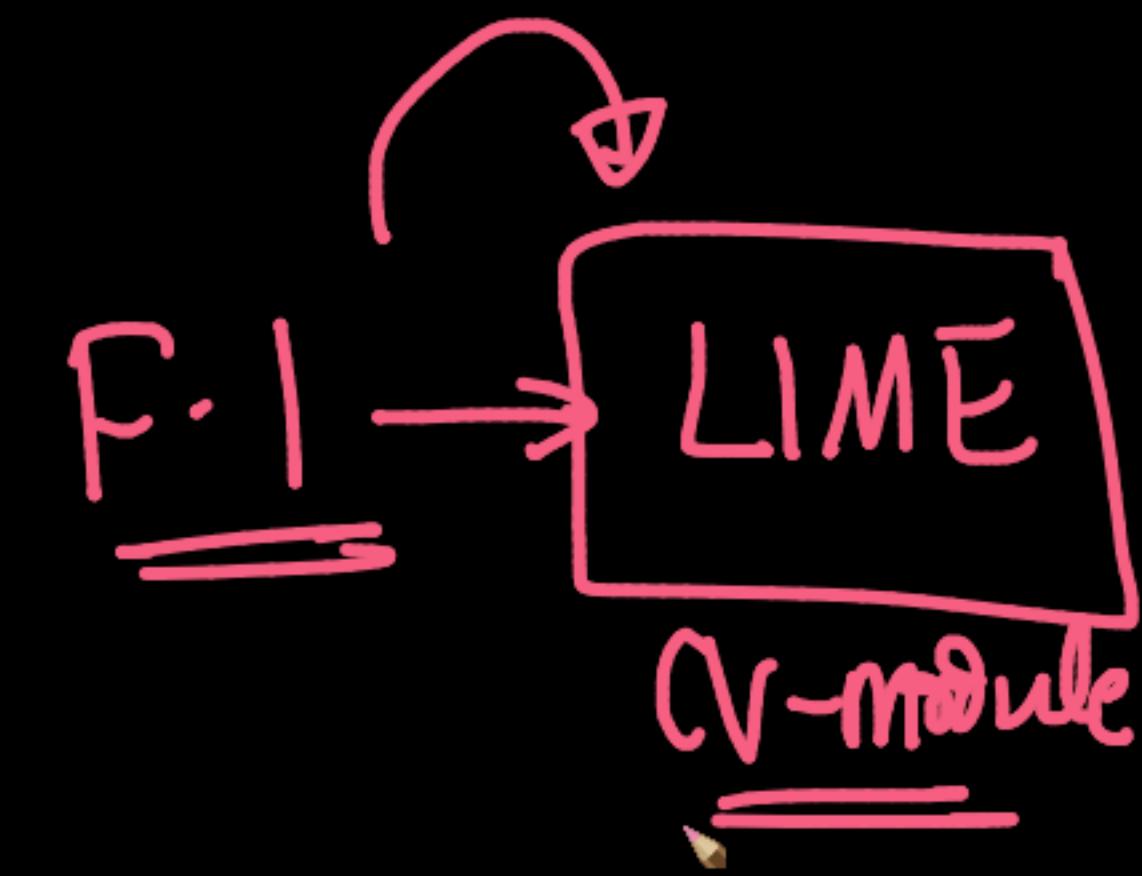
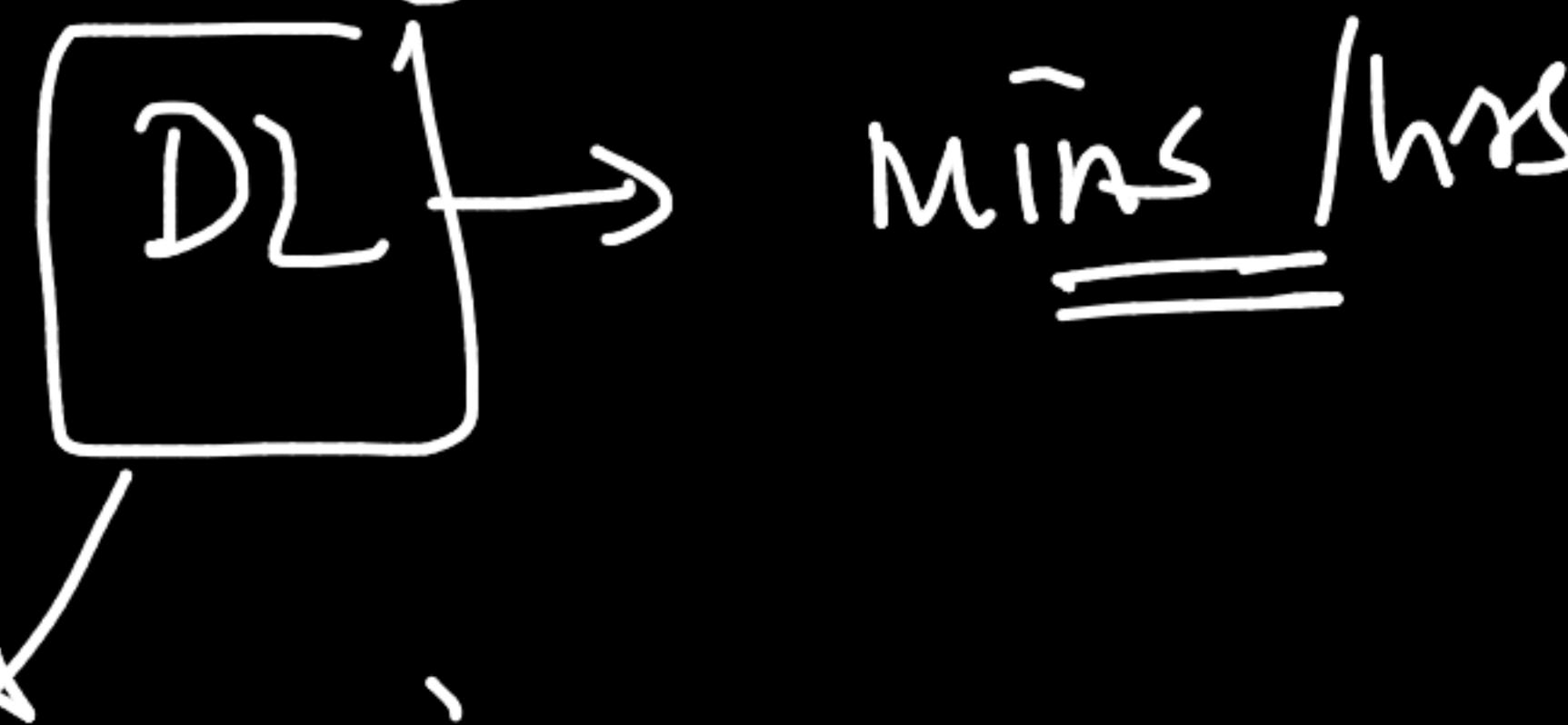


Convex:



Q

Lr. seg → few sec (1 min)



{
 V.V. expensive
 → hard way (GPUs)
 → lot more data → lots of params ✓

2 hrs + 30 min
(Q&A)

