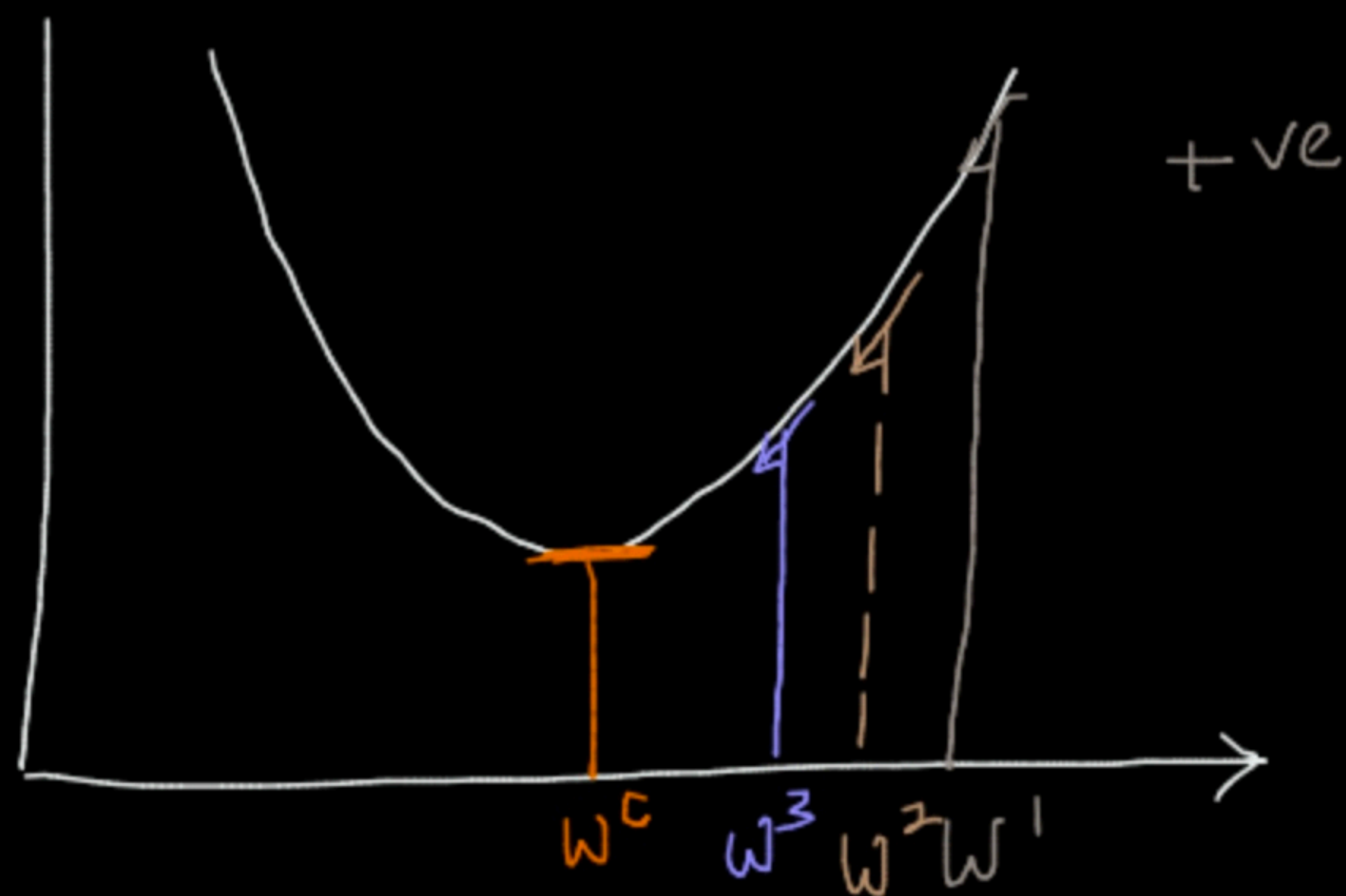


Gradient Descent Algorithm



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

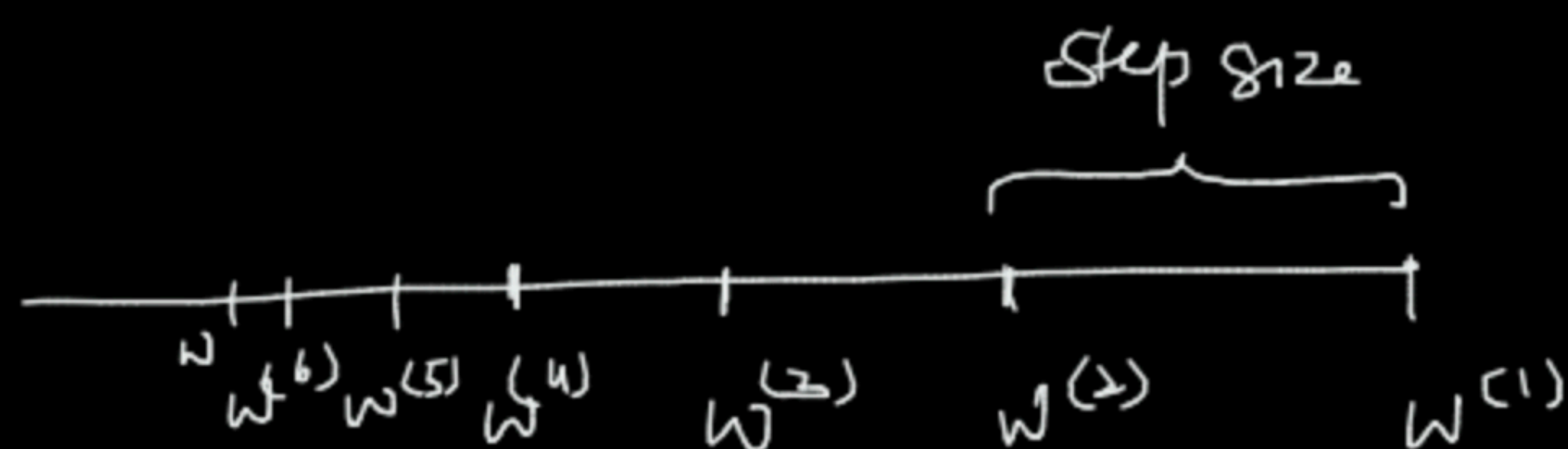
$$\hat{y} = w_1 x_1 + w_2 x_2 + b$$

$$\frac{\partial L}{\partial w_1} \quad \frac{\partial L}{\partial w_2} \quad \frac{\partial L}{\partial w_3}$$

$$w_1^{new} = w_1^{old} + \lambda \left[-\frac{\partial L}{\partial w_1} \right] w_1^{old}$$

$$w_2^{new} = w_2^{old} + \lambda \left[-\frac{\partial L}{\partial w_2} \right] w_2^{old}$$

$$w^{new} = w^{old} + \underbrace{\lambda}_{\text{Step size} \rightarrow 0}$$



λ - learning rate

Step size $\left[\lambda \frac{\partial L}{\partial w} \right]$ } optimizer

optimizers

1 Momentum — Gradient



$$W^{new} = W^{old} + \left(\frac{\partial L}{\partial W} \right)$$

2. Nesterov Momentum

$$W^5 \rightarrow$$

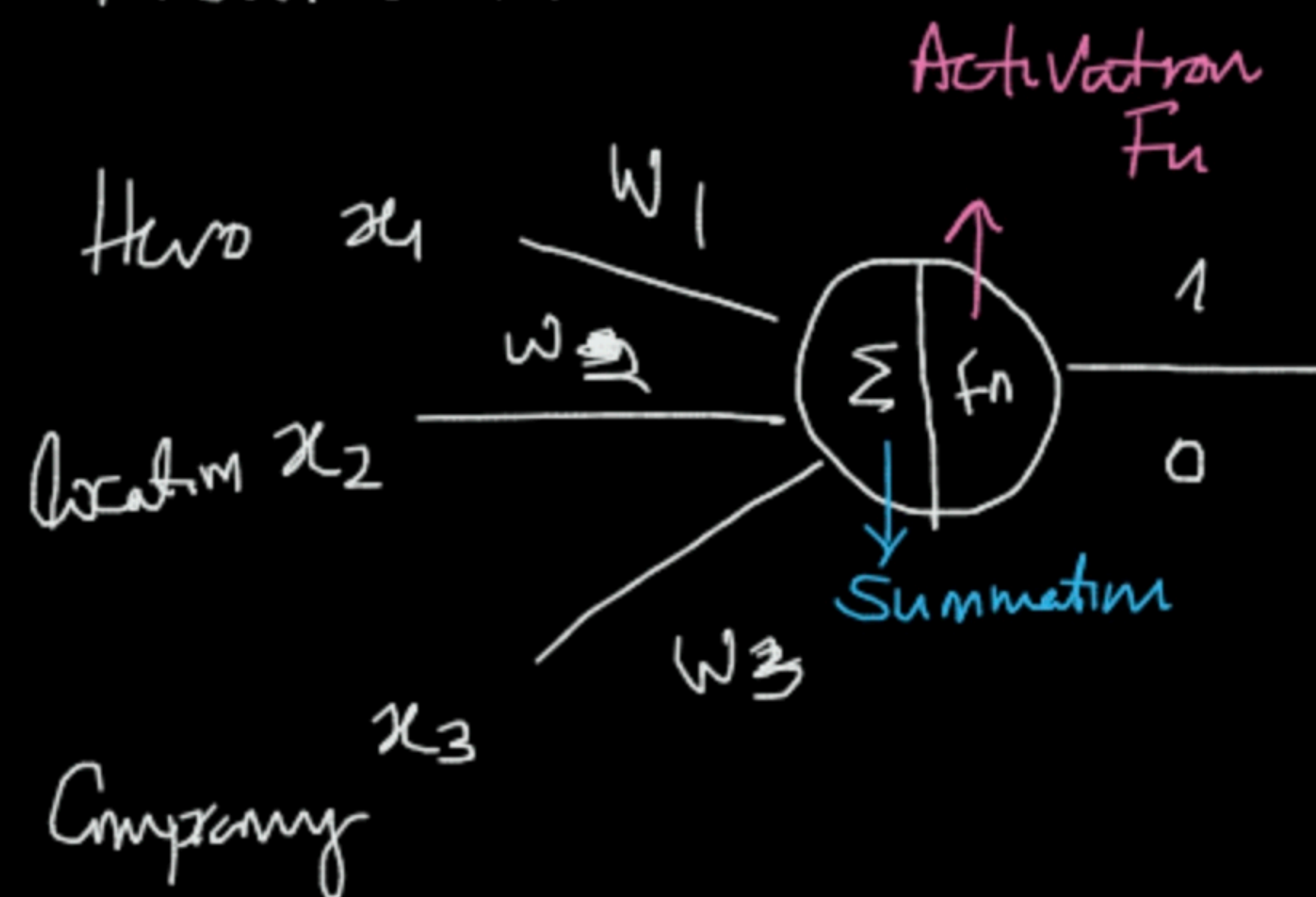
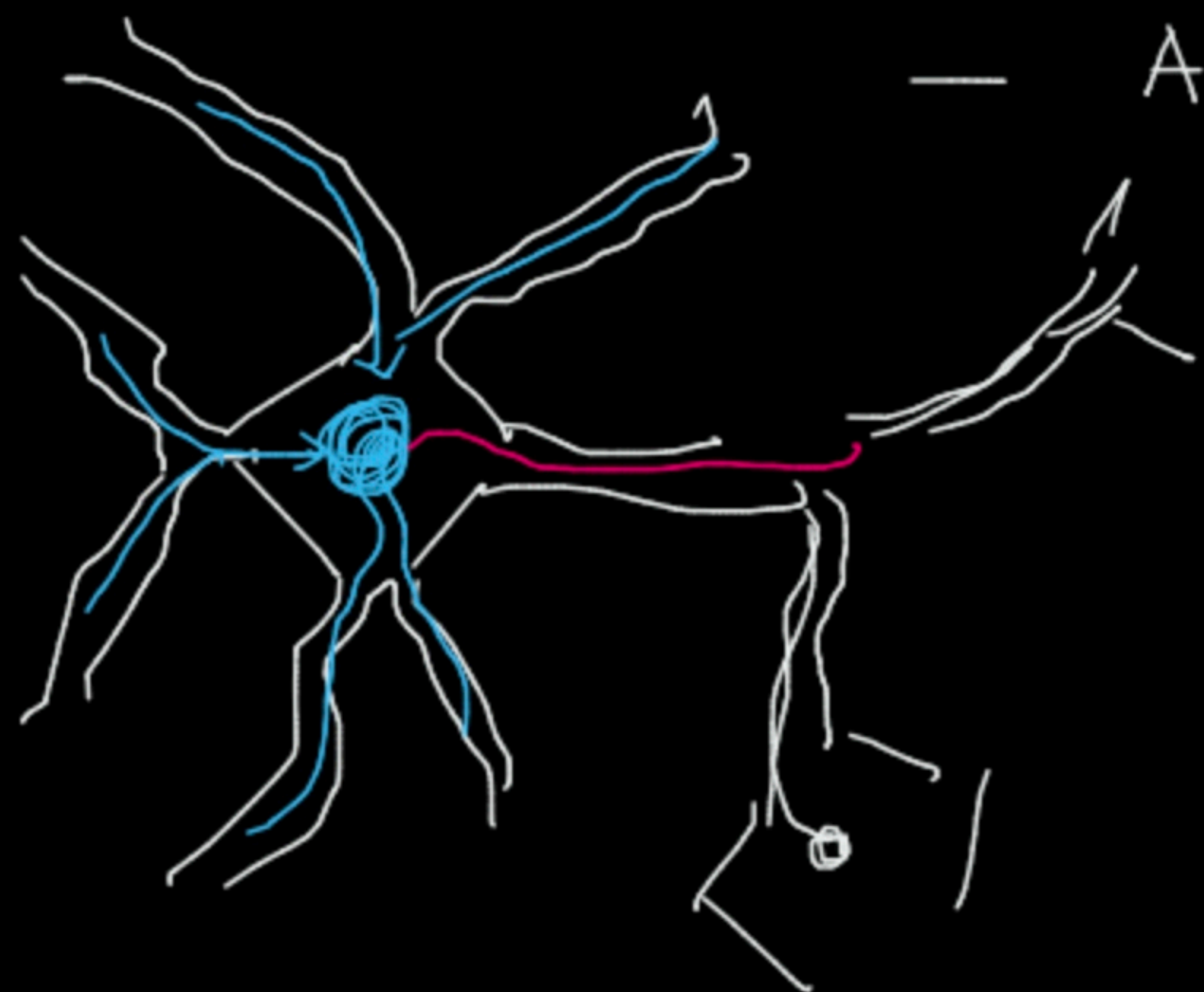
$$W^6 \rightarrow$$



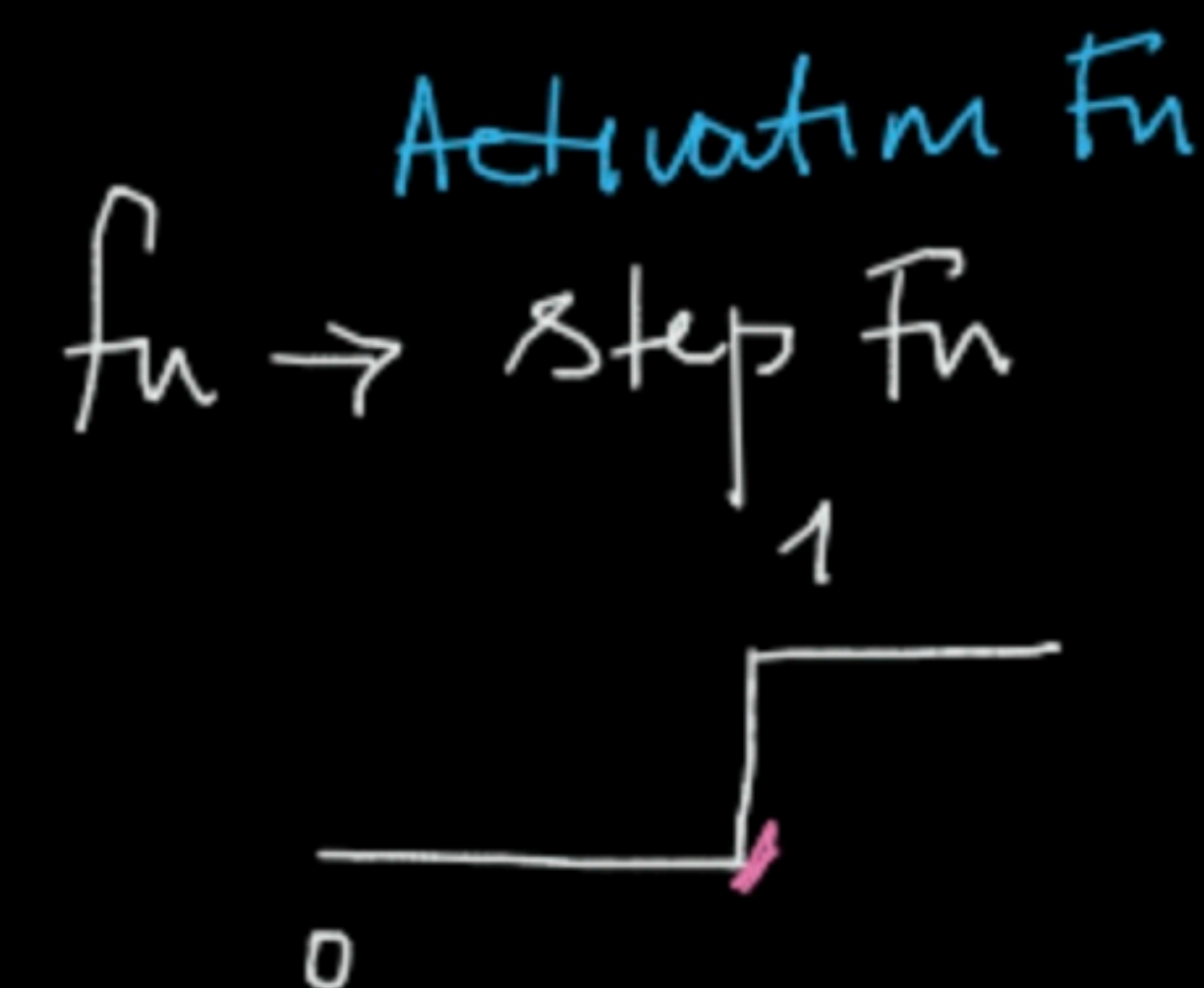
$$W^5 = W^4 - \lambda \underbrace{\left[\frac{\partial L}{\partial W^1} + \frac{\partial L}{\partial W^2} + \frac{\partial L}{\partial W^3} \right]}_{\text{Layer step size}}$$

Neural Networks

- Artificial Neurons
- Artificial Neural Networks



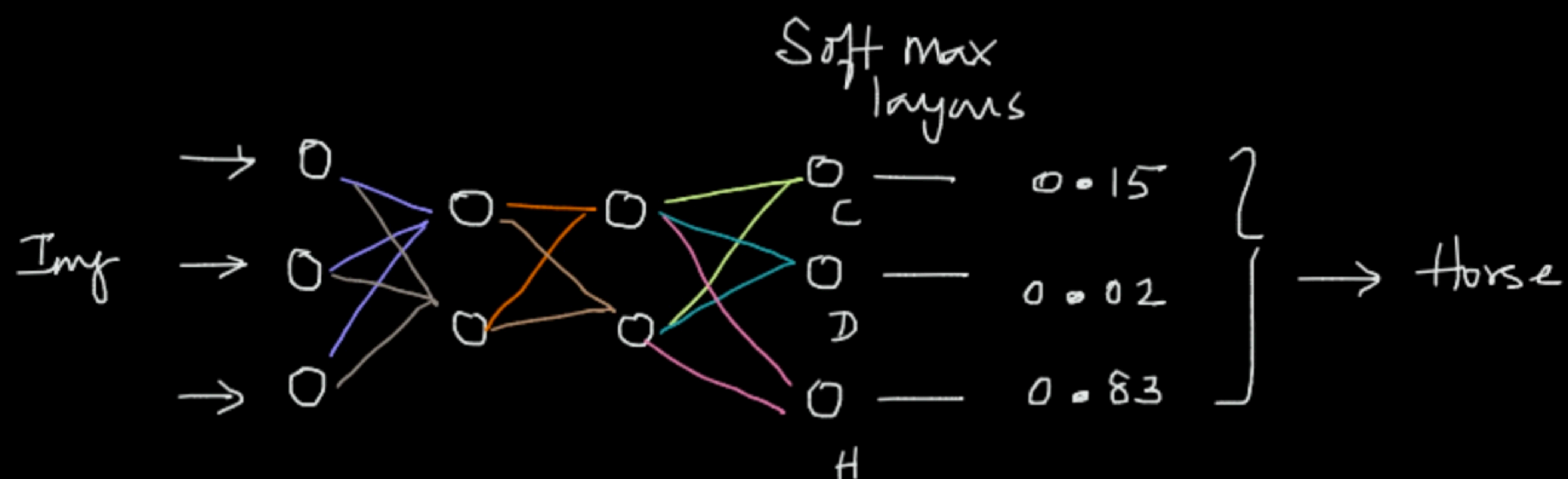
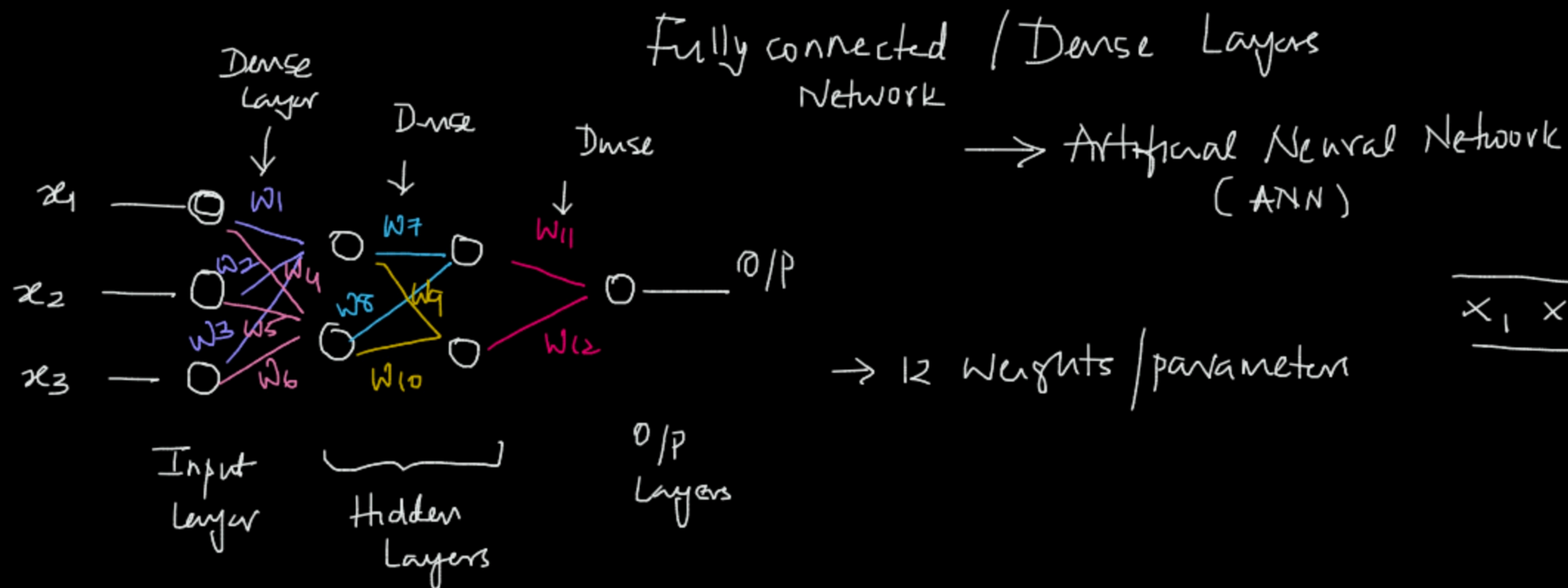
Perceptron



$$f_n(w_1x_1 + w_2x_2 + w_3x_3) \rightarrow \begin{matrix} 1 \\ \downarrow \\ \text{Activation} \\ f_n \end{matrix} \rightarrow \begin{matrix} 1 \\ \rightarrow 0 \end{matrix}$$

$$w_1x_1 + w_2x_2 + w_3x_3 \geq b \rightarrow \begin{matrix} 1 \\ \rightarrow 0 \end{matrix}$$

Activation $f_n \rightarrow$ Step f_n



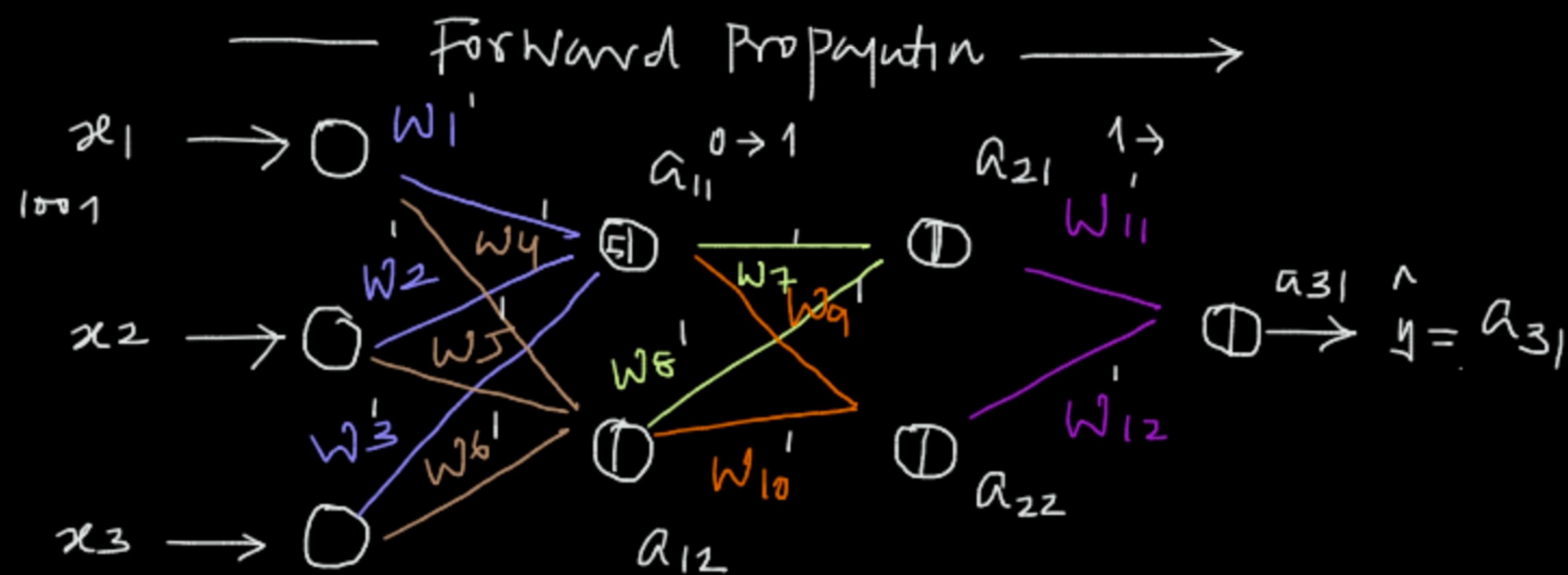
→ Dog
 → cat
 → Horse

} 3 o/p

→ 0
 → 1
 → 2

 → 9

} 10 o/p's



$$a_{11} = f_n(w_1'x_1 + w_2'x_2 + w_3'x_3)$$

$$a_{12} = f_n(w_4'x_1 + w_5'x_2 + w_6'x_3)$$

$$a_{21} = f_n(w_7'a_{11} + w_8'a_{12})$$

$$a_{22} = f_n(w_9'a_{11} + w_{10}'a_{12})$$

$$a_{31} = f_n(w_{11}'a_{21} + w_{12}'a_{22})$$

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

$$L = (y - a_{31})^2 \downarrow$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial a_{31}} \frac{\partial a_{31}}{\partial w_{11}}$$

$$\frac{\partial L}{\partial w_7} = \frac{\partial L}{\partial a_{31}} \frac{\partial a_{31}}{\partial a_{21}} \frac{\partial a_{21}}{\partial w_7}$$

- Gradient clipping
→ Exploding gradient

$$\left(\frac{\partial L}{\partial w_1} \right) = \left(\frac{\partial L}{\partial a_{31}} \right) \left(\frac{\partial a_{31}}{\partial a_{21}} \right) \left(\frac{\partial a_{21}}{\partial a_{11}} \right) \left(\frac{\partial a_{11}}{\partial w_1} \right)$$

0.2 0.1 0.01 0.02

$$= 0.000004$$

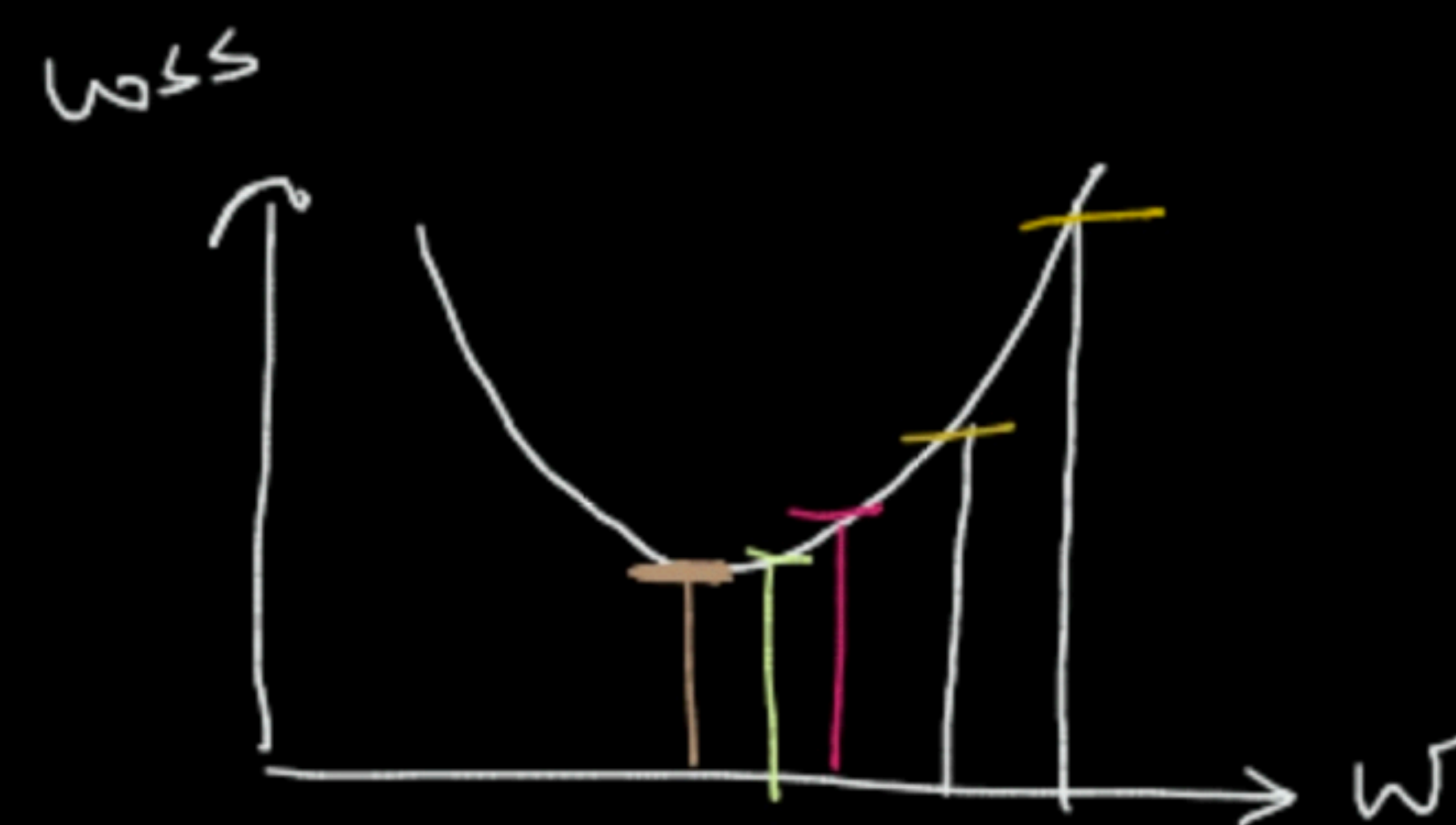
$w_1^{(1)} \rightarrow w_1^{(2)}$ → Vanishing gradients
- Never reach the min loss



Step 1 Randomly Choose 12 Weights

Step 2 Find the gradient of the loss fn w.r.t the weights

$$\frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}, \frac{\partial L}{\partial w_3}, \dots, \frac{\partial L}{\partial w_{12}}$$



Step 3 Simultaneously update weights

$$w_1^{new} = w_1^{old} - \lambda \left[\frac{\partial L}{\partial w_1} \right] w_1^{old}$$

$$w_2^{new} = w_2^{old} - \lambda \left[\frac{\partial L}{\partial w_2} \right] w_2^{old}$$

$$w_{12}^{new} = w_{12}^{old} - \lambda \left[\frac{\partial L}{\partial w_{12}} \right] w_{12}^{old}$$

	x_1	x_2	x_3	y	\hat{y}	$(y - \hat{y})^2$	\hat{y}	$(y - \hat{y})^2$	\hat{y}	$(y - \hat{y})^2$	\hat{y}
1	1000	3	5	47	20	27^2	23	24^2	30	17^2	45

Step 4 Repeat Steps 2 & 3 until convergence

→ all weights are stable

3 Blue 1 Brown

Composite Fn $\rightarrow y = f(g(h(x)))$

$$\frac{dy}{dx} = \left(\frac{dy}{df}\right) \left(\frac{df}{dg}\right) \left(\frac{dg}{dh}\right) \left(\frac{dh}{dx}\right)$$

↓ ↓ ↓ ↓

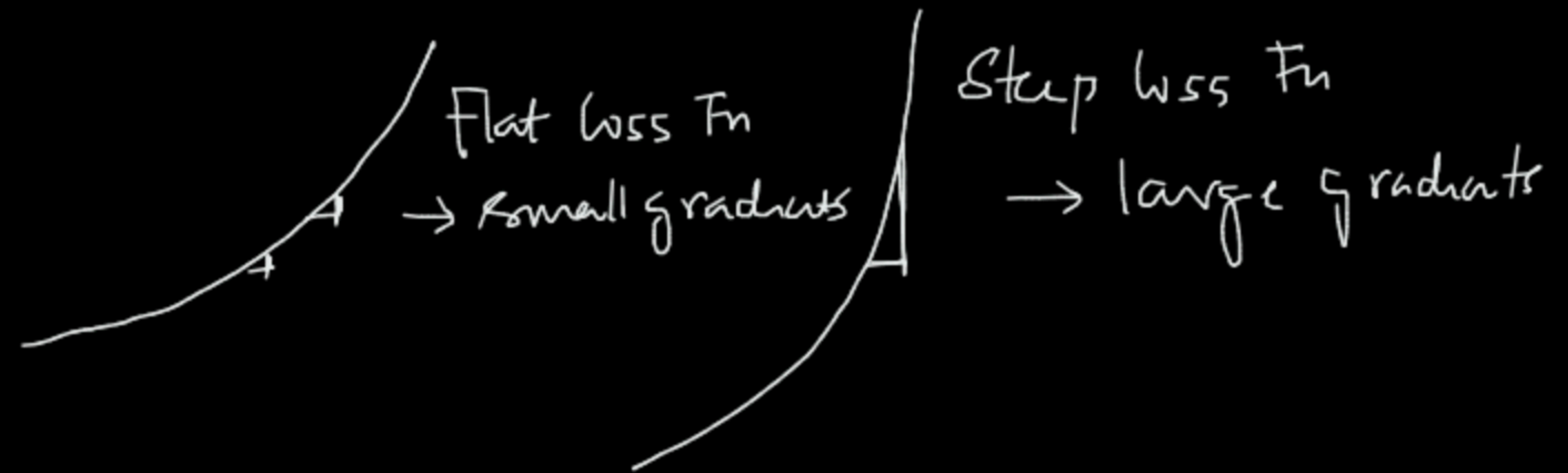
Chain Rule of differentiation

— product of multiple gradients

1/p $b - \Delta \rightarrow 0$
 $b \rightarrow 1$

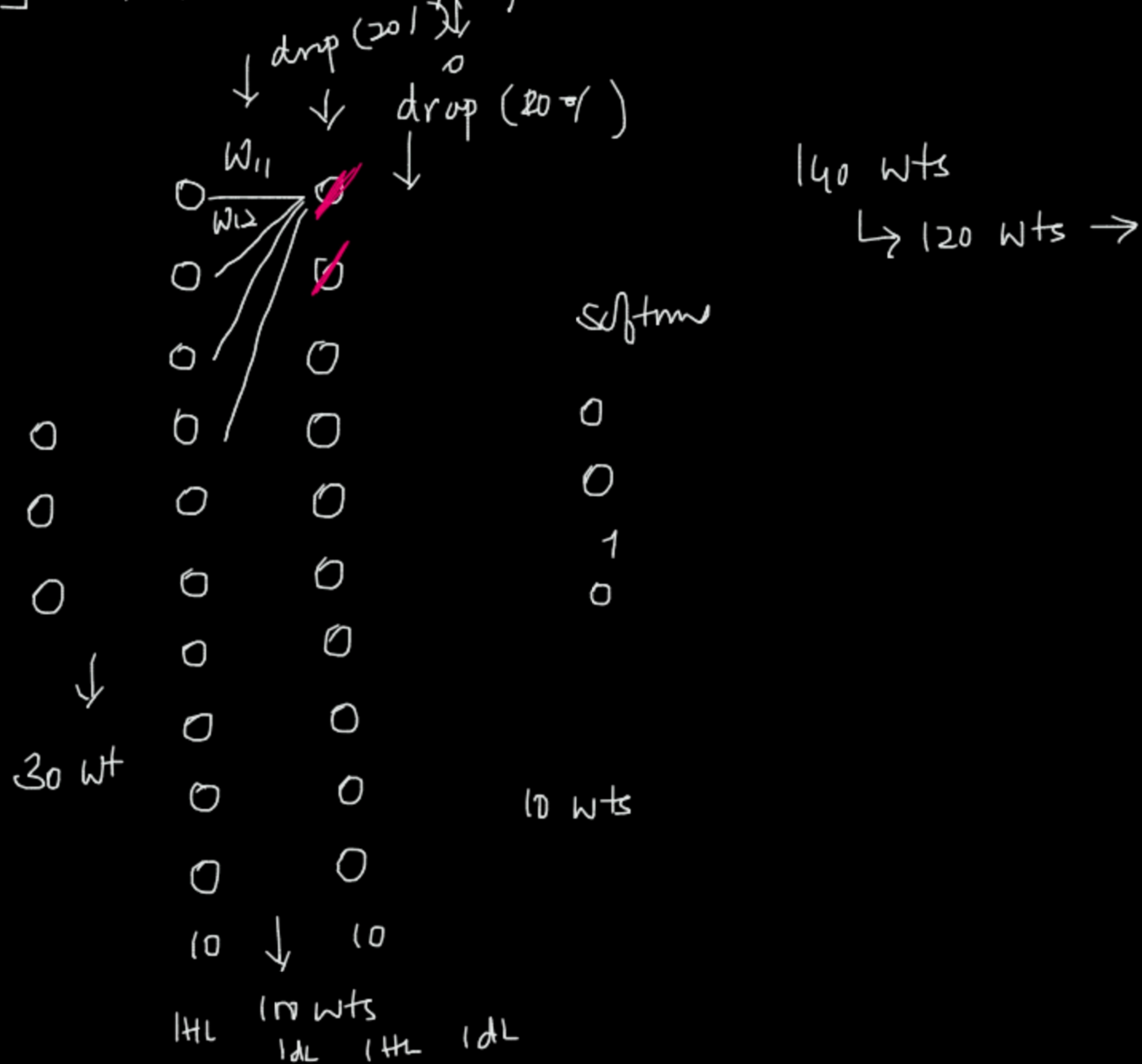


$$g(x) = \frac{1}{1 + e^{-x}}$$



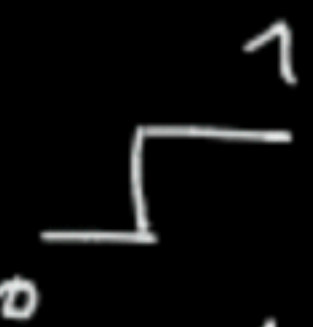
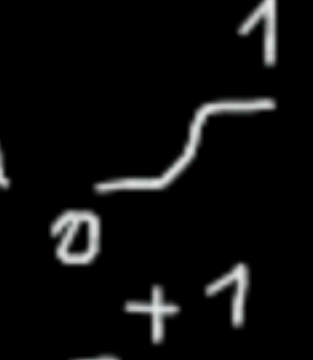

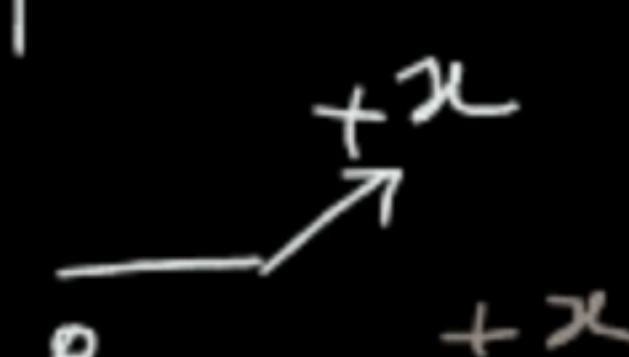

Overfitting \rightarrow Regularize

$$\hat{y} = w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_m x_m + b$$



140 Wts
 \rightarrow 120 Wts \rightarrow

• Activation Fn

- Step Fn 
- Sigmoid Fn 
- tanh 
- ReLu 
- Leaky ReLu 
- Linear $\rightarrow x$

• optimizers

- Momentum
- Nesterov Momentum
- Adam
- RMSprop
- Adagrad

• Weight Initialization

- Normal
- uniform
- Zero

• Network

- No 3 hidden layers
- No 4 units in HL

• Layers

- Dense layer
- Soft max
- dropout

• Regularization

L1 Ry

L2 Ry

• Loss Fn

Regression

— Squared loss

Classification

— Binary cross entropy (2 class)

— Categorical cross entropy (Multiclass)

Gradient Descent

- learning rate
- Batch size
- Epochs