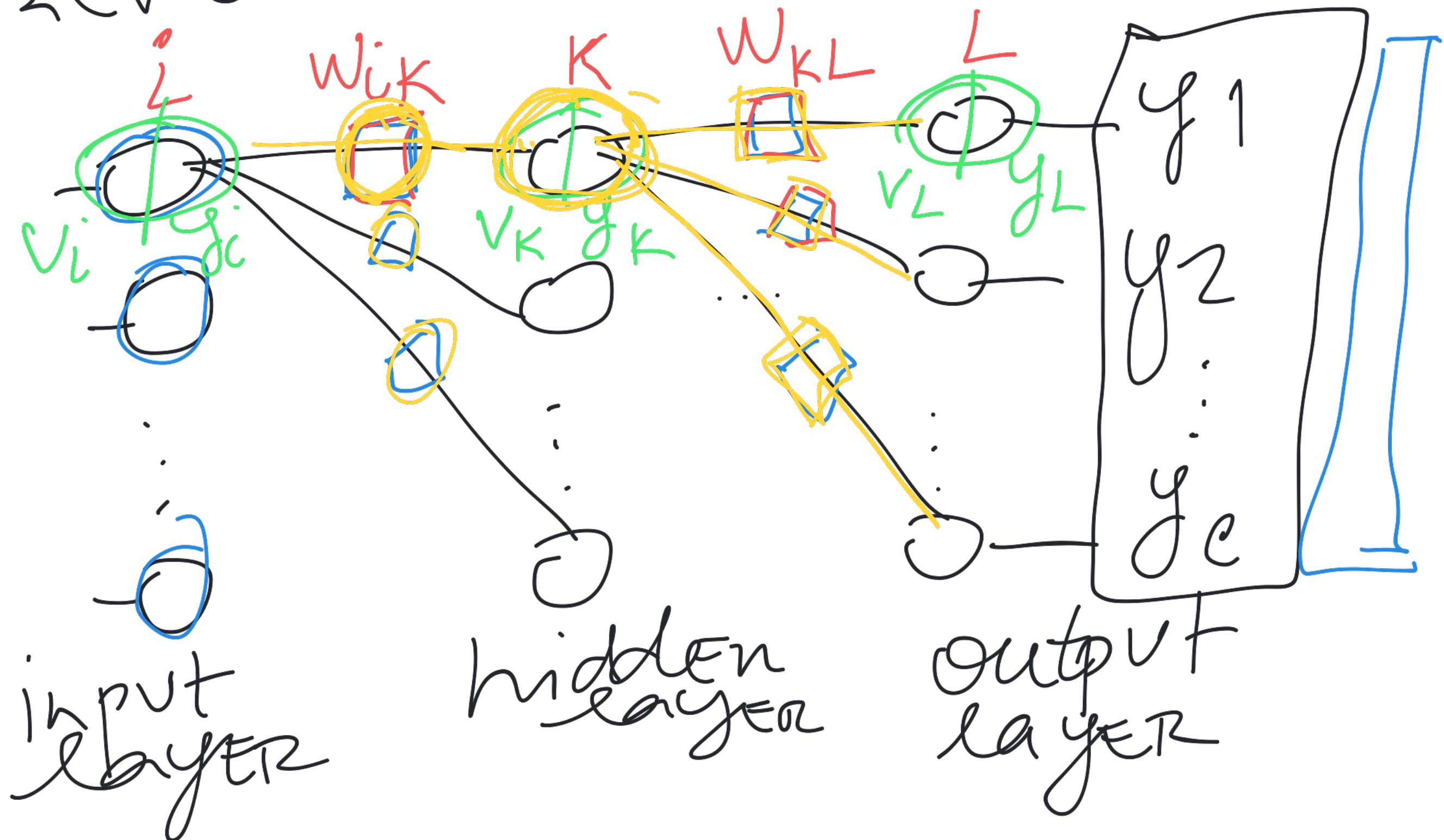


LECTURE 18

REVIEW BACK PROPAGATION



$$X = \{x_1, x_2, \dots, x_n\} \in \mathbb{R}^D$$

$$t = \{t_1, t_2, \dots, t_n\}$$

could have
ONE-hot
ENCODING

$$\begin{aligned} J &= \frac{1}{2} \sum_{i=1}^N e_i^2 \\ &= \frac{1}{2} \sum_{i=1}^N (t_i - y_i)^2 \end{aligned}$$

$t_1 \in C_2$

$t_1 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$

Chain Rule : $-S_L$

$$\frac{\partial J}{\partial w_{KL}} = \boxed{\frac{\partial J}{\partial e_L} \cdot \frac{\partial e_L}{\partial y_L} \cdot \frac{\partial y_L}{\partial v_L} \cdot \frac{\partial v_L}{\partial w_{KL}}}$$

$$e_L = t_L - y_L$$

$$\bar{J} = \frac{1}{2} \sum_{i=1}^n e_i^2$$

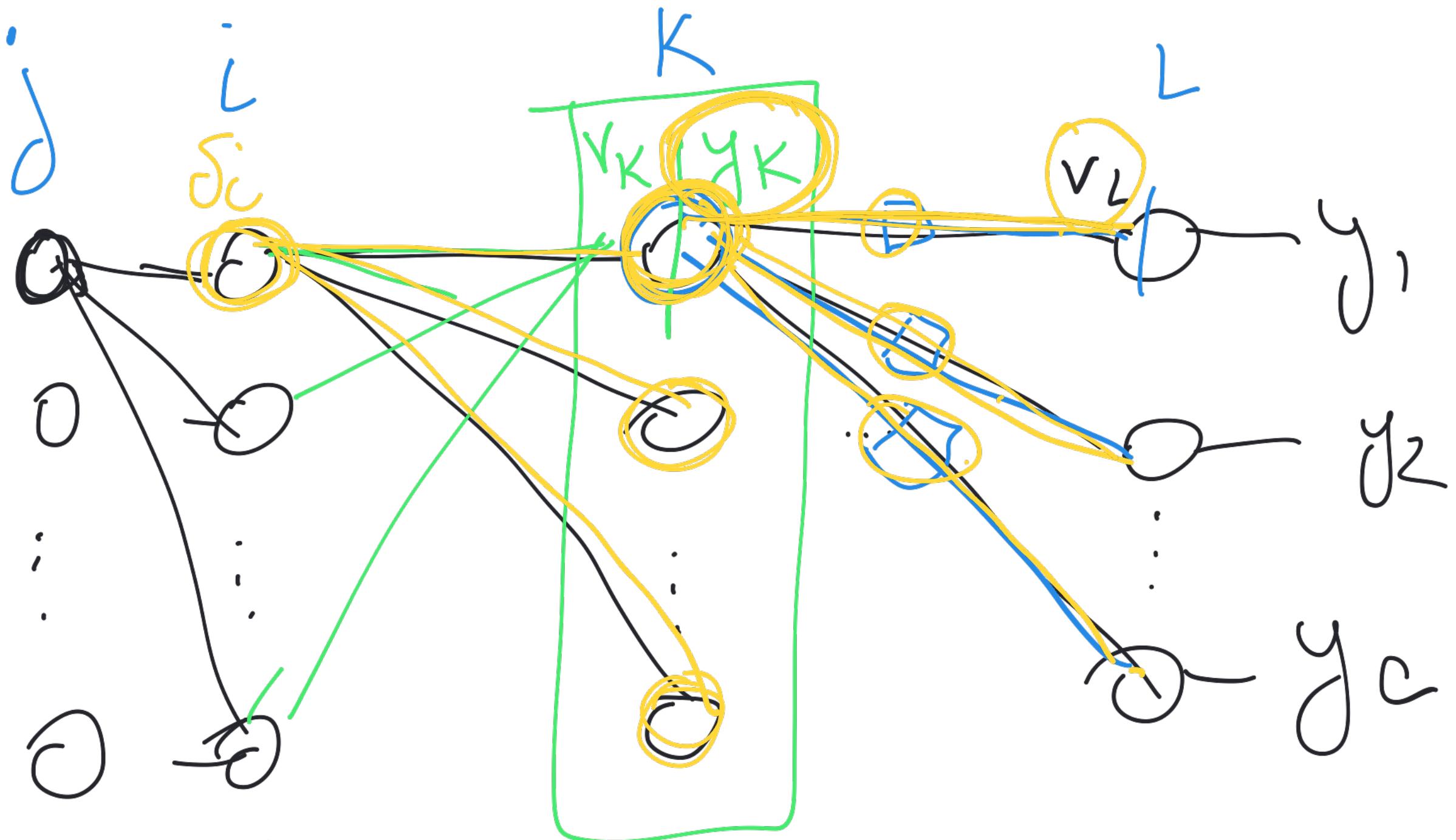
$$\underline{y}_L = \phi \left(\sum_K w_{KL} \cdot \underline{y}_K \right) = \phi(\underline{v}_L)$$

$$\frac{\partial J}{\partial w_{KL}} = e_L \cdot (-1) \cdot \phi'(v_L) \cdot y_K$$

$$w_{KL}^{(t+1)} \leftarrow w_{KL}^{(t)} - \eta \cdot \frac{\partial J}{\partial w_{KL}}$$

(t) \leftarrow iteration

learning
Rate



Local gradient:

$$\delta_i = \frac{\partial J}{\partial v_i}$$

$$\delta_K = \frac{\partial J}{\partial v_K}$$

$$\delta_L = -\frac{\partial J}{\partial v_L}$$

$$s_k = - \frac{\partial J}{\partial v_k}$$

$$= - \frac{\partial J}{\partial y_k} \cdot \frac{\partial y_k}{\partial v_k}$$

$\phi'(v_k)$

$$\begin{aligned}
 \frac{\partial J}{\partial y_k} &= \sum_L \left[\frac{\partial J}{\partial e_L} \cdot \frac{\partial e_L}{\partial y_L} \cdot \frac{\partial y_L}{\partial v_L} \cdot \frac{\partial v_L}{\partial y_k} \right] \\
 &= \sum_L (-\delta_L) \cdot w_{KL}
 \end{aligned}$$

$\frac{\partial J}{\partial v_L} = -\delta_L$

$$\begin{aligned}
 S_K &= - \frac{\partial J}{\partial y_K} \cdot \frac{\partial y_K}{\partial v_K} \\
 &= \left(\sum_L \delta_L \cdot w_{KL} \right) \cdot \phi'(v_K)
 \end{aligned}$$

Weight Update:

$$\Delta w_{ik}^{(t)} = -\gamma \cdot \delta_k \cdot y_i$$

$$w_{ik}^{(t+1)} = w_{ik}^{(t)} + \Delta w_{ik}^{(t)}$$

output Layer

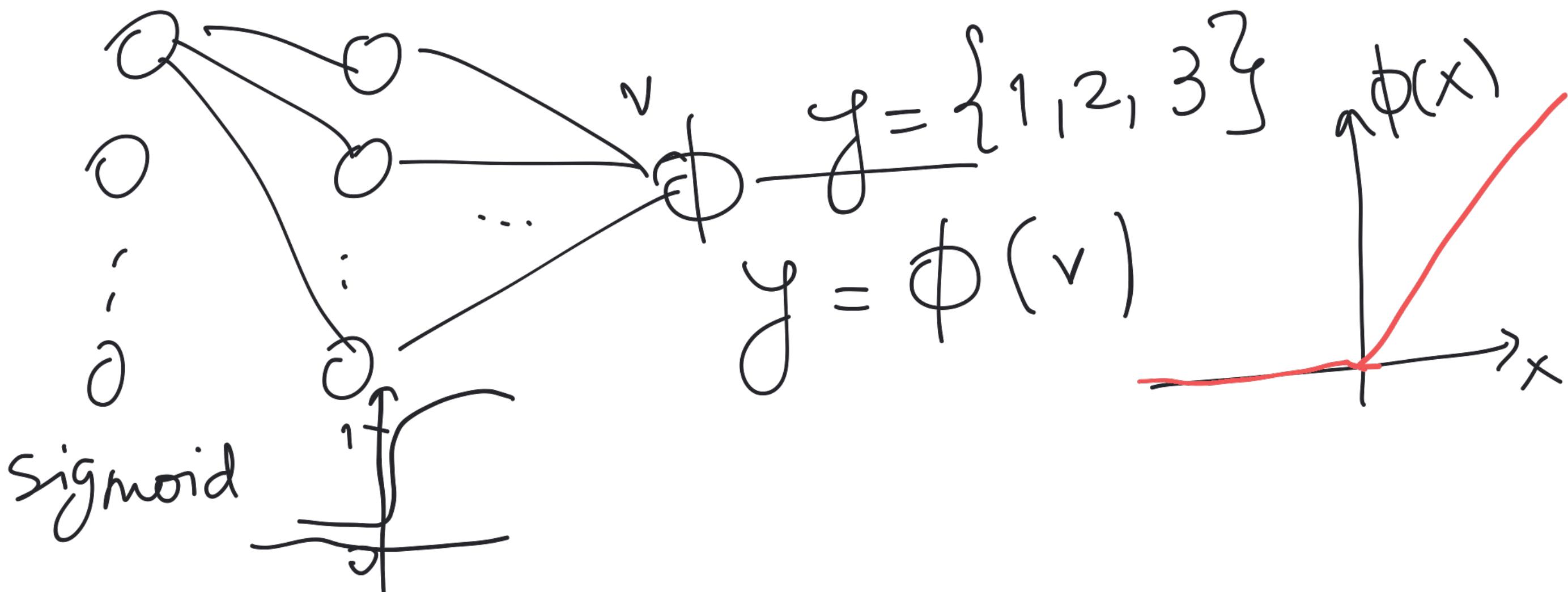
* Integer Encoding :

{Cat, dog, bird}

2 1 3

ReLU

$$\phi(x) = \begin{cases} 0, & x \leq 0 \\ x, & x > 0 \end{cases}$$

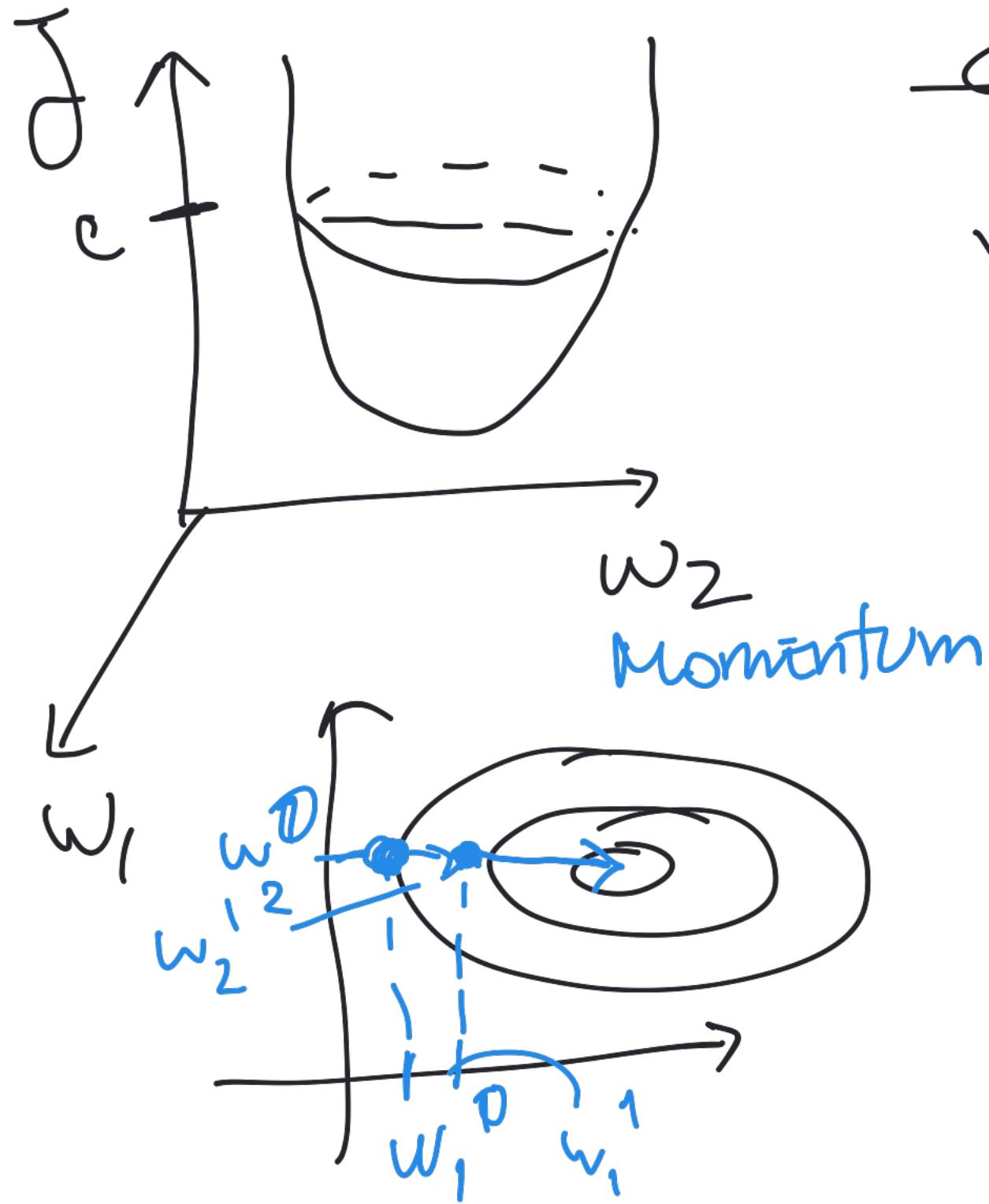


RELU : - handles
vanishing gradient
better than
Sigmoid or tanh.

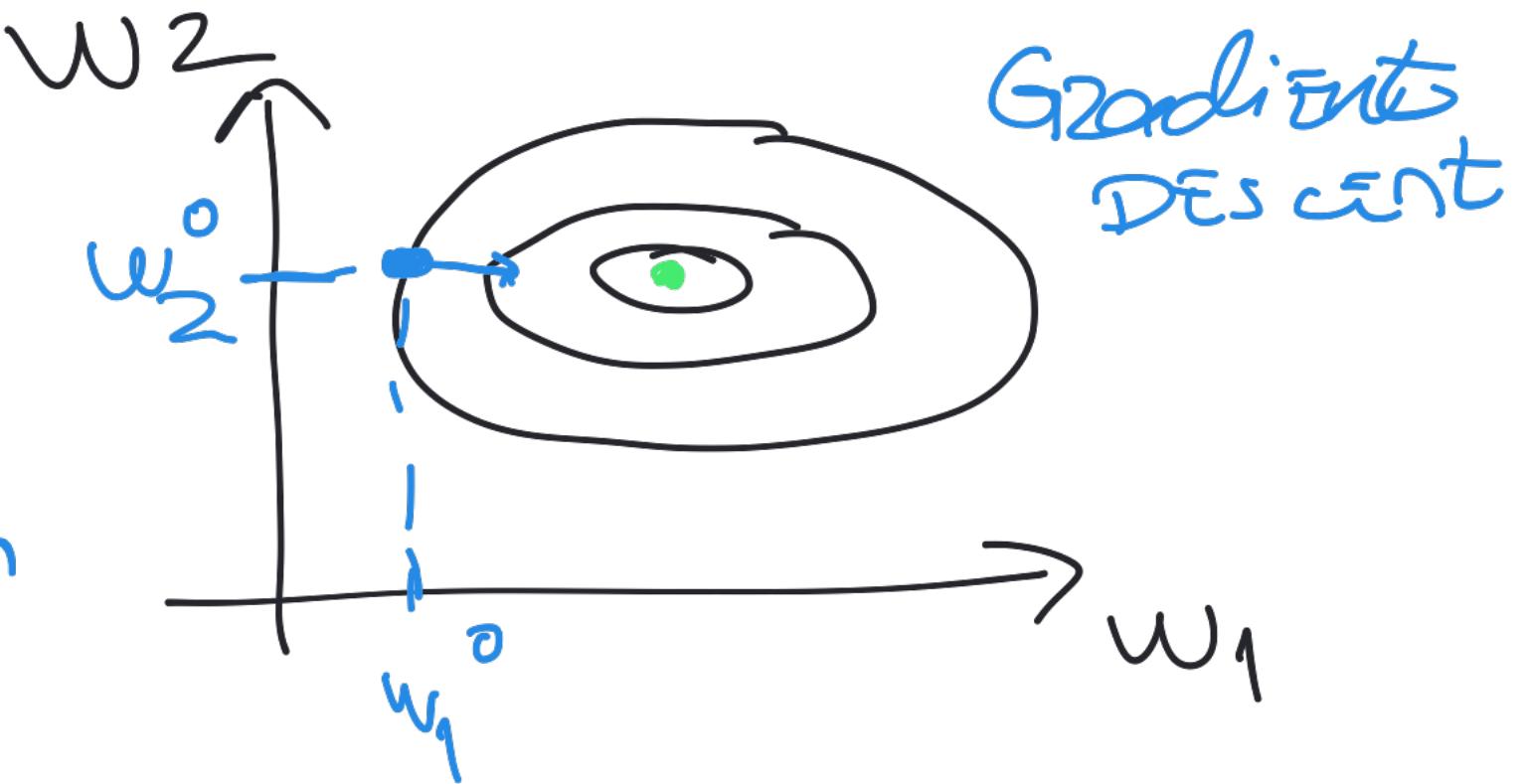
Common uses:

Sigmoid or tanh at hidden
layers followed by RELU at
output layer

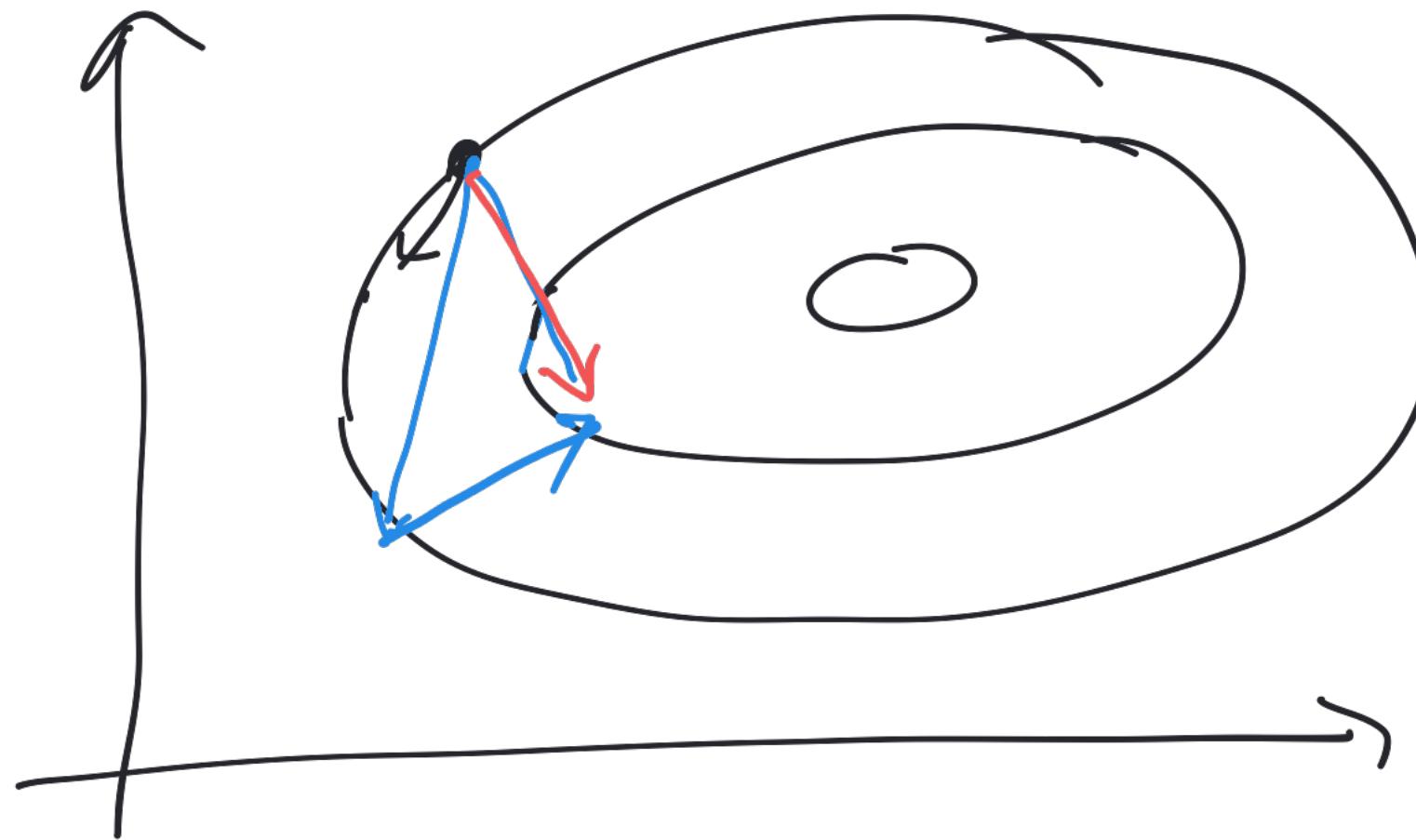
Accelerated Gradient



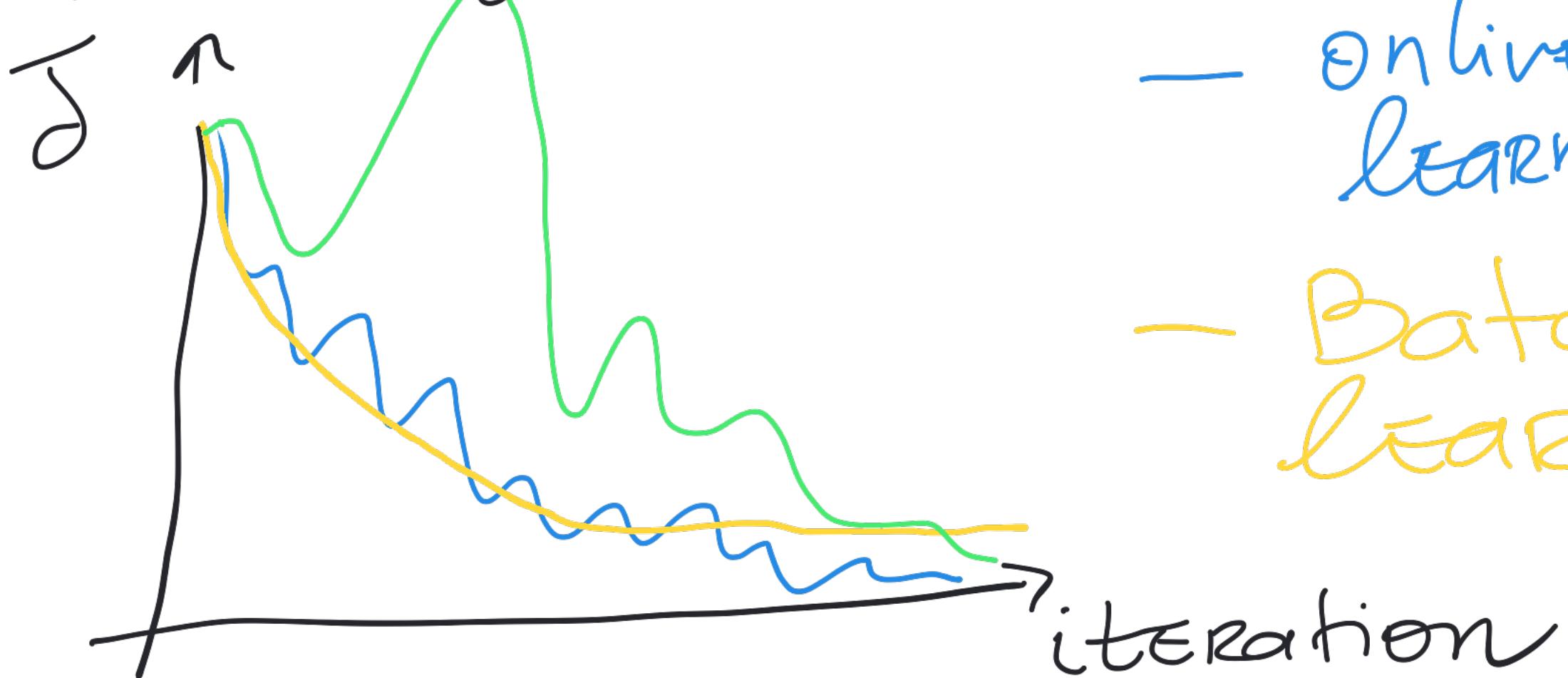
$J = c$: surface level



NAG (Nesterov's)



Learning CURVES



- online learning
- Batch learning

Online/Mini-Batch Learning

- shuffle data
- if using mini-batch,
partition data into batches
- for N epochs :
 - Compute local gradient
(iteration)