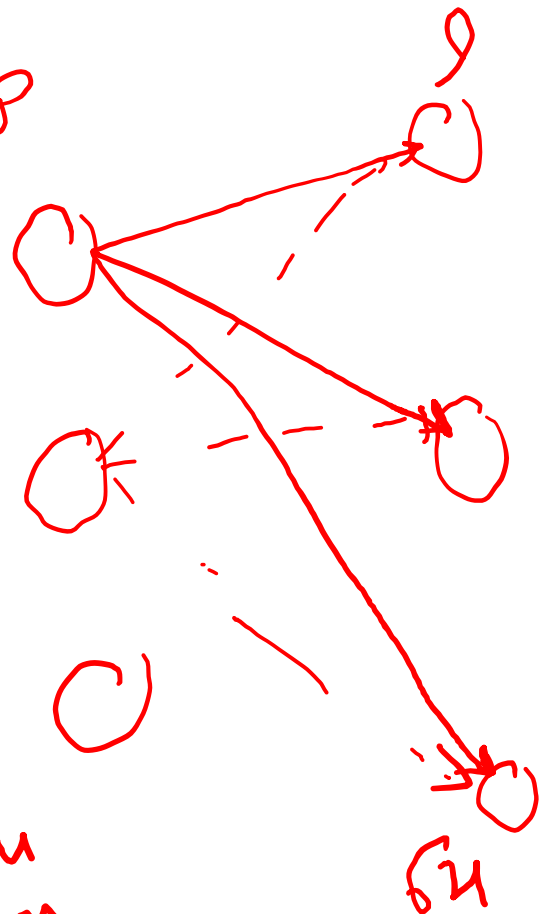


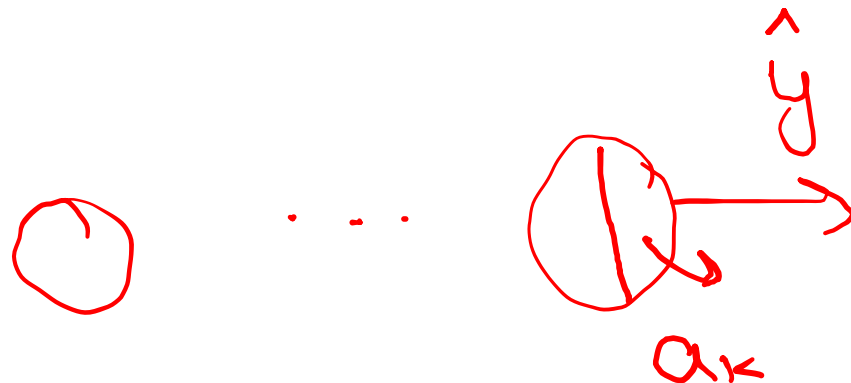
Lecture 17

* Back Propagation.

I/O



784 nodes

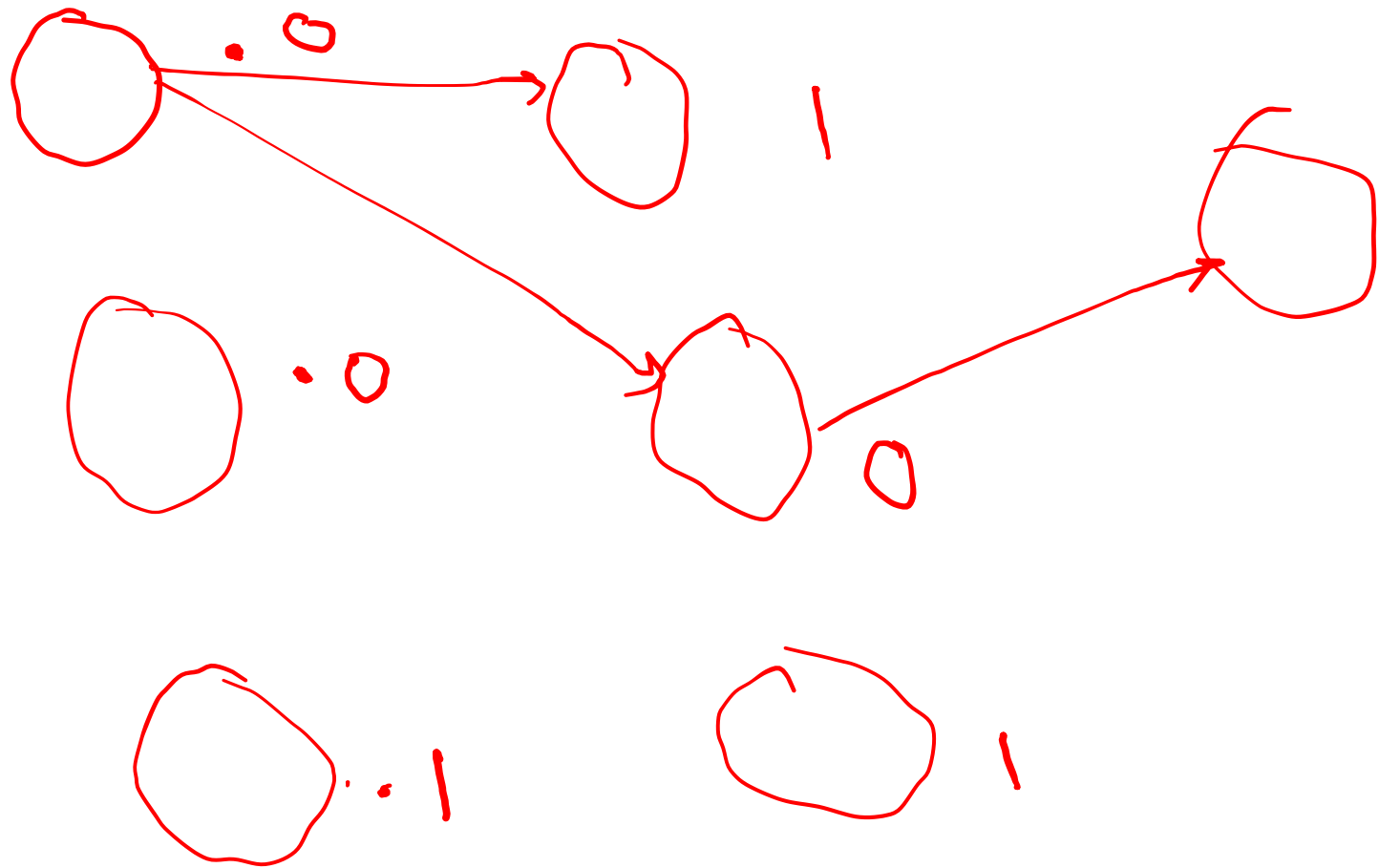


$$\delta_k = \frac{\partial E}{\partial a_k}$$

~~$\hat{y} = \sigma(a_k)$~~
 $\hat{y} = \sigma(a_k)$

$$\delta_i = \sum_j \delta_j U_{ji} \sigma'(a_i)$$

$$\delta_{kj} = \delta_k U_{kj} \sigma'(a_j)$$



$$\nabla_{\omega} J(\omega) = \begin{bmatrix} \partial J(\omega) / \partial \omega_1 \\ \partial J(\omega) / \partial \omega_2 \end{bmatrix}$$

Momentum

SGD with momentum remembers the update Δw at each iteration.

Each update is as a (convex) combination of the gradient and the previous update.

$$\Delta w := \eta \nabla Q_i(w) + \alpha \Delta w$$

$$w := w - \Delta w$$

$$w^{t+1} \leftarrow w^t - \underbrace{\eta [\nabla_{\omega} J(\omega)]^t}_{\text{gradient}} - \underbrace{\eta \alpha [\nabla_{\omega} J(\omega)]^{t-1}}_{\text{previous update}}$$

Rumelhart, David E.; Hinton, Geoffrey E.; Williams, Ronald J. (8 October 1986). "Learning representations by back-propagating errors". *Nature* 323 (6088): 535-536.

$$J(\omega) = (w_1 + w_2 x - y)^2$$

$$\nabla_{\omega} J(\omega) = \begin{bmatrix} 2(w_1 + w_2 x - y) \\ 2x(w_1 + w_2 x - y) \end{bmatrix}$$

at $t=0$, $\omega' \leftarrow \omega^0 - \eta [\nabla_{\omega} J(\omega)]^{t=0} - \cancel{\eta \alpha}$

$$\omega' \leftarrow \omega^0 - \eta \begin{bmatrix} 2(\omega_1 + \omega_2 x - y) \\ 2x(\omega_1 + \omega_2 x - y) \end{bmatrix} \quad \omega_1 = \omega_1^0, \omega_2 = \omega_2^0$$

at $t=1$, $\omega^2 \leftarrow \omega' - \eta [\nabla_{\omega} J(\omega)]^{t=1} - \eta \alpha \underbrace{[\nabla_{\omega} J(\omega)]^{t=0}}$
 $\omega_1 = \omega_1', \omega_2 = \omega_2'$

$W = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$ are the weights.

W at $t=0$ step

$$W^0 = \begin{bmatrix} w_1^0 \\ w_2^0 \end{bmatrix}$$

$W^0 \leftarrow$ initialized.

W at $t=1$ step

$$W^1 = \begin{bmatrix} w_1^1 \\ w_2^1 \end{bmatrix}$$

* Regression: $(y - \hat{y})^2 = E$

* Classification: Binary classification.
 Class labels $y \in \{0, 1\}$

$$-y(1 - \sigma(a_k)) + (1 - y)\sigma(a_k)$$

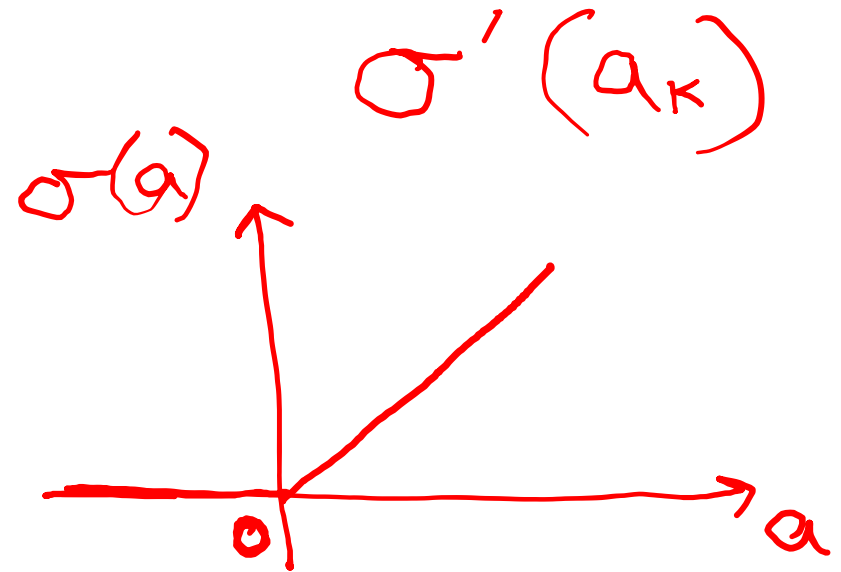
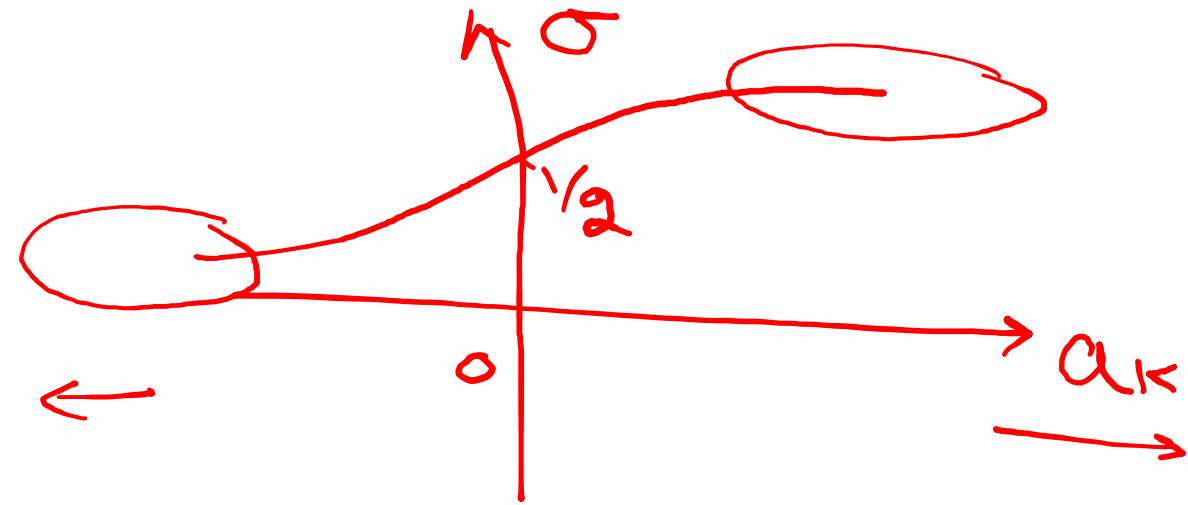
$$E = -[y \log \hat{y} + (1 - y) \log (1 - \hat{y})] \leftarrow \text{Binary cross-entropy}$$

$$\hat{y} \in \underline{[0, 1]}, (0, 1)$$

$$\partial_k: \frac{\partial E}{\partial a_k} = - \frac{\partial}{\partial a_k} [y \log \sigma(a_k) + (1 - y) \log (1 - \sigma(a_k))]$$

$$= - \frac{y \sigma(a_k)(1 - \sigma(a_k))}{\sigma(a_k)} - \frac{(1 - y)(-1)\sigma'(a_k)}{1 - \sigma(a_k)}$$

If $y = 0$, $\delta_k = \sigma(a_k)$



Vanishing gradient

$$\sigma(a) = \max(0, a)$$

ReLU \rightarrow Rectified linear unit.

$$E = -y \log \hat{y} - (1-y) \log (1-\hat{y}) \quad \Bigg| \quad E = e^{-y \hat{y}}$$

$$y = 0, \quad \hat{y} = 0.99$$

$$E = -\log(1-0.99) = \overset{2}{100} \rightarrow \text{error}$$

$$y = 1, \quad \hat{y} = 0.99 \Rightarrow$$

$$E = -1 \log(0.99) \approx 0$$

$$E = e^{-y \hat{y}}$$

$$y = 1$$

weight vector W

$$\|W^{t+1} - W^t\|$$

$$< \text{threshold.}$$

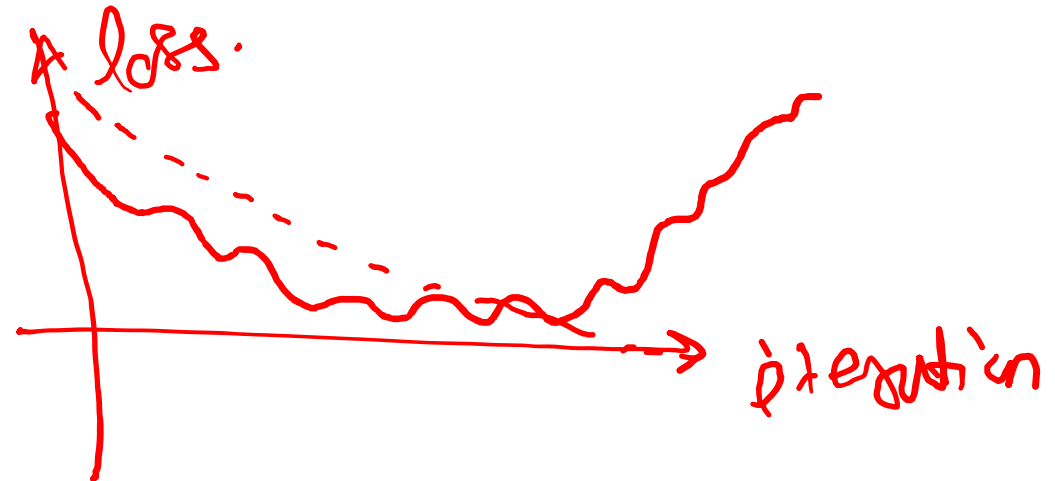
$$\text{threshold} = 10^{-3}$$

$$< \propto \|\nabla_w E(w)\|$$

$t+1^{\text{th}}$ iteration.

t^{th} iteration.

If weights do not change much b/w subsequent iterations,
exit the loop.



* multiclass cross-entropy.

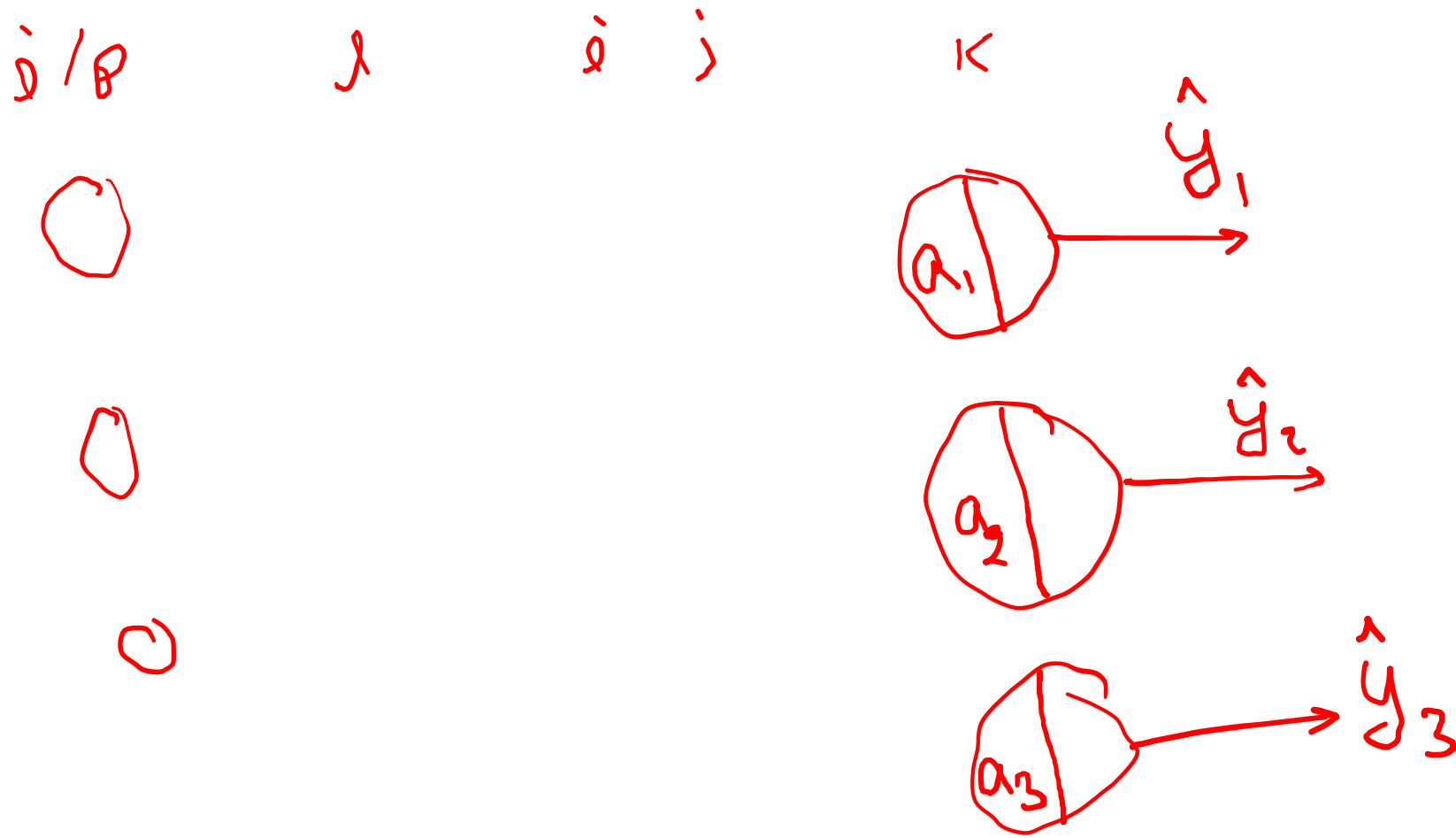
Classes $\rightarrow y \in \{0, 1, 2, \dots, C-1\}$ ~~$\leftarrow \log$~~

no. of classes, C

repr. of K^{th} class \rightarrow

$\theta = [0, 0, \dots, 1, \dots, 0] \rightarrow$ one hot encoding.
 \uparrow
 K^{th}

$$E = -\theta^T \log \hat{y}$$



o/p K^{th} layer will have C nodes for C -classes.
 let $C = 3$

activation \rightarrow softmax.

$$\hat{y}_1 = \frac{e^{a_1}}{\sum_{p=1}^C e^{a_p}}$$

$$\hat{y}_2 = \frac{e^{a_2}}{\sum_{p=1}^C e^{a_p}}$$

$$\hat{y}_K = \frac{e^{a_K}}{\sum_{p=1}^C e^{a_p}}$$

$\hat{y}_K \rightarrow$ probability