

www.menti.com

Code : 178406

CS230: Lecture 3

The mathematics of deep learning
Backpropagation, Initializations, Regularization

Kian Katanforoosh

I – Backpropagation

- derivatives

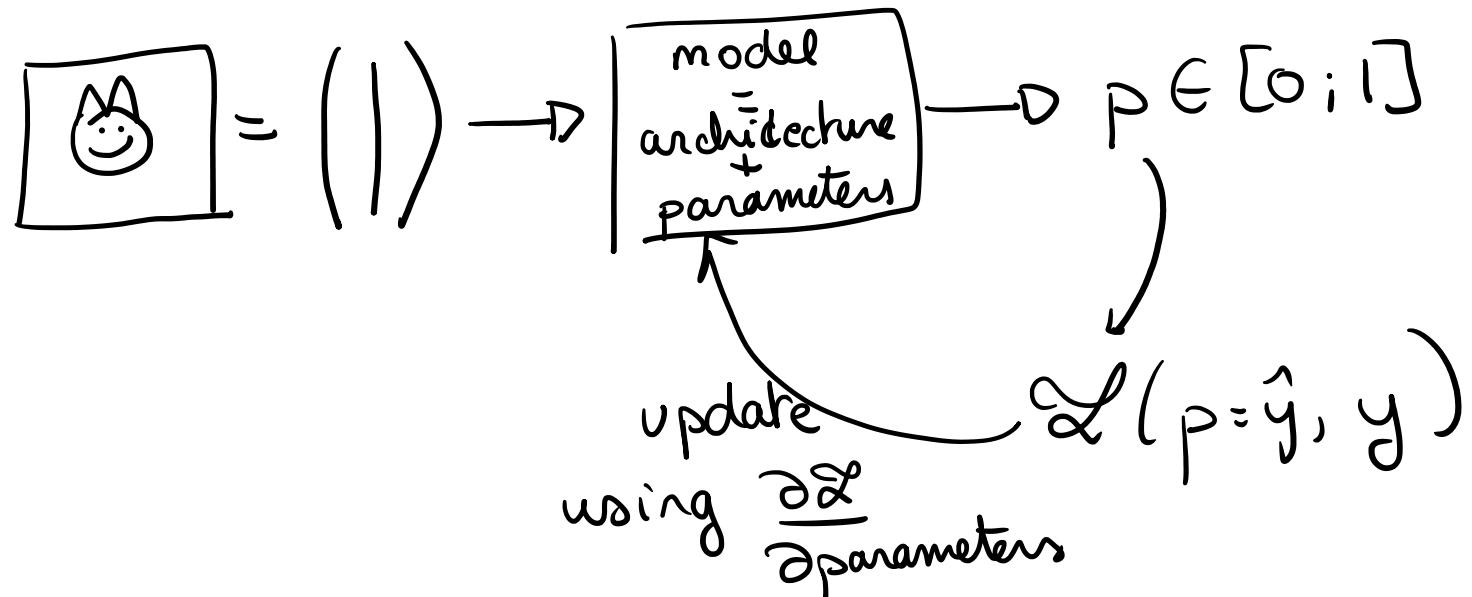
II – Initializations

- Variance
Mean

III – Regularization

I - Backpropagation

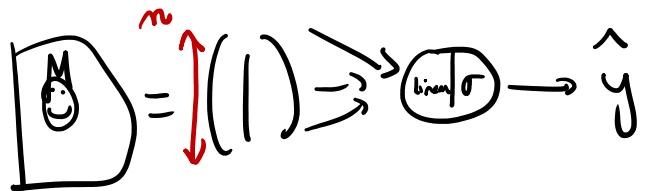
Problem statement



$$\text{GD : } w = w - \alpha \frac{\partial \mathcal{L}}{\partial w}$$
$$b = b - \alpha \frac{\partial \mathcal{L}}{\partial b}$$

A - Logistic Regression backpropagation for one training example

$$\frac{\partial \mathcal{L}}{\partial z} = \frac{\partial \mathcal{L}}{\partial a} \cdot \frac{\partial a}{\partial z}$$



forward propagation:

$$\begin{aligned} & \text{Input } x \in \mathbb{R}^{(1,n)} \\ & \text{Weights } w \in \mathbb{R}^{(n,1)} \\ & \text{Bias } b \in \mathbb{R}^{(1,1)} \\ & \text{Z} = w \cdot x + b \in \mathbb{R}^{(1,1)} \\ & \hat{y} = a = \sigma(z) \end{aligned}$$

$$\text{Loss: } \mathcal{L} \doteq - \left(y \log \hat{y} + (1-y) \log (1-\hat{y}) \right)$$

Backprop:

$$\frac{\partial \mathcal{L}}{\partial w} \quad \frac{\partial \mathcal{L}}{\partial b}$$

$$\frac{\partial z}{\partial w_1} = x_1$$

$$\frac{\partial z}{\partial w_2} = x_2$$

$$\begin{aligned} z &= w \cdot x + b = (w_1, \dots, w_n) \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} + b \\ &= \sum_{i=1}^n w_i x_i \end{aligned}$$

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

$$\frac{\partial \mathcal{L}}{\partial a} = - \left(y \frac{\partial \log a}{\partial a} + (1-y) \frac{\partial \log (1-a)}{\partial a} \right)$$

$$= - \left(y \frac{1}{a} + (1-y) \frac{1}{1-a} \cdot (-1) \right)$$

$$\frac{\partial \mathcal{L}}{\partial z} = - \left(y \frac{1}{a} \frac{\partial a}{\partial z} + (1-y) \frac{1}{1-a} \frac{\partial a}{\partial z} \right)$$

$$= - \left(y \frac{1}{a} a(1-a) + (1-y) \frac{-1}{1-a} a(1-a) \right)$$

$$(1,n) \quad = - y(1-a) - a(1-y) = \boxed{a-y}$$

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{\partial \mathcal{L}}{\partial z} \cdot \frac{\partial z}{\partial w} = \underbrace{(a-y)}_{(1,1)} \underbrace{x^T}_{(1,n)}$$

$$\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial z} \cdot \frac{\partial z}{\partial b} = (a-y) \cdot 1$$

$$\begin{cases} w = w - \alpha (a-y) x^T \\ b = b - \alpha (a-y) \cdot 1 \end{cases}$$

B - Logistic Regression backpropagation for a batch of m examples

* Forward propagation:

$$(1, m) \xrightarrow{x^{(n, m)}} z = Wx + b \xrightarrow{(n, m)} (1, 1) \text{ to } (1, m)$$

broadcasts

$b = \begin{pmatrix} b \\ b \\ \vdots \\ b \end{pmatrix}$

$$(1, m) \xrightarrow{A = \tau(z)} (1, m)$$

$$* \underline{\text{Loss}}: J = \frac{1}{m} \sum_{i=1}^m \mathcal{L}^{(i)}$$

binary cross entropy

$$* \underline{\text{Backprop}}: \frac{\partial J}{\partial w}, \frac{\partial J}{\partial b}$$

* Notations:

$$\omega = (\omega_1, \omega_2, \dots, \omega_n)$$

$$z = (z^{(1)}, z^{(2)}, \dots, z^{(m)})$$

$$A = (a^{(1)}, \dots, a^{(m)})$$

$$z^{(i)} = \sum_{k=1}^n \omega_k x_k^{(i)} + b$$

$$\forall k \in [1, n]$$

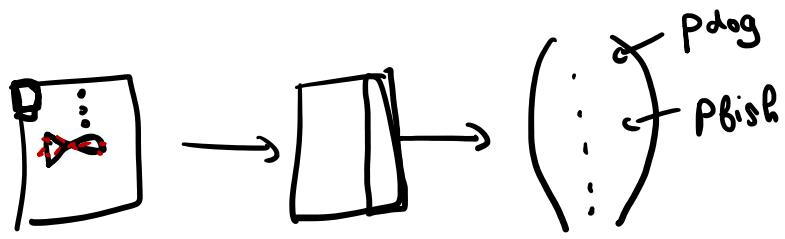
R = 1 ... n

$$\begin{aligned} \frac{\partial J}{\partial \omega_k} &= \frac{1}{m} \sum_{i=1}^m \frac{\partial \mathcal{L}^{(i)}}{\partial \omega_k} = \frac{1}{m} \sum_{i=1}^m y^{(i)} \frac{\partial \log a^{(i)}}{\partial \omega_k} + (1-y^{(i)}) \frac{\partial \log(1-a^{(i)})}{\partial \omega_k} \\ &= \frac{1}{m} \sum_{i=1}^m y^{(i)} \frac{\partial \log a^{(i)}}{\partial z^{(i)}} \cdot \frac{\partial z^{(i)}}{\partial \omega_k} + (1-y^{(i)}) \frac{\partial \log(1-a^{(i)})}{\partial z^{(i)}} \cdot \frac{\partial z^{(i)}}{\partial \omega_k} \\ &= \frac{1}{m} \sum_{i=1}^m (a^{(i)} - y^{(i)}) \cdot \frac{\partial z^{(i)}}{\partial \omega_k} \\ &= \frac{1}{m} \sum_{i=1}^m (a^{(i)} - y^{(i)}) \cdot \cancel{x_k^{(i)}} \underset{(1, 1)}{1} \end{aligned}$$

$$\frac{\partial J}{\partial \omega} = \left(\frac{\partial J}{\partial \omega_1}, \dots, \frac{\partial J}{\partial \omega_n} \right) = \frac{1}{m} (A - Y) \cdot X^T$$

$$\frac{\partial J}{\partial b} = \frac{1}{m} \underbrace{(A - Y)}_{(1, m)} \cdot \underbrace{\begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}}_{(m, 1)} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$$

Question: You have trained an animal classifier. Can you tell what part of the input led to this prediction?



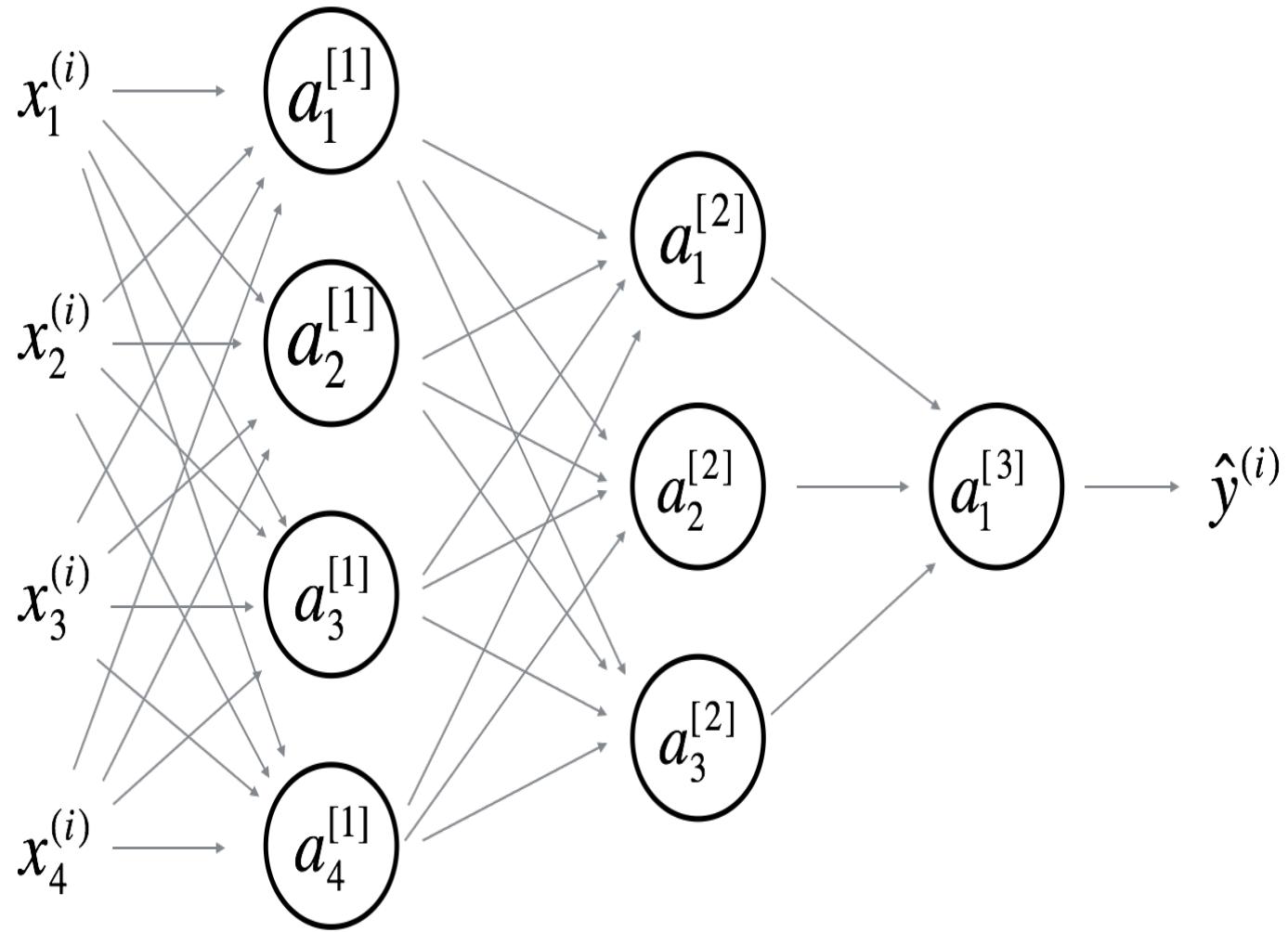
$$\frac{\partial \mathcal{L}}{\partial w}, \frac{\partial \mathcal{L}}{\partial b}$$

$\frac{\partial \mathcal{L}}{\partial x} =$ what is the value of
the impact of shifting x , on \mathcal{L}

$$\frac{\partial \mathcal{L}}{\partial x_{ii}}$$

II - Initializations

Problem statement



Forward propagation (L layers)

$$z^{[1]} = W^{[1]}x + b^{[1]}$$

$$a^{[1]} = g^{[1]}(z^{[1]})$$

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

$$a^{[2]} = g^{[2]}(z^{[2]})$$

...

$$z^{[L]} = W^{[L]}a^{[L-1]} + b^{[L]}$$

$$a^{[L]} = g^{[L]}(z^{[L]})$$

In class, you've seen that:

(i) - not too high / low

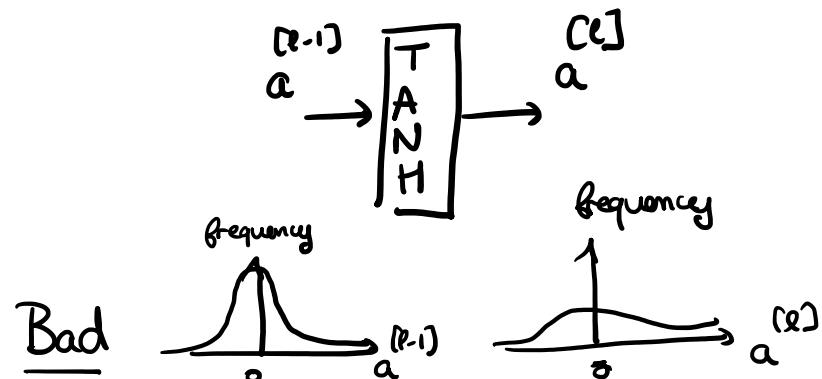
(ii) - breaking the symmetry

for a tanh activation:
at layer l : $\begin{cases} W^{[l]} \\ b^{[l]} = 0 \end{cases} = \text{np.random.randn(shape)} \cdot \text{np.sqrt}\left(\frac{1}{n^{[l-1]}}\right)$

\Leftrightarrow

$$W^{[l]} \sim \mathcal{N}\left(\mu = 0, \sigma^2 = \frac{1}{n^{[l-1]}}\right)$$

The goal: Given a layer ℓ , we'd like $\text{Var}(a^{[\ell-1]}) = \text{Var}(a^{[\ell]})$



$$\underbrace{\text{Var}(a^{[\ell-1]})}_{\Downarrow} = \text{Var}(a^{[\ell]})$$

$$\text{Var}(W^{[\ell]}) = \frac{1}{n^{[\ell-1]}}$$



vanishing and exploding
gradient

$$\text{Var}(a) = \text{Var}(a_1) = \text{Var}(a_2) = \dots$$

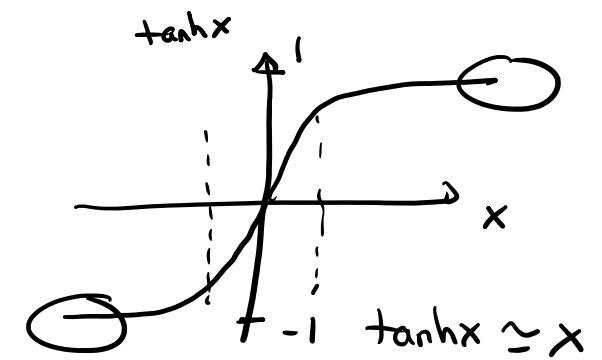
$$a = \begin{pmatrix} a_1 \\ \vdots \\ a_n \end{pmatrix}$$

Let's prove that: || $\text{Var}(a^{[l-1]}) = \text{Var}(a^{[l]}) \rightarrow \text{Var}(W^{[l]}) = \frac{1}{n^{[l-1]}}$

Forward prop $\bar{z}^{(l)} = W^{(l)} a^{(l-1)} + b^{(l)} \Rightarrow ?$

$a^{(l)} = \tanh(z^{(l)}) \Rightarrow a^{(l)} = \tanh(\bar{z}^{(l)})$

$\text{Var}(a^{(l)}) = \text{Var}(\tanh(\bar{z}^{(l)}))$

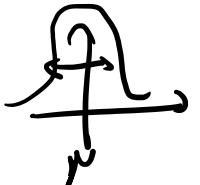


$$\bar{z}^{(l)} = W^{(l)} a^{(l-1)} + b^{(l)} = \begin{pmatrix} w_{11} & \dots & w_{1n}^{(l-1)} \\ \vdots & \ddots & \vdots \\ w_{n1}^{(l)} & \dots & w_{nn}^{(l)} \end{pmatrix} \begin{pmatrix} a_1^{(l-1)} \\ \vdots \\ a_n^{(l-1)} \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^n w_{1i}^{(l)} a_i^{(l-1)} \\ \vdots \\ \sum_{i=1}^n w_{ni}^{(l)} a_i^{(l-1)} \end{pmatrix}$$

$$\text{Var}(a^{(l)}) = \text{Var}(\bar{z}^{(l)}) = \text{Var}\left(\sum_{i=1}^n w_{ii}^{(l)} a_i^{(l-1)}\right) \leftarrow \underbrace{\dots}_{=1} \text{Var}(a_i^{(l-1)})$$

Let's prove that:

$$\text{Var}(a^{[l-1]}) = \text{Var}(a^{[l]}) \rightarrow \text{Var}(W^{[l]}) = \frac{1}{n^{[l-1]}}$$



Checkpoint:

$$\begin{aligned}\text{Var}(XY) &= E[X]^2 \text{Var}Y \\ &\quad + E[Y]^2 \text{Var}X \\ &\quad + \text{Var}X \text{Var}Y\end{aligned}$$

$$\begin{aligned}\text{Var}(a^{[l]}) &= \text{Var}(z^{[l]}) = \text{Var}\left(\sum_{i=1}^{n^{[l-1]}} W_{1,i}^{[l]} a_i^{[l-1]}\right) \\ &= \sum_{i=1}^m \text{Var}(W_{1,i}^{[l]} a_i^{[l-1]}) \\ &\xrightarrow{\substack{i=1 \\ \dots \\ n^{[l-1]}}} \sum_{i=1}^{n^{[l-1]}} \left[\underbrace{\mathbb{E}[W_{1,i}^{[l]}]}_0^2 \text{Var}(a_i^{[l-1]}) + \underbrace{\mathbb{E}[a_i^{[l-1]}]}_0^2 \text{Var}(W_{1,i}^{[l]}) + \text{Var}(W_{1,i}^{[l]}) \text{Var}(a_i^{[l-1]}) \right] \\ &= \sum_{i=1}^{n^{[l-1]}} \underbrace{\text{Var}(W_{1,i}^{[l]})}_{\text{II of } i} \underbrace{\text{Var}(a_i^{[l-1]})}_{\text{II of } i} \\ \underline{\text{Var}(a^{[l]})} &= \underline{n^{[l-1]}} \underline{\text{Var}(W^{[l]})} \underline{\text{Var}(a^{[l-1]})} \implies \text{Var}(W^{[l]}) = \frac{1}{n^{[l-1]}}\end{aligned}$$

Assumptions:

- weights are i.i.d
- inputs are i.i.d
- weights and inputs are mutually independent

We've shown that for every layer l:

$$\boxed{\text{Var}(a^{[l]}) = n^{[l-1]} \text{Var}(W^{[l]}) \text{Var}(a^{[l-1]})}$$

$$\begin{aligned}\text{Var}(a^{[L]}) &= n^{[L-1]} \text{Var}(W^{[L]}) \text{Var}(a^{[L-1]}) \\ &= n^{[L-1]} \text{Var}(W^{[L]}) \left(n^{[L-2]} \cdot \text{Var}(W^{[L-1]}) \cdot \text{Var}(a^{[L-2]}) \right)\end{aligned}$$

$$= \left[\prod_{l=1}^L n^{[l-1]} \cdot \text{Var}(W^{[l]}) \right] \text{Var}(a^{[0]}) = \beta^L \cdot \text{Var}(x)$$

$$\text{Var}(W^{[p]}) = \frac{1}{n^{[p-1]}} \cdot \beta$$

$$\beta > 1$$

