

BLG561E FALL 2021

Deep Learning

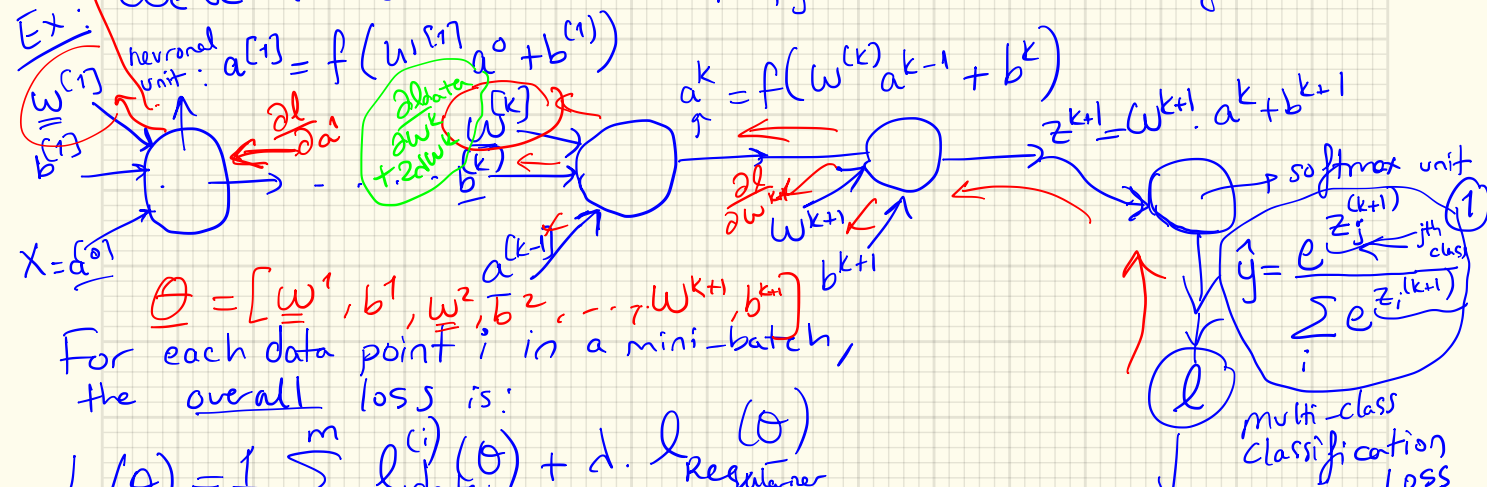
02.11.2021

Görde ÜNAL

$\frac{\partial L}{\partial w} + 2\lambda w = \nabla_{\underline{w}} L$
 Recap: We covered Backprop last time: $w^{(i)(t+1)} \leftarrow w^{(i)(t)} - LR \cdot \nabla_{\underline{w}} L_{\text{overall}}$

Today, we've seen computation graphs & how to backprop over the computation graph.

We've learned how to backpropagate the derivative of the loss fn.



$\theta = [w^1, b^1, w^2, b^2, \dots, w^{(k+1)}, b^{(k+1)}]$

For each data point i in a mini-batch, the overall loss is:

$$L(\theta) = \frac{1}{M} \sum_{i=1}^m l_{\text{data}}^{(i)}(\theta) + d \cdot l_{\text{regularizer}}(\theta)$$

$$d \cdot \|\underline{w}\|_2^2$$

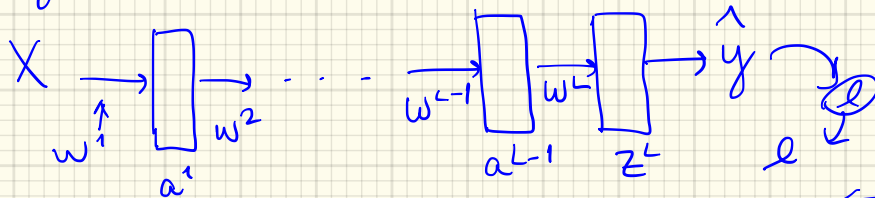
of: $\frac{\partial L_{\text{overall}}}{\partial w^1} = \frac{\partial l_{\text{data}}}{\partial w^1} + 2\lambda w^{(1)}$ *

$$l_{\text{data}}^{(i)} = \sum_{j=1}^c y_j^{(i)} \log \hat{y}_j^{(i)} \quad (2)$$

During backprop, only this gradient propagates back over the computation graph due to its dependence on the intermediate log activations. \therefore We add the regularizer derivative (*) at the update step.

Vanishing / Exploding Gradient Problem: in Deep Networks

Say L Layers in a FCN



Say activation

$$f(z) = z \text{ linear.}$$

Recall

← Forward pass uses \underline{W} multip.

Backward pass uses \underline{W}^T multip.

$$\hat{y} = \underline{W}^L \underline{W}^{L-1} \dots \underline{W}^2 \underline{W}^1 X$$

$$\frac{\partial l}{\partial \underline{W}^1} = \frac{\partial \underline{W}^1}{\partial \underline{W}^1} = (\underline{W}^1)^T (\underline{W}^2)^T \dots (\underline{W}^L)^T \frac{\partial l}{\partial a^L} f'(\cdot)$$

$$\underline{W}^1 \leftarrow \underline{W}^1 \propto d\underline{W}^L$$

Say $\underline{W} = 1.2 \underline{I}$ (say $L=200$ # layers) or $\underline{W} = 0.5 \underline{I}$

$$\hat{y} = (1.2)^{200} \underline{I} \rightarrow d\underline{W}^1 = (1.2)^{200} \underline{I} \text{ (other)}$$

$$d\underline{W}^1 = (0.5)^L \cdot \underline{I} \text{ (other)}$$

explosion $\nearrow \nearrow$

explosion 169

$(0.5) \cdot (0.5) \dots (0.5) \rightarrow 0$
200 vanishing gradient problem

To deal w/ the vanishing (exploding) gradients, we try to initialize our network weights carefully.

Q: How do we initialize Parameters?

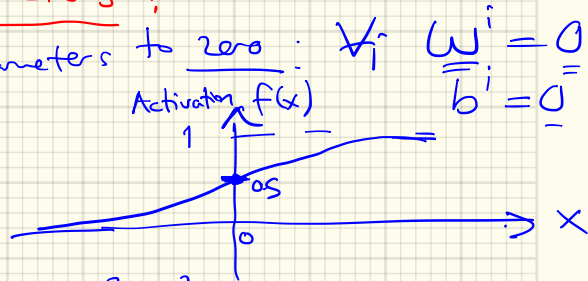
A. What if we initialize all parameters to zero: $\forall_i \underline{w^i} = 0$
 $\underline{b^i} = 0$

eg. 3 layer network

wrong \times

$$z^3 = \underline{w^3} a^2 + \underline{b^3}$$

$$a^3 = f(z^3)$$



$x \rightarrow z^1, a^1 \rightarrow z^2, a^2$

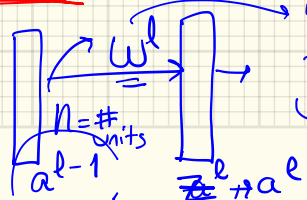
network outputs 0.5 no matter what the input is.

A. What if all parameters are set to some nonzero const. value.

\times

$$a^1 = f(z^1) = f(\underbrace{w^1 x}_{\sim} + \underbrace{b^1}_{\sim}) \rightarrow \begin{pmatrix} \cdot \\ \cdot \\ \cdot \end{pmatrix} \rightarrow \text{all the same const.}$$

① Random small weights w/ small variance:



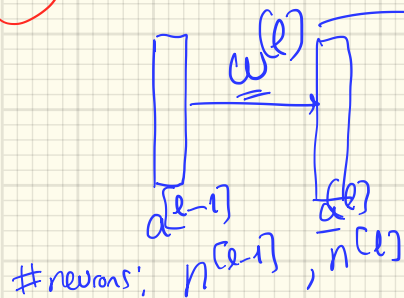
weights at the l^{th} layer:

$$z^l = w_1^l a_1^{l-1} + w_2^l a_2^{l-1} + \dots + w_n^l a_n^{l-1}$$

$$\text{Var}(w^l) \propto \frac{1}{n} \leftarrow \# \text{ units in } (l-1)^{\text{st}} \text{ layer.}$$

→ $\underline{w}^l \sim \mathcal{N}(0, \frac{1}{n^{(l-1)}})$ std $\propto \frac{1}{\sqrt{n}}$

② Xavier-He Initialization:



$\underline{w}^l \sim \mathcal{N}(0, \sqrt{\frac{2}{n^{(l)} + n^{(l-1)}}})$

Widely used in practice.

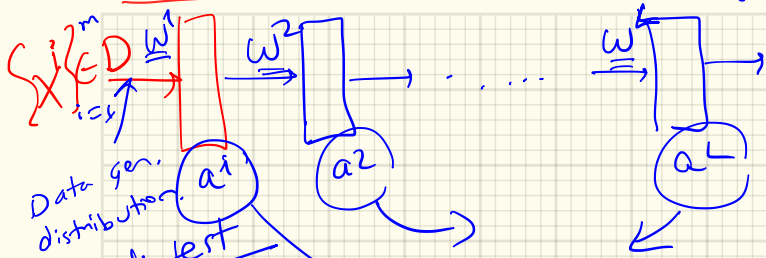
③ There are others such as Glorot initialization.

↑ check!
 network weight (exercise): Go & find out other initialization schemes/suggestions

④ We will talk about Transfer Learning:

then we will initialize our networks w/ an already trained (pretrained) model → their weights are used as our initialization.

BATCH NORMALIZATION (BN) applies normalization to deep layers of your NN.



COVARIATE SHIFT:
Learning on a shifting input distrib. for each layer.

Training & test data may not be coming from the same distrib.

Typically: $\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$, $\sigma^2 = \frac{1}{m-1} \sum_{i=1}^m (x_i - \mu)^2$

We normalize the training & test inputs to the network.

Problem: Intermediate layers in our networks shifting distributions.

BN addresses this problem: Normalization is added before activation, on $f(z)$'s

\therefore For an intermediate layer i , z^i are normalized as:

$$z_{\text{norm}}^i = \frac{z^i - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

\uparrow
 $\mu = \frac{1}{B} \sum_{i=1}^B z^i$; $\sigma^2 = \frac{1}{B} \sum_{i=1}^B (z^i - \mu)^2$
 \uparrow
 # batch samples

→ Our hidden units z_{norm}^i have $(0, \sigma^2)$

→ Now → we make them (γ, β)

$$z^{(i)} = \gamma^{(i)} z_{\text{norm}}^{(i)} + \beta^{(i)}$$

Now, we have 2 more learnable parameters for each layer.

Parameters of our model become:

$$\left. \begin{array}{l} \underline{W}^1, \beta^1, \gamma^1 \\ \underline{W}^2, \beta^2, \gamma^2 \\ \vdots \\ \underline{W}^L, \beta^L, \gamma^L \end{array} \right\} \begin{array}{l} \text{Training time} \\ \text{Add SGD / optimizer:} \\ \beta^i = \beta^i - \alpha d\beta^i \\ \gamma^i = \gamma^i - \alpha d\gamma^i \end{array}$$

Note; W / BN, bias is eliminated:

$$z^l = W^l a^{l-1} + \cancel{b^{l-1}} \quad \text{eliminated}$$

$$z_{\text{norm}}^l = \frac{z^l - \mu}{\sigma}$$

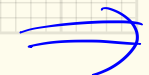
$$z^l = \gamma^i z_{\text{norm}}^l + \beta^l$$

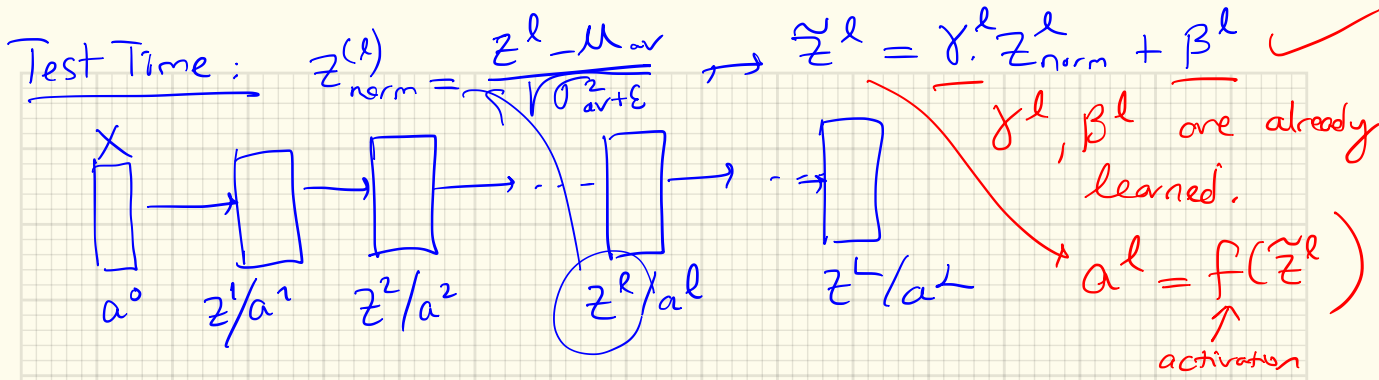
BN at Test Time: We need μ, σ^2 estimated through weighted averaging over mini-batches during training:

$$\left. \begin{array}{l} \mu_{\text{av}}^+ = \alpha \mu^{t-1} + (1-\alpha) \mu^t \\ \sigma_{\text{av}}^+ = \dots \end{array} \right\}$$

produce an overall

$$\mu_{\text{av}} \& \sigma_{\text{av}}^2$$





pytorch \rightarrow torch.nn.BatchNorm \leftarrow

Dropout: ✓ creates for us **Ensemble of Networks.**
 check, **MCDrop-out:** \rightarrow come up w/ confidence measures

w/ prob. p : Kill / randomly eliminate some nodes / hidden units w/ a dropout prob of p .

Helps w/ **overfitting**

Inverted Drop-out: In training: scale up; $1/p$ (to match the test time scales)

\rightarrow st. Output at test time = expected output at Training time

Test time: B/c all neurons are active
 we must scale the activations by p ,
 so that

Note: In today's lecture, we also went thru slides from

cs231n ~~X~~ some slides from
(Stanford)

Francois Fleuret's Deep learning course
at EPFL.

(You can study relevant material, if you want, from
those websites.)

Today,
we
covered:

- Batch Norm, Computation graphs,
- Dropout, vanishing gradients,
- activation functions,
- initializing NN weights.