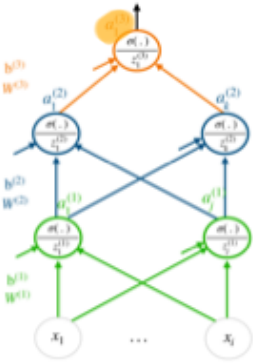


Lecture 21 - 14/11/24

EXAMPLE ANN WITH SEVERAL HIDDEN LAYERS:



It is for **BINARY CLASSIFICATION**

Its **LOSS FUNCTION** is:

$$L(z^{(3)}, y) = -(y \cdot \log(z) + (1-y) \cdot \log(1-z))$$

What we do now is:

Back Propagation!

$$L(z^{(3)}, y) = -(y \cdot \log(z^{(3)}) + (1-y) \cdot \log(1-z^{(3)}))$$

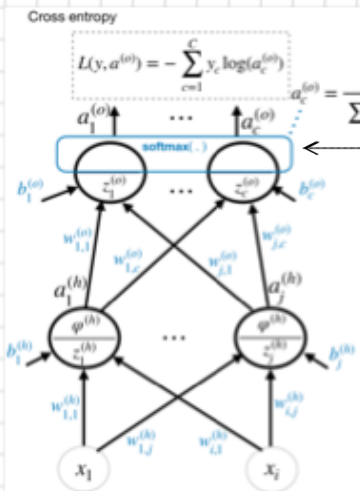
$$\begin{aligned} \frac{\partial L}{\partial z^{(3)}} &= - \left[y \cdot \frac{1}{z^{(3)}} + (1-y) \cdot \left(\frac{1}{1-z^{(3)}} \cdot -1 \right) \right] \\ &= - \left[\frac{y}{z^{(3)}} + \frac{(1-y)}{1-z^{(3)}} \right] = - \frac{y(1-z^{(3)}) + z^{(3)}(1-y)}{z^{(3)}(1-z^{(3)})} \\ &= - \frac{y - yz^{(3)} - z^{(3)} + z^{(3)}y}{z^{(3)}(1-z^{(3)})} = - \frac{-y + z^{(3)}}{z^{(3)}(1-z^{(3)})} \end{aligned}$$

$$\begin{aligned} \frac{\partial L}{\partial z^{(3)}} &= - \frac{-y + z^{(3)}}{z^{(3)}(1-z^{(3)})} \\ \frac{\partial L}{\partial z^{(3)}} &= \frac{\partial L}{\partial a^{(3)}} \cdot \frac{\partial a^{(3)}}{\partial z^{(3)}} \\ \frac{\partial L}{\partial a^{(3)}} &= \frac{\partial L}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial a^{(3)}} \\ \frac{\partial L}{\partial a^{(3)}} &= \frac{\partial L}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial a^{(3)}} \\ \frac{\partial L}{\partial a^{(3)}} &= \frac{\partial L}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial a^{(3)}} \\ \frac{\partial L}{\partial a^{(3)}} &= \frac{\partial L}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial a^{(3)}} \end{aligned}$$

EXAMPLE ANN WITH SOFTMAX OUTPUT LAYER

KNOW

SOFTMAX MAPS OUTPUT INTO PROBABILITIES!



$$\psi^{(1)}(.) = \text{softmax}(.)$$

$$\psi^{(h)}(.) = \sigma(.) \text{ (Sigmoid or ReLU)}$$

$$z^{(10)} = \frac{e^{z_c^{(10)}}}{\sum_{k=1}^C e^{z_k^{(10)}}}$$

$$L(y, z^{(10)}) = - \sum_{c=1}^C y_c \log \left(\frac{e^{z_c^{(10)}}}{\sum_{k=1}^C e^{z_k^{(10)}}} \right)$$

ALL RAW OUTPUTS

$$\frac{\partial L}{\partial z_c^{(10)}} = a_c^{(10)} - y_c$$

$$\begin{aligned} \frac{\partial L}{\partial w_{jc}^{(10)}} &= \frac{\partial L}{\partial z_c^{(10)}} \cdot a_j^{(10)} \\ \frac{\partial L}{\partial b_c^{(10)}} &= \frac{\partial L}{\partial z_c^{(10)}} \end{aligned}$$

$$\frac{\partial L}{\partial z_j^{(10)}} = \sum_c \left(\frac{\partial L}{\partial z_c^{(10)}} \cdot \frac{\partial z_c^{(10)}}{\partial z_j^{(10)}} \right)$$

$$\frac{\partial L}{\partial z_j^{(10)}} = \frac{\partial L}{\partial z_j^{(10)}} \cdot \frac{\partial z_j^{(10)}}{\partial z_j^{(10)}}$$

$$\begin{aligned} \frac{\partial L}{\partial w_{ij}^{(10)}} &= \frac{\partial L}{\partial z_j^{(10)}} \cdot x_i \\ \frac{\partial L}{\partial b_j^{(10)}} &= \frac{\partial L}{\partial z_j^{(10)}} \end{aligned}$$

REGULARIZATION

1. NEED TO TALK ABOUT **Bias & VARIANCE**

High Variance

• **Training Error** is much **lower** than **validation Error**.

• We have to **reduce model complexity**.

• We can do **bagging**, **under-training**, **add training data**.

ACCEPTABLE error

High Bias

• **Training Error** **higher** than **validation Error**.

• Use **more complex model**.

• **Add new features**, **boosting**.

1. We can apply a VARIETY OF REGULATION METHODS IN THE LEARNING OF ANN.
2. Play a VERY IMPORTANT ROLE SPECIALLY WHEN TRAINING SET IS SMALL.
 - WEIGHT DECAY
 - EARLY STOPPING
 - DROPOUT
 - DATA AUGMENTATION

WEIGHT DECAY

- ADD L^p NORM-REGULATION TERMS TO THE CRSECTIVE FUNCTION
- WITH L^2 NORM: INFORMALLY CALLED WEIGHT DECAY BECAUSE:
 - TRIES TO PUSH MAGNITUDE OF W TOWARDS ZERO!

WITHOUT REGULATION

Loss function: $Q(W)$

Gradient Descent update: $W := W - \alpha \frac{\partial Q(W)}{\partial W}$

WITH L^2 REGULATION

Loss Function: $Q(W) + \frac{\lambda}{2} \|W\|_2^2$

GD update: $W := W - \alpha \frac{\partial Q(W)}{\partial W} - \lambda \cdot W$

- THERE IS NO REGULATION FOR BIASES:
 - THEY DON'T CHANGE THE SHAPE!

DROPOUT REGULARIZATION

- Injecting NOISE IN LEARNING PROCESS
- \forall TRAINING SAMPLE (IN MIN BATCH):
 - MAKE A THINNED NETWORK BY RANDOMLY DROPPING OUT NODES
 - FORWARD & BACK-PROPAGATION ARE DONE ON THINNED NETWORK
- GRADIENTS FOR UPDATING EACH WEIGHT ARE AVERAGED OVER THE TRAINING SAMPLES.
 - THE DROPPED NODES REPORT A GRADIENT=0 FOR THEIR WEIGHTS.
- EACH NODE HAS THE CHANCE TO BE DROPPED EVERY TIME, FORCED TO REDUNDANT REPRESENTATION.
 - THIS PREVENTS OVERFITTING.
- LESS COMPLEX NETWORKS \rightarrow SUCH NOISE \rightarrow PREVENT OVERFITTING
- THE CONCEPT IS THAT FULL NETWORK LEARNS ON SCALED-DOWN VERSION OF ITSELF.
- FOR PREDICTION THE FULL NETWORK IS USED.

EARLY STOPPING

- WHEN PERFORMANCE ON VALIDATION SET IS BEING WORSE \rightarrow STOP!

DATA AUGMENTATION

- GENERATING MORE ALTERNATIVE TRAINING SAMPLES FROM THE RAW TRAINING DATA.
 - Inject small noise into RAW DATA — ALTERING



GD WITH MOMENTUM

- ADD A HISTORY OF THE GRADIENTS
- GIVE MORE WEIGHT TO MORE RECENT GRADIENTS

$$v_{t+1} = \beta v_t + \alpha \nabla L(w_t)$$

$$w_{t+1} = w_t - v_{t+1}$$

- GETS MOMENTUM & ACCELERATES IF THE GRADIENT WAS IN THE SAME DIRECTION
- DECREASES IF DIRECTION WAS MADE WITH LARGE STEPS.