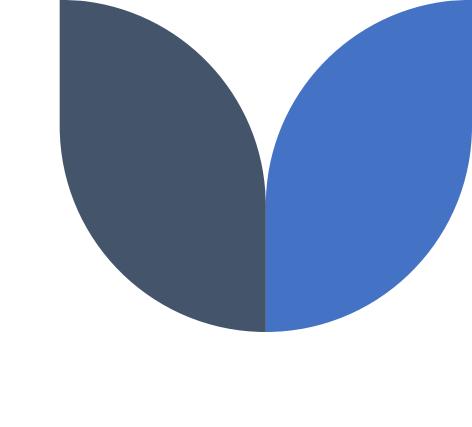
An Overview on Deep Neural Networks: Part 2



Backward Propagation

Calculating the Loss and Cost

$$A^{[0]} = X$$

$$Z^{[l]} = A^{[l-1]}W^{[l]} + b^{[l]}$$

$$A^{[l]} = \sigma^{[l]}(Z^{[l]})$$

$$\hat{Y} = A^{[L]}$$

$$C = \frac{1}{n}J(Y,\widehat{Y})$$

$$\forall l = 1, 2, ..., L$$

$$\forall l = 1,2,...,L$$

Gradient Descent and Chain Rule

$$A^{[l]} = \sigma^{[l]}(Z^{[l]})$$

$$\hat{Y} = A^{[L]}$$

$$C = \frac{1}{n}J(Y,\widehat{Y})$$

$$\frac{\partial C}{\partial \hat{Y}} = \frac{1}{n} \frac{\partial J}{\partial \hat{Y}}$$

$$\frac{\partial C}{\partial A^{[L]}} = \frac{\partial C}{\partial \hat{Y}}$$

$$\frac{\partial A^{[l]}}{\partial Z^{[l]}} = \sigma'^{[l]} (Z^{[l]})$$

$$\forall l = 1, 2, ..., L$$

$$\forall l = L, L - 1, ..., 1$$

Gradient Descent and Chain Rule

$$Z^{[l]} = A^{[l-1]}W^{[l]} + b^{[l]}$$

$$\forall l = 1, 2, ..., L$$

$$\frac{\partial Z^{[l]}}{\partial W^{[l]}} = A^{[l-1],T}$$

$$\forall l = L, L-1, ..., 1$$

$$\frac{\partial Z^{[l]}}{\partial b^{[l]}} = 1_{1 \times n}$$

$$\forall l = L, L-1, ..., 1$$

$$\frac{\partial Z^{[l]}}{\partial A^{[l-1]}} = W^{[l],T}$$

$$\forall l = L, L-1, \dots, 2$$

Putting It All Together

$$\bullet \quad \frac{\partial C}{\partial A^{[L]}} = \frac{1}{n} \frac{\partial J}{\partial \hat{Y}}$$

•
$$\frac{\partial C}{\partial W^{[l]}} = \frac{\partial Z^{[l]}}{\partial W^{[l]}} \left(\frac{\partial C}{\partial A^{[l]}} : \frac{\partial A^{[l]}}{\partial Z^{[l]}} \right) = A^{[l-1],T} \left(\frac{\partial C}{\partial A^{[l]}} : \sigma'^{[l]} (Z^{[l]}) \right)$$

•
$$\frac{\partial C}{\partial b^{[l]}} = \frac{\partial Z^{[l]}}{\partial b^{[l]}} \left(\frac{\partial C}{\partial A^{[l]}} : \frac{\partial A^{[l]}}{\partial Z^{[l]}} \right) = 1_{1 \times n} \left(\frac{\partial C}{\partial A^{[l]}} : \sigma'^{[l]} (Z^{[l]}) \right)$$

•
$$\frac{\partial C}{\partial A^{[l-1]}} = \left(\frac{\partial C}{\partial A^{[l]}} : \frac{\partial A^{[l]}}{\partial Z^{[l]}}\right) \frac{\partial Z^{[l]}}{\partial A^{[l-1]}} = \left(\frac{\partial C}{\partial A^{[l]}} : \sigma'^{[l]}(Z^{[l]})\right) W^{[l],T}$$

$$\forall l = L, L-1, ..., 1$$

$$\forall l = L, L-1, \dots, 1$$

$$\forall l = L, L - 1, ..., 2$$

Updating the Weights and the Biases

$$\Delta W^{[l]} = -\alpha \frac{\partial C}{\partial W^{[l]}}$$

$$\Delta b^{[l]} = -\alpha \frac{\partial C}{\partial b^{[l]}}$$

$$\forall l = L, L-1, \dots, 1$$

$$\forall l = L, L-1, \dots, 1$$

Loss Function Examples

Loss Function Examples

- Sum of Squared Errors (SSE)
- Binary Cross Entropy (BCE)
- Categorical Cross Entropy (CCE)

SSE

- Commonly used in regression.
- The sum is divided by two for easier derivative calculation.

$$J(Y, \hat{Y}) = \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{m_y} (\hat{Y}_{i,j} - Y_{i,j})^2$$

BCE

• Commonly used in binary or multi-label classification.

$$J(Y,\widehat{Y}) = -\sum_{i=1}^{n} \sum_{j=1}^{m_y} \left(Y_{i,j} \log \left(\widehat{Y}_{i,j} \right) + \left(1 - Y_{i,j} \right) \log \left(1 - \widehat{Y}_{i,j} \right) \right)$$

CCE

• Commonly used in multi-class classification.

$$J(Y,\widehat{Y}) = -\sum_{i=1}^{n} \sum_{j=1}^{m_{y}} (Y_{i,j} \log(\widehat{Y}_{i,j}))$$

Summary of Activation and Loss Functions

Summary of Activation and Loss Functions

Task	Output Layer Activation Function	Loss Function
Regression	ReLU	SSE
Binary of Multi-Label Classification	Sigmoid	BCE
Multi-Class Classification	Softmax	CCE