Activations, Problematic Gradients, Initialization, and Normalization

We want all values to be around zero. [-1,1]

Outline

- Activations (output of an activation function)
- Problematic gradients (exploding and vanishing)
- The benefits of depth in neural networks
- Proper parameter initialization
- Normalization (preprocessing) of inputs and activations (a type of layer)

Recap: Optimization Techniques

• Take five minutes to draw

Purpose of Assignments

I provide most (if not all code) on assignments for 2 reasons: Z = A W + b

- 1. This provides a nice set of working examples you can use for projects
- 2. It keeps assignments shorter so that you have more time for projects

You should end up spending more time reading my code and more time writing your own code.



Activation Functions, Two-Layer Network $A^{(2)} = A^{(0)} M + D$ Liner. Without Activation Functions **With Activation Functions** $Z^{[1]} = A^{[0]}W^{[1]T} + b^{[1]}$ $Z^{[1]} = A^{[0]} W^{[1]T} + h^{[1]}$ $A^{[1]} = Z^{[1]}$ $A^{[1]} = \sigma(Z^{[1]})$ $Z^{[2]} = A^{[1]} W^{[2]T} + h^{[2]}$ $Z^{[2]} = A^{[1]}W^{[2]T} + b^{[2]}$ $A^{[2]} = Z^{[2]}$ $A^{[2]} = \sigma(Z^{[2]})$ $= A^{CIJ}W^{C2JT} + h^{C2J}$ $= \mathcal{O}\left(\mathsf{Y}_{\mathbb{C},\mathbb{I}}^{\mathcal{I}}\mathsf{N}_{\mathbb{C},\mathbb{I}}^{\mathcal{I}} + \mathsf{P}_{\mathbb{C},\mathbb{I}}^{\mathcal{I}}\right)$ = 2 CD W CDT + h (2) $\mathcal{L}(\mathcal{Q}(5_{12})) \mathcal{M}_{1521} + \mathcal{P}_{152})$ $= \left(A^{(0)} W^{(1)T} + b^{(1)} \right) W^{(2)T} +$ $= \mathcal{O}\left(\mathcal{O}\left(A^{CO}\mathcal{W}^{CIT} + \mathcal{O}^{CI}\right)\right)$ $= A^{LOJ} W^{CIJT} L^{ZJT} + b^{LIJ} W^{C2JI} + b^{CIJ}$ 1, [2]T + 6[2]



Sigmoid and ReLU Activation Functions Most intersting around $O: \Delta, \Box$ Sigmoid ReLU $\sigma(z) = \frac{1}{1 + e^{-z}}$ ReLu'(z) = 0 if z < 0ReLu'(z) = 1 if $z \ge 0$ $\operatorname{ReLu}(z) = \max(0,z)$ $=\sigma(z)(1-\sigma(z))$ $\sigma(\mathbf{v})$ -6 -2 0 0 6 in E-1, 13 range Inputs



Fixes for Problematic Gradients







PyTorch Kaiming He Initialization $W = \frac{1}{3} + \frac{1}{3$

```
torch.nn.init.kaiming_uniform_(tensor, a=0, mode='fan_in',
nonlinearity='leaky_relu') [SOURCE]
```

Fills the input Tensor with values according to the method described in Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification - He, K. et al. (2015), using a uniform distribution. The resulting tensor will have values sampled from $\mathcal{U}(-\text{bound}, \text{bound})$ where

$$pound = gain \times \sqrt{fan mode}$$

gain = caiculate_gain(nonlinearity, a)
std = gain / math.sqrt(fan)
bound = math.sqrt(3.0) * std # Calculate uniform bounds from standard deviation
with torch.no_grad():
 return tensor.uniform_(-bound, bound)

Batch Normalization

- What is normalization?
 - Adjusting values to a different (common) scale.
 - For example, the standard-score normalization: $\bar{x} = \frac{x-\mu}{\sigma}$
- Normalize with respect to what?



Summary

- Activation functions behave nicely with inputs around zero
- Gradients behave nicely with values around zero (but not too small)
- Depth is good for extracting/finding complex features
- You should
 - Normalize input features
 - Initialize parameters properly
 - Prefer deeper over wider networks
 - Try/consider batch normalization