

1

Lecture 3 recap

Beyond linear

• Linear score function f = Wx



On CIFAR-10



Neural Network



Neural Network



Neural Network

• Linear score function f = Wx

• Neural network is a nesting of 'functions'

- 2-layers:
$$f = W_2 \max(0, W_1 x)$$

- 3-layers: $f = W_3 \max(0, W_2 \max(0, W_1 x))$
- 4-layers: $f = W_4 \tanh(W_3, \max(0, W_2 \max(0, W_1 x)))$
- 5-layers: $f = W_5 \sigma(W_4 \tanh(W_3, \max(0, W_2 \max(0, W_1 x))))$
- ... up to hundreds of layers

Computational Graphs

- Neural network is a computational graph
 - It has compute nodes
 - It has edges that connect nodes
 - It is directional
 - It is organized in 'layers'



Backprop con't

The importance of gradients

• All optimization schemes are based on computing gradients

$$\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$$

• One can compute gradients analytically but what if our function is too complex?

Break down gradient computation

Backpropagation

Backprop: Forward Pass

•
$$f(x, y, z) = (x + y) \cdot z$$

Initialization x = 1, y = -3, z = 4





What is
$$\frac{\partial f}{\partial x}$$
, $\frac{\partial f}{\partial y}$, $\frac{\partial f}{\partial z}$?











- x_k input variables
- $w_{l,m,n}$ network weights (note 3 indices)
 - I which layer
 - m which neuron in layer
 - n weights in neuron
- y_i computed output (i output dim; nout)
- *t_i* ground truth targets
- *L* is loss function



e.g., class label / regression target



We want to compute gradients w.r.t. all weights w



L2 Loss function

 $-t_0$

We want to compute gradients w.r.t. all weights w

Loss/cost

x*x



Input layer



 $L_i = (y_i - t_i)^2$





L2 loss -> simply sum up squares Energy to minimize is E=L

$$\frac{\partial L}{\partial w_{i,k}} = \frac{\partial L}{\partial y_i} \cdot \frac{\partial y_i}{\partial w_{i,k}}$$

-> use chain rule to compute partials

We want to compute gradients w.r.t. all weights w *AND* biases b



Gradient Descent

• From derivative to gradient

[x]

dx

Direction of greatest
 increase of the function

• Gradient steps in direction of negative gradient



For a given training pair {x,t}, we want to update all weights; i.e., we need to compute derivatives w.r.t. to all weights





Just go through layer by layer



Backpropagation

$$\frac{\partial L}{\partial w_{1,i,j}} = \frac{\partial L}{\partial y_i} \cdot \frac{\partial y_i}{\partial w_{1,i,j}}$$

$$\frac{\partial L}{\partial w_{0,j,k}} = \frac{\partial L}{\partial y_i} \cdot \frac{\partial y_i}{\partial h_j} \cdot \frac{\partial h_j}{\partial w_{0,j,k}}$$

$$\frac{\partial L_i}{\partial y_i} = 2(y_i - t_i) \qquad \dots$$

$$\frac{\partial y_i}{\partial w_{1,i,j}} = h_j \quad \text{if } > 0, \text{ else } 0$$



How many unknown weights?



Note that some activations have also weights







Derivatives of Neural Networks



Linear activation node + hinge loss + regularization

Can become quite complex...

input layer 1 hidden layer 2 hidden layer 3



Can become quite complex...

• These graphs can be huge!



Convolution AvgPool MaxPool Concat Dropout Fully connected Softmax

Another view of GoogLeNet's architecture.

The flow of the gradients



The flow of the gradients



Activation function

The flow of the gradients

Many many many many of these nodes form a neural network

NEURONS

• Each one has its own work to do

FORWARD AND BACKWARD PASS



Gradient descent



Gradient Descent

• From derivative to gradient

[x]

dx

Direction of greatest
 $\nabla_{\mathbf{x}} f(\mathbf{x})$ Direction of greatest
 the function

• Gradient steps in direction of negative gradient



Gradient Descent

• How to pick good learning rate?

• How to compute gradient for single training pair?

• How to compute gradient for large training set?

How to speed things up ☺

Next lecture

- This week:
 - Exercise 1 starts next Tuesday (Nov 20th)
 - Exercise session will introduce it!

- Next lecture on Nov 22nd:
 - Optimization of Neural Networks
 - In particular, introduction to SGD (our main method!)

See you next week!