CS 1678: Intro to Deep Learning Neural Network Basics

Prof. Adriana Kovashka University of Pittsburgh January 28, 2021

Plan for this lecture (next few classes)

- Definition
 - Architecture
 - Basic operations
 - Biological inspiration
- Goals
 - Loss functions
- Training
 - Gradient descent
 - Backpropagation
- Hands-on exercise

Definition

Neural network definition

Figure 5.1 Network diagram for the twolayer neural network corresponding to (5.7). The input, hidden, and output variables are represented by nodes, and the weight parameters are represented by links between the nodes, in which the bias parameters are denoted by links coming from additional input and hidden variables x_0 and z_0 . Arrows denote the direction of information flow through the network during forward propagation.



- Raw activations: $a_j = \sum_{i=1}^{j} w_{ji}^{(1)} x_i + w_{j0}^{(1)}$
- Nonlinear activation function h (e.g. sigmoid, tanh, RELU): $z_j = h(a_j)$ e.g. z = RELU(a) = max(0, a)



Activation functions



Leaky ReLU $\max(0.1x, x)$



 $\begin{array}{l} \textbf{Maxout} \\ \max(w_1^T x + b_1, w_2^T x + b_2) \end{array}$



A multi-layer neural network...



- Is a non-linear classifier
- Can approximate any continuous function to arbitrary accuracy given sufficiently many hidden units

Inspiration: Neuron cells

- Neurons
 - accept information from multiple inputs
 - transmit information to other neurons
- Multiply inputs by weights along edges
- Apply some function to the set of inputs at each node
- If output of function over threshold, neuron "fires"



Biological analog



Jia-bin Huang

Biological analog



- Cascade neurons together
- Output from one layer is the input to the next
- Each layer has its own sets of weights



• Inputs multiplied by initial set of weights



 Intermediate "predictions" computed at first hidden layer



- Intermediate predictions multiplied by second layer of weights
- Predictions are fed forward through the network to classify



Compute second set of intermediate predictions



• Multiply by final set of weights



• Compute output (e.g. probability of a particular class being present in the sample)



Deep neural networks

- Lots of hidden layers
- Depth = power (usually)



Goals

How do we train deep networks?

- No closed-form solution for the weights (can't set up a system A*w = b, solve for w)
- We will iteratively find such a set of weights that allow the outputs to match the desired outputs
- We want to minimize a loss function (a function of the weights in the network)
- For now, let's simplify and assume there's a single layer of weights in the network, and no activation function (i.e., output is a linear combination of the inputs)

Classification goal

airplane	the second second	-	X	*	+	2			-
automobile				-	The state			-	*
bird		12	X		-	1		- Se	4
cat			50		1		Å,	1 and a start	
deer	1 20	X	R		Y	Ŷ	N.	-	
dog	17 A.	X		1		-	C?	1	The second
frog				2			S.		5
horse		A	2	P	HCAL	1	No.	6	1
ship	2	diri:	-	144	-	2	12	10-1	
truck							1		ALL.

Example dataset: CIFAR-10 10 labels 50,000 training images each image is 32x32x3 10,000 test images.

Classification scores

$$f(x,W) = Wx$$

 $f(\mathbf{x},\mathbf{W})$



10 numbers, indicating class scores

[32x32x3] array of numbers 0...1 (3072 numbers total)

Andrej Karpathy

Linear classifier



Example with an image with 4 pixels, and 3 classes (cat/dog/ship)



Andrej Karpathy

Linear classifier

Going forward: Loss function/Optimization



TODO:

- 1. Define a **loss function** that quantifies our unhappiness with the scores across the training data.
- 2. Come up with a way of efficiently finding the parameters that minimize the loss function.
 (optimization)

Linear classifier

Suppose: 3 training examples, 3 classes. With some W the scores f(x, W) = Wx are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1

Suppose: 3 training examples, 3 classes. With some W the scores f(x, W) = Wx are:



cat	3.2	1.3	2.2		
car	5.1	4.9	2.5		
frog	-1.7	2.0	-3.1		



Hinge loss:

Given an example (x_i, y_i) where x_i is the image and where y_i is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the loss has the form:

$$L_i = \sum_{j
eq y_i} \max(0, s_j - s_{y_i} + 1)$$

Want:
$$S_{y_i} \ge S_j + 1$$

i.e. $S_j - S_{y_i} + 1 \le 0$

If true, loss is 0 If false, loss is magnitude of violation

Adapted from Andrej Karpathy

Suppose: 3 training examples, 3 classes. With some W the scores f(x, W) = Wx are:



Hinge loss:

Given an example (x_i, y_i) where x_i is the image and where y_i is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the loss has the form:

 $\begin{aligned} L_i &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \\ &= \max(0, 5.1 - 3.2 + 1) \\ &+ \max(0, -1.7 - 3.2 + 1) \\ &= \max(0, 2.9) + \max(0, -3.9) \\ &= 2.9 + 0 \\ &= 2.9 \end{aligned}$

Suppose: 3 training examples, 3 classes. With some W the scores f(x, W) = Wx are:



Hinge loss:

Given an example (x_i, y_i) where x_i is the image and where y_i is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the loss has the form:

 $\begin{aligned} L_i &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \\ &= \max(0, 1.3 - 4.9 + 1) \\ &+ \max(0, 2.0 - 4.9 + 1) \\ &= \max(0, -2.6) + \max(0, -1.9) \\ &= 0 + 0 \\ &= 0 \end{aligned}$

Suppose: 3 training examples, 3 classes. With some W the scores f(x, W) = Wx are:



Hinge loss:

Given an example (x_i, y_i) where x_i is the image and where y_i is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

= max(0, 2.2 - (-3.1) + 1)
+max(0, 2.5 - (-3.1) + 1)
= max(0, 5.3 + 1)
+ max(0, 5.6 + 1)
= 6.3 + 6.6
= 12.9

Suppose: 3 training examples, 3 classes. With some W the scores f(x, W) = Wx are:



cat	3.2	1.3	2.2	
car	5.1	4.9	2.5	
frog	-1.7	2.0	-3.1	
Losses:	2.9	0	12.9	-

Hinge loss:

Given an example (x_i, y_i) where x_i is the image and where y_i is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the loss has the form:

$$L_i = \sum_{j
eq y_i} \max(0, s_j - s_{y_i} + 1)$$

and the full training loss is the mean over all examples in the training data:

$$L = rac{1}{N} \sum_{i=1}^N L_i$$

L = (2.9 + 0 + 12.9)/3 = 15.8 / 3 = **5.3**

f(x,W) = Wx

 $L = rac{1}{N} \sum_{i=1}^{N} \sum_{j
eq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1)$

E.g. Suppose that we found a W such that L = 0. Is this W unique?

No! 2W is also has L = 0! How do we choose between W and 2W?

Slide from Fei-Fei, Johnson, Yeung

Weight Regularization

 λ = regularization strength (hyperparameter)

$$L(W) = \frac{1}{N} \sum_{i=1}^{N} L_i(f(x_i, W), y_i) + \lambda R(W)$$

Data loss: Model predictions should match training data

Regularization: Prevent the model from doing *too* well on training data

Simple examples

L2 regularization: $R(W) = \sum_{k} \sum_{l} W_{k,l}^2$ L1 regularization: $R(W) = \sum_{k} \sum_{l} |W_{k,l}|$ Elastic net (L1 + L2): $R(W) = \sum_{k} \sum_{l} \beta W_{k,l}^2 + |W_{k,l}|$

More complex:

Dropout

Batch normalization

|| Stochastic depth / pooling, etc

Why regularize?

- Express preferences over weights
- Make the model simple so it works on test data

Weight Regularization

Expressing preferences

$$x = [1, 1, 1, 1] \ w_1 = [1, 0, 0, 0]$$

$$R(W) = \sum_k \sum_l W_{k,l}^2$$

$$w_2 = \left[0.25, 0.25, 0.25, 0.25
ight]$$

L2 regularization likes to "spread out" the weights

$$w_1^T x = w_2^T x = 1$$

Slide from Fei-Fei, Johnson, Yeung

Weight Regularization

Preferring simple models



Regularization pushes against fitting the data *too* well so we don't fit noise in the data

Another loss: Cross-entropy



3.2

5.1

-1.7

scores = unnormalized log probabilities of the classes.

$$P(Y=k|X=x_i)=rac{e^{s_k}}{\sum_j e^{s_j}}$$
 where $oldsymbol{s}=f(x_i;W)$

Want to maximize the log likelihood, or (for a loss function) to minimize the negative log likelihood of the correct class:

$$L_i = -\log P(Y=y_i|X=x_i)$$

Andrej Karpathy

cat

car

frog
Another loss: Cross-entropy



unnormalized probabilities

Aside:

- This is multinomial logistic regression
- Choose weights to maximize the likelihood of the observed x/y data

Adapted from Fei-Fei, Johnson, Yeung

(Maximum Likelihood Estimation; more discussion in CS 1675)

Another loss: Cross-entropy



Other losses

• Triplet loss (Schroff, FaceNet, CVPR 2015)

$$\sum_{i=1}^{N} \left[\|f(x_{i}^{a}) - f(x_{i}^{p})\|_{2}^{2} - \|f(x_{i}^{a}) - f(x_{i}^{n})\|_{2}^{2} + \alpha \right]_{4}$$

a denotes anchor p denotes positive n denotes negative



Figure 3. The **Triplet Loss** minimizes the distance between an *an-chor* and a *positive*, both of which have the same identity, and maximizes the distance between the *anchor* and a *negative* of a different identity.

Anything you want! (almost)

Training

To minimize loss, use gradient descent



Andrej Karpathy

How to minimize the loss function?

In 1-dimension, the derivative of a function:

$$rac{df(x)}{dx} = \lim_{h o 0} rac{f(x+h) - f(x)}{h}$$

In multiple dimensions, the gradient is the vector of (partial derivatives).

Loss gradients

- Denoted as (diff notations): $\frac{\partial E}{\partial w_{ji}^{(1)}} = \nabla_W L$
- i.e. how does the loss change as a function of the weights
- We want to change the weights in such a way that makes the loss decrease as fast as possible <u>1</u>



current W:	
[0.34, -1.11, 0.78,	
0.12, 0.55, 2.81,	
-3.1, -1.5, 0.33,…] loss 1.25347	





current W:	W + h (first dim):	gradient dW:
[0.34,	[0.34 + 0.0001 ,	[?,
-1.11,	-1.11,	?,
0.78,	0.78,	?.
0.12,	0.12,	?
0.55,	0.55,	?
2.81,	2.81,	?
-3.1,	-3.1,	?
-1.5,	-1.5,	?
0.33,]	0.33,]	?]
loss 1.25347	loss 1.25322	

current W:	
[0.34,	
-1.11,	
0.78,	
0.12,	
0.55,	
2.81,	
-3.1,	
-1.5,	
0.33,]	
loss 1.25347	

W + h (first dim): [0.34 + **0.0001**, -1.11, 0.78, 0.12, 0.55, 2.81, -3.1, -1.5, 0.33,...] loss 1.25322



current W:	W + h (
[0.34,	[0.34,
-1.11,	-1.11 +
0.78,	0.78,
0.12,	0.12,
0.55,	0.55,
2.81,	2.81,
-3.1,	-3.1,
-1.5,	-1.5,
0.33,]	0.33,]
loss 1.25347	loss 1.2

second dim): 0.0001, 25353

gradient dW:

[-2.5, ?, ?, ?, ?, ?, ?, ?, ?, ?,

current W:	W + h
[0.34,	[0.34,
-1.11,	-1.11
0.78,	0.78,
0.12,	0.12,
0.55,	0.55,
2.81,	2.81,
-3.1,	-3.1,
-1.5,	-1.5,
0.33,]	0.33,.
loss 1.25347	loss '

n (second dim): + 0.0001, . . 1.25353



current W:	W + h (third dim):
[0.34,	[0.34,
-1.11,	-1.11,
0.78,	0.78 + 0.0001 ,
0.12,	0.12,
0.55,	0.55,
2.81,	2.81,
-3.1,	-3.1,
-1.5,	-1.5,
0.33,]	0.33,]
loss 1.25347	loss 1.25347

gradient dW:

[-2.5, 0.6, ?, ?, ?, ?, ?, ?, ?, ?,

This is silly. The loss is just a function of W:

$$L = \frac{1}{N} \sum_{i=1}^{N} L_i + \sum_k W_k^2$$

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$$s = f(x; W) = Wx$$
want $\nabla_W L$
Calculus
$$\nabla_W L = ...$$

.

.

Andrej Karpathy

. .

- -

current W:		gradient dW:
[0.34, -1.11, 0.78, 0.12, 0.55, 2.81, -3.1, -1.5, 0.33,]	dW = (some function data and W)	[-2.5, 0.6, 0, 0.2, 0.7, -0.5, 1.1, 1.3, -2.1,]

Gradient descent

- We'll update weights
- Move in direction opposite to gradient:

$$\mathbf{w}^{(\tau+1)}_{\uparrow} = \mathbf{w}^{(\tau)} - \eta \nabla E(\mathbf{w}^{(\tau)})_{\downarrow}$$

Time Learning rate



Gradient descent

- Iteratively *subtract* the gradient with respect to the model parameters (w)
- I.e., we're moving in a direction opposite to the gradient of the loss
- I.e., we're moving towards *smaller* loss

Learning rate selection



Andrej Karpathy

Comments on training algorithm

- Not guaranteed to converge to zero training error, may converge to local optima or oscillate indefinitely.
- However, in practice, does converge to low error for many large networks on real data.
- Local minima not a huge problem in practice for deep networks.
- Thousands of epochs (epoch = network sees all training data once) may be required, hours or days to train.
- May be hard to set learning rate and to select number of hidden units and layers.
- When in doubt, use validation set to decide on design/hyperparameters.
- Neural networks had fallen out of fashion in 90s, early 2000s; now significantly improved performance (deep networks trained with dropout and lots of data).

Gradient descent in multi-layer nets

- We'll update weights
- Move in direction opposite to gradient:

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E(\mathbf{w}^{(\tau)})$$

- How to update the weights at all layers?
- Answer: backpropagation of error from higher layers to lower layers

Gradient descent in multi-layer nets

- How to update the weights at all layers?
- Answer: backpropagation of error from higher layers to lower layers



Backpropagation: Graphic example

First calculate error of output units and use this to change the top layer of weights.



Backpropagation: Graphic example

Next calculate error for hidden units based on errors on the output units it feeds into.



Backpropagation: Graphic example

Finally update bottom layer of weights based on errors calculated for hidden units.



Computing gradient for each weight

• We need to move weights in direction opposite to gradient of loss wrt that weight:

 $w_{kj} = w_{kj} - \eta dL/dw_{kj}$ (output layer) $w_{ji} = w_{ji} - \eta dL/dw_{ji}$ (hidden layer)

 Loss depends on weights in an indirect way, so we'll use the chain rule and compute:

 $\frac{dL}{dw_{kj}} = \frac{dL}{dy_k} \frac{dy_k}{da_k} \frac{da_k}{dw_{kj}}$ $\frac{dL}{dw_{ji}} = \frac{dL}{dz_j} \frac{dz_j}{da_j} \frac{da_j}{dw_{ji}}$

$$L_{i} = -\log(\frac{e^{sy_{i}}}{\sum_{j} e^{s_{j}}})$$
$$y_{k} = \sigma(a_{k}) = \frac{1}{1 + \exp(-a_{k})} \qquad a_{k} = \sum_{j=1}^{M} w_{kj}^{(2)} z_{j} + w_{k0}^{(2)}$$

Gradient for output layer weights

 Loss depends on weights in an indirect way, so we'll use the chain rule and compute:

 $dL/dw_{kj} = dL/dy_k dy_k/da_k da_k/dw_{kj}$

- How to compute each of these?
- dL/dy_k : need to know form of error function
 - Example: if L = $(y_k y_k')^2$, where y_k' is the ground-truth label, then $dL/dy_k = 2(y_k y_k')$
- dy_k/da_k : need to know output layer activation
 - If $h(a_k) = \sigma(a_k)$, then $d(a_k) / d(a_k) = \sigma(a_k) (1 \sigma(a_k))$
- da_k/dw_{kj} : z_j since a_k is a linear combination

•
$$a_k = w_{k:}^T z = \Sigma_j w_{kj} z_j$$

Gradient for hidden layer weights

- We'll use the chain rule again and compute: $dL/dw_{ii} = dL/dz_i dz_i/da_i da_i/dw_{ii}$
- Unlike the previous case (weights for output layer), the error (dL/dz_j) is hard to compute (indirect, need chain rule again)
- We'll simplify the computation by doing it step by step via *backpropagation* of error
- You could directly compute this term— you will get the same result as with backprop (do as an exercise!)

Backprop – rough notation

- The following is a framework, slightly imprecise
- Let us denote the inputs at a layer i by *in_i*, the linear combination of inputs computed at that layer as *raw_i*, the activation as *act_i*
- We define a new quantity that will roughly correspond to accumulated error, err_i :

err_i = d L / d act_i * d act_i / d raw_i

• Then we can write the updates as

$$w = w - \eta * err_i * in_i$$

Backprop – formulation

- We'll write the weight updates as follows $\succ w_{kj} = w_{kj} - \eta \ \delta_k \ z_j$ for output units $\succ w_{ji} = w_{ji} - \eta \ \delta_j \ x_i$ for hidden units
- What are δ_k, δ_j?
 - They store error, gradient wrt raw activations (i.e. dL/da)
 - They're of the form dL/dz_i dz_i/da_i
 - The latter is easy to compute just use derivative of activation function
 - The former is easy for output e.g. $(y_k y_k')$
 - It is harder to compute for hidden layers
 - $dL/dz_j = \sum_k w_{kj} \delta_k$ (where did this come from?)



Deriving backprop (Bishop Eq. 5.56)

• In a neural network:

$$a_j = \sum_i w_{ji} z_i \qquad z_j = h(a_j)$$

• Gradient is (using chain rule):



• Denote the "errors" as:

$$\delta_j \equiv \frac{\partial E_n}{\partial a_j}$$

• Also:

$$\frac{\partial a_j}{\partial w_{ji}} = z_i$$

Equations from Bishop

Deriving backprop (Bishop Eq. 5.56)

- For output (identity output, L2 loss): $\delta_k = y_k t_k$
- For hidden units (using chain rule again):

$$\delta_j \equiv \frac{\partial E_n}{\partial a_j} = \sum_k \frac{\partial E_n}{\partial a_k} \frac{\partial a_k}{\partial a_j}$$

Backprop formula:



Putting it all together

 Example: use sigmoid at hidden layer and output layer, loss is L2 between true/predicted labels

Example algorithm for sigmoid, L2 error

- Initialize all weights to small random values
- Until convergence (e.g. all training examples' error small, or error stops decreasing) repeat:
 - For each (x, y'=class(x)) in training set:
 - Calculate network outputs: y_k
 - Compute errors (gradients wrt activations) for each unit:

$$\begin{split} & \gg \delta_k = y_k (1-y_k) (y_k - y_k') & \text{for output units} \\ & \gg \delta_j = z_j (1-z_j) \sum_k w_{kj} \delta_k & \text{for hidden units} \\ & - \text{Update weights:} \end{split}$$

 $w_{kj} = w_{kj} - \eta \ \delta_k \ z_j \qquad \text{for output units}$ $w_{ji} = w_{ji} - \eta \ \delta_j \ x_i \qquad \text{for hidden units}$

$$\frac{\text{Recall:}}{\delta_j = w_{ji}} = w_{ji} - \eta \, dE/dz_j \, dz_j/da_j \, da_j/dw_{ji}}$$

$$\delta_j = h'(a_j) \sum_k w_{kj} \delta_k$$

Adapted from R. Hwa, R. Mooney

Another example

- Two layer network w/ tanh at hidden layer: $h(a) \equiv \tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$
- Derivative: $h'(a) = 1 h(a)^2$
- Minimize: $E_n = \frac{1}{2} \sum_{k=1}^{K} (y_k t_k)^2$
- Forward propagation:

$$a_{j} = \sum_{i=0}^{D} w_{ji}^{(1)} x_{i}$$

$$z_{j} = \tanh(a_{j})$$

$$y_{k} = \sum_{j=0}^{M} w_{kj}^{(2)} z_{j}$$

Another example

• Errors at output (identity function at output):

$$\delta_k = y_k - t_k$$

• Errors at hidden units:

$$\delta_j = h'(a_j) \sum_k w_{kj} \delta_k$$

$$\delta_j = (1 - z_j^2) \sum_{k=1}^{\infty} w_{kj} \delta_k$$

K

$$\frac{\partial E_n}{\partial w_{ji}^{(1)}} = \delta_j x_i, \qquad \frac{\partial E_n}{\partial w_{kj}^{(2)}} = \delta_k z_j$$

Same example with graphic and math

First calculate error of output units and use this to change the top layer of weights.


Same example with graphic and math

Next calculate error for hidden units based on errors on the output units it feeds into.



Same example with graphic and math

Finally update bottom layer of weights based on errors calculated for hidden units.



Another way of keeping track of error

Computation graphs

- Accumulate upstream/downstream gradients at each node
- One set flows from inputs to outputs and can be computed without evaluating loss
- The other flows from outputs (loss) to inputs







Generic example



Generic example



Andrej Karpathy



$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4



$$f(x,y,z) = (x+y)z$$

e.g. x = -2, y = 5, z = -4
 $q = x + y$ $rac{\partial q}{\partial x} = 1, rac{\partial q}{\partial y} = 1$

$$f=qz$$
 $rac{\partial f}{\partial q}=z, rac{\partial f}{\partial z}=q$

Want:
$$\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$$



$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

$$x = -2, y = 5, z = -4$$

$$y = 5, z = -4$$

$$z = 4$$

$$y = 5, z = -4$$

$$z = 4$$

$$y = 5, z = -4$$

$$z = 4$$

Want:
$$\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$$

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

$$x = -2, y = 5, z = -4$$

$$y = 5, z = -4$$

$$z = 4$$

$$y = 5, z = -4$$

$$y = 5, z = -4$$

$$z = 4$$

$$y = 5, z = -4$$

$$y = 5, z = -4$$

$$y = 5, z = -4$$

$$z = 4$$

Want:
$$\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$$

$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

$$\frac{\partial f}{\partial z}$$

Want:
$$\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$$

$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

$$\frac{\partial f}{\partial z}$$

Want:
$$\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$$

$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

$$\frac{\partial f}{\partial q}$$

Want:
$$\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$$

$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

$$\frac{\partial f}{\partial q}$$

Want:
$$\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$$

$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

$$\frac{\partial f}{\partial y}$$

Want:
$$\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$$

$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

$$\overset{x -2}{y + \frac{q}{3}}$$

$$y = \frac{y}{4}$$

$$\frac{y - \frac{1}{4}}{3}$$

Chain rule: $\frac{\partial f}{\partial y}$

$$\frac{\partial f}{\partial y}$$

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

$$\frac{\partial f}{\partial x}$$

Want:
$$\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$$

Summary

- Feed-forward network architecture
- Training deep neural nets
 - We need an objective function that measures and guides us towards good performance
 - We need a way to minimize the loss function: gradient descent
 - We need backpropagation to propagate error towards all layers and change weights at those layers
- Next: Practices for preventing overfitting, training with little data, examining conditions for success, alternative optimization strategies