# CSCE 5218 & 4930 Deep Learning

## Neural Network Basics

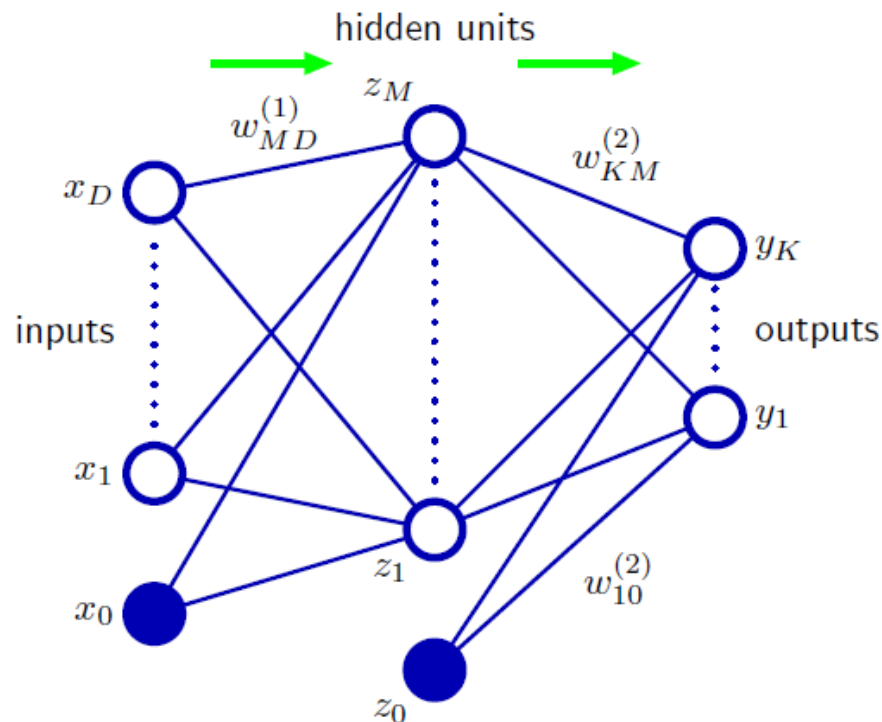Slides adapted from Adriana Kovashka

# Plan for this lecture (next few classes)

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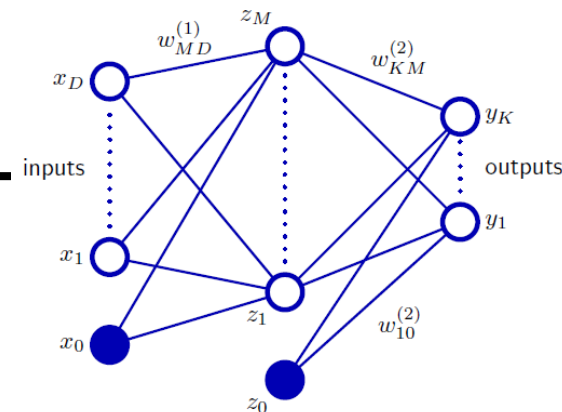  - Backpropagation
- Hands-on exercise

# Definition

# Neural network definition

**Figure 5.1** Network diagram for the two-layer 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: $$a_j = \sum_{i=1}^{D} w_{ji}^{(1)} x_i + w_{j0}^{(1)}$$

- Nonlinear activation function $h$ (e.g. sigmoid, tanh, RELU): $$z_j = h(a_j) = \text{RELU}(a) = \max(0, a)$$

# Neural network definition



- Layer 2

$$a_j = \sum_{i=1}^{D} w_{ji}^{(1)} x_i + w_{j0}^{(1)}$$

- Layer 3 (final)

$$a_k = \sum_{j=1}^{M} w_{kj}^{(2)} z_j + w_{k0}^{(2)}$$

- Outputs

(binary) $\quad y_k = \sigma(a_k) = \dfrac{1}{1 + \exp(-a_k)}$
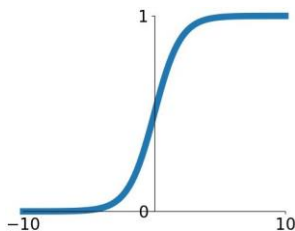
(multiclass)

$$y_k = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

- Finally:

(binary)

$$y_k(\mathbf{x}, \mathbf{w}) = \sigma \left( \sum_{j=1}^{M} w_{kj}^{(2)} h \left( \sum_{i=1}^{D} w_{ji}^{(1)} x_i + w_{j0}^{(1)} \right) + w_{k0}^{(2)} \right)$$
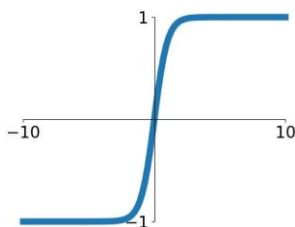
# Activation functions

**Sigmoid**

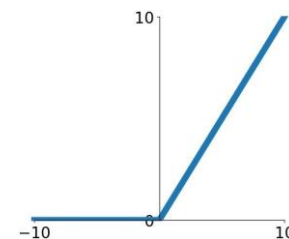$$\sigma(x) = \frac{1}{1+e^{-x}}$$

**tanh**

$$\tanh(x)$$

**Maxout**

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

**ReLU**
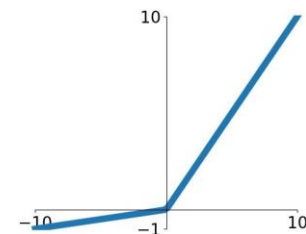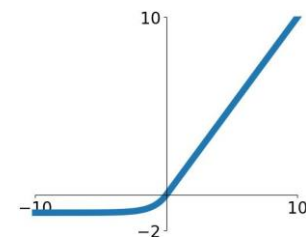
$$\max(0, x)$$

**Leaky ReLU**

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

**ELU**

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

# A multi-layer neural network…



Input Layer    Hidden Layer    Output Layer

Input #1 →
Input #2 →
Input #3 →
Input #4 →

Output

- 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



A biological neuron

An artificial neuron

**Input**

**Weights**

$x_1$   $w_1$

$x_2$   $w_2$

$x_3$   $w_3$

$x_d$   $w_d$

**Output:** $\sigma(\mathbf{w}\cdot\mathbf{x} + b)$

**Sigmoid function:**

$$\sigma(t) = \frac{1}{1 + e^{-t}}$$

# Biological analog



Hubel and Weisel's architecture



Multi-layer neural network

# Feed-forward networks

- Cascade neurons together

- Output from one layer is the input to the next

- Each layer has its own sets of weights

$x_0$

$x_1$

$x_2$

$x_P$

# Feed-forward networks

- Inputs multiplied by initial set of weights

# Feed-forward networks

- Intermediate "predictions" computed at first hidden layer

# Feed-forward networks

- Intermediate predictions multiplied by second layer of weights

- Predictions are fed forward through the network to classify

$x_0$
$x_1$
$x_2$
$x_P$

- Compute second set of intermediate predictions

# Feed-forward networks

- Multiply by final set of weights

# Feed-forward networks

- 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)

# Finding the optimal **w**: Example

- Suppose **w** is just a scalar, w, that can only take values 1, 2, 3

- Suppose $L(w=1) = 2$, $L(w=2) = 5$, $L(w=3) = 0.5$

- Find the optimal w as $\text{argmin}\_w\ L(w)$

- What is the optimal w?

- Now suppose $L^{\wedge}(w) = L(w) + ||w||$

- $L^{\wedge}(1) = 2+1 = 3$

- $L^{\wedge}(2) = ?\ L^{\wedge}(3) = ?$

- Now what is the optimal w?

# Classification goal



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($\textcolor{blue}{\mathbf{x}}$, $\textcolor{red}{\mathbf{W}}$)

**[32x32x3]**

array of numbers 0...1
(3072 numbers total)

**10** numbers,
indicating class
scores

# Linear classifier

$$f(x, W) = Wx \quad (+b)$$

10x1    10x3072    3072x1    10x1



**[32x32x3]**

array of numbers 0...1

**10** numbers, indicating class scores

parameters, or "weights"

# Linear classifier

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



stretch pixels into single column

| 0.2 | -0.5 | 0.1 | 2.0 |
| 1.5 | 1.3 | 2.1 | 0.0 |
| 0 | 0.25 | 0.2 | -0.3 |

$W$

| 56 |
| 231 |
| 24 |
| 2 |

$x_i$

$+$

| 1.1 |
| 3.2 |
| -1.2 |

$b$

$\longrightarrow$

| -96.8 | cat score |
| 437.9 | dog score |
| 61.95 | ship score |

$f(x_i; W, b)$

input image

# Linear classifier

Going forward: Loss function/Optimization



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

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

# Linear classifier: Hinge loss

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

|  | cat | car | frog |
|---|---|---|---|
| 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:

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

Want: $s_{y_i} >= s_j + 1$
i.e. $s_j - s_{y_i} + 1 <= 0$

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

# Linear classifier: Hinge loss

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



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

**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:

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

= 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

# Linear classifier: Hinge loss

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 | |

**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, 1.3 - 4.9 + 1)
  +max(0, 2.0 - 4.9 + 1)
= max(0, -2.6) + max(0, -1.9)
= 0 + 0
= 0

# Linear classifier: Hinge loss

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: $s = f(x_i, W)$

|  | cat | car | frog |
|---|---|---|---|
| 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 |

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

# Linear classifier: Hinge loss

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:

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

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

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

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

# Linear classifier: Hinge loss

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

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

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?**

# Weight Regularization

$\lambda = $ 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

$$x = [1, 1, 1, 1]$$

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

$$w_2 = \boxed{[0.25, 0.25, 0.25, 0.25]}$$

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

L2 Regularization

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

L2 regularization likes to "spread out" the weights

# 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

**scores = unnormalized log probabilities of the classes**

$$s = f(x_i; W)$$

cat      **3.2**

car      5.1

frog      -1.7

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

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

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

# Another loss: Cross-entropy

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$



Probabilities must be >= 0

Probabilities must sum to 1

|       | unnormalized log probabilities | exp | unnormalized probabilities | normalize | probabilities |
|-------|------|------|------|------|------|
| cat   | **3.2** | | **24.5** | | **0.13** |
| car   | 5.1 | | 164.0 | | 0.87 |
| frog  | -1.7 | | 0.18 | | 0.00 |

L_i = -log(0.13)
  = **0.89**

**Aside:**
- This is multinomial logistic regression
- Choose weights to maximize the likelihood of the observed x/y data (Maximum Likelihood Estimation)

# Other losses

- Triplet loss (Schroff, FaceNet, CVPR 2015)

$$\sum_i^N \left[ \|f(x_i^a) - f(x_i^p)\|_2^2 - \|f(x_i^a) - f(x_i^n)\|_2^2 + \alpha \right]_+$$

a denotes anchor
p denotes positive
n denotes negative



Figure 3. The **Triplet Loss** minimizes the distance between an *anchor* 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

# How to minimize the loss function?

In 1-dimension, the derivative of a function:

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

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

# Loss gradients

- Denoted as (diff notations):   $\dfrac{\partial E}{\partial w_{ji}^{(1)}}$   $\nabla_W L$

- 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



tangent line

slope= $f'(x)$

$x$

**current W:**

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

**gradient dW:**

[?,
?,
?,
?,
?,
?,
?,
?,
?,…]

**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

**gradient dW:**

[?,
?,
?,
?,
?,
?,
?,
?,
?,…]

**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**

**gradient dW:**

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

(1.25322 - 1.25347)/0.0001
= -2.5

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

?,
?,…]

**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** (second dim)**:**

[0.34,
-1.11 + **0.0001**,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
loss 1.25353

**gradient dW:**

[-2.5,
?,
?,
?,
?,
?,
?,
?,
?,…]

Andrej Karpathy

**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** (second dim)**:**

[0.34,
-1.11 + **0.0001**,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
**loss 1.25353**

**gradient dW:**

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

(1.25353 - 1.25347)/0.0001
= 0.6

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

?,…]

**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** (third dim)**:**

[0.34,
-1.11,
0.78 + **0.0001**,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
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

$$\nabla_W L = ...$$

**current W:**

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

dW = …
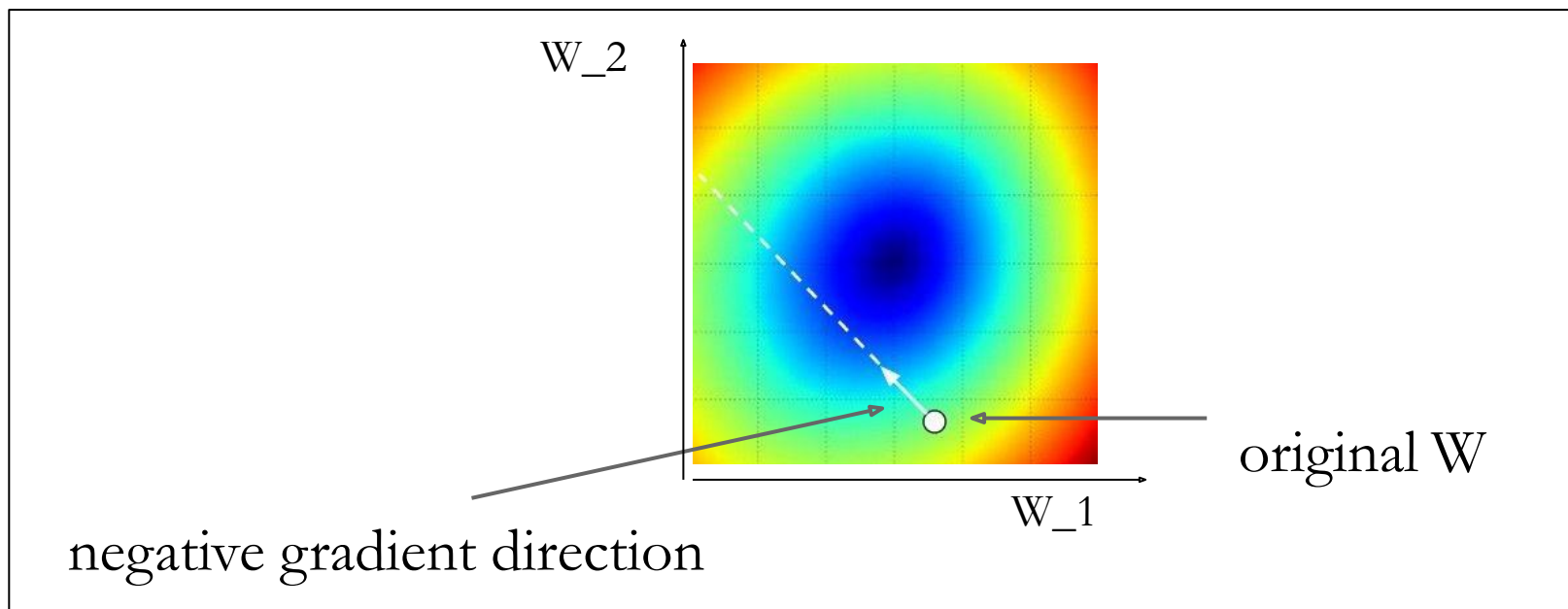(some function
data and W)

**gradient dW:**

[-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:

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

Time

Learning rate



W_2

W_1

original W

negative gradient direction

# 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

The effects of step size (or "learning rate")

# Comments on training algorithm

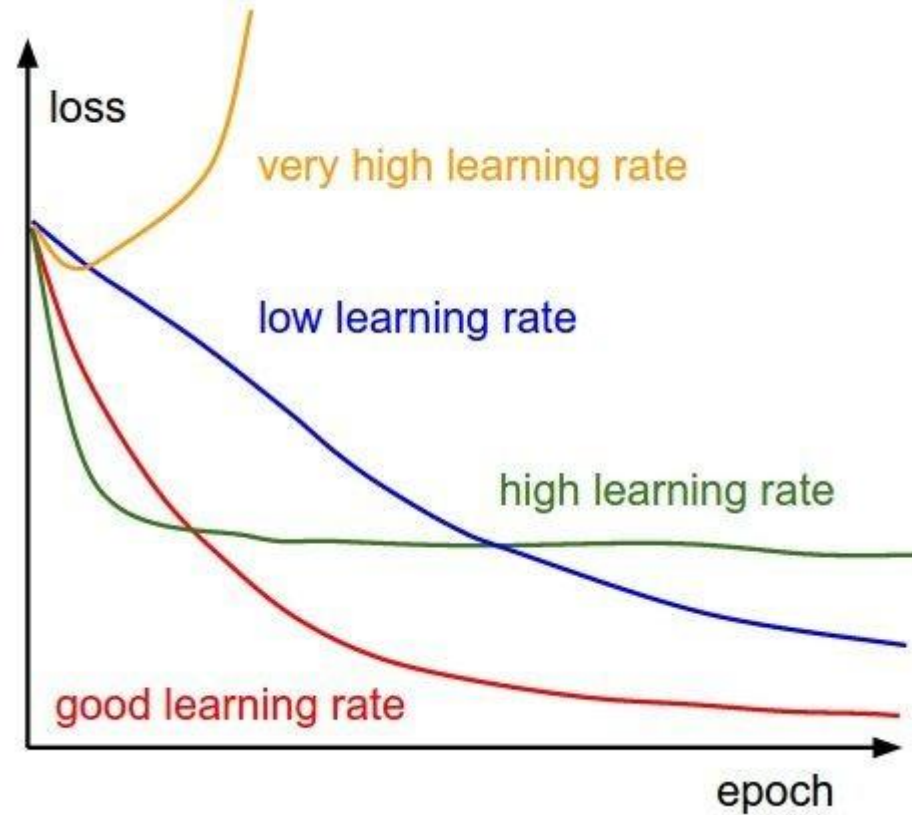- Not guaranteed to converge to zero training error, may converge to local optima or oscillate indefinitely.

- 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).
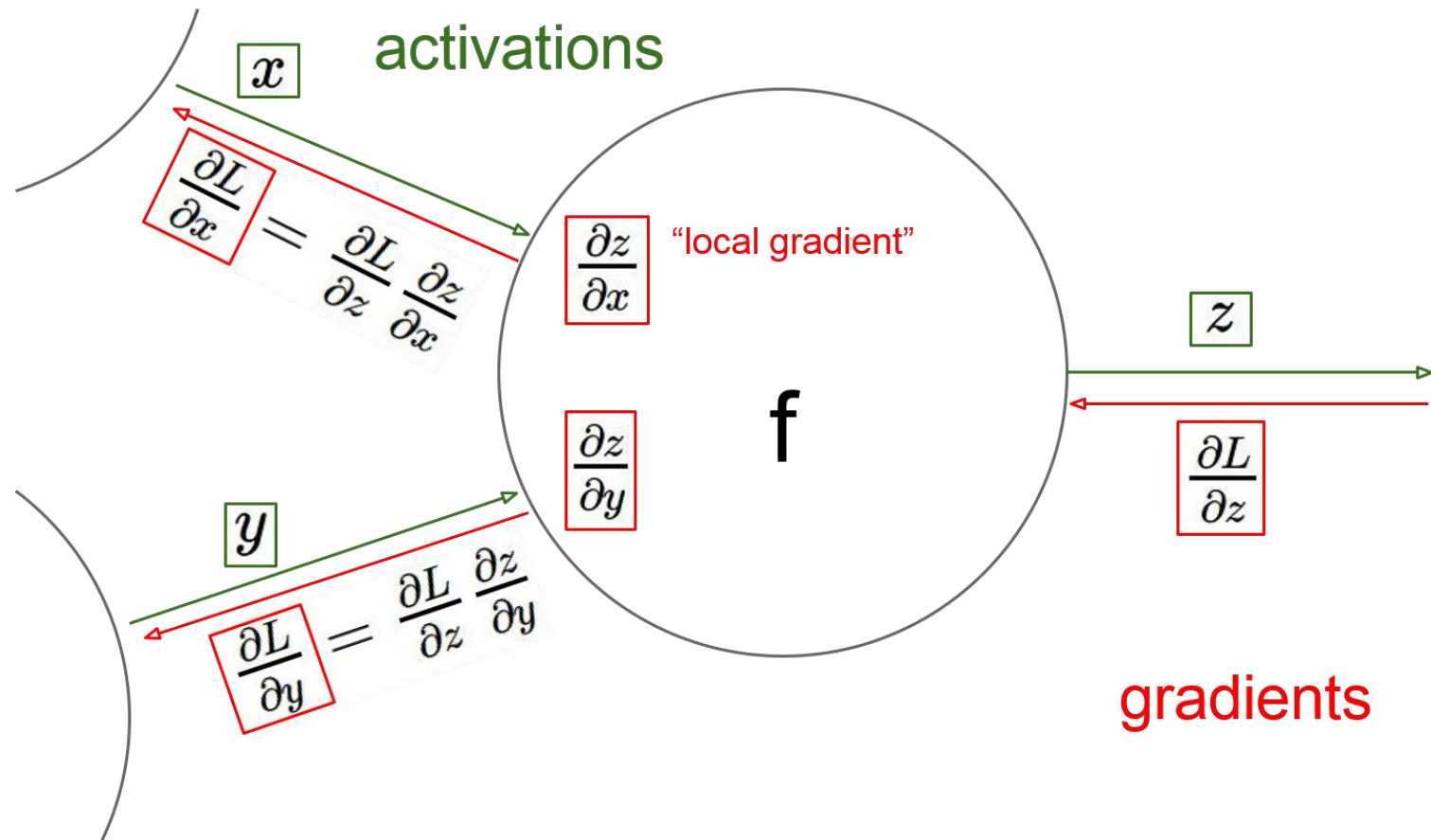
# Gradient descent in multi-layer nets

- 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



$x$ activations

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z}\frac{\partial z}{\partial x}$$

$\frac{\partial z}{\partial x}$ "local gradient"

$\frac{\partial z}{\partial y}$

$z$

$\frac{\partial L}{\partial z}$

f

$y$

$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial z}\frac{\partial z}{\partial y}$$

gradients

# Backpropagation: Graphic example

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

Update weights into $j$
(store diff, update @end)

output $k$

$w^{(2)}$

hidden $j$

$w^{(1)}$

input $i$

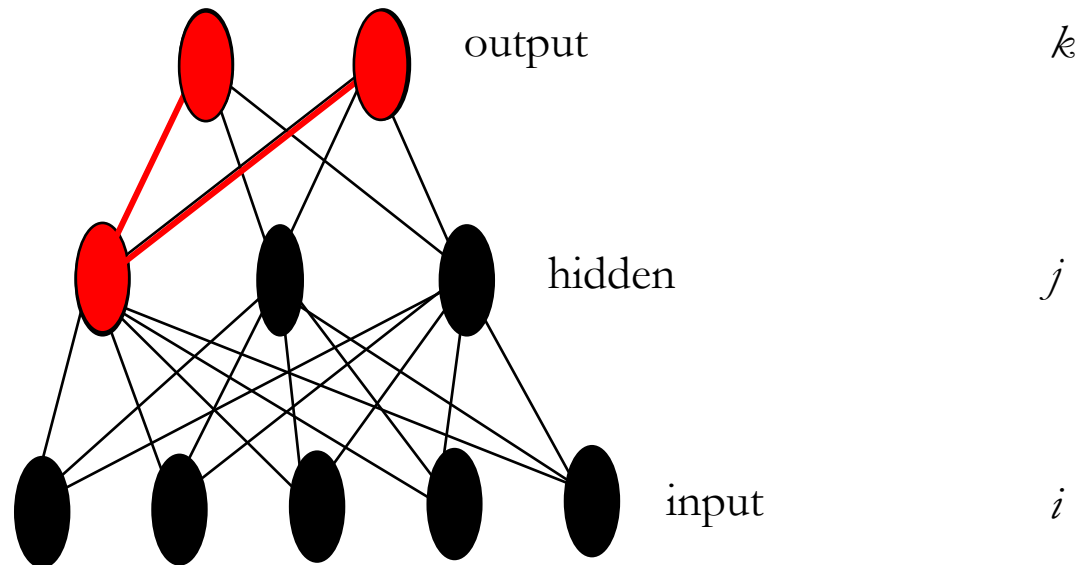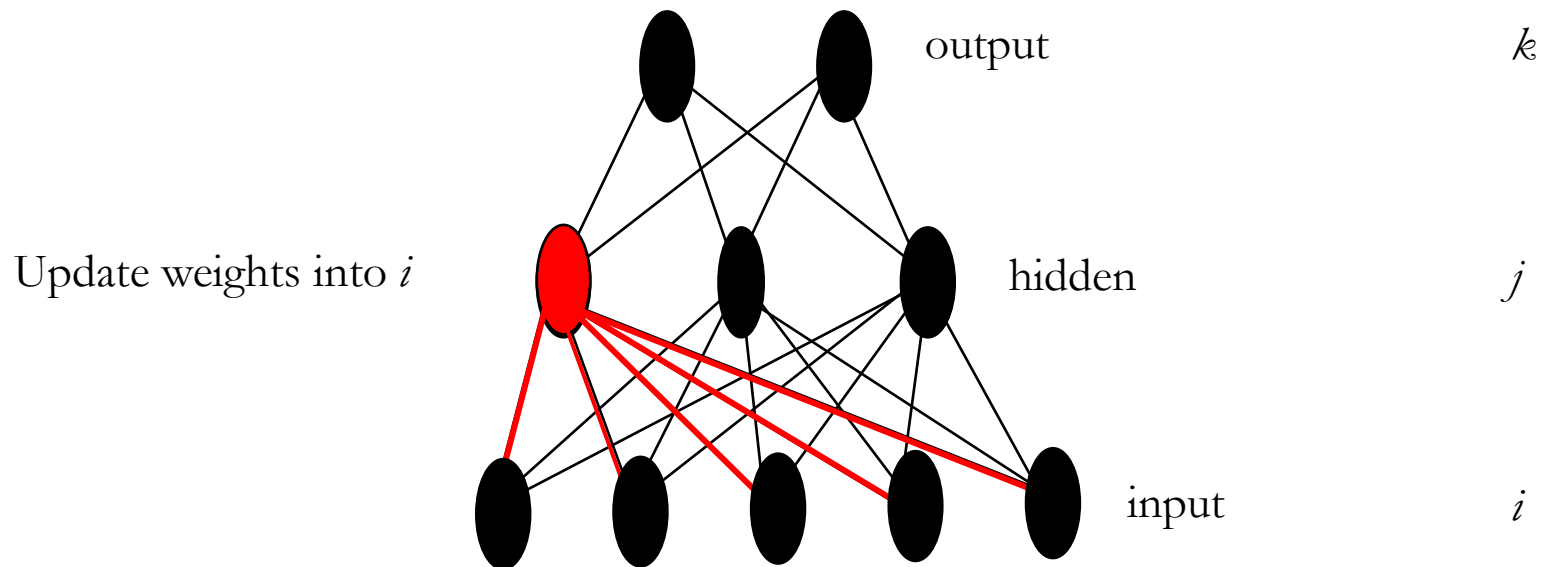# Backpropagation: Graphic example

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.

Update weights into $i$

output $\qquad$ $k$

hidden $\qquad$ $j$

input $\qquad$ $i$

# Computing gradient for each weight

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

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

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

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

$$\mathtt{dL/dw_{kj} = \textcolor{red}{dL/dy_k} \ \textcolor{green}{dy_k/da_k} \ \textcolor{blue}{da_k/dw_{kj}}}$$

$$\mathtt{dL/dw_{ji} = \textcolor{red}{dL/dz_j} \ \textcolor{green}{dz_j/da_j} \ \textcolor{blue}{da_j/dw_{ji}}}$$

$$\boxed{L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)}$$

$$\boxed{y_k = \sigma(a_k) = \frac{1}{1 + \exp(-a_k)}} \qquad \boxed{a_k = \sum_{j=1}^{M} w_{kj}^{(2)} z_j + w_{k0}^{(2)}}$$

# Gradient for output layer weights

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

$$\mathtt{dL/dw_{kj}} = \textcolor{red}{\mathtt{dL/dy_k}} \ \textcolor{green}{\mathtt{dy_k/da_k}} \ \textcolor{blue}{\mathtt{da_k/dw_{kj}}}$$

- How to compute each of these?

- $\textcolor{red}{\mathtt{dL/dy_k}}$ : need to know form of error function

  - Example: if $\mathtt{L = (y_k - y_k')^2}$, where $\mathtt{y_k'}$ is the ground-truth label, then $\mathtt{dL/dy_k = \cancel{2}(y_k - y_k')}$

- $\textcolor{green}{\mathtt{dy_k/da_k}}$ : need to know output layer activation

  - If $\mathtt{h(a_k)=\sigma(a_k)}$, then $\mathtt{d\ h(a_k)/d\ a_k = \sigma(a_k)(1-\sigma(a_k))}$

- $\textcolor{blue}{\mathtt{da_k/dw_{kj}}}$ :

  - $\mathtt{z_j}$ since $\mathtt{a_k}$ is a linear combination

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

# Gradient for hidden layer weights

- We'll use the chain rule again and compute:

$$\texttt{dL/dw}_{ji} \texttt{ = dL/dz}_j \texttt{ dz}_j\texttt{/da}_j \texttt{ da}_j\texttt{/dw}_{ji}$$
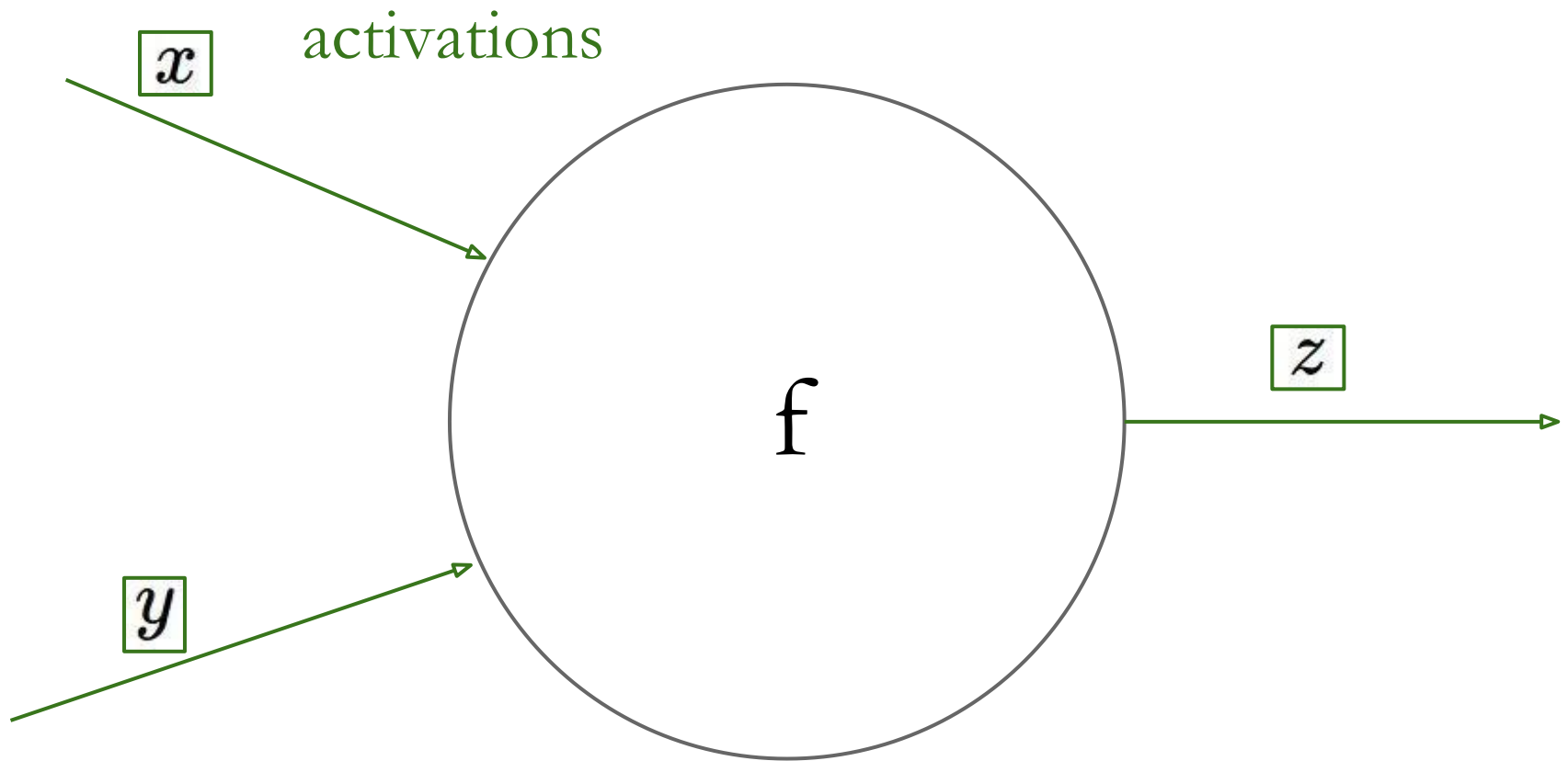
- Unlike the previous case (weights for output layer), the error ($\texttt{dL/dz}_j$) is hard to compute (indirect, need chain rule again)
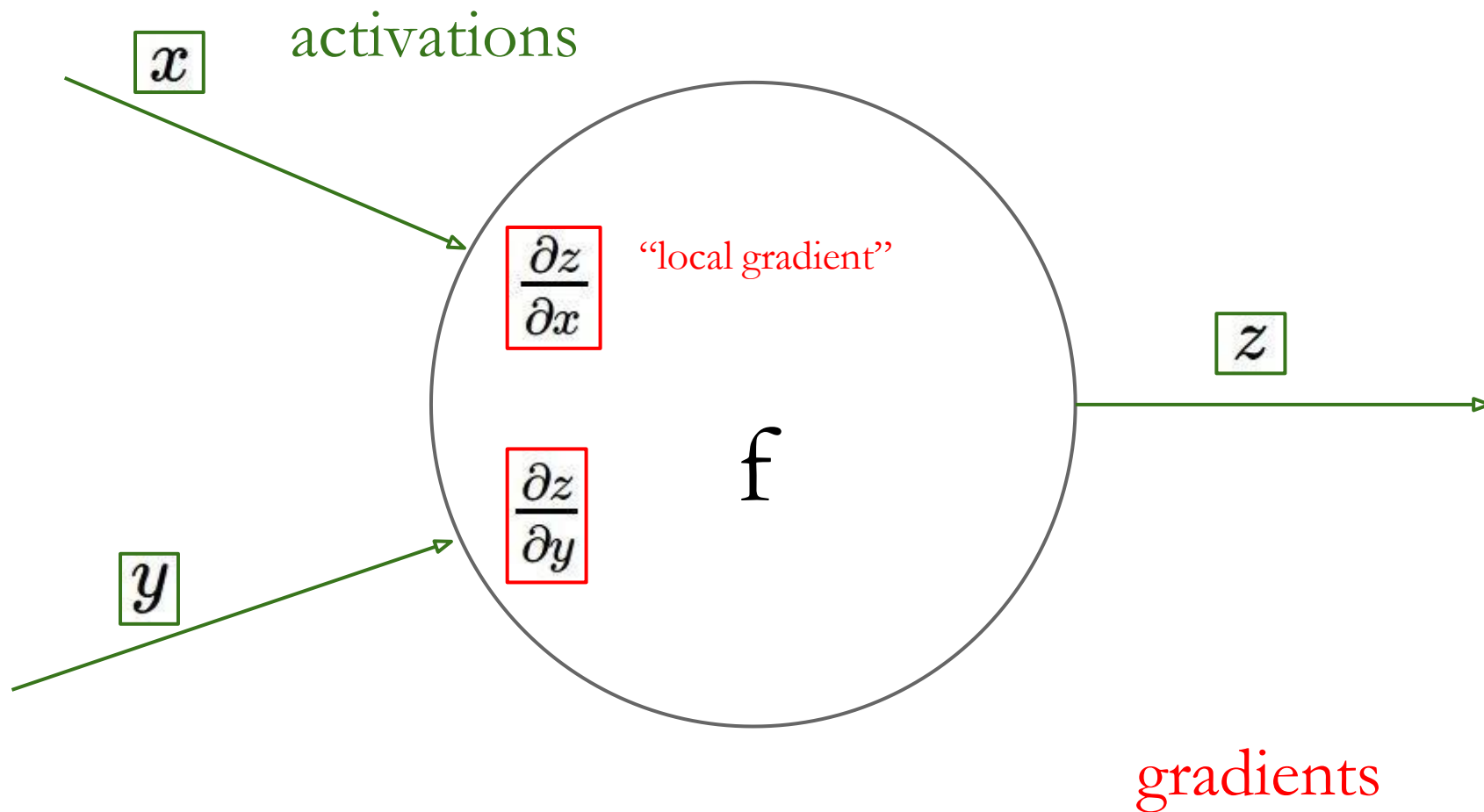
# Another way of keeping track of error

Computation graphs

- Accumulate upstream/downstream gradients at each node

- One set flows from inputs to outputs and can be computed without evaluating loss

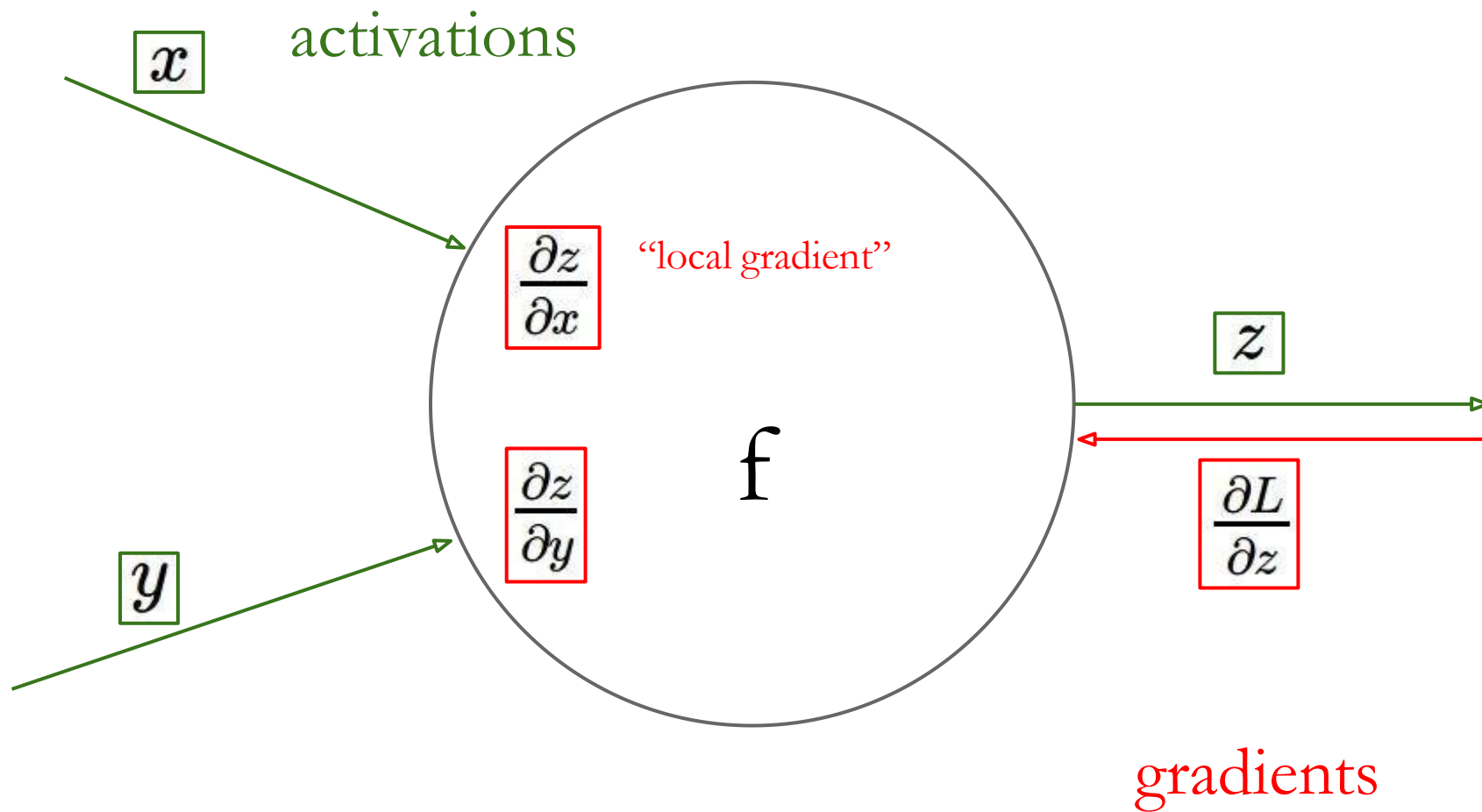- The other flows from outputs (loss) to inputs
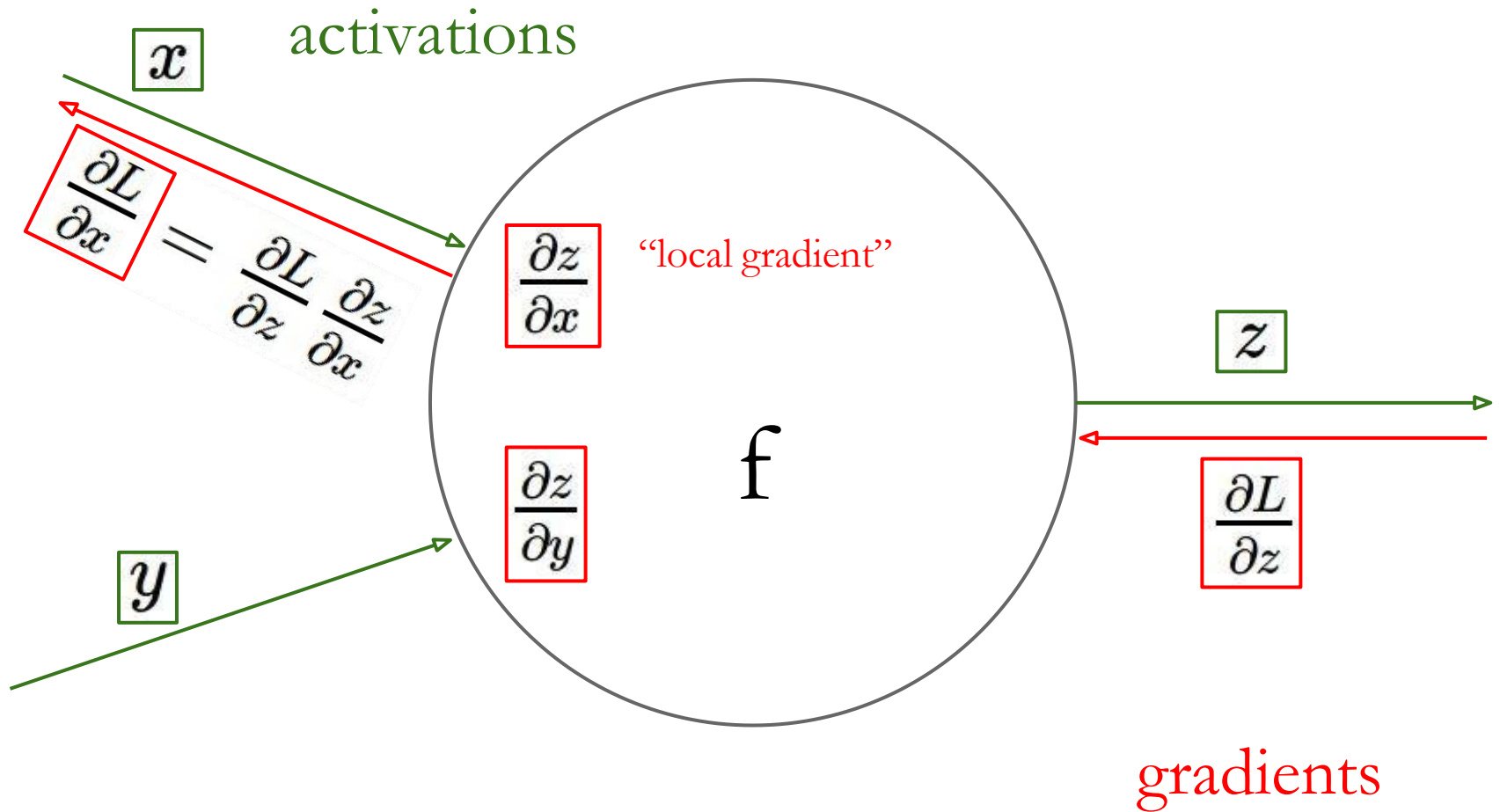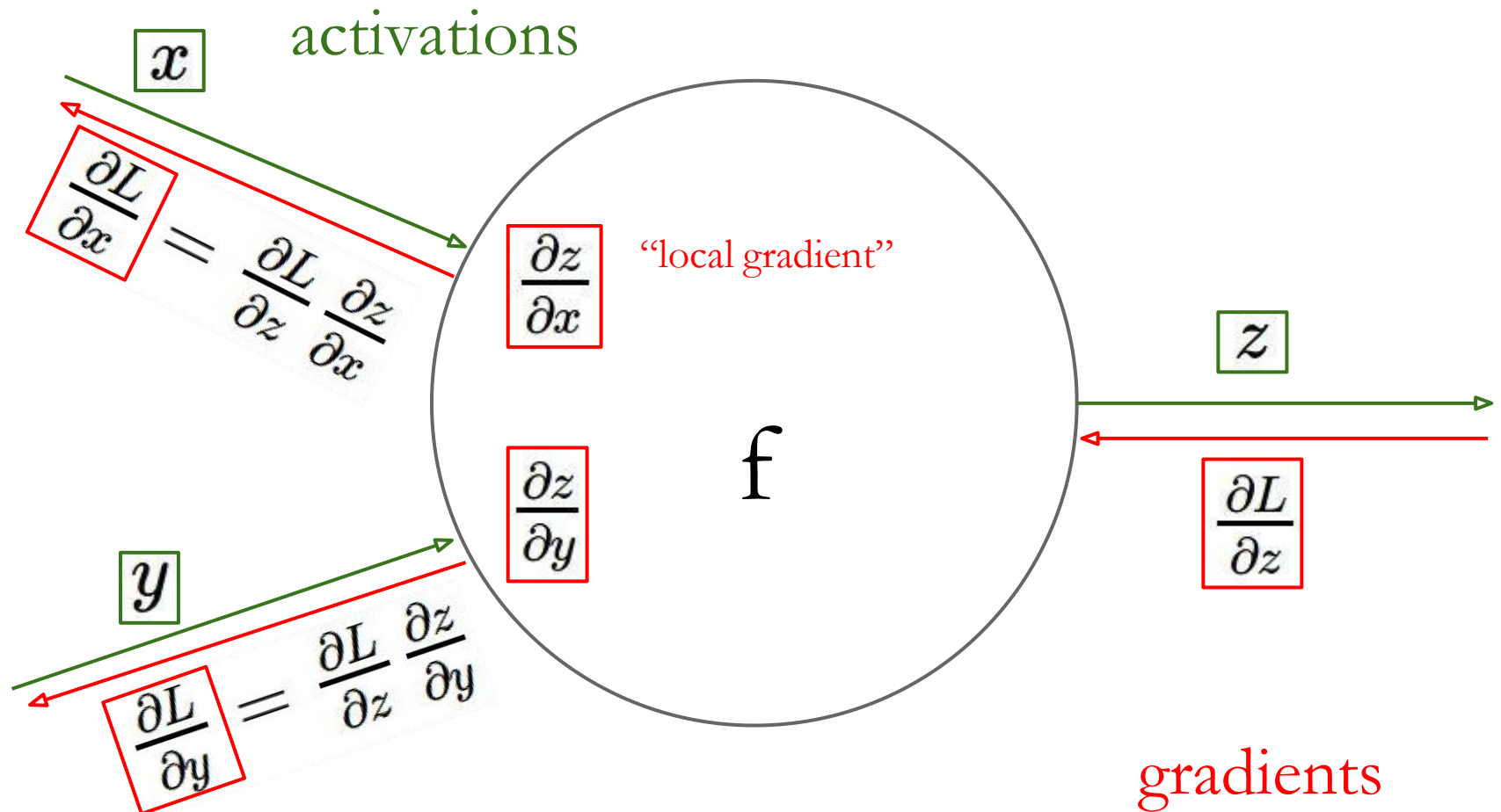
# Generic example

activations



$x$

$y$

f

$z$

# Generic example



activations

$x$

$$\frac{\partial z}{\partial x}$$ "local gradient"

$f$

$$\frac{\partial z}{\partial y}$$

$y$

$z$

gradients

# Generic example

activations



$x$

$\dfrac{\partial z}{\partial x}$   "local gradient"

$\dfrac{\partial z}{\partial y}$

$y$

f

$z$

$\dfrac{\partial L}{\partial z}$

gradients

# Generic example

activations

$x$

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial x}$$

$\frac{\partial z}{\partial x}$  "local gradient"

$\frac{\partial z}{\partial y}$

f

$y$

$z$

$\frac{\partial L}{\partial z}$

gradients

# Generic example



activations

$x$

$$\boxed{\frac{\partial L}{\partial x}} = \frac{\partial L}{\partial z}\frac{\partial z}{\partial x}$$

"local gradient"

$$\boxed{\frac{\partial z}{\partial x}}$$

$$\boxed{\frac{\partial z}{\partial y}}$$

f

$z$

$$\boxed{\frac{\partial L}{\partial z}}$$

$y$

$$\boxed{\frac{\partial L}{\partial y}} = \frac{\partial L}{\partial z}\frac{\partial z}{\partial y}$$

gradients

Andrej Karpathy

# Generic example

activations

$$x$$

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial x}$$

$$\frac{\partial z}{\partial x}$$   "local gradient"

$$z$$

$$\frac{\partial L}{\partial z}$$

f

$$y$$

$$\frac{\partial z}{\partial y}$$

$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial y}$$

gradients

Andrej Karpathy

# Another generic example

$$f(x, y, z) = (x + y)z$$
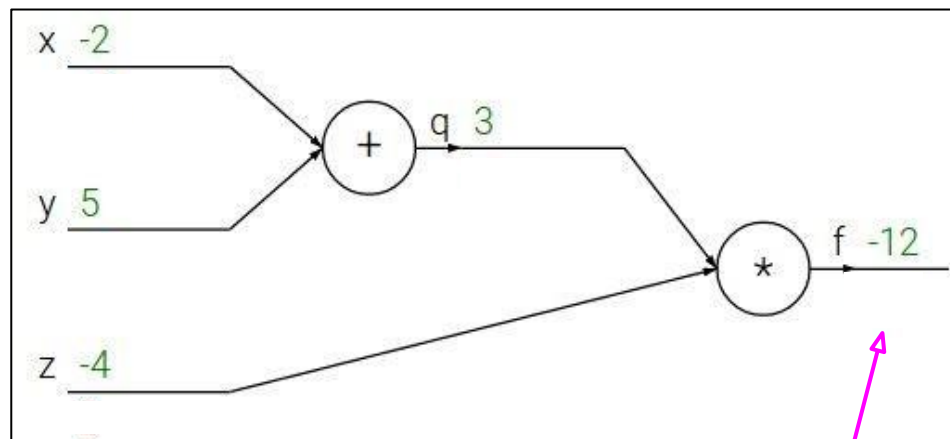
e.g. x = -2, y = 5, z = -4

# Another generic example

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

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



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

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

Want: $\quad \dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$
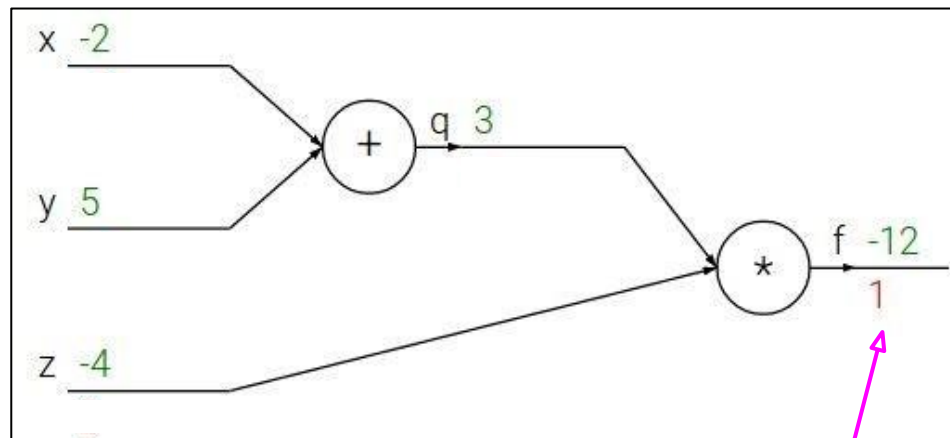
# Another generic example

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

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

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

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



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

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$
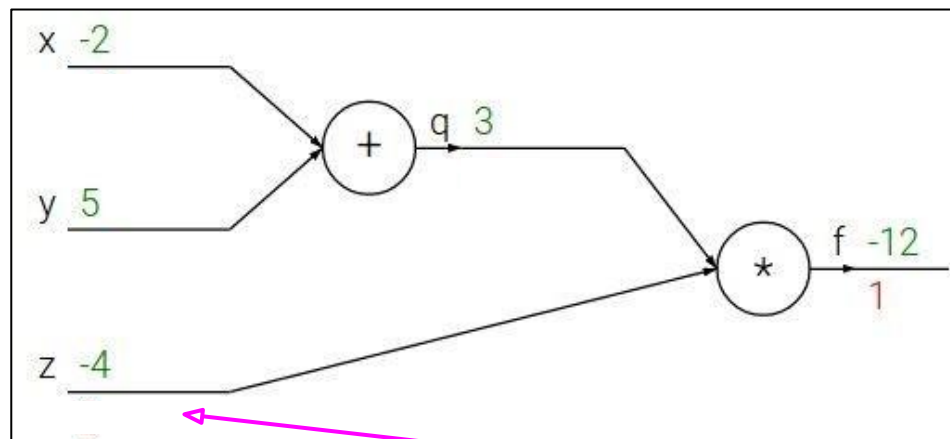
# Another generic example

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

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

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

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



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

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$
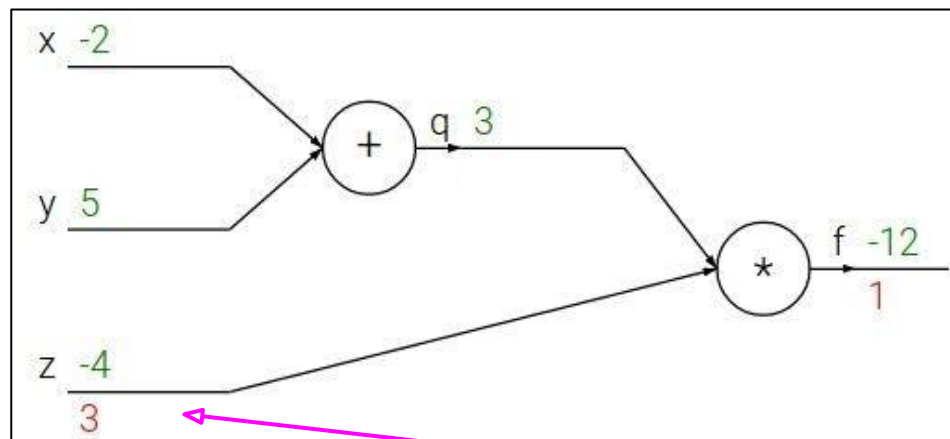
# Another generic example

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

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: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$
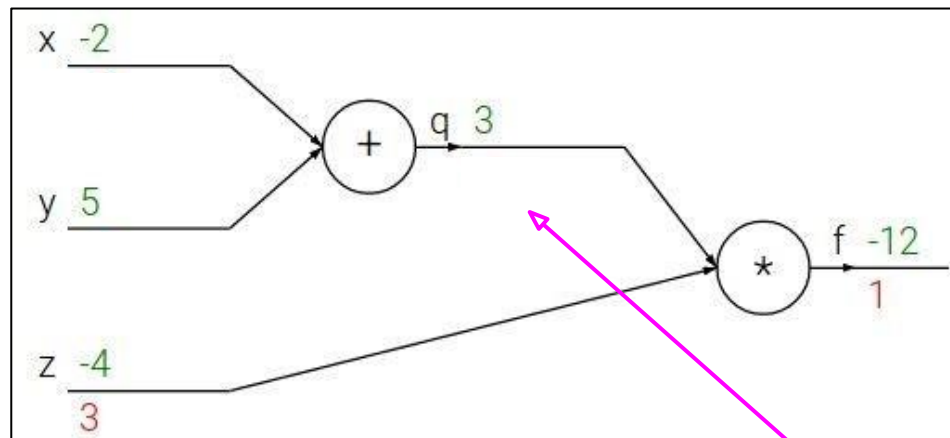
# Another generic example

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

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

# Another generic example

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

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: $$\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$$
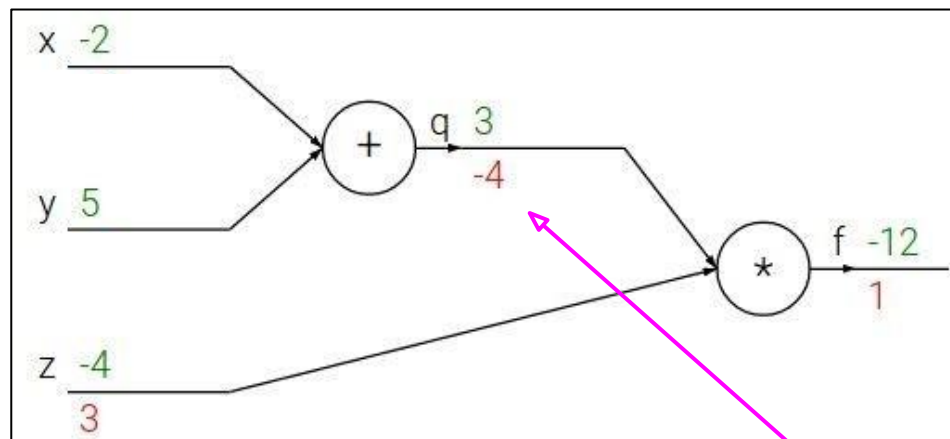
# Another generic example

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

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: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$
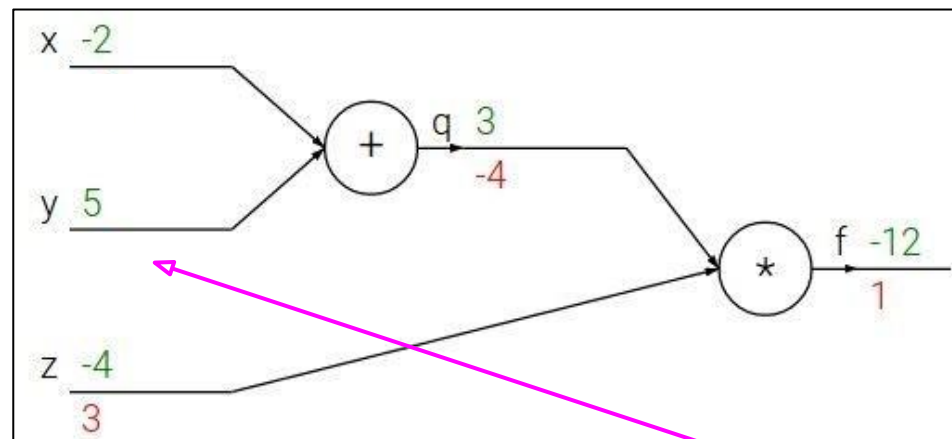
# Another generic example

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

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: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$
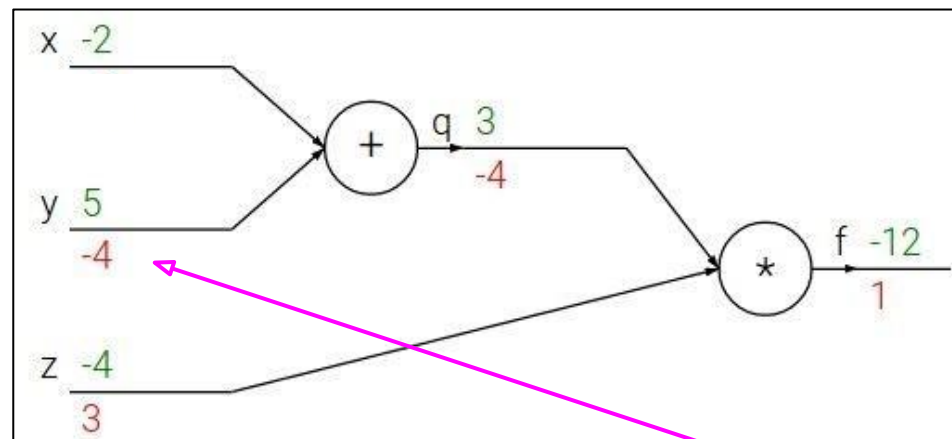
Andrej Karpathy

# Another generic example

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

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

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

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

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



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

Chain rule:

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y}$$
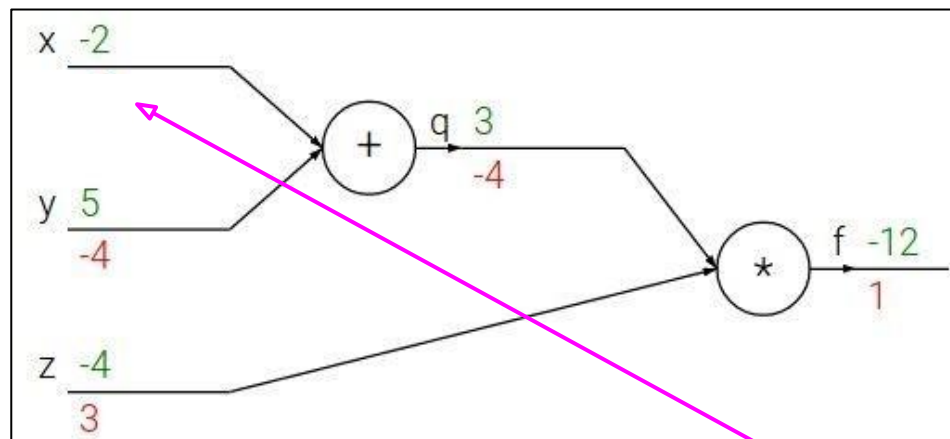
# Another generic example

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

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: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$
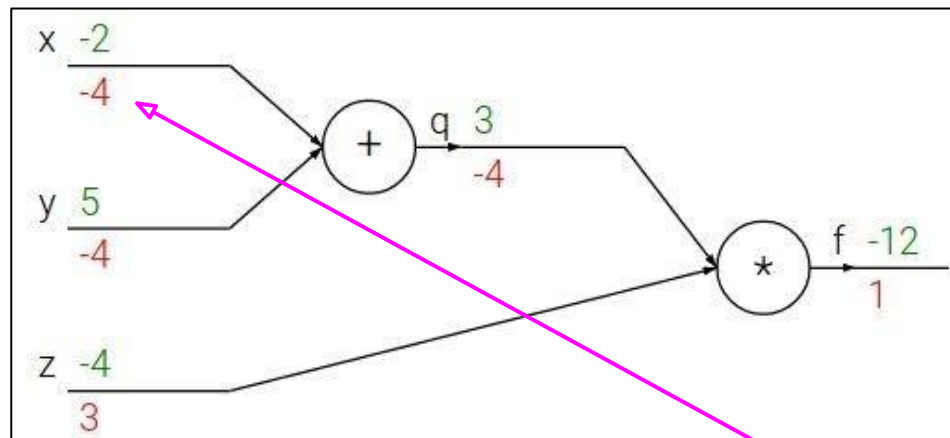
# Another generic example

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

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

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

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

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



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

Chain rule:

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

# 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