

# Backpropagation

TA: Yi Wen

April 17, 2020

CS231n Discussion Section

Slides credits: Barak Oshri, Vincent Chen, Nish Khandwala, Yi Wen

# Agenda

- Motivation
- Backprop Tips & Tricks
- Matrix calculus primer

# Agenda

- **Motivation**
- Backprop Tips & Tricks
- Matrix calculus primer

# Motivation

Recall: Optimization objective is minimize loss

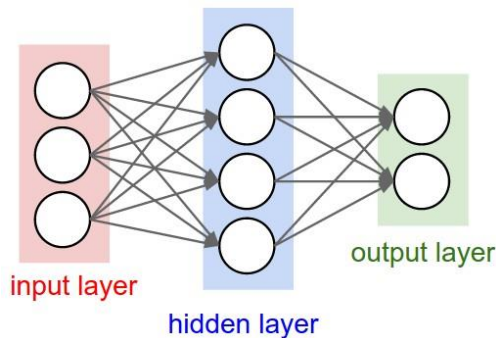
$$Loss = f(x, y; \theta)$$

# Motivation

Recall: Optimization objective is minimize loss

$$Loss = f(x, y; \theta)$$

**Goal: how should we tweak the parameters to decrease the loss?**



# Agenda

- Motivation
- **Backprop Tips & Tricks**
- Matrix calculus primer

# A Simple Example

Loss

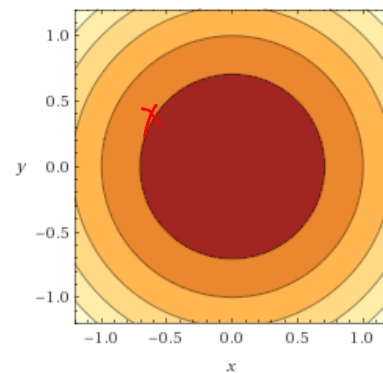
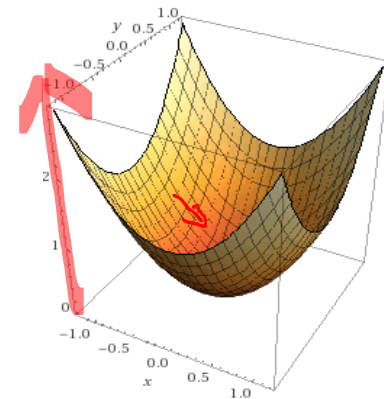
$$Loss = f(\overset{\text{dataset}}{\underbrace{x, y}}; \theta)$$

Goal: Tweak the parameters to minimize loss

=> **minimize a multivariable function** in parameter space

# A Simple Example

=> minimize a multivariable function



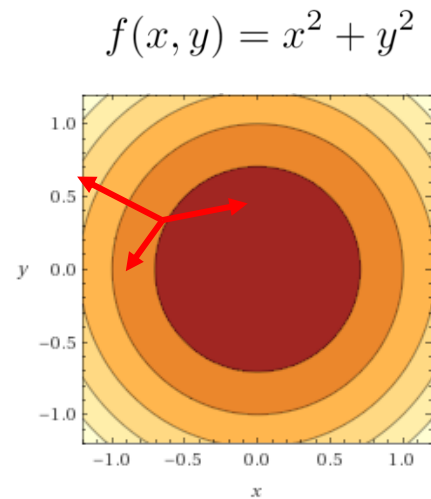
$$f(x, y) = x^2 + y^2$$

Plotted on WolframAlpha



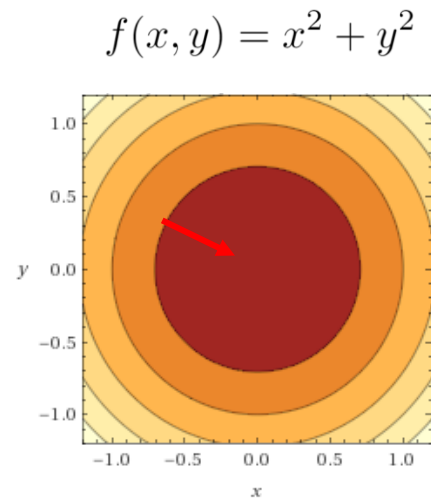
# Approach #1: Random Search

**Intuition:** the *step* we take in the domain of function



## Approach #2: Numerical Gradient

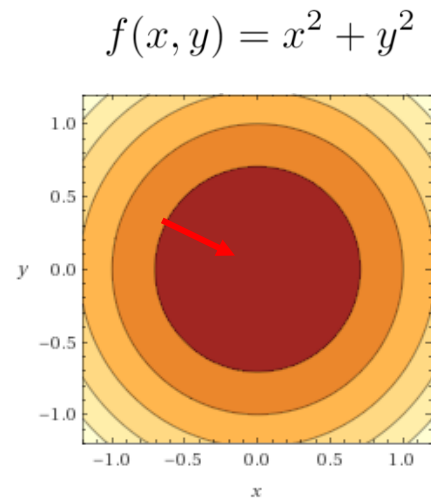
**Intuition:** rate of change of a function with respect to a variable surrounding a small region



## Approach #2: Numerical Gradient

**Intuition:** rate of change of a function with respect to a variable surrounding a small region

**Finite Differences:** 
$$\frac{f(x + h, y) - f(x, y)}{h}$$



# Approach #3: Analytical Gradient

**Recall:** partial derivative by limit definition

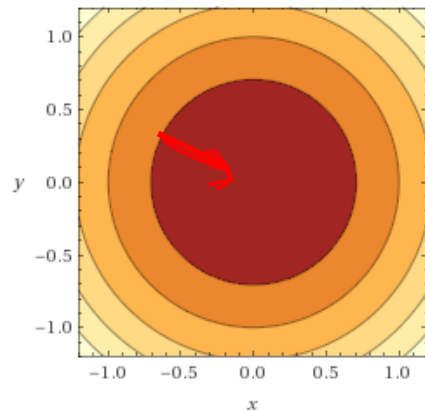
$$\frac{\partial f}{\partial x}(x, y) = \lim_{h \rightarrow 0} \frac{f(x + h, y) - f(x, y)}{h}$$

$$\frac{\partial f}{\partial x} = 2x$$

$$\frac{\partial f}{\partial y} = 2y$$

$$= \begin{bmatrix} 2x \\ 2y \end{bmatrix}$$

$$f(x, y) = x^2 + y^2$$



## Approach #3: Analytical Gradient

**Recall:** chain rule  $\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx}$ .

# Approach #3: Analytical Gradient

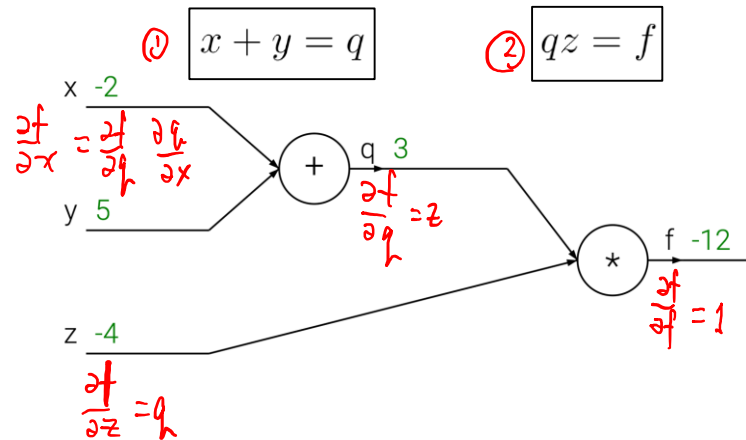
**Recall:** chain rule  $\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx}$ .

E.g.  $f(x, y, z) = (x + y)z$

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

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

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



# Approach #3: Analytical Gradient

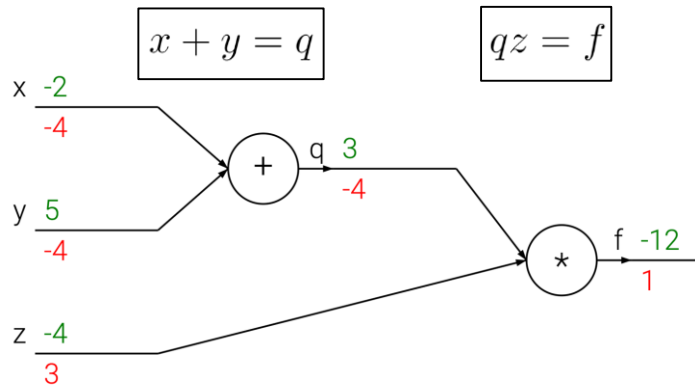
**Recall:** chain rule  $\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx}$ .

E.g.  $f(x, y, z) = (x + y)z$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x} = z \cdot 1 = z = -4$$

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y} = z \cdot 1 = z = -4$$

$$\frac{\partial f}{\partial z} = q = x + y = -2 + 5 = 3$$



## Approach #3: Analytical Gradient

**Recall:** chain rule  $\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx}$ .

**Intuition:** upstream gradient values propagate backwards -- we can reuse them!



# Gradient

$$f : R^n \rightarrow R$$

$$\nabla f : R^n \rightarrow R^n$$

$$\mathbf{x} = [x_1, \dots, x_n]^T \in R^n$$

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(\mathbf{x}) \\ \vdots \\ \frac{\partial f}{\partial x_n}(\mathbf{x}) \end{bmatrix}$$

“direction and rate of fastest increase”

Numerical Gradient vs Analytical Gradient

# What about Autograd?

- Deep learning frameworks can automatically perform backprop!
- Problems might surface related to underlying gradients when debugging your models

**“Yes You Should Understand Backprop”**

<https://medium.com/@karpathy/yes-you-should-understand-backprop-e2f06eab496b>

# Problem Statement: Backpropagation

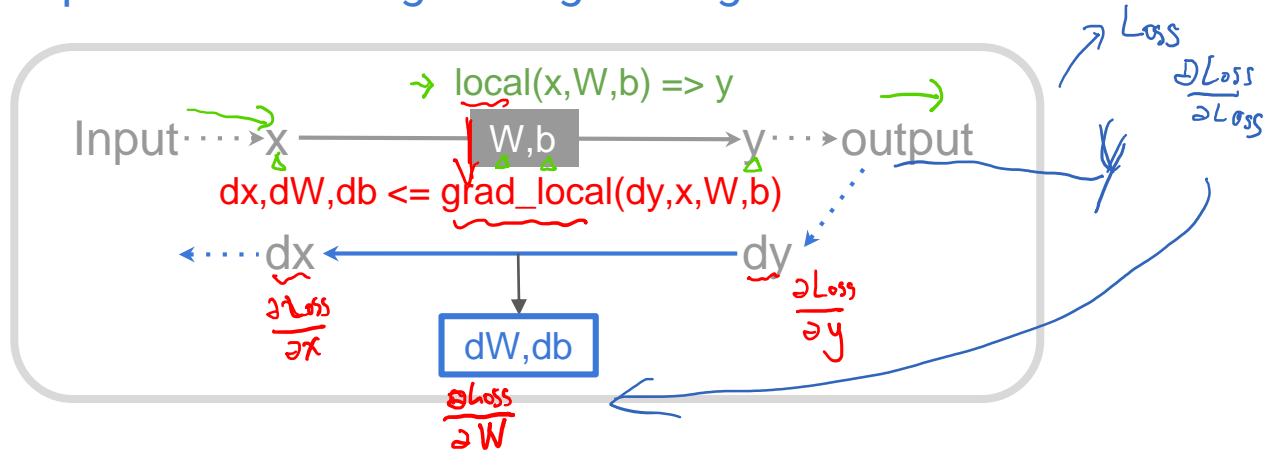
$$Loss = \frac{1}{N} \sum_i L_i(f(x_i, \theta), y_i)$$

Given a function ***f*** with respect to inputs ***x***, labels ***y***, and parameters ***θ***  
compute the gradient of **Loss** with respect to ***θ***

# Problem Statement: Backpropagation

An algorithm for computing the gradient of a **compound** function as a series of **local, intermediate gradients**:

1. Identify intermediate functions (forward prop)
2. Compute local gradients (chain rule)
3. Combine with upstream error signal to get full gradient



# Modularity: Previous Example

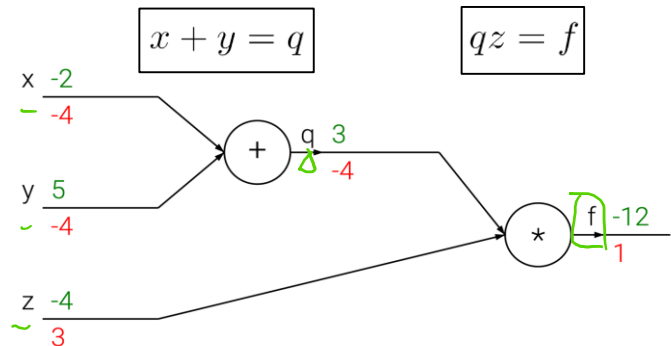
Compound function

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

Intermediate Variables  
(forward propagation)

$$q = x + y$$

$$f = qz$$



# Modularity: 2-Layer Neural Network

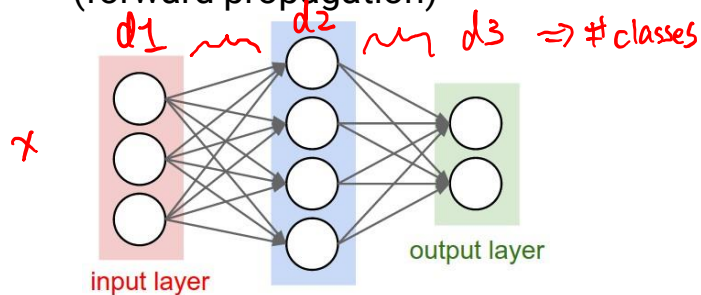
Compound function

$$Loss = L\left(\sigma(xW_1 + b_1)W_2 + b_2, y\right)$$

Intermediate Variables

(forward propagation)

*classifier*  
 $\Rightarrow \# \text{ classes}$

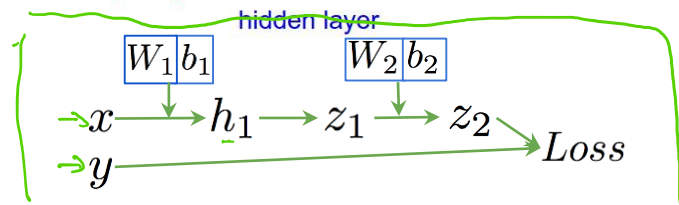


$$h_1 = xW_1 + b_1$$

$$z_1 = \sigma(h_1)$$

$$z_2 = z_1W_2 + b_2$$

$Loss \Rightarrow$  Squared Euclidean Distance  
between  $z_2$  and  $y$



# Intermediate Variables

(forward propagation)

coding, Batch size = N

$N \times d_1$

$d_1 \times d_2$

$$h_1 = x \underline{W_1} + b_1$$

↓

$N \times d_2$

↑

$$z_1 = \sigma(h_1)$$

$N \times d_2$   $d_2 \times d_3$

$$z_2 = z_1 W_2 + b_2$$

$N \times d_3$

only one image

$$\rightarrow ? f(x; W, b) = \underline{W}x + b ?$$

$$\Rightarrow [ ]_{d_2 \times 1}$$

(↑lecture note) Input one feature **vector**

(←here) Input a batch of data (**matrix**)

## Intermediate **Variables**

(forward propagation)

1. intermediate functions

2. local gradients

3. full gradients

## Intermediate **Gradients**

(backward propagation)

$$h_1 = xW_1 + b_1$$

$$z_1 = \sigma(h_1)$$

$$z_2 = z_1W_2 + b_2$$

$$Loss = \frac{1}{N} \|z_2 - y\|_F^2$$

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}$$

$$\begin{array}{l} \frac{\partial h_1}{\partial W_1}, \frac{\partial h_1}{\partial b_1} \quad \frac{\partial h_1}{\partial x} \text{ ? ? ? } W_1 \\ \frac{\partial z_1}{\partial h_1} \text{ ? ? } z_1(1 - z_1) \\ \frac{\partial z_2}{\partial W_2}, \frac{\partial z_2}{\partial b_2} \quad \frac{\partial z_2}{\partial z_1} \text{ ? ? ? } W_2 \end{array}$$

not scalar

$$\frac{\partial Loss}{\partial z_2} = \frac{2}{N} (z_2 - y)$$



# Agenda

- Motivation
- Backprop Tips & Tricks
- **Matrix calculus primer**

Previous :  $\frac{\partial \text{Scalar1}}{\partial \text{Scalar2}}$

## Derivative w.r.t. Vector

Scalar-by-Vector

$$\frac{\partial y}{\partial \mathbf{x}_n} = \left[ \frac{\partial y}{\partial x_1} \quad \frac{\partial y}{\partial x_2} \quad \cdots \quad \frac{\partial y}{\partial x_n} \right]$$

shape change

col  $\begin{bmatrix} \end{bmatrix}_{n \times 1}$   $\begin{bmatrix} \end{bmatrix}$  row  $\begin{bmatrix} \end{bmatrix}_{1 \times n}$

Vector-by-Vector

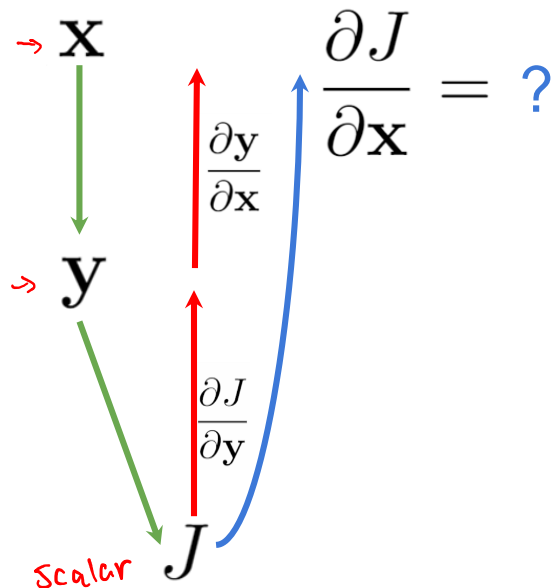
$$\frac{\partial \mathbf{y}_m}{\partial \mathbf{x}_n} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}_{m \times n}$$

# Derivative w.r.t. Vector: Chain Rule

1. intermediate functions

2. local gradients

3. full gradients



$$\begin{aligned}\frac{\partial J}{\partial x_j} &= \sum_i \frac{\partial J}{\partial y_i} \frac{\partial y_i}{\partial x_j} \\ &= \sum_i \left( \frac{\partial J}{\partial \mathbf{y}} \right)_{1i} \left( \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)_{ij}\end{aligned}$$

$$\boxed{\frac{\partial J}{\partial \mathbf{x}} = \frac{\partial J}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{x}}}$$

$1 \times n \quad 1 \times m \quad m \times n$

## Derivative w.r.t. Vector: Takeaway

$$\mathbf{y} = A_{m \times n} \mathbf{x} \quad \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = A$$
$$\frac{\partial J}{\partial \mathbf{x}} = \frac{\partial J}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \frac{\partial J}{\partial \mathbf{y}} A$$

---

elementwise function

$$\mathbf{y} = \omega(\mathbf{x})$$

$= \begin{bmatrix} \omega(x_1) \\ \vdots \\ \omega(x_n) \end{bmatrix}$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \omega'(x_1) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \omega'(x_n) \end{bmatrix}$$
$$\frac{\partial J}{\partial \mathbf{x}} = \frac{\partial J}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \left[ \frac{\partial J}{\partial y_1} \omega'(x_1), \dots, \frac{\partial J}{\partial y_n} \omega'(x_n) \right] = \frac{\partial J}{\partial \mathbf{y}} \circ \omega'(\mathbf{x})^T$$

# Derivative w.r.t. Matrix

Loss      input/param  
Scalar-by-Matrix

scalar

$$\frac{\partial y}{\partial A}_{m \times n} = \begin{bmatrix} \frac{\partial y}{\partial A_{11}} & \frac{\partial y}{\partial A_{12}} & \cdots & \frac{\partial y}{\partial A_{1n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y}{\partial A_{m1}} & \frac{\partial y}{\partial A_{m2}} & \cdots & \frac{\partial y}{\partial A_{mn}} \end{bmatrix}_{m \times n}$$

matrix

1D      2D  
Vector-by-Matrix → 3D

e.g.  $y = \sum_{i,j} A_{ij} = \boxed{A_{11}} + A_{12} + \cdots + A_{1n} + A_{21} + \cdots + A_{m1} + \cdots + A_{mn}$

?

$$\frac{\partial y}{\partial A_{ij}} = 1$$

$$\frac{\partial y}{\partial A_{11}} = \frac{\partial (A_{11} + \cdots)}{\partial A_{11}} = 1$$

# Derivative w.r.t. Matrix: Dimension Balancing

When you take **scalar-by-matrix** gradients



The gradient has **shape of denominator**

- Dimension balancing is the “cheap” but **efficient** approach to gradient calculations in most practical settings



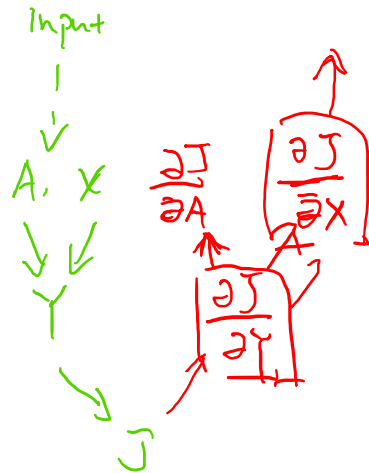
# Derivative w.r.t. Matrix: Takeaway

known  $\frac{\partial J}{\partial Y} \leftarrow \text{scalar}$

$$Y_{m \times l} = A_{m \times n} X_{n \times l}$$

$$\frac{\partial J}{\partial X}_{n \times l} = A^T_{n \times m} \frac{\partial J}{\partial Y}_{m \times l}$$

$$\frac{\partial J}{\partial A}_{m \times n} = \frac{\partial J}{\partial Y}_{m \times l} X^T_{l \times n}$$



elementwise function : activation function

$$Y_{m \times n} = \omega(X_{m \times n})$$

$$\frac{\partial J}{\partial X}_{m \times n} = \frac{\partial J}{\partial Y}_{m \times n} \overset{\text{elementwise mult}}{\odot} \omega'(X)_{m \times n}$$

## Intermediate **Variables**

(forward propagation)

Batch size:  $N$

$$\begin{aligned} & \overset{N \times d_1}{h_1} = xW_1 + b_1 \\ & \overset{N \times d_2}{z_1} = \sigma(h_1) \\ & \overset{N \times d_3}{z_2} = z_1W_2 + b_2 \quad \frac{\partial \text{Loss}}{\partial W_2} = z_1^T \frac{\partial \text{Loss}}{\partial z_2} \\ & \text{Loss} = \frac{1}{N} \|z_2 - y\|_F^2 \end{aligned}$$

1. intermediate functions

2. local gradients

3. full gradients

## Intermediate **Gradients**

(backward propagation)

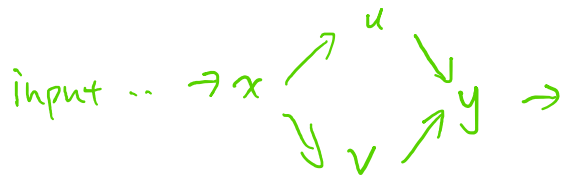
$$\begin{aligned} & \left[ \frac{\partial \text{Loss}}{\partial W_1}, \frac{\partial \text{Loss}}{\partial b_1} \right] \quad \frac{\partial \text{Loss}}{\partial x} \\ & \quad \uparrow \\ & \quad \frac{\partial \text{Loss}}{\partial h_1} \\ & \left[ \frac{\partial \text{Loss}}{\partial W_2}, \frac{\partial \text{Loss}}{\partial b_2} \right] \quad \frac{\partial \text{Loss}}{\partial z_1} = \frac{\partial \text{Loss}}{\partial z_2} \frac{\partial z_2}{\partial z_1} \\ & \quad \uparrow \\ & \quad \frac{\partial \text{Loss}}{\partial z_2} \end{aligned}$$

Red arrows indicate the flow of gradients from the loss back through the intermediate variables. A red circle highlights the gradient  $\frac{\partial \text{Loss}}{\partial z_2}$  and its contribution to the gradient of  $z_1$ .



# Backprop Menu for Success

1. Write down variable graph
2. Keep track of error signals
3. Compute derivative of loss function
4. Enforce shape rule on error signals, especially when deriving over a linear transformation



$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial u} \frac{\partial u}{\partial x} \boxed{+} \frac{\partial y}{\partial v} \frac{\partial v}{\partial x}$$



## Vector-by-vector

$$\mathbf{y} = A_{m \times n} \mathbf{x}$$

$$y_i = \sum_{j=1}^n A_{ij} x_j$$

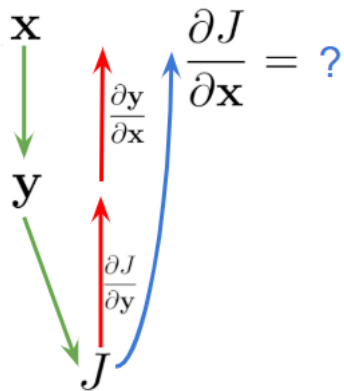
$$\frac{\partial y_i}{\partial x_j} = A_{ij}$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = A$$



## Vector-by-vector

$$\mathbf{y} = A_{m \times n} \mathbf{x}$$
$$\lambda^T = \frac{\partial J}{\partial \mathbf{y}}$$



$$y_i = \sum_{j=1}^n A_{ij} x_j$$
$$\lambda_i = \frac{\partial J}{\partial y_i}$$
$$\frac{\partial y_i}{\partial x_j} = A_{ij}$$

$$\frac{\partial J}{\partial x_j} = \sum_i \frac{\partial J}{\partial y_i} \frac{\partial y_i}{\partial x_j} = \sum_{i=1}^m \lambda_i A_{ij}$$

$$\frac{\partial J}{\partial \mathbf{x}} = \lambda^T A = \frac{\partial J}{\partial \mathbf{y}} A$$



## Vector-by-vector

$$\mathbf{y} = \omega(\mathbf{x})$$

$$y_i = \omega(x_i)$$

$$\frac{\partial y_i}{\partial x_i} = \omega'(x_i)$$

$$\frac{\partial y_i}{\partial x_j} = 0 \quad (i \neq j)$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \omega'(x_1) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \omega'(x_n) \end{bmatrix}$$



## Vector-by-vector

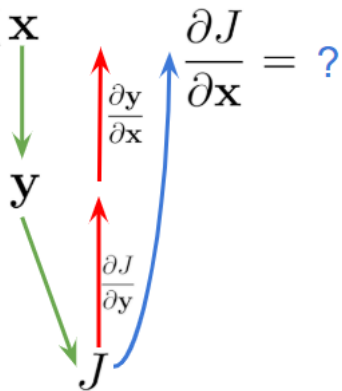
$$\mathbf{y} = \omega(\mathbf{x})$$
$$\lambda^T = \frac{\partial J}{\partial \mathbf{y}}$$

$$y_i = \omega(x_i)$$

$$\lambda_i = \frac{\partial J}{\partial y_i}$$

$$\frac{\partial y_i}{\partial x_i} = \omega'(x_i)$$

$$\frac{\partial y_i}{\partial x_j} = 0 \quad (i \neq j)$$



$$\frac{\partial J}{\partial x_j} = \sum_i \frac{\partial J}{\partial y_i} \frac{\partial y_i}{\partial x_j} = \frac{\partial J}{\partial y_j} \frac{\partial y_j}{\partial x_j} = \lambda_j \omega'(x_j)$$

$$\frac{\partial J}{\partial \mathbf{x}} = \lambda^T \circ \omega'(\mathbf{x})^T = \frac{\partial J}{\partial \mathbf{y}} \circ \omega'(\mathbf{x})^T$$



$$\mathbf{y} = A_{m \times n} \mathbf{x}$$

$$\frac{\partial J}{\partial \mathbf{x}} = \frac{\partial J}{\partial \mathbf{y}} A$$

$$\frac{\partial J}{\partial \mathbf{x}} = A^T \frac{\partial J}{\partial \mathbf{y}}$$

## Matrix multiplication [Backprop]

$$Y_{m \times l} = A_{m \times n} X_{n \times l}$$

$$\frac{\partial J}{\partial Y} = \Lambda_{m \times l}$$

$$Y_{ij} = \sum_{k=1}^n A_{ik} X_{kj}$$

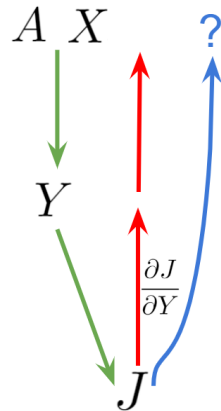
$$\Lambda_{ij} = \frac{\partial J}{\partial Y_{ij}}$$

$$\begin{aligned} \frac{\partial J}{\partial X_{ij}} &= \sum_{i',j'} \frac{\partial J}{\partial Y_{i'j'}} \frac{\partial Y_{i'j'}}{\partial X_{ij}} = \sum_{i'=1}^m \frac{\partial J}{\partial Y_{i'j}} \frac{\partial Y_{i'j}}{\partial X_{ij}} \\ &= \sum_{i'=1}^m A_{i'i} \Lambda_{i'j} = \sum_{k=1}^m A_{ki} \Lambda_{kj} \end{aligned}$$

$$\begin{aligned} \frac{\partial J}{\partial A_{ij}} &= \sum_{i',j'} \frac{\partial J}{\partial Y_{i'j'}} \frac{\partial Y_{i'j'}}{\partial A_{ij}} = \sum_{j'=1}^l \frac{\partial J}{\partial Y_{ij'}} \frac{\partial Y_{ij'}}{\partial A_{ij}} \\ &= \sum_{j'=1}^l \Lambda_{ij'} X_{jj'} = \sum_{k=1}^l \Lambda_{ik} X_{jk} \end{aligned}$$

$$\frac{\partial J}{\partial X} = A^T \Lambda = A^T \frac{\partial J}{\partial Y}$$

$$\frac{\partial J}{\partial A} = \Lambda X^T = \frac{\partial J}{\partial Y} X^T$$





$$\mathbf{y} = \omega(\mathbf{x}) \quad \frac{\partial J}{\partial \mathbf{x}} = \frac{\partial J}{\partial \mathbf{y}} \circ \omega'(\mathbf{x})^T$$
$$\frac{\partial J^T}{\partial \mathbf{x}} = \frac{\partial J^T}{\partial \mathbf{y}} \circ \omega'(\mathbf{x})$$

## Elementwise function [Backprop]

$$Y_{m \times n} = \omega(X_{m \times n})$$

$$\frac{\partial J}{\partial Y} = \Lambda_{m \times n}$$

$$Y_{ij} = \omega(X_{ij})$$

$$\Lambda_{ij} = \frac{\partial J}{\partial Y_{ij}}$$

$$\begin{aligned} \frac{\partial J}{\partial X_{ij}} &= \sum_{i', j'} \frac{\partial J}{\partial Y_{i' j'}} \frac{\partial Y_{i' j'}}{\partial X_{ij}} \\ &= \frac{\partial J}{\partial Y_{ij}} \frac{\partial Y_{ij}}{\partial X_{ij}} = \Lambda_{ij} \omega'(X_{ij}) \end{aligned}$$

$$\frac{\partial J}{\partial X} = \Lambda \circ \omega'(X) = \frac{\partial J}{\partial Y} \circ \omega'(X)$$

