

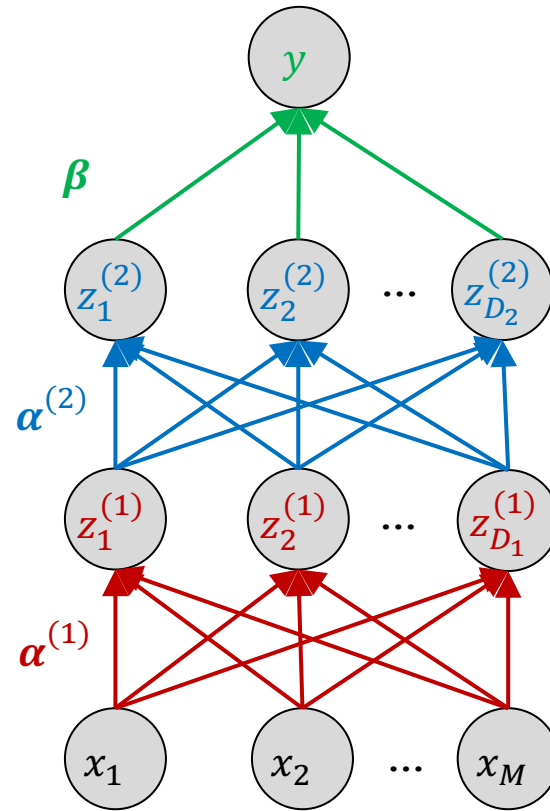
10-301/601: Introduction to Machine Learning

Lecture 13 – Differentiation

Henry Chai & Matt Gormley & Hoda Heidari

2/26/24

Recall: Neural Networks (Matrix Form)



$$\beta \in \mathbb{R}^{D_2}$$

$$\beta_0 \in \mathbb{R}$$

$$\alpha^{(2)} \in \mathbb{R}^{M \times D_2}$$

$$\mathbf{b}^{(2)} \in \mathbb{R}^{D_2}$$

$$\alpha^{(1)} \in \mathbb{R}^{M \times D_1}$$

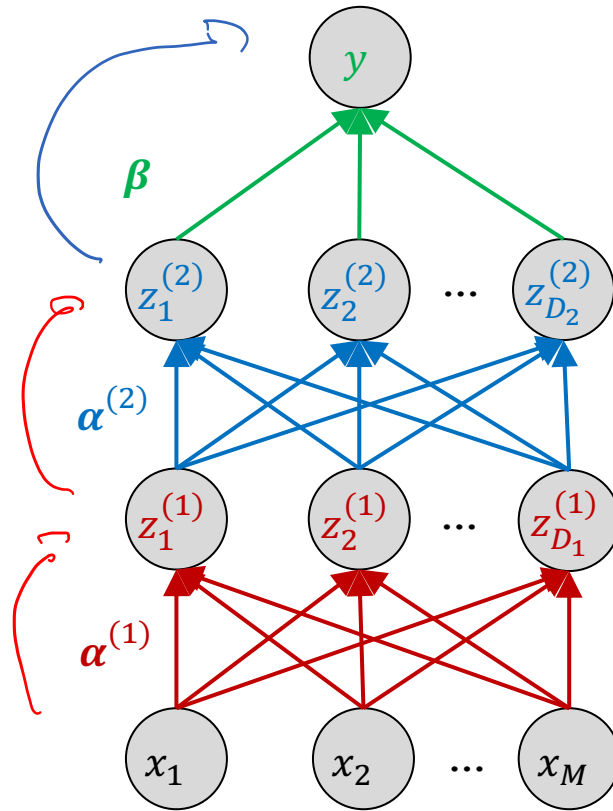
$$\mathbf{b}^{(1)} \in \mathbb{R}^{D_1}$$

$$y = \sigma((\beta)^T \mathbf{z}^{(2)} + \beta_0)$$

$$\mathbf{z}^{(2)} = \sigma((\alpha^{(2)})^T \mathbf{z}^{(1)} + \mathbf{b}^{(2)})$$

$$\mathbf{z}^{(1)} = \sigma((\alpha^{(1)})^T \mathbf{x} + \mathbf{b}^{(1)})$$

Recall: Neural Networks (Matrix Form)



$$y = \sigma \left(\boldsymbol{\beta}'^T \begin{bmatrix} 1 \\ \mathbf{z}^{(2)} \end{bmatrix} \right)$$

$$\boldsymbol{\beta}' = \begin{bmatrix} \beta_0 \\ \boldsymbol{\beta} \end{bmatrix} \in \mathbb{R}^{D_2+1}$$

$$\mathbf{z}^{(2)} = \sigma \left(\boldsymbol{\alpha}^{(2)'}^T \begin{bmatrix} 1 \\ \mathbf{z}^{(1)} \end{bmatrix} \right)$$

$$\boldsymbol{\alpha}^{(2)'} = \begin{bmatrix} \mathbf{b}^{(2)T} \\ \boldsymbol{\alpha}^{(2)} \end{bmatrix} \in \mathbb{R}^{(D_1+1) \times D_2}$$

$$\mathbf{z}^{(1)} = \sigma \left(\boldsymbol{\alpha}^{(1)'}^T \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix} \right)$$

$$\boldsymbol{\alpha}^{(1)'} = \begin{bmatrix} \mathbf{b}^{(1)T} \\ \boldsymbol{\alpha}^{(1)} \end{bmatrix} \in \mathbb{R}^{(M+1) \times D_1}$$

Forward Propagation for Making Predictions

- Inputs: weights $\alpha^{(1)}, \dots, \alpha^{(L)}, \beta$ and a query data point \mathbf{x}'
- Initialize $\mathbf{z}^{(0)} = \mathbf{x}'$
- For $l = 1, \dots, L$
 - $\mathbf{a}^{(l)} = \alpha^{(l)T} \mathbf{z}^{(l-1)}$
 - $\mathbf{z}^{(l)} = \sigma(\mathbf{a}^{(l)})$
- $\hat{y} = \sigma(\beta^T \mathbf{z}^{(L)})$
- Output: the prediction \hat{y}

Stochastic Gradient Descent for Learning

- Input: $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N, \gamma$
- Initialize all weights $\boldsymbol{\alpha}^{(1)}, \dots, \boldsymbol{\alpha}^{(L)}, \boldsymbol{\beta}$
- While TERMINATION CRITERION is not satisfied
 - For $i \in \text{shuffle}(\{1, \dots, N\})$
 - Compute $\mathbf{g}_{\boldsymbol{\beta}} = \nabla_{\boldsymbol{\beta}} J^{(i)}(\boldsymbol{\alpha}^{(1)}, \dots, \boldsymbol{\alpha}^{(L)}, \boldsymbol{\beta})$
 - For $l = 1, \dots, L$
 - Compute $\mathbf{g}_{\boldsymbol{\alpha}^{(l)}} = \nabla_{\boldsymbol{\alpha}^{(l)}} J^{(i)}(\boldsymbol{\alpha}^{(1)}, \dots, \boldsymbol{\alpha}^{(L)}, \boldsymbol{\beta})$
 - Update $\boldsymbol{\beta} = \boldsymbol{\beta} - \gamma \mathbf{g}_{\boldsymbol{\beta}}$
 - For $l = 1, \dots, L$
 - Update $\boldsymbol{\alpha}^{(l)} = \boldsymbol{\alpha}^{(l)} - \gamma \mathbf{g}_{\boldsymbol{\alpha}^{(l)}}$
- Output: $\boldsymbol{\alpha}^{(1)}, \dots, \boldsymbol{\alpha}^{(L)}, \boldsymbol{\beta}$

Two questions:

1. What is this loss function $J^{(i)}$?

2. How on earth do we take these gradients?

- Input: $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N, \gamma$
- Initialize all weights $\boldsymbol{\alpha}^{(1)}, \dots, \boldsymbol{\alpha}^{(L)}, \boldsymbol{\beta}$
- While TERMINATION CRITERION is not satisfied
 - For $i \in \text{shuffle}(\{1, \dots, N\})$
 - Compute $\mathbf{g}_{\boldsymbol{\beta}} = \nabla_{\boldsymbol{\beta}} J^{(i)}(\boldsymbol{\alpha}^{(1)}, \dots, \boldsymbol{\alpha}^{(L)}, \boldsymbol{\beta})$
 - For $l = 1, \dots, L$
 - Compute $\mathbf{g}_{\boldsymbol{\alpha}^{(l)}} = \nabla_{\boldsymbol{\alpha}^{(l)}} J^{(i)}(\boldsymbol{\alpha}^{(1)}, \dots, \boldsymbol{\alpha}^{(L)}, \boldsymbol{\beta})$
 - Update $\boldsymbol{\beta} = \boldsymbol{\beta} - \gamma \mathbf{g}_{\boldsymbol{\beta}}$
 - For $l = 1, \dots, L$
 - Update $\boldsymbol{\alpha}^{(l)} = \boldsymbol{\alpha}^{(l)} - \gamma \mathbf{g}_{\boldsymbol{\alpha}^{(l)}}$
- Output: $\boldsymbol{\alpha}^{(1)}, \dots, \boldsymbol{\alpha}^{(L)}, \boldsymbol{\beta}$

Two questions:

1. What is this loss function $J^{(i)}$?

2. How on earth do we take these gradients?

- Input: $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N, \gamma$
- Initialize all weights $\boldsymbol{\alpha}^{(1)}, \dots, \boldsymbol{\alpha}^{(L)}, \boldsymbol{\beta}$
- While TERMINATION CRITERION is not satisfied
 - For $i \in \text{shuffle}(\{1, \dots, N\})$
 - Compute $\mathbf{g}_{\boldsymbol{\beta}} = \nabla_{\boldsymbol{\beta}} J^{(i)}(\boldsymbol{\alpha}^{(1)}, \dots, \boldsymbol{\alpha}^{(L)}, \boldsymbol{\beta})$
 - For $l = 1, \dots, L$
 - Compute $\mathbf{g}_{\boldsymbol{\alpha}^{(l)}} = \nabla_{\boldsymbol{\alpha}^{(l)}} J^{(i)}(\boldsymbol{\alpha}^{(1)}, \dots, \boldsymbol{\alpha}^{(L)}, \boldsymbol{\beta})$
 - Update $\boldsymbol{\beta} = \boldsymbol{\beta} - \gamma \mathbf{g}_{\boldsymbol{\beta}}$
 - For $l = 1, \dots, L$
 - Update $\boldsymbol{\alpha}^{(l)} = \boldsymbol{\alpha}^{(l)} - \gamma \mathbf{g}_{\boldsymbol{\alpha}^{(l)}}$
- Output: $\boldsymbol{\alpha}^{(1)}, \dots, \boldsymbol{\alpha}^{(L)}, \boldsymbol{\beta}$

Loss Functions for Neural Networks

- Let $\Theta = \{\alpha^{(1)}, \dots, \alpha^{(L)}, \beta\}$ be the parameters of our neural network
- Regression - squared error (same as linear regression!)

$$J^{(i)}(\Theta) = (\hat{y}_{\Theta}(\mathbf{x}^{(i)}) - y^{(i)})^2$$

- Binary classification - cross-entropy loss (same as logistic regression!)
 - Assume $Y \in \{0,1\}$ and $P(Y = 1|\mathbf{x}, \Theta) = \hat{y}_{\Theta}(\mathbf{x})$

$$J^{(i)}(\Theta) = -\log P(y^{(i)}|\mathbf{x}^{(i)}, \Theta)$$

Loss Functions for Neural Networks

- Let $\Theta = \{\alpha^{(1)}, \dots, \alpha^{(L)}, \beta\}$ be the parameters of our neural network
- Multi-class classification - cross-entropy loss again!

- Express the label as a one-hot or one-of- C vector e.g.,

$$\mathbf{y} = [0 \quad 0 \quad 1 \quad 0 \quad \dots \quad 0]$$

- Assume the neural network output is also a vector of length C , $\hat{\mathbf{y}}_{\Theta}$

$$P(\mathbf{y}[c] = 1 | \mathbf{x}, \Theta) = \hat{\mathbf{y}}_{\Theta}(\mathbf{x}^{(i)})[c]$$

Okay but
how do
we get
our
network
to output
this
vector?

- Let $\Theta = \{\alpha^{(1)}, \dots, \alpha^{(L)}, \beta\}$ be the parameters of our neural network

- Multi-class classification - cross-entropy loss

- Express the label as a one-hot or one-of- C vector e.g.,

$$\mathbf{y} = [0 \quad 0 \quad 1 \quad 0 \quad \dots \quad 0]$$

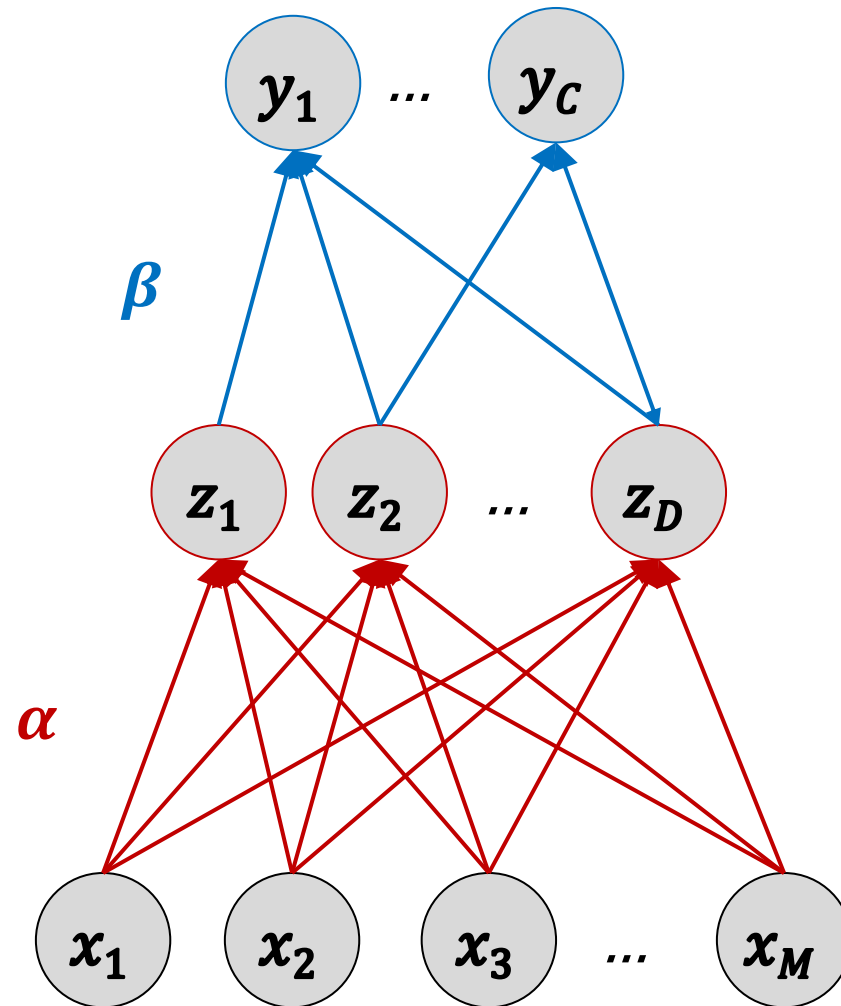
- Assume the neural network output is also a vector of length C , $\hat{\mathbf{y}}_{\Theta}$

$$P(\mathbf{y}[c] = 1 | \mathbf{x}, \Theta) = \hat{\mathbf{y}}_{\Theta}(\mathbf{x}^{(i)})[c]$$

- Then the cross-entropy loss is

$$\begin{aligned} J^{(i)}(\Theta) &= -\log P(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}, \Theta) \\ &= -\sum_{c=1}^C \mathbf{y}^{(i)}[c] \log(\hat{\mathbf{y}}_{\Theta}(\mathbf{x}^{(i)})[c]) \end{aligned}$$

Softmax



$$y_c = \frac{\exp b_c}{\sum_{k=1}^C \exp b_k}$$

$$b_c = \sum_{j=0}^D \beta_{c,j} z_j$$

$$z_d = \sigma(a_d)$$

$$a_d = \sum_{i=0}^M \alpha_{d,i} x_i$$

Two questions:

1. What is this loss function $J^{(i)}$?

2. How on earth do we take these gradients?

- Input: $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N, \gamma$
- Initialize all weights $\boldsymbol{\alpha}^{(1)}, \dots, \boldsymbol{\alpha}^{(L)}, \boldsymbol{\beta}$
- While TERMINATION CRITERION is not satisfied
 - For $i \in \text{shuffle}(\{1, \dots, N\})$
 - Compute $\mathbf{g}_{\boldsymbol{\beta}} = \nabla_{\boldsymbol{\beta}} J^{(i)}(\boldsymbol{\alpha}^{(1)}, \dots, \boldsymbol{\alpha}^{(L)}, \boldsymbol{\beta})$
 - For $l = 1, \dots, L$
 - Compute $\mathbf{g}_{\boldsymbol{\alpha}^{(l)}} = \nabla_{\boldsymbol{\alpha}^{(l)}} J^{(i)}(\boldsymbol{\alpha}^{(1)}, \dots, \boldsymbol{\alpha}^{(L)}, \boldsymbol{\beta})$
 - Update $\boldsymbol{\beta} = \boldsymbol{\beta} - \gamma \mathbf{g}_{\boldsymbol{\beta}}$
 - For $l = 1, \dots, L$
 - Update $\boldsymbol{\alpha}^{(l)} = \boldsymbol{\alpha}^{(l)} - \gamma \mathbf{g}_{\boldsymbol{\alpha}^{(l)}}$
- Output: $\boldsymbol{\alpha}^{(1)}, \dots, \boldsymbol{\alpha}^{(L)}, \boldsymbol{\beta}$

Matrix Calculus

		Numerator		
		scalar	vector	matrix
Denominator	Types of Derivatives			
	scalar	$\frac{\partial y}{\partial x}$	$\frac{\partial \mathbf{y}}{\partial x}$	$\frac{\partial \mathbf{Y}}{\partial x}$
	vector	$\frac{\partial y}{\partial \mathbf{x}}$	$\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$	$\frac{\partial \mathbf{Y}}{\partial \mathbf{x}}$
matrix	$\frac{\partial y}{\partial \mathbf{X}}$	$\frac{\partial \mathbf{y}}{\partial \mathbf{X}}$	$\frac{\partial \mathbf{Y}}{\partial \mathbf{X}}$	

Matrix Calculus: Denominator Layout

- Derivatives of a scalar always have the *same shape* as the entity that the derivative is being taken with respect to.

<i>Types of Derivatives</i>	scalar
scalar	$\frac{\partial y}{\partial x} = \left[\frac{\partial y}{\partial x} \right]$
vector	$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_P} \end{bmatrix}$
matrix	$\frac{\partial y}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial y}{\partial X_{11}} & \frac{\partial y}{\partial X_{12}} & \cdots & \frac{\partial y}{\partial X_{1Q}} \\ \frac{\partial y}{\partial X_{21}} & \frac{\partial y}{\partial X_{22}} & \cdots & \frac{\partial y}{\partial X_{2Q}} \\ \vdots & & & \vdots \\ \frac{\partial y}{\partial X_{P1}} & \frac{\partial y}{\partial X_{P2}} & \cdots & \frac{\partial y}{\partial X_{PQ}} \end{bmatrix}$

Matrix Calculus: Denominator Layout

<i>Types of Derivatives</i>	scalar	vector
scalar	$\frac{\partial y}{\partial x} = \left[\frac{\partial y}{\partial x} \right]$	$\frac{\partial \mathbf{y}}{\partial x} = \left[\frac{\partial y_1}{\partial x} \quad \frac{\partial y_2}{\partial x} \quad \dots \quad \frac{\partial y_N}{\partial x} \right]$
vector	$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_P} \end{bmatrix}$	$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} & \dots & \frac{\partial y_N}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_N}{\partial x_2} \\ \vdots & \vdots & \dots & \vdots \\ \frac{\partial y_1}{\partial x_P} & \frac{\partial y_2}{\partial x_P} & \dots & \frac{\partial y_N}{\partial x_P} \end{bmatrix}$

Three Approaches to Differentiation

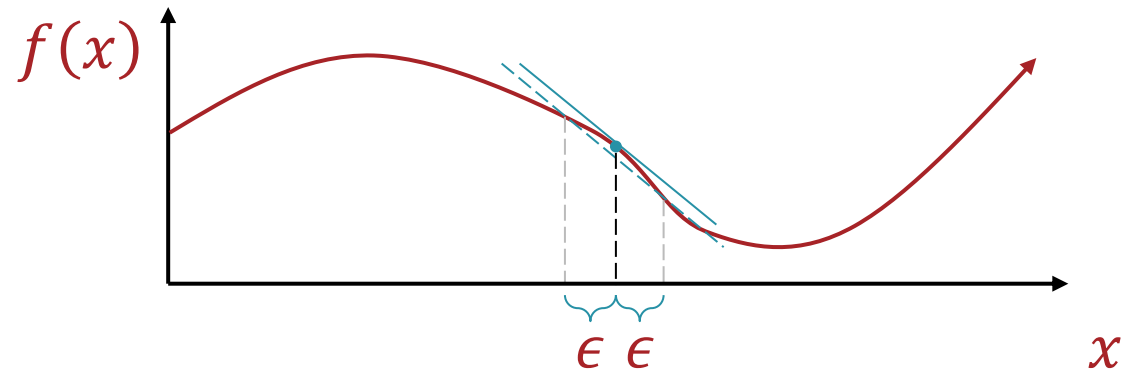
• Given $f: \mathbb{R}^D \rightarrow \mathbb{R}$, compute $\nabla_{\mathbf{x}} f(\mathbf{x}) = \partial f(\mathbf{x}) / \partial \mathbf{x}$

1. Finite difference method
2. Symbolic differentiation
3. Automatic differentiation (reverse mode)

Approach 1: Finite Difference Method

- Given $f: \mathbb{R}^D \rightarrow \mathbb{R}$, compute $\nabla_x f(\mathbf{x}) = \partial f(\mathbf{x}) / \partial \mathbf{x}$
$$\frac{\partial f(\mathbf{x})}{\partial x_i} \approx \frac{f(\mathbf{x} + \epsilon \mathbf{d}_i) - f(\mathbf{x} - \epsilon \mathbf{d}_i)}{2\epsilon}$$

where \mathbf{d}_i is a one-hot vector with a 1 in the i^{th} position



- We want ϵ to be small to get a good approximation but we run into floating point issues when ϵ is too small
- Getting the full gradient requires computing the above approximation for each dimension of the input

Approach 1: Finite Difference Method Example

- Given

$$y = f(x, z) = e^{xz} + \frac{xz}{\ln(x)} + \frac{\sin(\ln(x))}{xz}$$

what are $\frac{\partial y}{\partial x}$ and $\frac{\partial y}{\partial z}$ at $x = 2, z = 3$?

Three Approaches to Differentiation

- Given $f: \mathbb{R}^D \rightarrow \mathbb{R}$, compute $\nabla_{\mathbf{x}} f(\mathbf{x}) = \partial f(\mathbf{x}) / \partial \mathbf{x}$
- 1. Finite difference method
 - Requires the ability to call $f(\mathbf{x})$
 - Great for checking accuracy of implementations of more complex differentiation methods
 - Computationally expensive for high-dimensional inputs
- 2. Symbolic differentiation
- 3. Automatic differentiation (reverse mode)

Approach 2: Symbolic Differentiation

- Given

$$y = f(x, z) = e^{xz} + \frac{xz}{\ln(x)} + \frac{\sin(\ln(x))}{xz}$$

what are $\frac{\partial y}{\partial x}$ and $\frac{\partial y}{\partial z}$ at $x = 2, z = 3$?

The Chain Rule of Calculus

- If $y = f(z)$ and $z = g(x)$ then the corresponding computation graph is
- If $y = f(z_1, z_2)$ and $z_1 = g_1(x), z_2 = g_2(x)$ then
- If $y = f(\mathbf{z})$ and $\mathbf{z} = g(x)$ then

Poll Question 1

- If $y = f(\mathbf{z})$, $\mathbf{z} = g(\mathbf{w})$ and $\mathbf{w} = h(x)$, does the equation

$$\frac{\partial y}{\partial x} = \sum_{d=1}^D \frac{\partial y}{\partial z_d} \frac{\partial z_d}{\partial x}$$

still hold?

- A. Yes
- B. No
- C. Only on Fridays (TOXIC)

Approach 2: Symbolic Differentiation

- Given

$$y = f(x, z) = e^{xz} + \frac{xz}{\ln(x)} + \frac{\sin(\ln(x))}{xz}$$

what are $\frac{\partial y}{\partial x}$ and $\frac{\partial y}{\partial z}$ at $x = 2, z = 3$?

Three Approaches to Differentiation

- Given $f: \mathbb{R}^D \rightarrow \mathbb{R}$, compute $\nabla_{\mathbf{x}} f(\mathbf{x}) = \partial f(\mathbf{x}) / \partial \mathbf{x}$
1. Finite difference method
 - Requires the ability to call $f(\mathbf{x})$
 - Great for checking accuracy of implementations of more complex differentiation methods
 - Computationally expensive for high-dimensional inputs
 2. Symbolic differentiation
 - Requires systematic knowledge of derivatives
 - Can be computationally expensive if poorly implemented
 3. Automatic differentiation (reverse mode)

Approach 3: Automatic Differentiation (reverse mode)

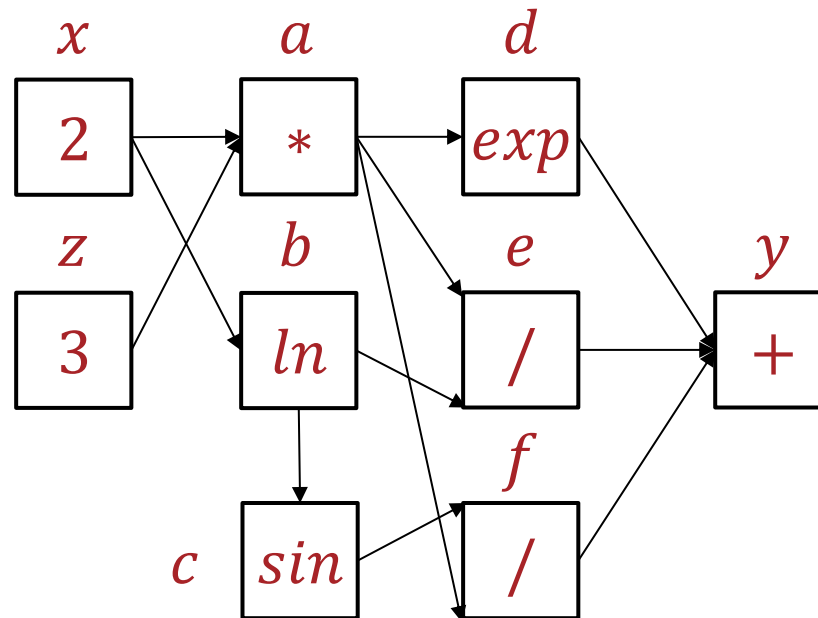
- Given

$$y = f(x, z) = e^{xz} + \frac{xz}{\ln(x)} + \frac{\sin(\ln(x))}{xz}$$

what are $\frac{\partial y}{\partial x}$ and $\frac{\partial y}{\partial z}$ at $x = 2, z = 3$?

- First define some intermediate quantities, draw the computation graph and run the “forward” computation

$$\begin{aligned} a &= xz \\ b &= \ln(x) \\ c &= \sin(b) \\ d &= e^a \\ e &= a/b \\ f &= c/a \\ y &= d + e + f \end{aligned}$$



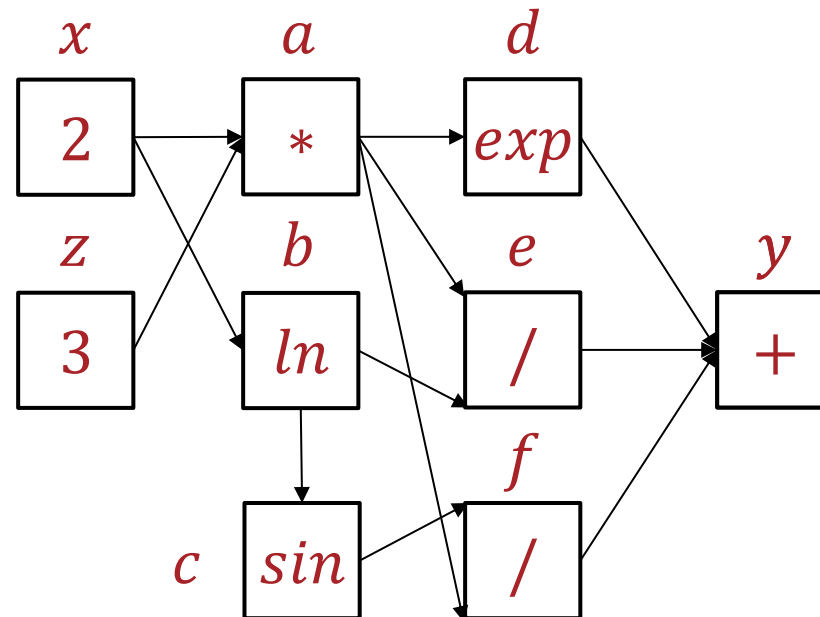
Approach 3: Automatic Differentiation (reverse mode)

- Given

$$y = f(x, z) = e^{xz} + \frac{xz}{\ln(x)} + \frac{\sin(\ln(x))}{xz}$$

what are $\frac{\partial y}{\partial x}$ and $\frac{\partial y}{\partial z}$ at $x = 2, z = 3$?

- Then compute partial derivatives, starting from y and working back



Three Approaches to Differentiation

- Given $f: \mathbb{R}^D \rightarrow \mathbb{R}$, compute $\nabla_{\mathbf{x}} f(\mathbf{x}) = \partial f(\mathbf{x}) / \partial \mathbf{x}$
- 1. Finite difference method
 - Requires the ability to call $f(\mathbf{x})$
 - Great for checking accuracy of implementations of more complex differentiation methods
 - Computationally expensive for high-dimensional inputs
- 2. Symbolic differentiation
 - Requires systematic knowledge of derivatives
 - Can be computationally expensive if poorly implemented
- 3. Automatic differentiation (reverse mode)
 - Requires systematic knowledge of derivatives *and* an algorithm for computing $f(\mathbf{x})$
 - Computational cost of computing $\partial f(\mathbf{x}) / \partial \mathbf{x}$ is proportional to the cost of computing $f(\mathbf{x})$

Automatic Differentiation

- Given $f: \mathbb{R}^D \rightarrow \mathbb{R}^C$, compute $\nabla_{\mathbf{x}} f(\mathbf{x}) = \partial f(\mathbf{x}) / \partial \mathbf{x}$
3. Automatic differentiation (reverse mode)
 - Requires systematic knowledge of derivatives *and* an algorithm for computing $f(\mathbf{x})$
 - Computational cost of computing $\nabla_{\mathbf{x}} f(\mathbf{x})_c = \partial f(\mathbf{x})_c / \partial \mathbf{x}$ is proportional to the cost of computing $f(\mathbf{x})$
 - Great for high-dimensional inputs and low-dimensional outputs ($D \gg C$)
 4. Automatic differentiation (forward mode)
 - Requires systematic knowledge of derivatives *and* an algorithm for computing $f(\mathbf{x})$
 - Computational cost of computing $\partial f(\mathbf{x}) / \partial x_d$ is proportional to the cost of computing $f(\mathbf{x})$
 - Great for low-dimensional inputs and high-dimensional outputs ($D \ll C$)

Computation Graph: 10-301/601 Conventions

- The diagram represents *an algorithm*
- Nodes are rectangles with one node per intermediate variable in the algorithm
- Each node is labeled with the function that it computes (inside the box) and the variable name (outside the box)
- Edges are directed and do not have labels
- For neural networks:
 - Each weight, feature value, label and *bias term* appears as a node
 - *We can* include the loss function

Neural Network Diagram Conventions

- The diagram represents a *neural network*
- Nodes are circles with one node per hidden unit
- Each node is labeled with the variable corresponding to the hidden unit
- Edges are directed and each edge is labeled with its weight
- Following standard convention, the bias term is typically *not* shown as a node, but rather is assumed to be part of the activation function i.e., its weight does not appear in the picture anywhere.
- The diagram typically does *not* include any nodes related to the loss computation

Backprop Learning Objectives

You should be able to...

- Differentiate between a neural network diagram and a computation graph
- Construct a computation graph for a function as specified by an algorithm
- Carry out the backpropagation on an arbitrary computation graph
- Construct a computation graph for a neural network, identifying all the given and intermediate quantities that are relevant
- Instantiate the backpropagation algorithm for a neural network
- Instantiate an optimization method (e.g. SGD) and a regularizer (e.g. L2) when the parameters of a model are comprised of several matrices corresponding to different layers of a neural network
- Use the finite difference method to evaluate the gradient of a function
- Identify when the gradient of a function can be computed at all and when it can be computed efficiently
- Employ basic matrix calculus to compute vector/matrix/tensor derivatives.