

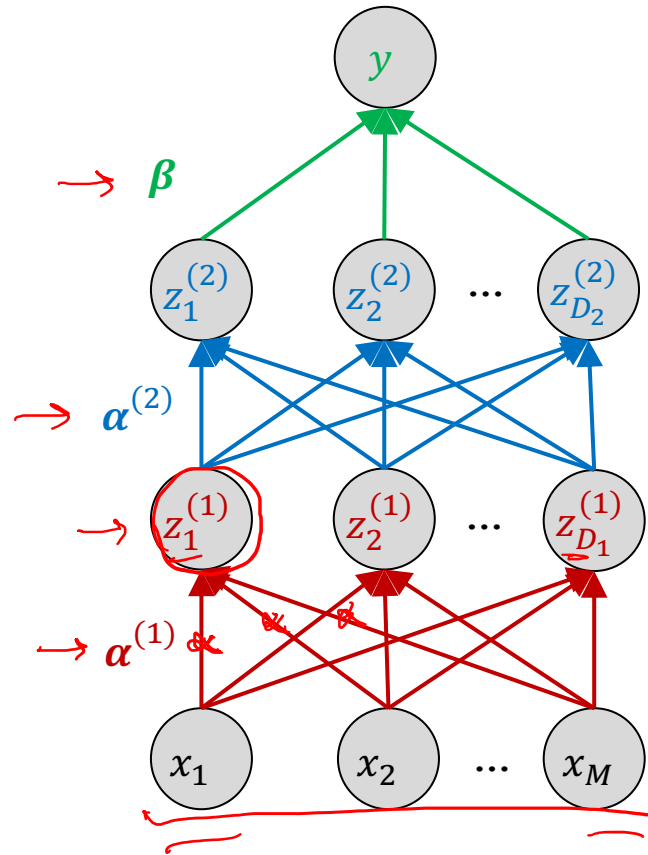
10-301/601: Introduction to Machine Learning

Lecture 13 – Differentiation

Henry Chai & Matt Gormley & Hoda Heidari

2/26/24

Recall: Neural Networks (Matrix Form)



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

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

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

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

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

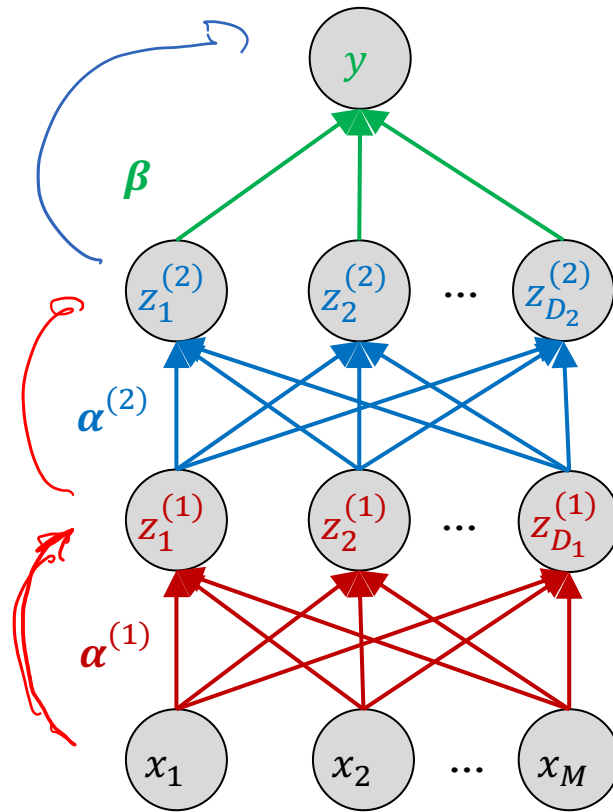
$$\mathbf{b}^{(1)} \in \mathbb{R}^{D_1}$$

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

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

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

Recall: Neural Networks (Matrix Form)



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

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

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

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

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

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

Forward Propagation for Making Predictions

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

Stochastic Gradient Descent for Learning

- Input: $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N, \gamma$ $J = \text{obj. fu.}$

- Initialize all weights $\alpha^{(1)}, \dots, \alpha^{(L)}, \beta$

- While TERMINATION CRITERION is not satisfied

- For $i \in \text{shuffle}(\{1, \dots, N\})$

gradient
computation

- Compute $\underline{g}_{\beta} = \nabla_{\underline{\beta}} J^{(i)}(\alpha^{(1)}, \dots, \alpha^{(L)}, \beta)$

- For $l = 1, \dots, L$

- Compute $\underline{g}_{\alpha^{(l)}} = \nabla_{\alpha^{(l)}} J^{(i)}(\alpha^{(1)}, \dots, \alpha^{(L)}, \beta)$

descent

- Update $\beta = \beta - \gamma g_{\beta}$

- For $l = 1, \dots, L$

- Update $\alpha^{(l)} = \alpha^{(l)} - \gamma g_{\alpha^{(l)}}$

- Output: $\alpha^{(1)}, \dots, \alpha^{(L)}, \beta$

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!)

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

\rightarrow The prediction of NN w. param Θ on $x^{(i)}$

- Binary classification - cross-entropy loss (same as logistic regression!)

- Assume $Y \in \{0,1\}$ and $P(Y=1|x, \Theta) = \hat{y}_{\Theta}(x)$

$$J^{(i)}(\Theta) = -\log P(y^{(i)} | x^{(i)}, \Theta) = \begin{cases} \text{if } y_i=1 \rightarrow \hat{y}_{\Theta}(x^{(i)}) \leftarrow \\ \text{if } y_i=0 \rightarrow 1 - \hat{y}_{\Theta}(x^{(i)}) \leftarrow \end{cases}$$

$$= -\log \left(\hat{y}_{\Theta}(x^{(i)})^{y_i} (1 - \hat{y}_{\Theta}(x^{(i)}))^{(1-y_i)} \right)$$

$$= - \left(\underline{y_i} \log \hat{y}_{\Theta}(x^{(i)}) + \underline{(1-y_i)} \log (1 - \hat{y}_{\Theta}(x^{(i)})) \right)$$

$$\log b^a = a \log b$$

Loss Functions for Neural Networks

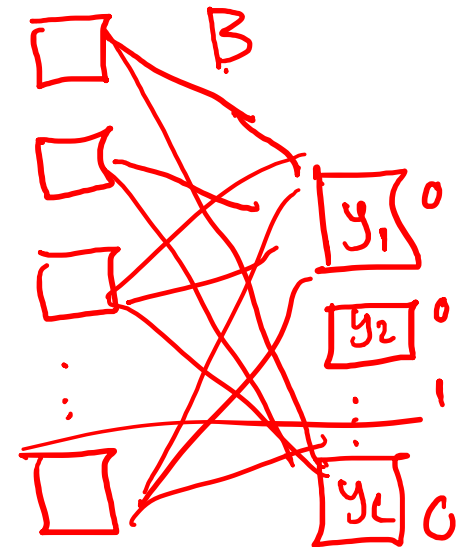
$$y = \{1, 2, \dots, C\} \quad (\vec{x}, \mathbf{z}) \rightarrow (\vec{x}, \underbrace{[\overset{t}{0}, \overset{t}{0}, \overset{t}{1}, \dots, \overset{t}{0}]}_C)$$

- 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.,

$$y = [0 \quad 0 \quad 1 \quad 0 \quad \dots \quad 0]$$

- Assume the neural network output is also a vector of length C , \hat{y}_Θ

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



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 , $\hat{\mathbf{y}}_{\Theta}$

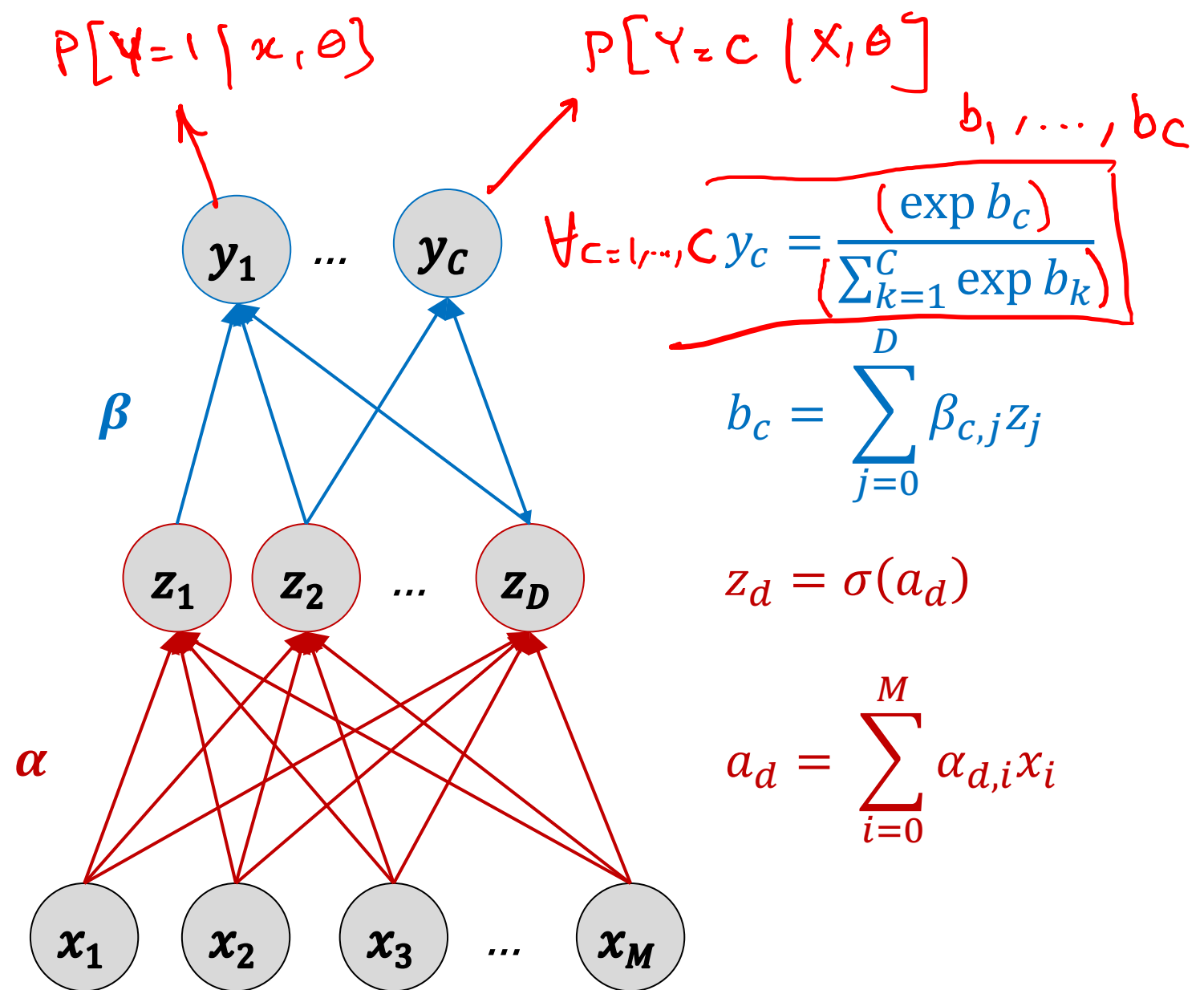
$$P(\mathbf{y}[c] = 1 | \mathbf{x}, \Theta) = \hat{\mathbf{y}}_{\Theta}(\mathbf{x}^{(i)})[c]$$

- Then the cross-entropy loss is

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

$$= -\sum_{c=1}^C \underbrace{\mathbf{y}^{(i)}[c]}_{\text{circled}} \log(\underbrace{\hat{\mathbf{y}}_{\Theta}(\mathbf{x}^{(i)})[c]}_{\text{boxed}})$$

Softmax



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 $\alpha^{(1)}, \dots, \alpha^{(L)}, \beta$
- While TERMINATION CRITERION is not satisfied
 - For $i \in \text{shuffle}(\{1, \dots, N\})$
 - Compute $g_{\beta} = \nabla_{\beta} J^{(i)}(\alpha^{(1)}, \dots, \alpha^{(L)}, \beta)$
 - For $l = 1, \dots, L$
 - Compute $g_{\alpha^{(l)}} = \nabla_{\alpha^{(l)}} J^{(i)}(\alpha^{(1)}, \dots, \alpha^{(L)}, \beta)$
 - Update $\beta = \beta - \gamma g_{\beta}$
 - For $l = 1, \dots, L$
 - Update $\alpha^{(l)} = \alpha^{(l)} - \gamma g_{\alpha^{(l)}}$
- Output: $\alpha^{(1)}, \dots, \alpha^{(L)}, \beta$

Matrix Calculus

$$\frac{\partial \square}{\Delta \triangle}$$

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

Table courtesy of Matt Gormley

Matrix Calculus: Denominator Layout

$$\frac{\partial y}{\partial \Delta}$$

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

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

Matrix Calculus: Denominator Layout

Types of Derivatives	scalar	vector
scalar	$\frac{\partial y}{\partial x} = \left[\frac{\partial y}{\partial x} \right]$	$\mathbf{y} = (y_1, \dots, y_N) \downarrow$ $\frac{\partial \mathbf{y}}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x} & \frac{\partial y_2}{\partial x} & \dots & \frac{\partial y_N}{\partial x} \end{bmatrix}$
vector	$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_P} \end{bmatrix}$	$\mathbf{y} \in \mathbb{R}^N$ $\mathbf{x} \in \mathbb{R}^P$ $\in \mathbb{R}^{P \times N}$ $\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} & \dots & \frac{\partial y_N}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_N}{\partial x_2} \\ \vdots & \vdots & \dots & \vdots \\ \frac{\partial y_1}{\partial x_P} & \frac{\partial y_2}{\partial x_P} & \dots & \frac{\partial y_N}{\partial x_P} \end{bmatrix}$

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

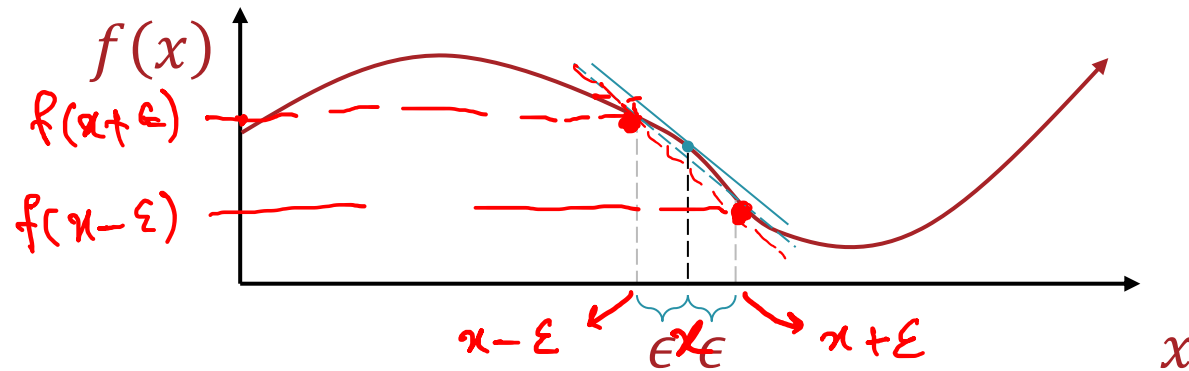
3. Automatic differentiation (reverse mode)



Approach 1: Finite Difference Method

- Given $f: \mathbb{R}^D \rightarrow \mathbb{R}$, compute $\nabla_x f(x) = \partial f(x) / \partial x$
$$\frac{\partial f(x)}{\partial x_i} \approx \frac{f(x + \epsilon d_i) - f(x - \epsilon d_i)}{2\epsilon}$$
 where $d_i = (0, 0, \dots, \overset{i}{1}, \dots, 0)$

where 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

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

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

$$\frac{\partial y}{\partial x} \approx \frac{f(x+\varepsilon, z) - f(x-\varepsilon, z)}{2\varepsilon} \quad \left\{ \begin{array}{l} f(x+\varepsilon, z) \\ f(x-\varepsilon, z) \end{array} \right.$$

$$\frac{\partial y}{\partial z} \approx \frac{f(x, z+\varepsilon) - f(x, z-\varepsilon)}{2\varepsilon} \quad \left\{ \begin{array}{l} f(x, z+\varepsilon) \\ f(x, z-\varepsilon) \end{array} \right.$$

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

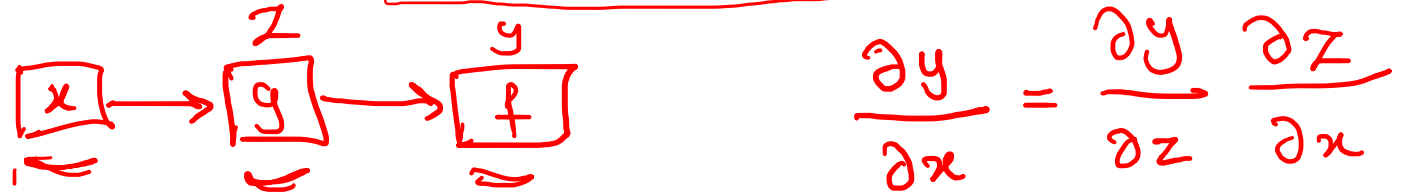
- Given

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

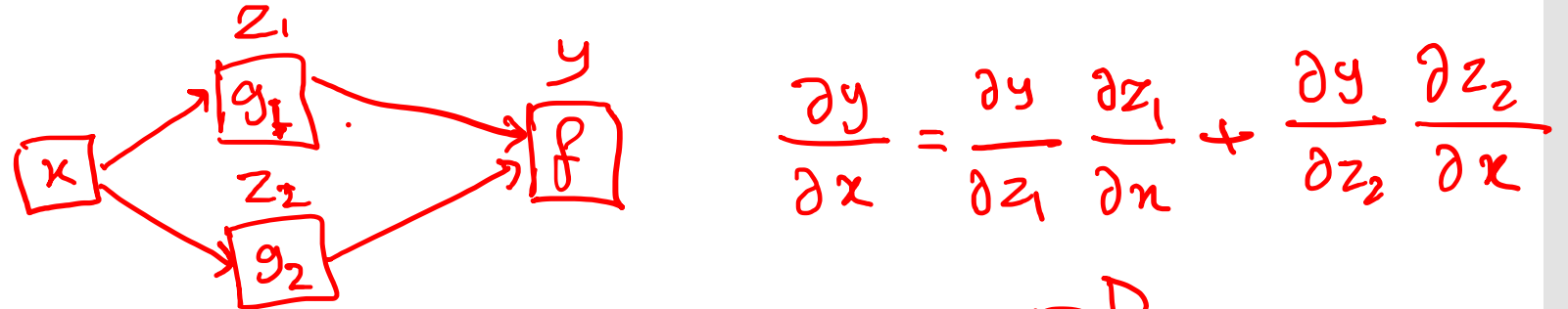
what are $\frac{\partial y}{\partial x}$ and $\frac{\partial y}{\partial z}$ at $x = 2, z = 3$?

The Chain Rule of Calculus

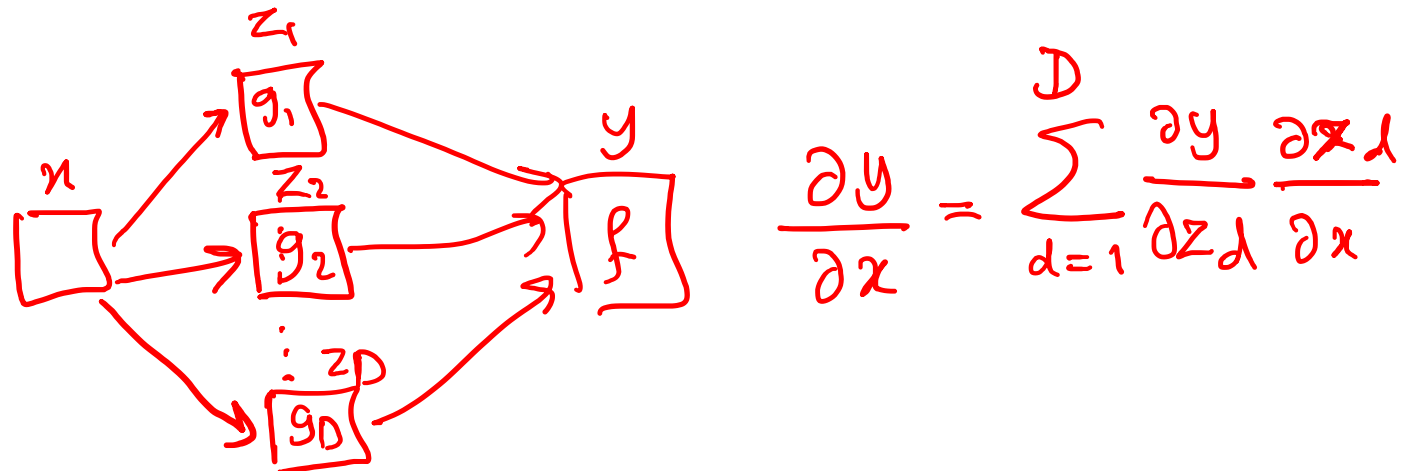
- If $y = f(z)$ and $z = g(x)$ then the corresponding computation graph is



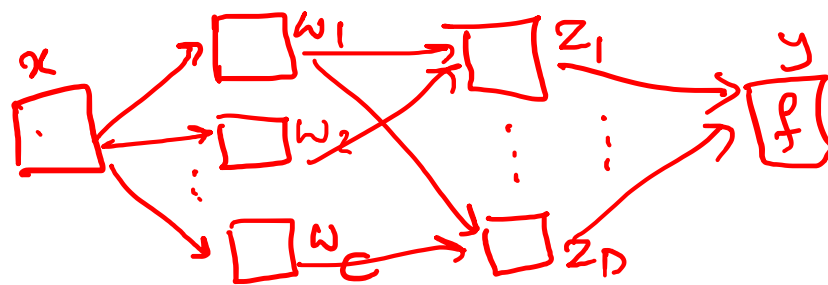
- If $y = f(z_1, z_2)$ and $z_1 = g_1(x), z_2 = g_2(x)$ then



- If $y = f(\mathbf{z})$ and $\mathbf{z} = g(x)$ then $\mathbf{z} \in \mathbb{R}^D$



Poll Question 1



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

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

$\sum_{c=1}^C \frac{\partial z_d}{\partial w_c} \frac{\partial w_c}{\partial x}$

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 $x = 2, z = 3$?

$$\frac{\partial y}{\partial x} = e^{xz} z + \frac{z}{\log x} + \frac{xz}{(\log x)^2} \left(-\frac{1}{x}\right) + \frac{\cos(\log x)}{x^2 z} + \frac{-\sin \log x}{x^2 z}$$

$$= 3e^6 + \frac{3}{\log 2} + \frac{3}{(\log 2)^2} + \frac{\cos(\log 2)}{12} + \frac{-\sin(\log 2)}{12}$$

$$\frac{\partial y}{\partial z} \Big|_{x=2, z=3} = 2e^6 + \frac{2}{\log(2)} - \frac{\sin(\log(2))}{18}$$

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)

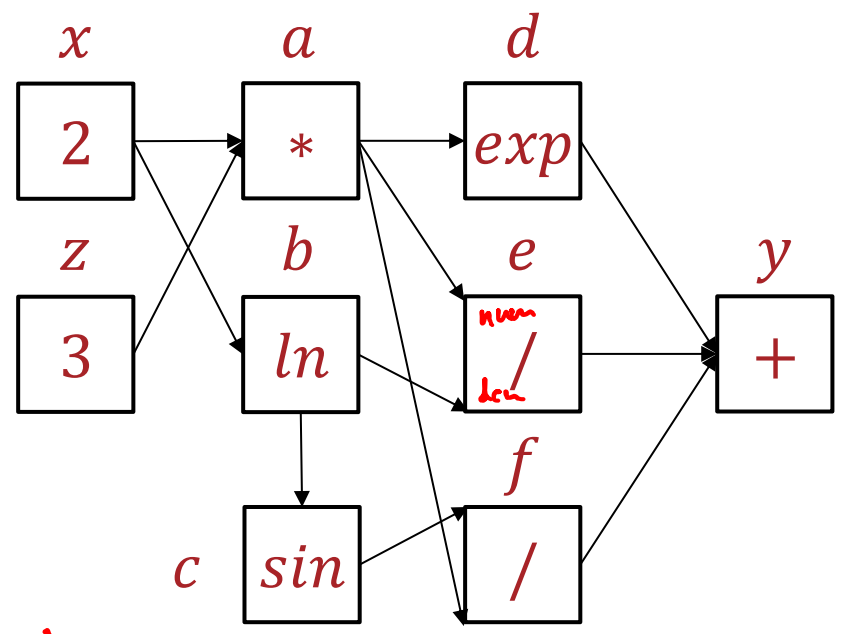
- Given

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

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

- $\rightarrow a = xz$
- $\rightarrow b = \ln(x)$
- $\rightarrow c = \sin(b)$
- $\rightarrow d = e^a$
- $\rightarrow e = a/b$
- $\rightarrow f = c/a$
- $y = d + e + f$



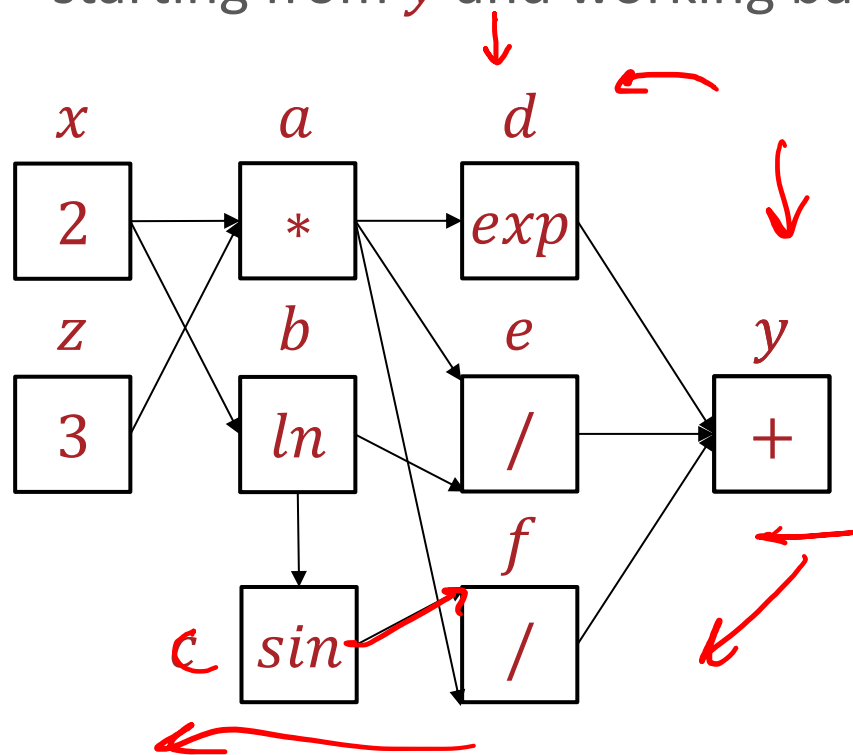
Approach 3: Automatic Differentiation (reverse mode)

- 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



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

$$g_z = \frac{\partial y}{\partial z} = 1 = \frac{\partial y}{\partial e} = \frac{\partial y}{\partial f}$$

$$g_c = \frac{\partial y}{\partial c} = g_f \frac{\partial f}{\partial c} = g_f \left(\frac{1}{a} \right)$$

g_b

g_a

g_x, g_z

Approach 3: Automatic Differentiation (reverse mode)

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.