

# 10-301/601: Introduction to Machine Learning

## Lecture 13 – Differentiation

Henry Chai

6/13/23

# Front Matter

- Announcements
  - PA3 released 6/8, due 6/15 at 11:59 PM
  - Quiz 4: Neural Networks on 6/20 (next Tuesday)
  - No lecture **or OH** on 6/19 (next Monday) for Juneteenth
  - Midterm on 6/23, one week from Friday
- Recommended Readings
  - None

# Midterm Logistics

- Time and place: *(tentative)*
  - Friday, 6/23 from 12:00PM to 3:00 PM in **DH 2302**
- Closed book/notes
  - 1-page cheatsheet allowed, both back and front; can be typeset or handwritten

# Midterm Coverage

- Lectures: 1 – 14 (through tomorrow's lecture)
  - Foundations: probability, linear algebra, calculus
  - Important concepts: inductive bias, overfitting, model selection/hyperparameter optimization, regularization
  - Models: decision trees, kNN, Perceptron, linear regression, logistic regression, neural networks
  - Methods: (stochastic) gradient descent, closed-form optimization, backpropagation, MLE/MAP

# Midterm Preparation

- Review midterm practice problems, posted to the course website (under [Recitations](#))
- Attend the exam review recitation on 6/20 (after the quiz)
- Review this semester's quizzes and study guides
- Consider whether you understand the “Key Takeaways” for each lecture / section
- Write your cheat sheet

# Recall: Forward Propagation for Making Predictions

- Input: weights  $W^{(1)}, \dots, W^{(L)}$  and a query data point  $\mathbf{x}$
- Initialize  $\mathbf{o}^{(0)} = [1, \mathbf{x}]^T$
- For  $l = 1, \dots, L$ 
  - $\mathbf{s}^{(l)} = W^{(l)} \mathbf{o}^{(l-1)}$
  - $\mathbf{o}^{(l)} = [1, \theta(\mathbf{s}^{(l)})]^T$
- Output:  $h_{W^{(1)}, \dots, W^{(L)}}(\mathbf{x}) = \mathbf{o}^{(L)}$

# Recall: Gradient Descent for Learning

- Input:  $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N, \eta^{(0)}$
- Initialize all weights  $W_{(0)}^{(1)}, \dots, W_{(0)}^{(L)}$  to small, random numbers and set  $t = 0$  (???)
- While TERMINATION CRITERION is not satisfied (???)
  - For  $l = 1, \dots, L$
  - Compute  $G^{(l)} = \nabla_{W^{(l)}} \ell_{\mathcal{D}} \left( W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)} \right)$  (???)
  - Update  $W^{(l)}$ :  $W_{(t+1)}^{(l)} = W_{(t)}^{(l)} - \eta_0 G^{(l)}$
  - Increment  $t$ :  $t = t + 1$
- Output:  $W_{(t)}^{(1)}, \dots, W_{(t)}^{(L)}$

# 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	$\frac{\partial y}{\partial \mathbf{X}}$	$\frac{\partial \mathbf{y}}{\partial \mathbf{X}}$	$\frac{\partial \mathbf{Y}}{\partial \mathbf{X}}$	

Table courtesy of Matt Gormley



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



# 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

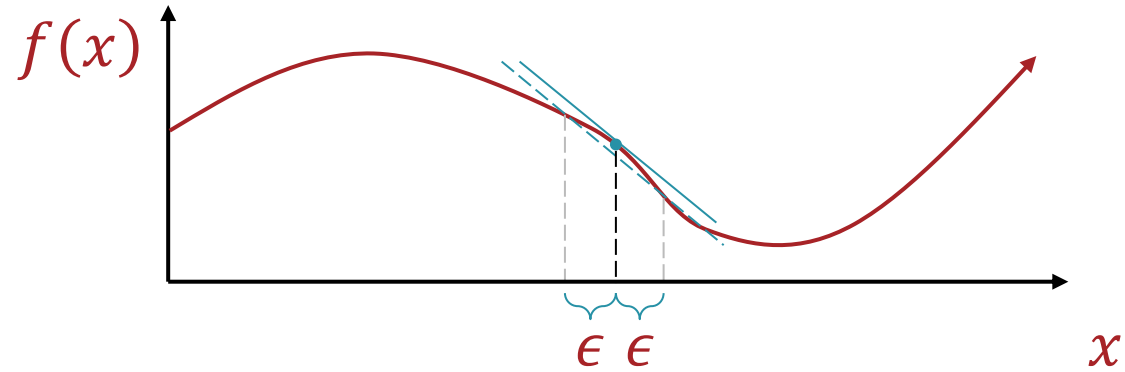
• Given  $f: \mathbb{R}^D \rightarrow \mathbb{R}$ , compute  $\nabla_{\mathbf{x}} f(\mathbf{x}) = \partial f(\mathbf{x}) / \partial \mathbf{x}$

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  $i^{\text{th}}$  position



- We want  $\epsilon$  to be small to get a good approximation but we run into floating point issues when  $\epsilon$  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

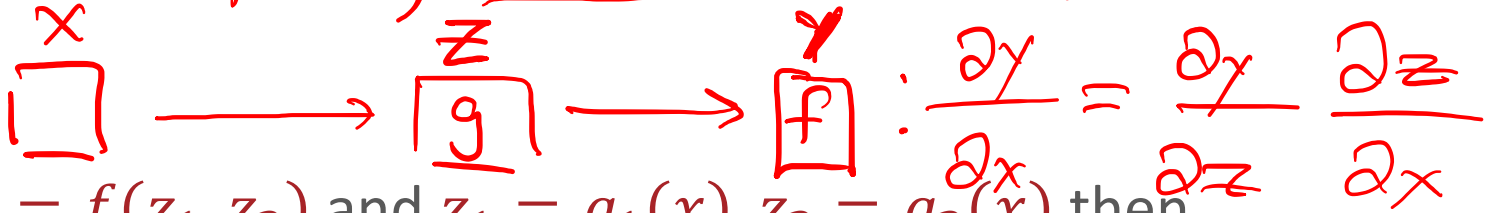
$$y = f(x, z) = e^{xz} + \frac{xz}{\ln(x)} + \frac{\sin(\ln(x))}{xz}$$

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

# The Chain Rule of Calculus

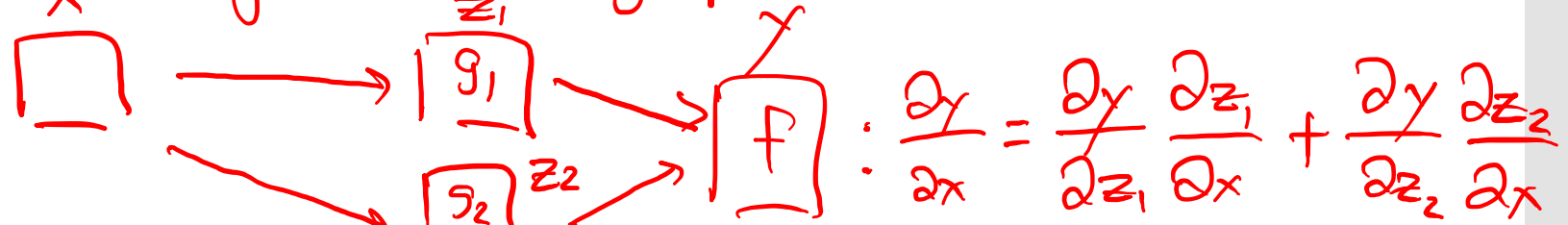
- If  $y = f(z)$  and  $z = g(x)$  then

the corresponding computational graph



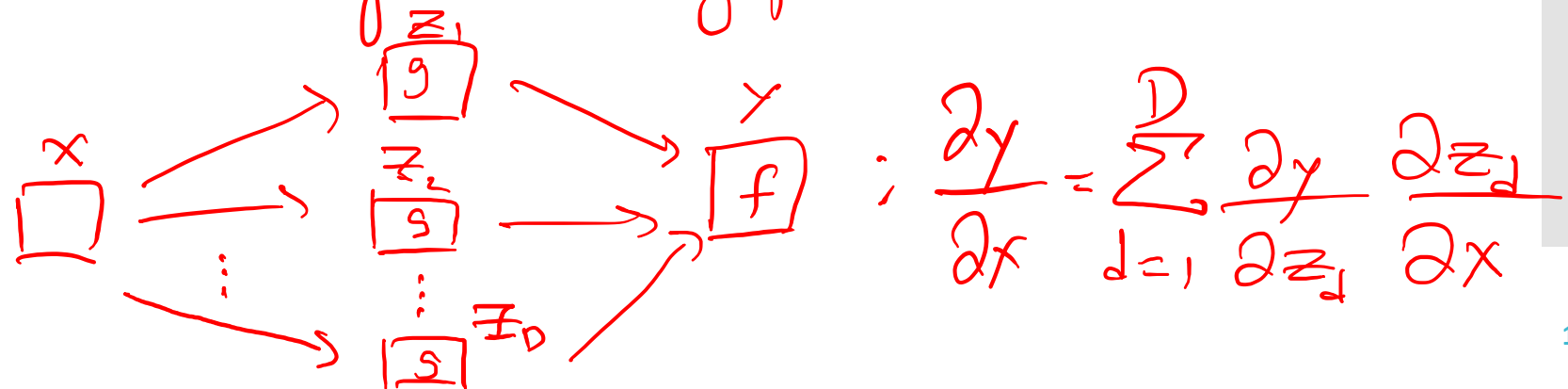
- If  $y = f(z_1, z_2)$  and  $z_1 = g_1(x), z_2 = g_2(x)$  then

computational graph:



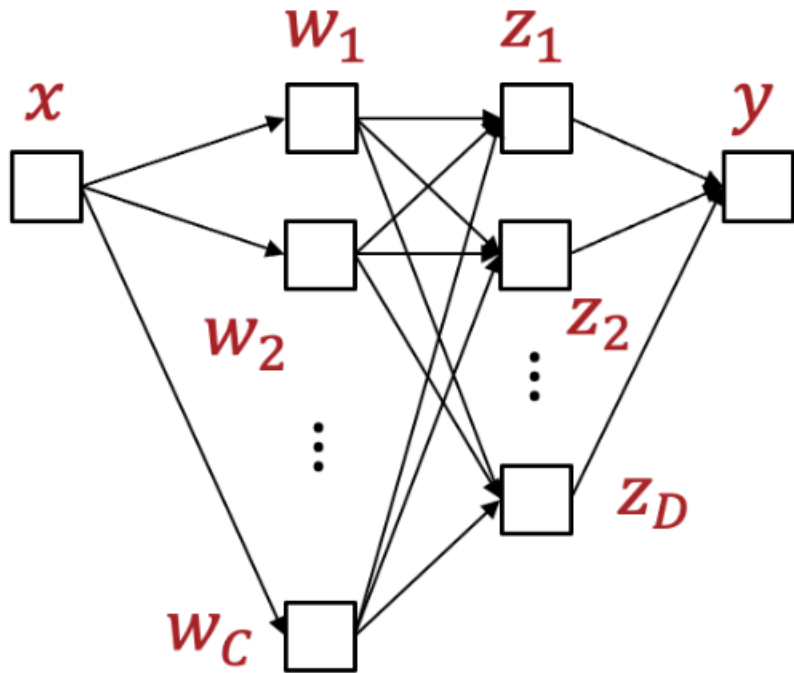
- If  $y = f(\mathbf{z})$  and  $\mathbf{z} = g(x)$  then ( $\mathbf{z} \in \mathbb{R}^D$ )

computational graph





Suppose  $y = f(\mathbf{z})$ ,  $\mathbf{z} = h(\mathbf{w})$  and  $\mathbf{w} = g(x)$ . 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?



Yes

No

Unsure

# 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  $x = 2, z = 3$ ?

$$\begin{aligned} \frac{\partial y}{\partial x} &= \frac{\partial}{\partial x} (e^{xz}) + \frac{\partial}{\partial x} \left( \frac{xz}{\ln(x)} \right) + \frac{\partial}{\partial x} \left( \frac{\sin(\ln(x))}{xz} \right) \\ &= ze^{xz} + \frac{z}{\ln(x)} - \frac{xz}{\ln(x)^2} \left( \frac{1}{x} \right) + \frac{\cos(\ln(x))}{x^2 z} - \frac{\sin(\ln(x))}{x^2 z} \\ &= 3e^6 + \frac{3}{\ln(2)} - \frac{3}{\ln(2)^2} + \frac{\cos(\ln(2))}{12} - \frac{\sin(\ln(2))}{12} \end{aligned}$$

$$\begin{aligned} \frac{\partial y}{\partial z} &= \frac{\partial}{\partial z} (e^{xz}) + \frac{\partial}{\partial z} \left( \frac{xz}{\ln(x)} \right) + \frac{\partial}{\partial z} \left( \frac{\sin(\ln(x))}{xz} \right) \\ &= xe^{xz} + \frac{x}{\ln(x)} - \frac{\sin(\ln(x))}{x^2 z^2} \\ &= 2e^6 + \frac{2}{\ln(2)} - \frac{\sin(\ln(2))}{18} \end{aligned}$$

# 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)

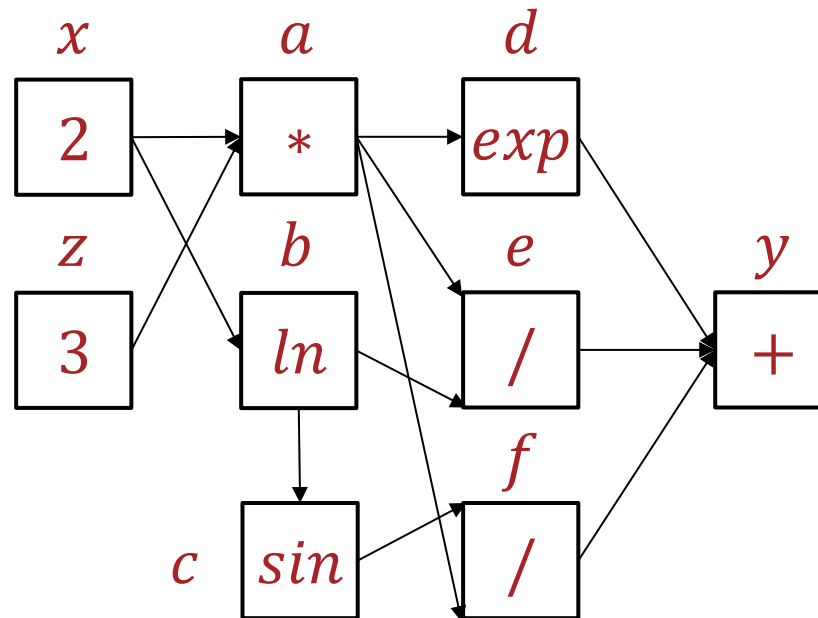
- Given

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

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

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

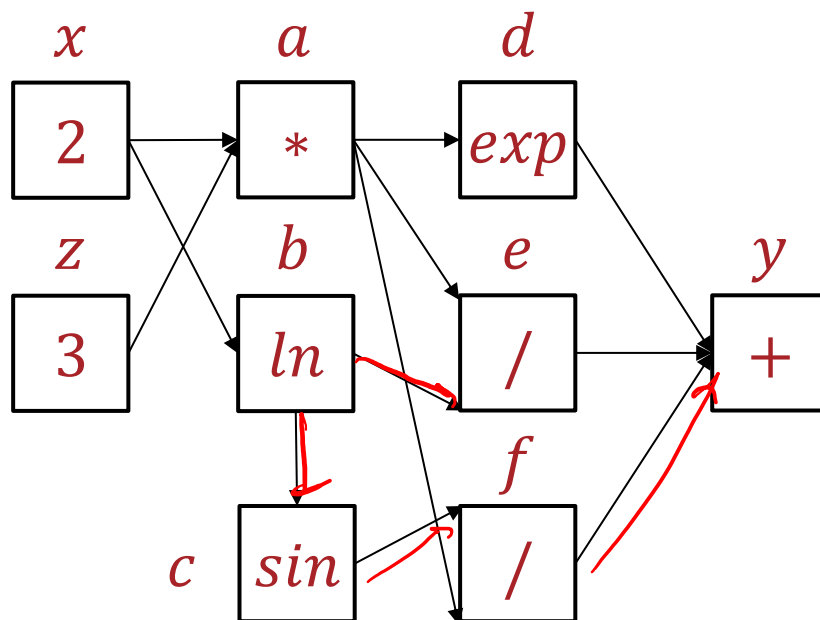


- Given

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

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



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

$$g_d = \frac{\partial y}{\partial d} = g_e = g_f = 1$$

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

$$g_b = \frac{\partial y}{\partial b} = g_e \frac{\partial e}{\partial b} + g_f \frac{\partial f}{\partial b}$$

$$= (1) \left( \frac{-a}{b^2} \right) + \left( \frac{1}{a} \right) \cos(b)$$

$$g_a = \frac{\partial y}{\partial a} = g_d \frac{\partial d}{\partial a} + g_e \frac{\partial e}{\partial a} + g_f \frac{\partial f}{\partial a}$$

$$= (1) \exp(a) + (1) \left( \frac{1}{b} \right) + (1) \left( \frac{-c}{a^2} \right)$$

## 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  $f(\mathbf{x})$

# 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



# Key Takeaways

- Denominator layout for matrix calculus
- Finite difference method is a simple but computationally expensive approximation technique
  - ***You should use this to unit test your implementation of backpropagation!***
- Symbolic differentiation is the “default” differentiation method but can also be computationally expensive
- Automatic differentiation (reverse mode) saves intermediate quantities for computational efficiency
  - Backpropagation is an instance of automatic differentiation in the reverse mode