

10-301/10-601 Introduction to Machine Learning

Machine Learning Department School of Computer Science Carnegie Mellon University

Neural Networks

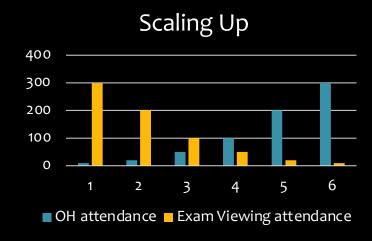
+

Backpropagation

Matt Gormley & Henry Chai Lecture 12 Feb. 19, 2025

Reminders

- Post-Exam Followup:
 - Exam Viewing
 - Exit Poll: Exam 1
 - Grade Summary 1
- Homework 4: Logistic Regression
 - Out: Mon, Feb 17
 - Due: Wed, Feb 26 at 11:59pm



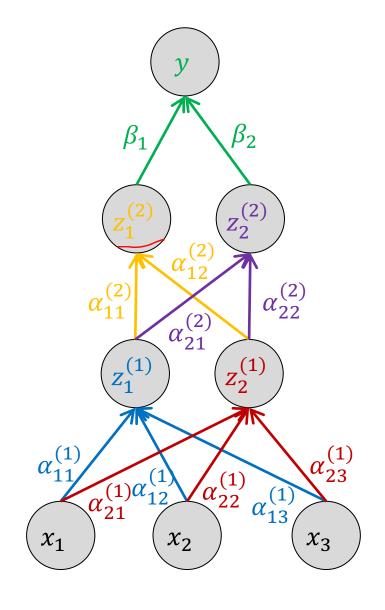
Neural Network Architectures

Even for a basic Neural Network, there are many design decisions to make:

- # of hidden layers (depth)
- 2. # of units per hidden layer (width)
- 3. Type of activation function (nonlinearity)
- 4. Form of objective function
- 5. How to initialize the parameters

Recall

Neural Network



Example: Neural Network with 2 Hidden Layers and 2 Hidden Units

$$z_{1}^{(1)} = \underline{\sigma}(\alpha_{11}^{(1)}x_{1} + \alpha_{12}^{(1)}x_{2} + \alpha_{13}^{(1)}x_{3} + \underline{\alpha_{10}^{(1)}})$$

$$z_{2}^{(1)} = \underline{\sigma}(\alpha_{21}^{(1)}x_{1} + \alpha_{22}^{(1)}x_{2} + \alpha_{23}^{(1)}x_{3} + \underline{\alpha_{20}^{(1)}})$$

$$z_{1}^{(2)} = \underline{\sigma}(\alpha_{21}^{(2)}z_{1}^{(1)} + \alpha_{12}^{(2)}z_{2}^{(1)} + \alpha_{10}^{(2)})$$

$$z_{2}^{(2)} = \underline{\sigma}(\alpha_{21}^{(2)}z_{1}^{(1)} + \alpha_{22}^{(2)}z_{2}^{(1)} + \alpha_{20}^{(2)})$$

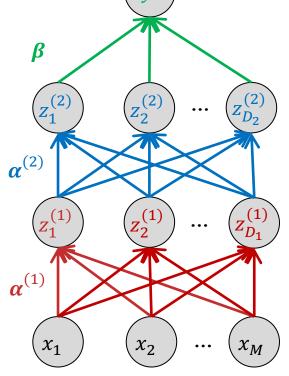
$$y = \underline{\sigma}(\beta_{1} z_{1}^{(2)} + \beta_{2} z_{2}^{(2)} + \beta_{0})$$

Recall

Neural Network (Matrix Form)

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

Example: Arbitrary Feed-forward Neural Network



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

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

$$\boldsymbol{\gamma} = \sigma((\boldsymbol{\beta})^T \boldsymbol{z}^{(2)} + \beta_0)$$

$$\boldsymbol{\alpha}^{(2)} \in \mathbb{R}^{D_1 \times D_2}$$

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

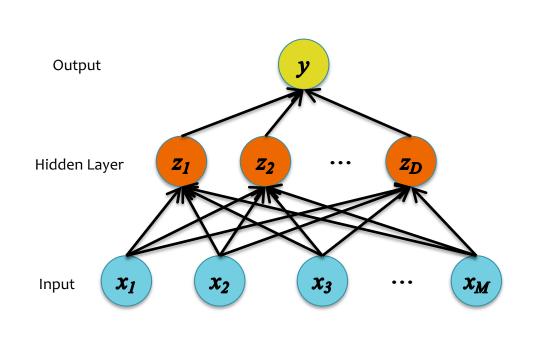
$$\boldsymbol{b}^{(2)} \in \mathbb{R}^{D_2}$$

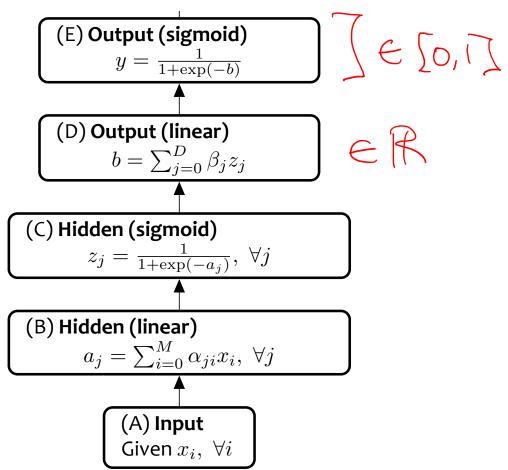
$$\boldsymbol{z}^{(1)} = \sigma((\boldsymbol{\alpha}^{(1)})^T \boldsymbol{x} + \boldsymbol{b}^{(1)})$$

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

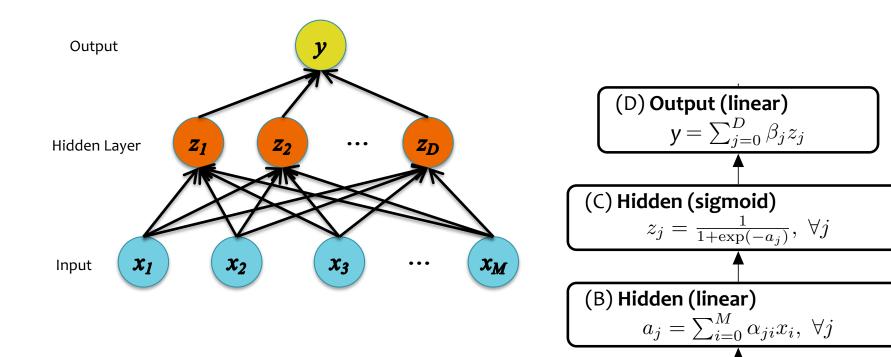
LOSS FUNCTIONS & OUTPUT LAYERS

Neural Network for Classification





Neural Network for Regression



(A) Input

Given $x_i, \ \forall i$

Objective Functions for NNs

Quadratic Loss:

- the same objective as Linear Regression
- i.e. mean squared error

$$J = \ell_Q(y, y^{(i)}) = \frac{1}{2}(y - y^{(i)})^2$$

$$\frac{dJ}{dy} = y - y^{(i)}$$
 Predicted

Binary Cross-Entropy:

- the same objective as Binary Logistic Regression
- i.e. negative log likelihood
 This requires our output y to be a probability in [0,1], y (3) ∈ 50,13

$$J = \ell_{CE}(y, y^{(i)}) = -(y^{(i)} \log(y) + (1 - y^{(i)}) \log(1 - y))$$

$$\frac{dJ}{dy} = -\left(y^{(i)} \frac{1}{y} + (1 - y^{(i)}) \frac{1}{y - 1}\right)$$

Cross-entropy vs. Quadratic loss

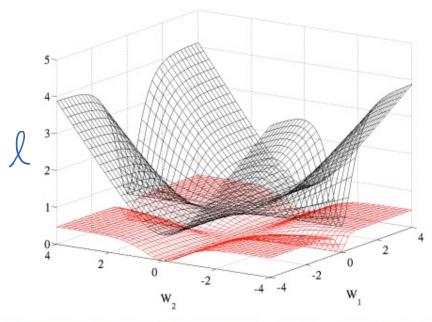
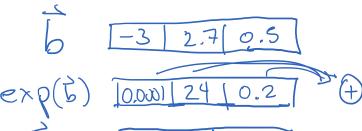
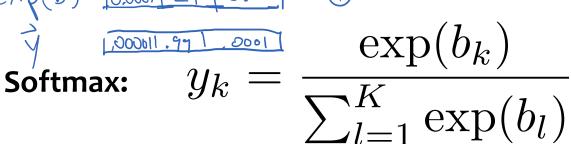
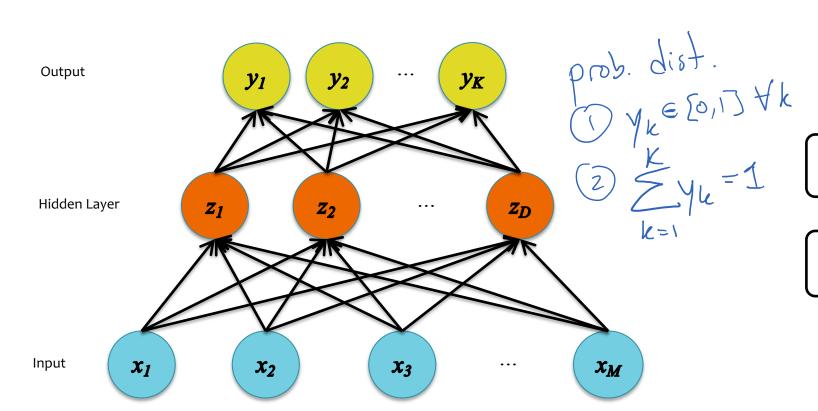


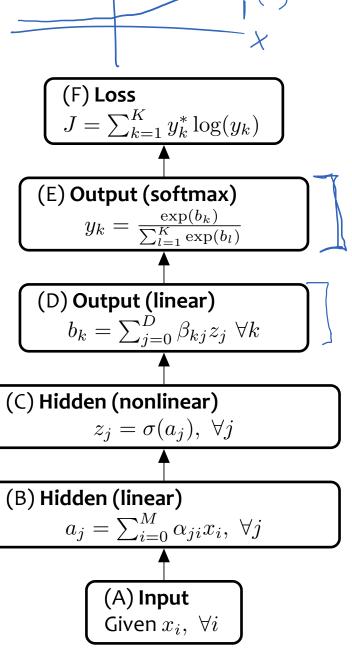
Figure 5: Cross entropy (black, surface on top) and quadratic (red, bottom surface) cost as a function of two weights (one at each layer) of a network with two layers, W_1 respectively on the first layer and W_2 on the second, output layer.



Multiclass Output







Objective Functions for NNs

- Cross-Entropy for Multiclass Outputs:
 - i.e. negative log likelihood for multiclass outputs
 - Suppose output is a random variable Y that takes one of K values
 - Let $\mathbf{y}^{(i)}$ represent our true label as a one-hot vector:

$$\mathbf{y}^{(i)} = \begin{bmatrix} 0 & 0 & 0 & 1 & 0 & 0 & \dots & 0 \\ & 1 & 2 & 3 & 4 & 5 & 6 & \dots & K \end{bmatrix}$$

Assume our model outputs a length K vector of probabilities:

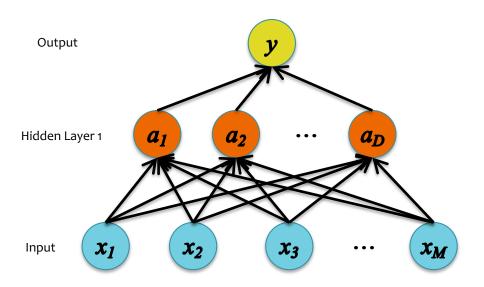
$$y = softmax(f_{scores}(x, \theta))$$

Then we can write the log-likelihood of a single training example $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$ as:

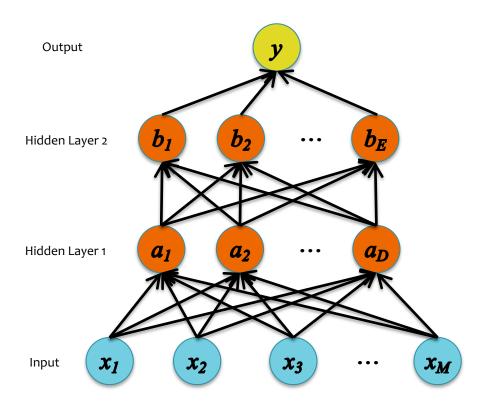
$$J = \ell_{CE}(\mathbf{y}, \mathbf{y}^{(i)}) = -\sum_{k=1}^{K} y_k^{(i)} \log(y_k)$$

$$= -\left(0 | \circ_{\mathcal{S}}(\gamma_1) + 0 | \circ_{\mathcal{S}}(\gamma_2) + 0 | \circ_{\mathcal{S}}(\gamma_3) + 1 | \circ_{\mathcal{S}}(\gamma_4) + \cdots | \circ_{\mathcal{S}(\gamma_4) + \cdots | \circ_{\mathcal{S}}(\gamma_4) + \cdots | \circ_{\mathcal{S}}(\gamma_4) + \cdots | \circ_{\mathcal{S}}(\gamma_4)$$

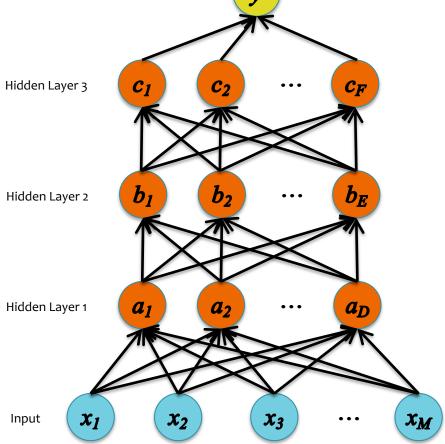
Q: How many layers should we use?



Q: How many layers should we use?



Q: How many layers should we use?



Q: How many layers should we use?

Theoretical answer:

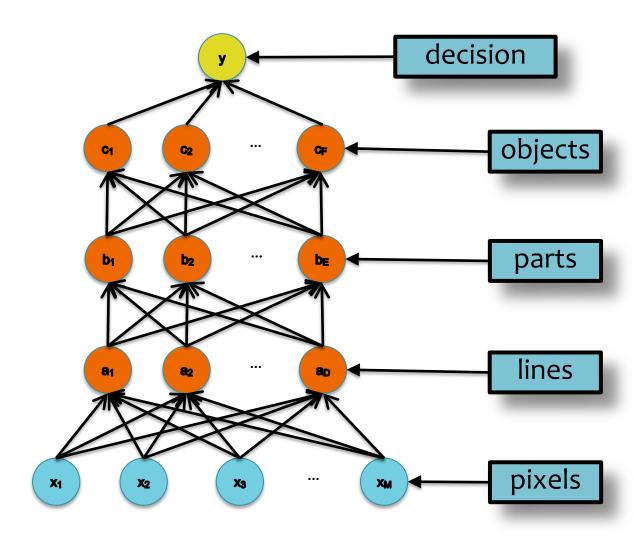
- A neural network with 1 hidden layer is a universal function approximator
- Cybenko (1989): For any continuous function g(x), there exists a 1-hidden-layer neural net $h_{\theta}(x)$ s.t. $|h_{\theta}(x) g(x)| < \epsilon$ for all x, assuming sigmoid activation functions

Empirical answer:

- Before 2006: "Deep networks (e.g. 3 or more hidden layers) are too hard to train"
- After 2006: "Deep networks are easier to train than shallow networks (e.g. 2 or fewer layers) for many problems"

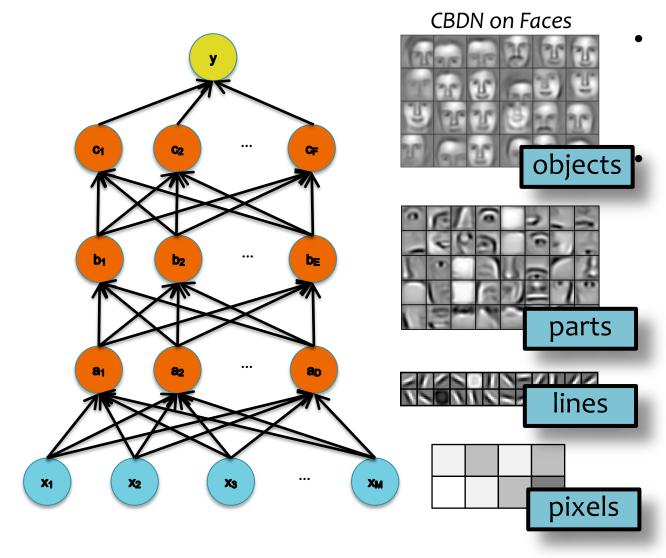
Big caveat: You need to know and use the right tricks.

Feature Learning



- Traditional feature engineering: build up levels of abstraction by hand
- Deep networks (e.g. convolution networks): learn the increasingly higher levels of abstraction from data
 - each layer is a learned feature representation
 - sophistication increases in higher layers

Feature Learning



Traditional feature engineering: build up levels of abstraction by hand

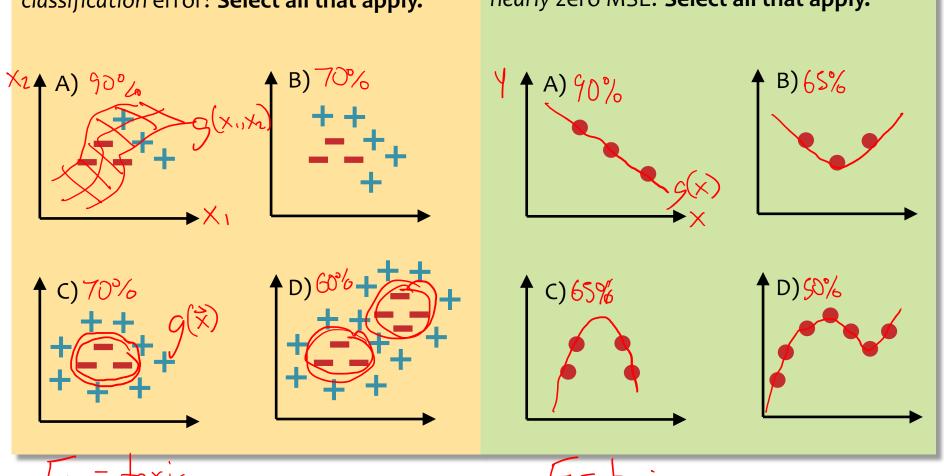
Deep networks (e.g. convolution networks): learn the increasingly higher levels of abstraction from data

- each layer is a learned feature representation
- sophistication increases in higher layers

Neural Network Errors

Poll Question 1: For which of the datasets below does there exist a one-hidden layer neural network that achieves zero classification error? **Select all that apply.**

Poll Question 2: For which of the datasets below does there exist a one-hidden layer neural network for *regression* that achieves *nearly* zero MSE? **Select all that apply.**



Neural Networks Objectives

You should be able to...

- Explain the biological motivations for a neural network
- Combine simpler models (e.g. linear regression, binary logistic regression, multinomial logistic regression) as components to build up feed-forward neural network architectures
- Explain the reasons why a neural network can model nonlinear decision boundaries for classification
- Compare and contrast feature engineering with learning features
- Identify (some of) the options available when designing the architecture of a neural network
- Implement a feed-forward neural network

Computing Gradients

APPROACHES TO DIFFERENTIATION

Background

A Recipe for Machine Learning

1. Given training data:

$$\{oldsymbol{x}_i, oldsymbol{y}_i\}_{i=1}^N$$

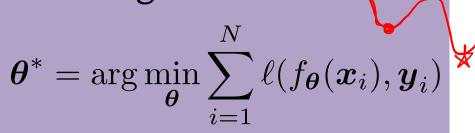
- 2. Choose each of these:
 - Decision function

$$\hat{\boldsymbol{y}} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

Loss function

$$\ell(\hat{m{y}},m{y}_i)\in\mathbb{R}$$

3. Define goal:



4. Train with SGD:

(take small steps opposite the gradient)

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

Background

A Recipe for

Gradients

1. Given training dat

$$\{oldsymbol{x}_i, oldsymbol{y}_i\}_{i=1}^N$$
 gradient! And it's a

- 2. Choose each of tl
 - Decision function

$$\hat{\boldsymbol{y}} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

Loss function

$$\ell(\hat{m{y}}, m{y}_i) \in \mathbb{R}$$

Backpropagation can compute this gradient!

And it's a special case of a more general algorithm called reversemode automatic differentiation that can compute the gradient of any differentiable function efficiently!

opposite the gradient)
$$\boldsymbol{\theta}^{(t)} = \boldsymbol{\eta}_t \nabla \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

Approaches to Differentiation

Question 1:

When can we compute the gradients for an arbitrary neural network?

Question 2:

When can we make the gradient computation efficient?

Given
$$f: \mathbb{R}^A \to \mathbb{R}^B, f(\mathbf{x})$$
Compute $\frac{\partial f(\mathbf{x})_i}{\partial x_j} \forall i, j$

Approaches to Differentiation



1. Finite Difference Method

- Pro: Great for testing implementations of backpropagation
- Con: Slow for high dimensional inputs / outputs
- Required: Ability to call the function f(x) on any input x

2. Symbolic Differentiation

- Note: The method you learned in highschool
- Note: Used by Mathematica / Wolfram Alpha / Maple
- Pro: Yields easily interpretable derivatives
- Con: Leads to exponential computation time if not carefully implemented
- Required: Mathematical expression that defines f(x)

Given
$$f: \mathbb{R}^A \to \mathbb{R}^B, f(\mathbf{x})$$
Compute $\frac{\partial f(\mathbf{x})_i}{\partial x_j} \forall i, j$

Approaches to Differentiation

3. Automatic Differentiation – Reverse Mode

- Note: Called Backpropagation when applied to Neural Nets
- Pro: Computes partial derivatives of one output f(x)_i with respect to all inputs x_j in time polynomial in the computation time of f(x)
- Con: Slow for high dimensional outputs (e.g. vector-valued functions)
- Required: Algorithm for computing f(x)

4. Automatic Differentiation - Forward Mode

- Note: Easy to implement. Uses dual numbers.
- Pro: Computes partial derivatives of all outputs f(x)_i with respect to one input x_j in time polynomial in the computation time of f(x)
- Con: Slow for high dimensional inputs (e.g. vector-valued x)
- Required: Algorithm for computing f(x)

THE FINITE DIFFERENCE METHOD

Finite Difference Method

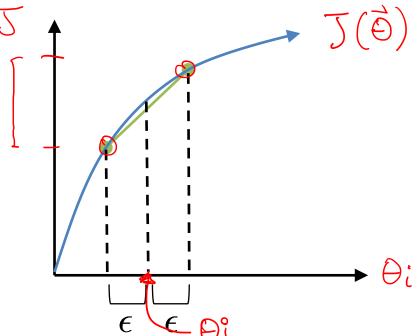
The centered finite difference approximation is:

$$\frac{\partial}{\partial \theta_i} J(\boldsymbol{\theta}) pprox \frac{(J(\boldsymbol{\theta} + \epsilon \cdot \boldsymbol{d}_i) - J(\boldsymbol{\theta} - \epsilon \cdot \boldsymbol{d}_i))}{2\epsilon}$$
 (1)

where d_i is a 1-hot vector consisting of all zeros except for the ith entry of d_i , which has value 1.

Notes:

- Suffers from issues of floating point precision, in practice
- Typically only appropriate to use on small examples with an appropriately chosen epsilon



Differentiation Quiz

Poll Question 3: Differentiation Quiz #1

Suppose x = 2 and z = 3, what are dy/dx and dy/dz for the function below? Round your answer to the nearest integer.

$$y = \exp(xz) + \frac{xz}{\log(x)} + \frac{\sin(\log(x))}{xz}$$

Answer: Answers below are in the form [dy/dx, dy/dz]

Differentiation Quiz

Differentiation Quiz #2:

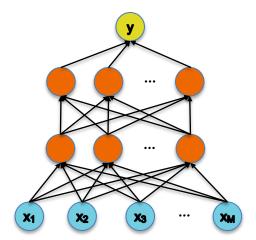
A neural network with 2 hidden layers can be written as:

$$y = \sigma(\boldsymbol{\beta}^T \sigma((\boldsymbol{\alpha}^{(2)})^T \sigma((\boldsymbol{\alpha}^{(1)})^T \mathbf{x}))$$

where $y \in \mathbb{R}$, $\mathbf{x} \in \mathbb{R}^{D^{(0)}}$, $\boldsymbol{\beta} \in \mathbb{R}^{D^{(2)}}$ and $\boldsymbol{\alpha}^{(i)}$ is a $D^{(i)} \times D^{(i-1)}$ matrix. Nonlinear functions are applied elementwise:

$$\sigma(\mathbf{a}) = [\sigma(a_1), \dots, \sigma(a_K)]^T$$

Let σ be sigmoid: $\sigma(a)=\frac{1}{1+exp-a}$ What is $\frac{\partial y}{\partial \beta_j}$ and $\frac{\partial y}{\partial \alpha_j^{(i)}}$ for all i,j.



THE CHAIN RULE OF CALCULUS

Chain Rule

Definition 1:

$$y = f(\mathbf{u})$$

$$v = g(x)$$

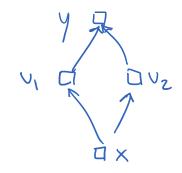
$$\frac{9x}{24} = \frac{90}{90} \frac{9x}{90}$$

Definition 2:

$$y = f(v_1, v_2)$$

$$v_2 = g_2(x)$$

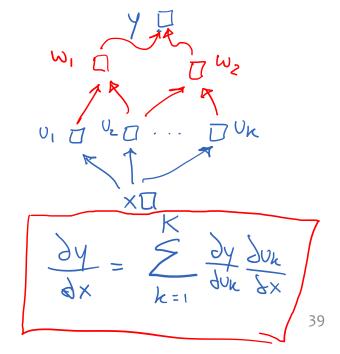
$$v_1 = g_1(x)$$



$$\frac{\partial x}{\partial y} = \frac{\partial v}{\partial y} \frac{\partial x}{\partial v_1} + \frac{\partial v}{\partial y} \frac{\partial x}{\partial v_2}$$

Definition 3:

$$y = f(\vec{v})$$
 $y \in \mathbb{R}$
 $\vec{v} = g(x)$ $\vec{v} \in \mathbb{R}^k$
 $x \in \mathbb{R}$

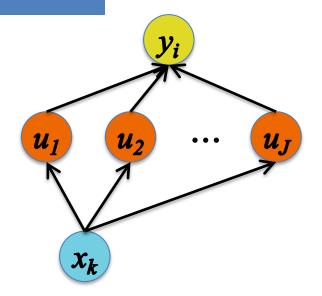


Chain Rule

Given: y = g(u) and u = h(x).

Chain Rule:

$$\frac{dy_i}{dx_k} = \sum_{j=1}^{J} \frac{dy_i}{du_j} \frac{du_j}{dx_k}, \quad \forall i, k$$



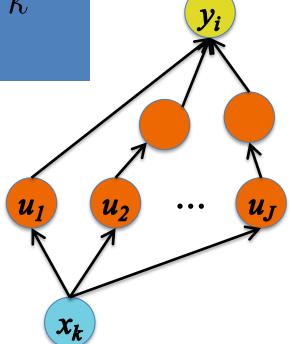
Chain Rule

Given: y = g(u) and u = h(x).

Chain Rule:

$$\frac{dy_i}{dx_k} = \sum_{j=1}^{J} \frac{dy_i}{du_j} \frac{du_j}{dx_k}, \quad \forall i, k$$

Backpropagation is just repeated application of the chain rule from Calculus 101.



Algorithm

FORWARD COMPUTATION FOR A COMPUTATION GRAPH

Backpropagation

Whiteboard

- From equation to forward computation
- Representing a simple function as a computation graph

Differentiation Quiz #1:

Suppose x = 2 and z = 3, what are dy/dx and dy/dz for the function below? Round your answer to the nearest integer.

$$y = \exp(xz) + \frac{xz}{\log(x)} + \frac{\sin(\log(x))}{xz}$$

Backpropagation

Differentiation Quiz #1:

Suppose x = 2 and z = 3, what are dy/dx and dy/dz for the function below? Round your answer to the nearest integer.

$$y = \exp(xz) + \frac{xz}{\log(x)} + \frac{\sin(\log(x))}{xz}$$

Now let's solve this in a different way!

Given:
$$y = \exp(xz) + \frac{xz}{\log(x)} + \frac{\sin(\log(x))}{xz}$$

Forward Computation:

Computation Graph:

Backward Computation

Given
$$x = 2$$
, $z = 3$

$$Q = XZ$$

$$b = log(x)$$

$$C = sin(b)$$

$$d = exp(a)$$

$$e = q/b$$

$$f = c/a$$

$$y = d + e + f$$

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

$$g_{\beta} = \frac{\partial y}{\partial y} = 1, \quad g_{\theta} = \frac{\partial y}{\partial z} = 1$$

$$g_{c} = \frac{\partial y}{\partial z} = \frac{\partial y}{\partial z} = (g_{\beta})(\frac{1}{a}) \quad -\frac{a}{b^{2}}$$

$$g_{b} = \frac{\partial y}{\partial z} = \frac{\partial y}{\partial z} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{e})(\frac{1}{a})$$

$$g_{a} = \frac{\partial y}{\partial z} = \frac{\partial z}{\partial z} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{e})(\frac{1}{a})$$

$$g_{a} = \frac{\partial y}{\partial z} = \frac{\partial z}{\partial z} + \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{e})(\frac{1}{a})$$

$$g_{a} = \frac{\partial z}{\partial z} = \frac{\partial z}{\partial z} = (g_{a})(z) + (g_{b})(\frac{1}{a})$$

$$g_{b} = \frac{\partial z}{\partial z} = (g_{a})(z) + (g_{b})(\frac{1}{a})$$

$$g_{c} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{c})(\frac{1}{a})$$

$$g_{c} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{c})(\frac{1}{a})$$

$$g_{c} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{c})(\frac{1}{a})$$

$$g_{c} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{c})(\frac{1}{a})$$

$$g_{c} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{c})(\frac{1}{a})$$

$$g_{c} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{c})(\frac{1}{a})$$

$$g_{c} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{c})(\frac{1}{a})$$

$$g_{c} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{c})(\frac{1}{a})$$

$$g_{c} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{c})(\frac{1}{a})$$

$$g_{c} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{c})(\frac{1}{a})$$

$$g_{c} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{c})(\frac{1}{a})$$

$$g_{c} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{c})(\frac{1}{a})$$

$$g_{c} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{c})(\frac{1}{a})$$

$$g_{c} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{c})(\frac{1}{a})$$

$$g_{c} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{c})(\frac{1}{a})$$

$$g_{c} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{c})(\cos(b))$$

$$g_{c} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{c})(\cos(b))$$

$$g_{c} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{c})(\cos(b))$$

$$g_{c} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{c})(\cos(b))$$

$$g_{c} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{c})(\cos(b))$$

$$g_{c} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{c})(\cos(b))$$

$$g_{c} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{c})(\cos(b))$$

$$g_{c} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b)) + (g_{c})(\cos(b))$$

$$g_{c} = \frac{\partial z}{\partial z} = (g_{c})(\cos(b))$$

$$g_{$$

Backpropagation

Updates for Backpropagation:

$$g_x = \frac{\partial y}{\partial x} = \sum_{k=1}^K \frac{\partial y}{\partial u_k} \frac{\partial u_k}{\partial x}$$
$$= \sum_{k=1}^K g_{u_k} \frac{\partial u_k}{\partial x}$$

Backprop is efficient b/c of reuse in the forward pass and the backward pass.