

# Lecture 2: Deep Learning - Introduction and Applications

Isaiah Hull<sup>1,2</sup>

<sup>1</sup>BI Norwegian Business School

<sup>2</sup>CogniFrame

October 31, 2023



# Introduction

## Lecture 2: Overview

1. Introduction to Deep Learning.
2. Deep Learning Applications.

# 1. Introduction to Deep Learning

# Introduction to Deep Learning

## Overview

- ▶ Based on Ng and Ma (2023).
  - ▶ [https://cs229.stanford.edu/main\\_notes.pdf](https://cs229.stanford.edu/main_notes.pdf)
- ▶ Overview of neural networks.
- ▶ Vectorization.
- ▶ Backpropagation.

# Introduction to Deep Learning

## Supervised Learning and Model Types

- ▶ Supervised learning: predicting  $y$  from input  $x$ .
- ▶ Model:  $h_{\theta}(x)$ .
- ▶ Familiar examples:
  - ▶ Linear regression:  $h_{\theta}(x) = \theta^{\top} x$ .
  - ▶ With feature map:  $h_{\theta}(x) = \theta^{\top} \phi(x)$ .

# Introduction to Deep Learning

## Supervised Learning and Model Types

- ▶ Common feature is linearity in parameters  $\theta$ .
- ▶ Eventually consider non-linear models in  $\theta$  and  $x$ .
- ▶ Will focus on neural networks specifically.
- ▶ Initially, consider  $h_{\theta}(x)$  as an abstract non-linear model.

# Introduction to Deep Learning

## Training and Regression with Non-Linear Models

- ▶ Terminology differs from econometrics.
- ▶ Training examples (observations):  $\{ (x^{(i)}, y^{(i)}) \}_{i=1}^n$ .
- ▶ Define a nonlinear model and an associated loss or cost function.
- ▶ Regression problem: continuous target (dependent variable).

# Introduction to Deep Learning

## Training and Regression with Non-Linear Models

- ▶ The output is a real number ( $y^{(i)} \in \mathbb{R}$ ).
- ▶ Model also outputs a real number:  $h_{\theta}(x) \in \mathbb{R}$ .
- ▶ Loss function example: least squares.



# Introduction to Deep Learning

## Regression and Mean-Square Cost Function

- ▶  $i$ -th example (observation) cost:

$$J^{(i)}(\theta) = \frac{1}{2} \left( h_{\theta} \left( x^{(i)} \right) - y^{(i)} \right)^2$$

- ▶ Mean-square cost for entire dataset (sample):

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n J^{(i)}(\theta)$$

# Introduction to Deep Learning

## Regression and Mean-Square Cost Function

- ▶ Similar to OLS but:
  - ▶ Average, rather than sum of squared residuals.
  - ▶ Doesn't change local/global minima.
  - ▶ Different parameterization for  $h_{\theta}(x)$ .
- ▶ "Loss" and "cost" used interchangeably.

# Introduction to Deep Learning

## Binary Classification and Loss Function

- ▶ Inputs:  $x \in \mathbb{R}^d$ .
- ▶ Parameterized model:  $\bar{h}_\theta : \mathbb{R}^d \rightarrow \mathbb{R}$ .
- ▶ Output (logit):  $\bar{h}_\theta(x) \in \mathbb{R}$ .
- ▶ Logistic function  $g(\cdot)$  to convert logit:

$$h_\theta(x) = g(\bar{h}_\theta(x)) = \frac{1}{1 + \exp(-\bar{h}_\theta(x))}$$

# Introduction to Deep Learning

## Binary Classification and Loss Function

- ▶ Conditional distribution of  $y$ :

$$P(y = 1 \mid x; \theta) = h_{\theta}(x)$$

$$P(y = 0 \mid x; \theta) = 1 - h_{\theta}(x)$$

- ▶ Negative likelihood loss function:

$$J^{(i)}(\theta) = -\log p\left(y^{(i)} \mid x^{(i)}; \theta\right) = \ell_{\text{logistic}}\left(\bar{h}_{\theta}\left(x^{(i)}\right), y^{(i)}\right)$$

# Introduction to Deep Learning

## Loss Functions

- ▶ Total loss:

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n J^{(i)}(\theta)$$

- ▶ Multi-class classification:

- ▶ Response variable  $y$  can be  $k$  values:  $y \in \{1, 2, \dots, k\}$ .
- ▶ Parameterized model:  $\bar{h}_{\theta} : \mathbb{R}^d \rightarrow \mathbb{R}^k$ .

# Introduction to Deep Learning

## Multi-Class Classification

- ▶ Outputs (logits):  $\bar{h}_\theta(x) \in \mathbb{R}^k$ .
- ▶ Use softmax for probability vector:

$$P(y = j \mid x; \theta) = \frac{\exp(\bar{h}_\theta(x)_j)}{\sum_{s=1}^k \exp(\bar{h}_\theta(x)_s)}$$

# Introduction to Deep Learning

## Multi-Class Classification

- ▶  $i$ -th example loss (negative log-likelihood):

$$J^{(i)}(\theta) = \ell_{\text{ce}} \left( \bar{h}_{\theta} \left( x^{(i)} \right), y^{(i)} \right)$$

- ▶ Average of loss for training examples:

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n J^{(i)}(\theta)$$

- ▶ Exponential family and non-linear parameterization.

# Introduction to Deep Learning

## Optimizers and SGD

- ▶ Common Optimizers: gradient descent (GD), stochastic gradient descent (SGD), and variants.

- ▶ Gradient descent update rule:

$$\theta := \theta - \alpha \nabla_{\theta} J(\theta)$$

- ▶ Learning Rate ( $\alpha$ ): controls step size.
- ▶ Modification of gradient descent: stochastic gradient descent.



# Introduction to Deep Learning

---

**Algorithm 1** Stochastic Gradient Descent

---

- 1: Hyperparameter: learning rate  $\alpha$ , number of total iteration  $n_{\text{iter}}$ .
- 2: Initialize  $\theta$  randomly.
- 3: **for**  $i = 1$  to  $n_{\text{iter}}$  **do**
- 4:     Sample  $j$  uniformly from  $\{1, \dots, n\}$ , and update  $\theta$  by

$$\theta := \theta - \alpha \nabla_{\theta} J^{(j)}(\theta) \tag{7.9}$$

---

Source: Ng and Ma (2023).

# Introduction to Deep Learning

## Hardware Parallelization and Mini-batch SGD

- ▶ Computing gradient of  $B$  (batch size) examples simultaneously often faster due to hardware parallelization.
- ▶ Mini-batch version of SGD prevalent in deep learning.
- ▶ Variants of SGD and mini-batch SGD exist with different sampling schemes.

# Introduction to Deep Learning

---

**Algorithm 2** Mini-batch Stochastic Gradient Descent

---

- 1: Hyperparameters: learning rate  $\alpha$ , batch size  $B$ , # iterations  $n_{\text{iter}}$ .
- 2: Initialize  $\theta$  randomly
- 3: **for**  $i = 1$  to  $n_{\text{iter}}$  **do**
- 4:     Sample  $B$  examples  $j_1, \dots, j_B$  (without replacement) uniformly from  $\{1, \dots, n\}$ , and update  $\theta$  by

$$\theta := \theta - \frac{\alpha}{B} \sum_{k=1}^B \nabla_{\theta} J^{(j_k)}(\theta) \quad (7.10)$$

---

Source: Ng and Ma (2023).

# Introduction to Deep Learning

## Steps in Training a Deep Learning Model

1. Define a neural network parametrization  $h_{\theta}(x)$ .
2. Implement the backpropagation algorithm for efficient gradient computation of  $J^{(j)}(\theta)$ .
3. Execute SGD, mini-batch SGD, or other gradient-based optimizers with  $J(\theta)$ .

# Introduction to Deep Learning

## Neural Networks: Introduction

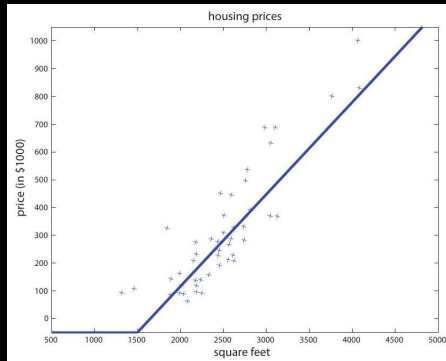
- ▶ Broad type of non-linear models:  $\bar{h}_\theta(x)$ .
- ▶ Combines matrix multiplications with element-wise non-linear operations.
- ▶ Regression:  $h_\theta(x) = \bar{h}_\theta(x)$ .
- ▶ Classification:
  - ▶ Binary:  $h_\theta(x) = 1 / (1 + \exp(-\bar{h}_\theta(x)))$
  - ▶ Multi-class:  $h_\theta(x) = \text{softmax}(\bar{h}_\theta(x))$

# Introduction to Deep Learning

## Neural Network with a Single Neuron: The Housing Example

- ▶ Example in this unit: predict house price from its size.
- ▶ Starting point: fit a straight line.
- ▶ New requirement: disallow negative prices, introducing a “kink.”
- ▶ How can we represent a function with a kink as  $\bar{h}_\theta(x)$ ?

# Introduction to Deep Learning



Source: Ng and Ma (2023).

# Introduction to Deep Learning

## Parameterization and Activation Functions

► Define  $h_{\theta}(x)$ :

$$\bar{h}_{\theta}(x) = \max(wx + b, 0),$$

where  $\theta = (w, b) \in \mathbb{R}^2$



# Introduction to Deep Learning

## Parameterization and Activation Functions

- ▶ Rectified linear unit (ReLU):

$$\text{ReLU}(t) \equiv \max\{t, 0\}$$

- ▶ ReLU is an example of an activation function.
- ▶ Single neuron due to single non-linear activation function.

# Introduction to Deep Learning

## Basic Terminology

- ▶  $b$ : often referred to as the “bias.”
- ▶  $w$ : known as the weight vector.
- ▶ Can extend to multi-neuron case.

# Introduction to Deep Learning

## Neural Network with a Single Neuron for Multidimensional Input

$$\bar{h}_{\theta}(x) = \text{ReLU} \left( w^{\top} x + b \right) \quad (1)$$

► Where:  $w \in \mathbb{R}^d$ ,  $b \in \mathbb{R}$ , and  $\theta = (w, b)$ .

# Introduction to Deep Learning

## Stacking Neurons

- ▶ Neural networks can be made more complex by stacking neurons.
- ▶ Each neuron passes its output as input to the next neuron.
- ▶ Allows for flexibility in modeling complex functions.

# Introduction to Deep Learning

## Deepening the Housing Example

- ▶ New features: size, number of bedrooms, zip code, wealth of neighborhood.
- ▶ Neural network have modular structure, which allows us to stack them to create complex structures.
- ▶ Individual neurons are fundamental building block.

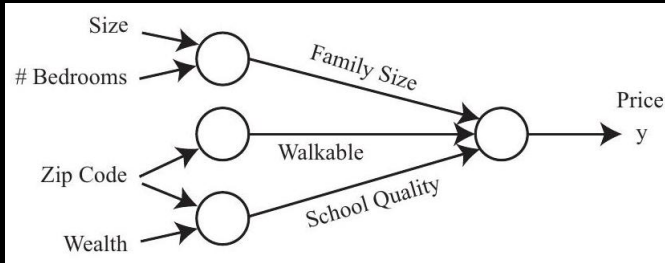
# Introduction to Deep Learning

## Deriving Features for Housing Price Prediction

- ▶ Family size: function of size of the house and number of bedrooms.
- ▶ Walkability: inferred from zip code.
- ▶ School quality: based on zip code and neighborhood wealth.
- ▶ House price: depends on family size, walkability, and school quality.

# Introduction to Deep Learning

## Small Neural Network for Predicting Housing Prices



Source: Ng and Ma (2023).

# Introduction to Deep Learning

## Input and Hidden Units

- ▶ Input features:  $x_1, x_2, x_3, x_4$ .
- ▶ Intermediate variables (hidden units):  $a_1$  (family size),  $a_2$  (walkable),  $a_3$  (school quality)

$$a_1 = \text{ReLU}(\theta_1 x_1 + \theta_2 x_2 + \theta_3)$$

$$a_2 = \text{ReLU}(\theta_4 x_3 + \theta_5)$$

$$a_3 = \text{ReLU}(\theta_6 x_3 + \theta_7 x_4 + \theta_8)$$

- ▶ Parameters:  $(\theta_1, \dots, \theta_8)$ .



# Introduction to Deep Learning

## Output Parameterization

- ▶ Output representation:  $\bar{h}_\theta(x)$ .

$$\bar{h}_\theta(x) = \theta_9 a_1 + \theta_{10} a_2 + \theta_{11} a_3 + \theta_{12}$$

- ▶ Parameters:  $\theta = (\theta_1, \dots, \theta_{12})$ .
- ▶ Complex function of  $x$  with parameters  $\theta$ .
- ▶ Try to learn  $\theta$ .

# Introduction to Deep Learning

## Inspiration from Biological Neural Networks

- ▶ Artificial neural networks inspired by biological ones.
- ▶ Hidden units: analogous to biological neurons.
- ▶ Parameters  $\theta_i$ : correspond to synapses.
- ▶ Uncertain similarity: deep artificial networks vs. biological ones.
- ▶ Open question: human brain learning mechanism vs. backpropagation in artificial networks.

# Introduction to Deep Learning

## Two-layer Fully-Connected Neural Networks

- ▶ Intermediate variables are functions of all  $x_1, \dots, x_4$ .
- ▶ This setup avoids the need for prior assumptions:

$$a_1 = \text{ReLU} \left( w_1^\top x + b_1 \right)$$

$$a_2 = \text{ReLU} \left( w_2^\top x + b_2 \right)$$

$$a_3 = \text{ReLU} \left( w_3^\top x + b_3 \right)$$

- ▶ The network is fully-connected: all  $a_i$ 's depend on all  $x_i$ 's.

# Introduction to Deep Learning

## General Two-layer Fully-Connected Neural Network

- ▶ For input  $x \in \mathbb{R}^d$  and  $m$  hidden units.
- ▶ Connection to previous slide: we generalize to  $d$  dimensions.

$$\begin{aligned}\forall j \in [1, \dots, m], \quad z_j &= w_j^{[1]\top} x + b_j^{[1]} \\ a_j &= \text{ReLU}(z_j), \\ a &= [a_1, \dots, a_m]^\top \\ \bar{h}_\theta(x) &= w^{[2]\top} a + b^{[2]}\end{aligned}$$

- ▶ Notation:  $a$  is a column vector, and indices  $^{[1]}$  and  $^{[2]}$  distinguish parameters.

# Introduction to Deep Learning

## Vectorization in Neural Networks

- ▶ Simplify neural network expressions using matrix and vector notations.
- ▶ Speed up implementation by avoiding for loops.
- ▶ Utilize matrix algebra and optimized numerical packages.
- ▶ Modern networks need vectorization for efficiency.

# Introduction to Deep Learning

## Vectorizing the Neural Network: Weight Matrix

- ▶ Define weight matrix  $W^{[1]}$  by concatenating all vectors  $w_j^{[1]}$ .
- ▶ Dimension of  $W^{[1]}$ :  $\mathbb{R}^{m \times d}$ :

$$W^{[1]} = \begin{bmatrix} -w_1^{[1]\top} & - \\ -w_2^{[1]\top} & - \\ \vdots & \\ -w_m^{[1]\top} & - \end{bmatrix}$$

# Introduction to Deep Learning

## Computing $z$ in Vectorized Form

- ▶ Use matrix-vector multiplication.
- ▶ Define  $z$ :  $[z_1, \dots, z_m]^T \in \mathbb{R}^m$ .

$$z = W^{[1]}x + b^{[1]}$$

# Introduction to Deep Learning

## Vectors as Column Vectors

- ▶ By default, vectors in  $\mathbb{R}^d$  are viewed as column vectors.
- ▶ Vectors can also be treated as  $d \times 1$  matrices.



# Introduction to Deep Learning

## Computing Activations

- ▶ Activations are computed using the ReLU function.
- ▶ Element-wise non-linear application.
- ▶ Element-wise ReLU is parallelized efficiently.

$$a = \text{ReLU}(z)$$

# Introduction to Deep Learning

## Summarizing the Model

- ▶ Define  $W^{[2]}$ :  $\begin{bmatrix} w^{[2]\top} \end{bmatrix} \in \mathbb{R}^{1 \times m}$ .
- ▶ Use weight matrices for efficient computation.

$$a = \text{ReLU} \left( W^{[1]}x + b^{[1]} \right)$$
$$\bar{h}_{\theta}(x) = W^{[2]}a + b^{[2]}$$

# Introduction to Deep Learning

## Structure of a Neural Network

- ▶  $W^{[1]}$  and  $W^{[2]}$ : weight matrices for first and second layers.
- ▶  $b^{[1]}$  and  $b^{[2]}$ : biases for the first and second layers.
- ▶ Activation  $a$  is the hidden layer.
- ▶ A two-layer neural network is a one-hidden-layer neural network.

# Introduction to Deep Learning

## Multi-layer Neural Networks

- ▶ Generalizes concept to multiple layers.
- ▶ Activation:  $a^{[k]} = \text{ReLU}(W^{[k]}a^{[k-1]} + b^{[k]})$ .
- ▶ Output:  $\bar{h}_{\theta}(x) = W^{[r]}a^{[r-1]} + b^{[r]}$ .

# Introduction to Deep Learning

## Activation Functions

- ▶ Functions applied element-wise introducing non-linearity.
- ▶ Example functions:
  - ▶ Sigmoid
  - ▶ Tanh
  - ▶ Leaky ReLU
  - ▶ GELU
  - ▶ Softplus

# Introduction to Deep Learning

## Linearity vs Non-Linearity

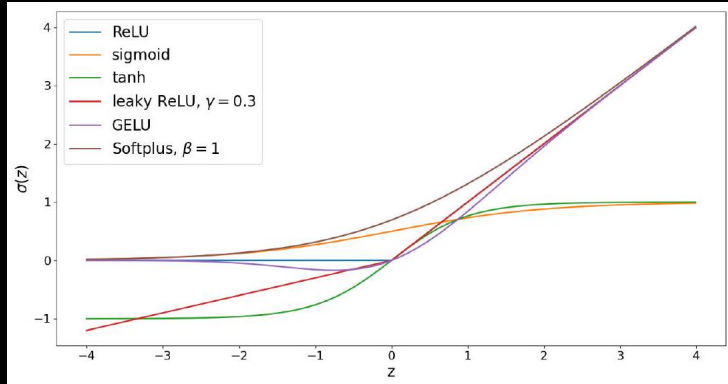
- ▶ Why not  $\sigma(z) = z$ ?
- ▶ With linear activation, multiple layers collapse to a single matrix product:

$$\tilde{W}x \text{ where } \tilde{W} = W^{[2]} W^{[1]}.$$

- ▶ Importance of non-linear activation functions: represent complex relationships.

# Introduction to Deep Learning

## Activation Functions



Source: Ng and Ma (2023).

# Introduction to Deep Learning

## Feature Maps in Traditional Machine Learning

- ▶ Represent non-linear functions:  $\theta^\top \phi(x)$ .
  - ▶  $\theta$ : parameters;  $\phi(x)$ : feature map.
- ▶ Performance depends on the feature map choice.
- ▶ Process of choosing: feature engineering.



# Introduction to Deep Learning

## Deep Learning as Feature Learning

- ▶ Feature map is learned, not handcrafted.
- ▶ Let  $\beta$  be parameters (excluding final layer).
- ▶ Abstract:  $a^{[r-1]} = \phi_{\beta}(x)$ .
- ▶ Model:  $\bar{h}_{\theta}(x) = W^{[r]}\phi_{\beta}(x) + b^{[r]}$ .

# Introduction to Deep Learning

## Training and Optimization

- ▶ When  $\beta$  is fixed,  $\phi_{\beta}(\cdot)$  is a feature map.
- ▶ During training,  $\beta$ ,  $W^{[r]}$ , and  $b^{[r]}$  are optimized.
- ▶ Learn the model and feature map simultaneously.
- ▶ Less dependence on domain knowledge and feature engineering.

# Introduction to Deep Learning

## Deep Learning Representations

- ▶ Second to last layer  $a^{[l]}$  is the “learned feature.”
  - ▶ E.g., in house pricing, neural network might learn “family size.”
- ▶ Learned representations can be transferred across datasets.
  - ▶ Might be complex and difficult to interpret.

# Introduction to Deep Learning

## Interpretability Challenge

- ▶ Neural networks can be viewed as a “black box.”
- ▶ Difficult to discern the features they discover.
- ▶ This poses challenges for interpretability.

# Introduction to Deep Learning

## Modules in Modern Neural Networks

- ▶ Modern neural networks are more sophisticated.
- ▶ Networks incorporate layers and building blocks.
- ▶ Will examine building blocks.

# Introduction to Deep Learning

## Matrix Multiplication as a Building Block

- ▶ Matrix multiplication serves as a primary building block:

$$\text{MM}_{W,b}(z) = Wz + b$$

- ▶  $W$  denotes the weight matrix.
- ▶  $b$  is the bias vector.
- ▶  $z$  is the input.

# Introduction to Deep Learning

## MLP: Composition of Modules

- ▶ MLP is a sequence of matrix multiplication modules interwoven with nonlinear activation modules:

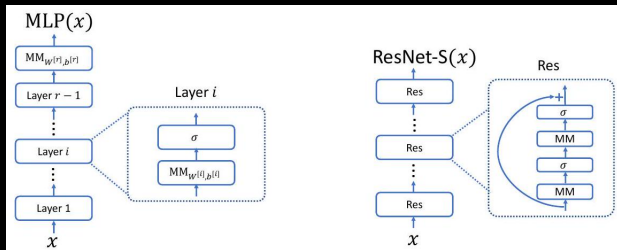
$$\text{MLP}(x) = \text{MM}(\sigma(\text{MM}(\sigma(\cdots \text{MM}(x))))$$

- ▶ Each module typically has its unique set of parameters.

# Introduction to Deep Learning

## Larger Modules from Smaller Ones

- ▶ A single 'layer' often encompasses an activation function,  $\sigma$ , and a matrix multiplication module, MM.



Source: Ng and Ma (2023).



# Residual Connections

- ▶ ResNets have revolutionized vision applications
- ▶ A residual block is given as follows:

$$\text{Res}(z) = z + \sigma(\text{MM}(\sigma(\text{MM}(z))))$$

- ▶ A simplified ResNet:

$$\text{ResNet-S}(x) = \text{MM}(\text{Res}(\text{Res}(\cdots \text{Res}(x))))$$

# Introduction to Deep Learning

## Layer Normalization

- ▶ Layer normalization, LN, is pivotal post-nonlinear activations. It transforms:

$$\text{LN-S}(z) = \begin{bmatrix} \frac{z_1 - \hat{\mu}}{\hat{\sigma}} \\ \vdots \\ \frac{z_m - \hat{\mu}}{\hat{\sigma}} \end{bmatrix}$$

- ▶ However, merely having zero mean and standard deviation 1 might not be ideal. Hence:

$$\text{LN}(z) = \beta + \gamma \cdot \text{LN-S}(z)$$

# Introduction to Deep Learning

## Computational Graph Representation

- ▶ Each computation visualized as a computational graph.
- ▶ Nodes represent operations, variables, or functions.
- ▶ Example:  $y = \sigma(Wx + b)$ 
  - ▶  $x$ : input
  - ▶  $W$ : weight matrix
  - ▶  $b$ : bias
  - ▶  $\sigma$ : activation function
- ▶ Equations can be broken down into smaller operations.

# Introduction to Deep Learning

## Layer Normalization

► Relies on chain rule.

► If  $f(g(h(x)))$ , then:

$$\frac{df}{dx} = \frac{df}{dg} \times \frac{dg}{dh} \times \frac{dh}{dx}$$

► Compute local gradient and multiply with gradient from next operation.

# Introduction to Deep Learning

## Forward and Backward Passes

### 1. Forward Pass:

- ▶ Compute output layer by layer.
- ▶ Standard computation pass.

### 2. Backward Pass:

- ▶ Compute gradient of loss with respect to each neuron's output.
- ▶ Continue layer-by-layer in reverse.

# Introduction to Deep Learning

## Weight Update

- ▶ Update weights using optimization algorithm.
- ▶ Adjust weights to decrease error:

$$W_{new} = W_{old} - \alpha \times \frac{\partial L}{\partial W}$$

- ▶ where  $L$  = loss and  $\alpha$  = learning rate.

# Introduction to Deep Learning

## Challenges with Backpropagation

- ▶ Vanishing gradient: gradients diminish to near zero.
- ▶ Exploding gradient: gradients become too large.
- ▶ Local minima: possible to get stuck in suboptimal solutions.

# Introduction to Deep Learning

## Improvements and Variants

- ▶ Stochastic and mini-batch GD: update after a small set.
- ▶ Momentum: use previous weight change for current update.
- ▶ Advanced Optimizers: Adam, RMSprop.
- ▶ Regularization techniques: dropout, L1/L2 regularization, early stopping.



# Introduction to Deep Learning

## Overview

- ▶ Review of the basic chain rule.
- ▶ Introduction to the general strategy for backpropagation.
- ▶ Computation of the backward function for basic modules in neural networks.
- ▶ Concrete backprop algorithm for MLPs.

# Introduction to Deep Learning

## Partial Derivatives in NNs

- ▶ Scalar variable  $J$  depending on variables  $z$ .
- ▶ Notation:  $\frac{\partial J}{\partial z}$ .
- ▶ Dimension of  $\frac{\partial J}{\partial z}$  is same as  $z$ .
- ▶ Example: If  $z \in \mathbb{R}^{m \times n}$ , then  $\frac{\partial J}{\partial z} \in \mathbb{R}^{m \times n}$ .

# Introduction to Deep Learning

## Remark

- ▶ When both  $J$  and  $z$  are not scalars, partial derivatives can become a matrix or tensor.
- ▶ Computationally expensive and rarely explicitly constructed.
- ▶ Focus is on derivatives of scalar function w.r.t vectors, matrices, or tensors.

# Introduction to Deep Learning

## Chain Rule

- ▶ Reviewing the chain rule:

$$z \in \mathbb{R}^m$$

$$u = g(z) \in \mathbb{R}^n$$

$$J = f(u) \in \mathbb{R}.$$

- ▶ Final variable  $J$  is a scalar.

# Introduction to Deep Learning

## Chain Rule (cont.)

$$u = (u_1, \dots, u_n)$$

$$g(z) = (g_1(z), \dots, g_n(z))$$

► Chain rule provides:

$$\forall i \in \{1, \dots, m\}, \quad \frac{\partial J}{\partial z_i} = \sum_{j=1}^n \frac{\partial J}{\partial u_j} \cdot \frac{\partial g_j}{\partial z_i}$$

# Introduction to Deep Learning

## Vectorized Notation

► When  $\mathbf{z}$  and  $\mathbf{u}$  are vectors:

$$\frac{\partial J}{\partial \mathbf{z}} = \begin{bmatrix} \frac{\partial g_1}{\partial z_1} & \dots & \frac{\partial g_n}{\partial z_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial g_1}{\partial z_m} & \dots & \frac{\partial g_n}{\partial z_m} \end{bmatrix} \cdot \frac{\partial J}{\partial \mathbf{u}}$$

# Introduction to Deep Learning

## Backward Function and Linear Map

- ▶ The backward function is always a linear map from  $\frac{\partial J}{\partial u}$  to  $\frac{\partial J}{\partial z}$ .
- ▶ The mapping can depend on  $z$  in complex ways.
- ▶ The matrix on the RHS is the transpose of the Jacobian matrix of  $g$ .

# Introduction to Deep Learning

## Advantages

- ▶ Avoids in-depth discussion about Jacobian matrices.
- ▶ Convenient in cases like  $\mathbf{z} \in \mathbb{R}^{r \times s}$ , giving:

$$\forall i, k, \quad \frac{\partial J}{\partial z_{ik}} = \sum_{j=1}^n \frac{\partial J}{\partial u_j} \cdot \frac{\partial g_j}{\partial z_{ik}}.$$



# Introduction to Deep Learning

## Key Interpretation of the Chain Rule

- ▶ Chain rule provides a way to compute  $\frac{\partial J}{\partial z}$  from  $\frac{\partial J}{\partial u}$ .
- ▶ Only requires knowledge about  $g$ .

# Introduction to Deep Learning

## Backward Function Notation

- ▶ Denoted as  $\mathcal{B}[g, z]$ .
- ▶ For fixed  $z$ ,  $\mathcal{B}[g, z]$  is a linear map from  $\mathbb{R}^n$  to  $\mathbb{R}^m$ .

# Introduction to Deep Learningx

## Backward Function as Matrix

- ▶ Viewed as a matrix, but depends on changing  $z$ .
- ▶ Inputs:  $z$  (input to  $g$ ) and  $v$  (gradient w.r.t to output of  $g$ ).
- ▶ Outputs: gradient of  $J$  w.r.t  $z$ .

# Introduction to Deep Learning

## General Strategy of Backpropagation

- ▶ Neural networks: complex compositions of building blocks.
- ▶ Modules: MM,  $\sigma$ , Conv2D, LN, and others.
- ▶ Losses are also viewed as modules.

# Introduction to Deep Learning

## Example of Modular Composition

- ▶ Example:  $J = M_k (M_{k-1} (\cdots M_1(x)))$ .
- ▶ Modules involve parameters or fixed operations.
- ▶ Each  $M_i$  involves a set of parameters  $\theta^{[i]}$ .

# Introduction to Deep Learning

## Intermediate Variables

- ▶ Introduce intermediate variables.

$$u^{[0]} = x$$

$$u^{[1]} = M_1 \left( u^{[0]} \right)$$

$$u^{[2]} = M_2 \left( u^{[1]} \right)$$

$$\vdots$$

$$J = u^{[k]} = M_k \left( u^{[k-1]} \right) .$$

# Introduction to Deep Learning

## Overview of Backpropagation

- ▶ Backpropagation consists of two passes:
  - ▶ Forward pass.
  - ▶ Backward pass.

# Introduction to Deep Learning

## Intermediate Variables

- ▶ Compute  $u^{[1]}, \dots, u^{[k]}$  from  $i = 1, \dots, k$ , sequentially using the definition in (F).
- ▶ Save all intermediate variables  $u^{[i]}$ 's in memory.



# Introduction to Deep Learning

## Backward Pass

- ▶ Compute derivatives w.r.t intermediate variables:  $\frac{\partial J}{\partial u^{[k]}}, \dots, \frac{\partial J}{\partial u^{[1]}}$ .
- ▶ Compute derivatives of parameters:  $\frac{\partial J}{\partial \theta^{[l]}}$ .
- ▶ Interleaved computations are possible.

# Introduction to Deep Learning

## Chain Rule in Backpropagation

Efficient computation of  $\frac{\partial J}{\partial u^{[i-1]}}$  from  $\frac{\partial J}{\partial u^{[i]}}$  and  $u^{[i-1]}$ :

$$\frac{\partial J}{\partial u^{[i]}} \Longrightarrow \frac{\partial J}{\partial u^{[i-1]}}.$$

$$\frac{\partial J}{\partial u^{[i-1]}} = \mathcal{B} \left[ M_i, u^{[i-1]} \right] \left( \frac{\partial J}{\partial u^{[i]}} \right).$$

# Introduction to Deep Learning

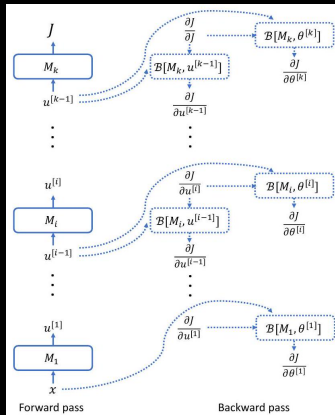
## Further Chain Rule Application

- Instantiating with  $z = \theta^{[l]}$  and  $u = u^{[l]}$ :

$$\frac{\partial J}{\partial \theta^{[l]}} = \mathcal{B} \left[ M_i, \theta^{[l]} \right] \left( \frac{\partial J}{\partial u^{[l]}} \right).$$

# Introduction to Deep Learning

## Further Chain Rule Application



Source: Ng and Ma (2023)

# Introduction to Deep Learning

## Computational Efficiency and Granularity

- ▶ Small modules have efficient backward functions.
- ▶ Backward functions of atomic modules like addition, multiplication, ReLU are computationally efficient.
- ▶ Neural networks as compositions of atomic operations.
- ▶ Modularize using matrix multiplication, layer norm, etc. for practicality.

# Introduction to Deep Learning

## Backward Functions for Basic Modules

- ▶ Compute the backward function for:
  - ▶ Basic module MM.
  - ▶ Activations  $\sigma$ .
  - ▶ Loss functions.

# Introduction to Deep Learning

$$\text{MM}_{W,b}(z) = Wz + b \quad (2)$$

$$\mathcal{B}[\text{MM}, z](v) = W^\top v \in \mathbb{R}^m \quad (3)$$

$$\mathcal{B}[\text{MM}, W](v) = vz^\top \in \mathbb{R}^{n \times m} \quad (4)$$

$$\mathcal{B}[\text{MM}, b](v) = v \quad (5)$$

# Introduction to Deep Learning

## Computational Efficiency

- ▶ Computational efficiency for computing the backward function is  $O(mn)$ .
- ▶ Equivalent to evaluating the result of matrix multiplication up to a constant factor.



# Introduction to Deep Learning

## Backward Function for Activations

- ▶ Let  $M(z) = \sigma(z)$  where  $\sigma$  is an element-wise activation function and  $z \in \mathbb{R}^m$ .
- ▶ We have:

$$\mathcal{B}[\sigma, z](v) = \sigma'(z) \odot v \in \mathbb{R}^m.$$

# Introduction to Deep Learning

## Notes on Activation Backward Function

- ▶  $\frac{\partial \sigma(z_j)}{\partial z_i} = 0$  when  $j \neq i$ .
- ▶ The backward function looks like  $O(m^2)$  but is  $O(m)$  in practice.
- ▶ Using smaller modules simplifies the process.

# Introduction to Deep Learning

## Backward Function for Loss Functions

- ▶ Given a module  $M$  taking vector  $z$  and outputting a scalar, we have:

$$\mathcal{B}[M, z](v) = \frac{\partial M}{\partial z} \cdot v$$

- ▶ Examples:

- ▶ Squared Loss:  $\mathcal{B}[\ell_{\text{MSE}}, z](v) = (z - y) \cdot v$

- ▶ Logistics Loss:  $\mathcal{B}[\ell_{\text{logistic}}, t](v) = (1/(1 + \exp(-t)) - y) \cdot v$

- ▶ Cross-Entropy Loss:  $\mathcal{B}[\ell_{\text{ce}}, t](v) = (\phi - e_y) \cdot v$

# Introduction to Deep Learning

## Back-propagation for MLPs

- If we are given backward functions for modules to evaluate loss, can compute gradient of loss with respect to hidden activations and parameters.

# Introduction to Deep Learning

## Back-propagation for MLPs: Forward Pass

- ▶ An  $r$ -layer MLP with logistic loss:

$$\begin{aligned} z^{[1]} &= \text{MM}_{W^{[1]}, b^{[1]}}(x), \\ &\vdots \\ J &= \ell_{\text{logistic}} \left( z^{[r]}, y \right). \end{aligned}$$

# Introduction to Deep Learning

## Back-propagation for MLPs: Backward Pass

- Compute the gradient of loss  $J$  w.r.t  $z^{[r]}$ :

$$\frac{\partial J}{\partial z^{[r]}} = \left( 1 / \left( 1 + \exp \left( -z^{[r]} \right) \right) - y \right)$$

- Then, compute the gradient with respect to parameters  $W^{[k+1]}$  and  $b^{[k+1]}$  and others iteratively.

# Introduction to Deep Learning

## Algorithm 3: Back-propagation for Multi-layer Neural Networks

1. Forward pass. Compute and store the values of  $a^{[k]}$  's,  $z^{[k]}$  's, and  $J$ .
2. Backward pass. Compute the gradient of loss  $J$  with respect to  $z^{[r]}$ .
3. Iterate from  $k = r - 1$  to 0 for the rest.

# Introduction to Deep Learning

## Backpropagation Equations

$$\begin{aligned}\frac{\partial J}{\partial \mathbf{a}^{[k]}} &= \mathcal{B} \left[ \sigma, \mathbf{a}^{[k]} \right] \left( \frac{\partial J}{\partial \mathbf{z}^{[k+1]}} \right) \\ &= \mathbf{W}^{[k+1]\top} \frac{\partial J}{\partial \mathbf{z}^{[k+1]}} \\ \frac{\partial J}{\partial \mathbf{z}^{[k]}} &= \mathcal{B} \left[ \sigma, \mathbf{z}^{[k]} \right] \left( \frac{\partial J}{\partial \mathbf{a}^{[k]}} \right) \\ &= \sigma' \left( \mathbf{z}^{[k]} \right) \odot \frac{\partial J}{\partial \mathbf{a}^{[k]}}\end{aligned}$$



# Introduction to Deep Learning

## Vectorization Over Training Examples

- ▶ Avoid explicit loops by stacking examples in matrices.
- ▶ Each column in a matrix represents a training example.
- ▶ Efficient computation on hardware optimized for matrix operations.

# Introduction to Deep Learning

## Column Vector Notation

- ▶ Use broadcasting to add bias  $b^{[1]}$  to each column of  $W^{[1]}X$ .

$$Z^{[1]} = W^{[1]}X + \tilde{b}^{[1]}$$

- ▶ Not necessary to construct  $\tilde{b}^{[1]}$  explicitly.

# Introduction to Deep Learning

## Data Representation Mismatch

- ▶ Mismatch between deep learning packages (row vectors) and notation in papers (column vectors).
- ▶ Adjustments:
  - ▶ Columns become row vectors and vice-versa.
  - ▶ Matrices are transposed.
  - ▶ Order of matrix multiplication is flipped.
- ▶ Possible reason: natural inclination to multiply a matrix to a vector from the left.

## 2. Deep Learning Applications

# Introduction to TensorFlow

## Colab Tutorial

- ▶ Introduction to Deep Learning

# References I

Ng, Andrew and Tengyu Ma (2023) “CS229 Lecture Notes,” June.