

CS540 Introduction to Artificial Intelligence **Neural Networks: Review**

University of Wisconsin-Madison Spring 2025



Announcements

Homework: HW7 is due on Monday Apr. 7th at 11:59 PM _

Class roadmap:

•

Deep Learning and Neural Network's Summary

search

- Search I: Un-Informed
- Search II: Informed search

How to classify

Cats vs. dogs?

Neural networks can also be used for regression. - Typically, no activation on outputs, mean squared error loss function.



Single-layer Perceptron

Multi-layer Perceptron

Training of neural networks

Convolutional neural networks









Inspiration from neuroscience

- Inspirations from human brains - Networks of simple and homogenous units (a.k.a neuron)



(soma)

(wikipedia)





Perceptron • Given input x, weight w and bias b, perceptron outputs: $o = \sigma(\mathbf{w}^{\mathsf{T}}\mathbf{x} + b)$ $\sigma(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$ Activation function

Cats vs. dogs?











Perceptron

• Goal: learn parameters $\mathbf{W} = \{w_1, w_2, \dots, w_d\}$ and b to minimize the classification error

Cats vs. dogs?









Example 2: Predict whether a user likes a song or not









Example 2: Predict whether a user likes a song or not using Perceptron





User Sharon









Learning logic functions using perceptron The perceptron can learn an AND function





 $w_1 = 1, w_2 = 1, b = -1.5$













XOR Problem (Minsky & Papert, 1969)

The perceptron cannot learn an XOR function (neurons can only generate linear separators)



This contributed to the first AI winter





Quiz break

Which one of the following is NOT true about perceptron?

- A. Perceptron only works if the data is linearly separable.
- B. Perceptron can learn AND function
- C. Perceptron can learn XOR function
- D. Perceptron is a supervised learning algorithm

Quiz break

Which one of the following is NOT true about perceptron?

- A. Perceptron only works if the data is linearly separable.
- B. Perceptron can learn AND function
- C. Perceptron can learn XOR function
- D. Perceptron is a supervised learning algorithm

Multilayer Perceptron



Single Hidden Layer

How to classify Cats vs. dogs?









Single Hidden Layer

- Input $\mathbf{x} \in \mathbb{R}^d$
- Hidden $\mathbf{W} \in \mathbb{R}^{m \times d}$, $\mathbf{b} \in \mathbb{R}^m$
- Intermediate output

$\mathbf{h} = \sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$

 σ is an element-wise activation function



Neural networks with one hidden layer



W



Single Hidden Layer

- Output $f = \mathbf{w}_2^\mathsf{T}\mathbf{h} + b_2$
- Normalize the output into probability using sigmoid

$$p(y = 1 | \mathbf{x}) = \frac{1}{1 + e^{-f}}$$

Output

Multi-class classification

Turns outputs f into k probabilities (sum up to 1 across k classes)

$p(y|\mathbf{x}) = softmax(\mathbf{f})$ $\frac{\exp f_y(x)}{\sum_{i}^{k} \exp f_i(x)}$

Deep neural networks (DNNs)

$\mathbf{h}_1 = \sigma(\mathbf{W}_1\mathbf{x} + \mathbf{b}_1)$ $\mathbf{h}_2 = \sigma(\mathbf{W}_2\mathbf{h}_1 + \mathbf{b}_2)$ $h_3 = \sigma(W_3h_2 + b_3)$ $f = W_4 h_3 + b_4$ $\mathbf{y} = \operatorname{softmax}(\mathbf{f})$

NNs are composition of nonlinear functions

Classify MNIST handwritten digits

How to train a neural network?

Loss function: $\frac{1}{|D|} \sum_{i}^{l} \ell(\mathbf{x}_i, y_i)$

Per-sample loss:

 $\ell(\mathbf{x}, y) = \sum_{j=1}^{K} - y_j \log p_j$

Also known as cross-entropy loss or softmax loss

Cross-Entropy Loss

$L_{CE} = \sum_{i} - y_{i} \log(p_{i})$

$= -\log(0.8)$

Goal: push p and Y to be identical

How to train a neural network?

Update the weights W to minimize the loss function

$$L = \frac{1}{|D|} \sum_{i}^{l} \ell(\mathbf{x}_i, y_i)$$

Use gradient descent!

Gradient Descent

- Choose a learning rate $\alpha > 0$
- Initialize the model parameters W_0
- For *t* =1, 2, ...
 - Update parameters:

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \alpha_{\overline{\partial}}$$

 $= w_{t-1} - \alpha_{-1}$

Repeat until converges

Minibatch Stochastic Gradient Descent

- Choose a learning rate $\alpha > 0$
- Initialize the model parameters W_0
- For *t* =1, 2, ...
 - Randomly sample a subset (mini-batch) $B \subset D$ Update parameters:

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \alpha$$

Repeat

Calculate gradient: backpropagation with chain rule • Define a loss function L, must compute $\frac{\partial L}{\partial W}$, $\frac{\partial L}{\partial h}$ for all

- weights and biases.
- Gradient to a variable = X Z_1 *

W

Non-convex Optimization

[Gao and Li et al., 2018]

How to classify Cats vs. dogs?

Dual 12202 wide-angle and telephoto cameras

36M floats in a RGB image!

Fully Connected Networks

Cats vs. dogs?

~ 36M elements x 100 = \sim 3.6B parameters!

Convolutions come to rescue!

Where is Waldo?

Why Convolution?

- Translation
 Invariance
- Locality

2-D Convolution

 $0 \times 0 + 1 \times 1 + 3 \times 2 + 4 \times 3 = 19,$ $1 \times 0 + 2 \times 1 + 4 \times 2 + 5 \times 3 = 25,$ $3 \times 0 + 4 \times 1 + 6 \times 2 + 7 \times 3 = 37,$ $4 \times 0 + 5 \times 1 + 7 \times 2 + 8 \times 3 = 43.$

19	25
37	43

(vdumoulin@ Github)

2-D Convolution Layer

- **X**: $n_h \times n_w$ input matrix
- W: $k_h \times k_w$ kernel matrix
- b: scalar bias
- **Y**: $(n_h k_h + 1) \times (n_w k_w + 1)$ output matrix

$\mathbf{Y} = \mathbf{X} \star \mathbf{W} + b$

• W and b are learnable parameters

25 19 37 43

2-D Convolution Layer with Stride and Padding Stride is the #rows/#columns per slide Padding adds rows/columns around input Kernel Input Output 2 0 8 5 * 3 4 **Kernel/filter size** 6 2 3 8 6 8 0 1 0 1 0 1

- Output shape

$$\begin{bmatrix} n_h - k_h + p_h + s_h \\ \uparrow & \uparrow & \uparrow \\ \end{bmatrix}$$
Input size Pad Stride

Pad

 $| \times |(n_w - k_w + p_w + s_w)/s_w|$

Multiple Input Channels

- Input and kernel can be 3D, e.g., an RGB image have 3 channels
- channels

Have a kernel for each channel, and then sum results over

Multiple Input Channels

- Input and kernel can be 3D, e.g., an RGB image have 3 channels
- channels

Have a 2D kernel for each channel, and then sum results over

Multiple Input Channels Input and kernel can be 3D, e.g., an RGB image have 3 channels

output channel (due to summation over channels)

Also call each 3D kernel a "filter", which produce only one

RGB (3 input channels)

Multiple filters (in one layer) • Apply multiple filters on the input Each filter may learn different features about the input Each filter (3D kernel) produces one output channel

RGB (3 input channels)

Conv1 Filters in AlexNet

- 96 filters (each of size 11x11x3)
- Gabor filters

Figures from Visualizing and Understanding Convolutional Networks by M. Zeiler and R. Fergus

Multiple Output Channels

- The # of output channels = # of filters
- Input X: $c_i \times n_h \times n_w$
- Kernel W: $c_o \times c_i \times k_h \times k_w$
- Output Y: $c_o \times m_h \times m_w$

 $\mathbf{Y}_{i\ldots} = \mathbf{X} \star \mathbf{W}_{i\ldots}$ $for i = 1, ..., c_{0}$

Convolutional Neural Networks

LeNet Architecture

gluon-cv.mxnet.io

answer: 0

Y. LeCun, L. Bottou, Y. Bengio, P. Haffner, 1998 Gradient-based learning applied to document recognition

Quiz break

Which one of the following is NOT true?

- A. LeNet has two convolutional layers
- input
- C. Pooling is performed right after convolution
- D. Pooling layer does not have learnable parameters

B. The first convolutional layer in LeNet has 5x5x6x3 parameters, in case of RGB

Quiz break

Which one of the following is NOT true?

- A. LeNet has two convolutional layers
- input
- C. Pooling is performed right after convolution
- D. Pooling layer does not have learnable parameters

B. The first convolutional layer in LeNet has 5x5x6x3 parameters, in case of RGB

Pooling is performed after ReLU: conv -> relu -> pooling

Evolution of neural net architectures

Deng et al. 2009

[Krizhevsky et al. 2012]

AlexNet vs LeNet Architecture

Larger pool size, change to max pooling

Larger kernel size, stride because of the increased image size, and more output channels.

AlexNet Architecture

ResNet: Going deeper in depth

ImageNet Top-5 error%

152 layers

[He et al. 2015]

Other neural network architectures

- Convolutional neural networks are one of many special types of layers.
 - Main use is for processing images.
 - Also can be useful for handling time series.
- Other common architectures:
 - Recurrent neural networks: hidden activations are a function of input and activations from previous inputs. Designed for sequential data such as text.
 - Graph neural networks: take graph data as input.
 - Transformers: take sequences as input and learn what parts of input to pay attention to.

Brief history of neural networks

What we've learned today...

- Modeling a single neuron
 - Linear perceptron
 - Limited power of a single neuron
- Multi-layer perceptron
- Training of neural networks
 - Loss function (cross entropy)
 - Backpropagation and SGD
- Convolutional neural networks
 - Convolution, pooling, stride, padding
 - Basic architectures (LeNet etc.)
 - More advanced architectures (AlexNet, ResNet etc)

Thank you!

Some of the slides in these lectures have been adapted from materials developed by Alex Smola and Mu Li: <u>https://courses.d21.ai/berkeley-stat-157/index.html</u>