



# Deep Neural Networks II

#### Sen Wang

UDRC Co-Investigator Associate Professor in Robotics and Autonomous Systems Institute of Signals, Sensors and Systems Heriot-Watt University

UDRC-EURASIP Summer School

30th June 2021

Slides adapted from Andrej Karpathy, Kaiming He

#### **Learning features for machines to solve problems**

1

- Convolutional Neural Networks (CNNs)
- CNN Architectures learning features
- Some Deep Learning Applications problems
  - Object detection (image, radar, sonar)
  - Semantic segmentation
  - Object detection and tracking using radar and sonar
  - Visual odometry
  - 3D reconstruction
  - Semantic mapping
  - Robot navigation
  - Manipulation and grasping
  - 0 .....

## Deep Learning

Deep Learning: a learning technique combining layers of neural networks to **automatically identify features** that are relevant to the problem to solve



### Deep Learning in Robotics



https://youtu.be/l8zKZLqkfll https://youtu.be/qhUvQiKec2U https://youtu.be/9j2a1oAHDL8 https://youtu.be/2N\_wKXQ6MXA

2021 UDRC-EURASIP Summer School

-2-

# Feed-forward Neural Networks

 $\bullet$   $\bullet$   $\bullet$ 

#### From Biological Neuron to Neural Network



#### Feedforward Neural Networks

• Feedforward Neural Networks or Multi-Layer Perceptrons (MLPs): fully-connected layers



#### Feedforward Neural Networks



 $a_1^2 = f(w_{1,1}^2 I_1 + w_{1,2}^2 I_2 + w_{1,3}^2 I_3 + b_1^2) \quad f(.) \text{ is an activation function}$  $a_2^2 = f(w_{2,1}^2 I_1 + w_{2,2}^2 I_2 + w_{2,3}^2 I_3 + b_2^2)$ 

a<sup>3</sup><sub>2</sub>?
2021 UDRC-EURASIP Summer School

#### Activation Functions f(\*)

- activation function
  - o introduce non-linearity into the neural network



Leaky ReLU  $\max(0.1x, x)$ 



 $\begin{array}{l} \textbf{Maxout} \\ \max(w_1^T x + b_1, w_2^T x + b_2) \end{array}$ 



#### Feedforward Neural Networks: Vector Formulation





#### Feedforward Neural Networks: Vector Formulation



$$a^{2} = f(W^{2^{T}}I + b^{2}) \quad a^{3} = f(W^{3^{T}}a^{2} + b^{3}) \quad O = f(W^{4^{T}}a^{3} + b^{4})$$

# forward pass of a 3 layer feed-forward neural networks
f = lambda x: 1.0/(1.0 + np.exp(-x)) # activation function (sigmoid)
I = np.random.randn(3, 1) # random input vector
a2 = f(np.dot(W2, I) + b2) # first hidden layer activation
a3 = f(np.dot(W3, a2) + b3) # second hidden layer activation
0 = f(np.dot(W4, a3) + b4) # output

# Convolutional Neural Networks (CNNs)

 $\bullet$   $\bullet$   $\bullet$ 

#### From MLPs to CNNs

- Feed-forward Neural Networks or Multi-Layer Perceptrons (MLPs)
  - many multiplications
- CNNs are similar to Feed-forward Neural Networks

   convolution instead of general matrix multiplication

$$S(i,j) = (I * K)(i,j)$$
$$= \sum_{m} \sum_{n} I(i+m, j+n)K(m,n)$$



### CNNs

- 3 Main Types of Layers:

   convolutional layer
   activation layer
   pooling layer
- repeat many times





5x5x3 filter

**Convolve** the filter with the image i.e. "slide over the image spatially, computing dot products"

Slides courtesy of Andrej Karpathy



Filters always extend the full depth of the input volume

**Convolve** the filter with the image i.e. "slide over the image spatially, computing dot products"









#### consider a second, green filter



For example, if we had 6 of 5x5 filters, we'll get 6 separate activation maps:



We stack these up to get a "new image" of size 28x28x6

For example, if we had 6 of 5x5 filters, we'll get 6 separate activation maps:



We processed [32x32x3] volume into [28x28x6] volume.

Q: how many parameters and multiplies would this be if we used a fully connected layer instead?

A: (32\*32\*3)\*(28\*28\*6) = **14.5M parameters**, ~**14.5M multiplies** 

For example, if we had 6 of 5x5 filters, we'll get 6 separate activation maps:



## **CNNs: Activation Layer**

- 3 Main Types of Layers:

   convolutional layer
   activation layer
  - pooling layer





## CNNs: Pooling Layer

- 3 Main Types of Layers:

   convolutional layer
   activation layer
  - pooling layer



fully-connected layer

## **CNNs: A sequence of Convolutional Layers**



# **CNN** Architectures

 $\bullet$   $\bullet$   $\bullet$ 

#### Hand-Crafted Features by Human



#### Feature Engineering and Representation

#### **Pervasive Data**





#### Deep Learning: Representation Learning

**Pervasive Data** 



#### LeNet - 1998

- Convolution for neural networks
- Fully-connected outputs
- weight-sharing is a key to reduce number of parameters

#### Foundation of modern ConvNets



Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

"Gradient-based learning applied to document recognition", LeCun et al. 1998 2021 UDRC-EURASIP Summer School

## AlexNet – 2012

8 layers: 5 conv and max-pooling + 3 fully-connected

LeNet-style backbone, plus:

- ReLU
  - Accelerate training
  - o better gradprop (vs. tanh)
- Dropout
  - Reduce overfitting
- Data augmentation
  - Image transformation
  - Reduce overfitting



"ImageNet Classification with Deep Convolutional Neural Networks", Krizhevsky, Sutskever, Hinton. NIPS 2012 2021 UDRC-EURASIP Summer School

3x3 conv, 64

#### 3x3 conv, 64, pool/2

3x3 conv, 128

3x3 conv, 128, pool/2

3x3 conv, 256

3x3 conv, 256

3x3 conv, 256

3x3 conv, 256, pool/2

3x3 conv, 512

3x3 conv, 512

3x3 conv, 512

3x3 conv, 512, pool/2

3x3 conv, 512

3x3 conv, 512

3x3 conv, 512

3x3 conv, 512, pool/2

fc, 4096

fc, 4096

fc, 1000

"Very Deep Convolutional Networks for Large-Scale Image Recognition", Simonyan & Zisserman. arXiv 2014 (ICLR 2015) 2021 UDRC-EURASIP Summer School 33

## VGG16/19 - 2014

Very deep ConvNet

Modularized design

- 3x3 Conv as the module
- Stack the same module
- Same computation for each module

Stage-wise training

• VGG-11 => VGG-13 => VGG-16

## GoogleNet/Inception - 2014

22 layers Multiple branches

- e.g., 1x1, 3x3, 5x5 convolutions and pooling
- merged by concatenation
- Reduce dimensionality by 1x1 conv before expensive 3x3/5x5 conv (change num of channels)



Szegedy et al. "Going deeper with convolutions". arXiv 2014 (CVPR 2015)

2021 UDRC-EURASIP Summer School



Going Deeper

## Simply stacking layers?



- Plain nets: stacking 3x3 conv layers
- 56-layer net has **higher training error and test error** than 20-layer net
- A deeper model should not have higher training error

## Going Deeper



## Cannot go deeper for deep neural networks

#### Problem:

deeper plain nets have higher training error on various datasets

#### Optimization difficulties:

- vanishing gradient
- solvers struggle to find the solution when going deeper

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

#### ResNets-2016



Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016. 2021 UDRC-EURASIP Summer School



#### ResNets-2016



- Deep ResNets can be trained easier
- Deeper ResNets have lower training error, and also lower test error

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

#### ImageNet experiments



top 5 error %

#### DenseNets - 2017

- simply connect every layer directly with each other
  - each layer has direct access to the gradients from the loss function and the original input image
- DenseNets concatenate the output feature maps of the layer with the incoming feature maps.



$$x_{l} = H_{l}(x_{l-1})$$

$$x_{l} = H_{l}(x_{l-1}) + x_{l-1}$$

$$x_{l} = H_{l}([x_{0}, x_{1}, \dots, x_{l-1}])$$

G. Huang, Z. Liu and L. van der Maaten, "Densely Connected Convolutional Networks," 2017.

#### MobileNets - 2017

#### better accuracy vs better efficiency

Light-weight ConvNets for mobile applications using depthwise convolutions



Howard. et. al. MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications 2017

#### EfficientNet - 2019

carefully design architecture to achieve top results with reasonable parameters

built upon a model (B0) from neural architecture search

Compound scaling: scale up network depth (more layers), width (more channels per layer), resolution (input image) simultaneously

#### EfficientNet-B7

- state-of-the-art 84.4% top-1 / 97.1% top-5 accuracy
- 8.4x smaller and 6.1x faster on inference than the best existing ConvNet.

Tan, Mingxing, and Quoc Le. "Efficientnet: Rethinking model scaling for convolutional neural networks." ICML 2019.



## **CNN** Timeline



# Deep Learning Applications

 $\bullet$   $\bullet$   $\bullet$ 

## Deep Learning Applications



#### **Object Detection and Recognition**

#### Vision based: RCNN, Fast RCNN, Faster RCNN, YOLO, SSD,.....







46

#### **Object Detection and Recognition**

#### Radar and polarized sensor based object detection and recognition



#### Vehicle detection and recognition using radar [Sheeny 2021]

Pape

Heriot-Watt RADIATE Dataset

att RADIATE Dataset Home Documentation Downloads

Welcome to Heriot-Watt RADIATE Dataset Website

Multi-Modality Radar Dataset in Adverse Weather with Object Annotation

#### http://pro.hw.ac.uk/radiate/

Choose to Start



2021 UDRC-EURASIP Summer School



vehicle detection using polarised infrared sensors [Sheeny 2018] 48

## Semantic Segmentation

#### FCN, SegNet, RefineNet, PSPNet, .....





#### **Robot** Navigation







#### Extra Resource

#### Book



#### **Online Codes and Practice**

- <u>https://codelabs.developers.google.com/</u>
- <u>TensorFlow, Keras and deep learning,</u> <u>without a PhD</u>







 <u>Build and deploy a custom object</u> <u>detection model with TensorFlow Lite</u> <u>(Android)</u>

#### Summary

- Deep Learning is a powerful tool
- Leveraging pre-trained models for feature extraction whenever possible (similar modality with smaller dataset)
- **Learning representation** is the key for Deep Learning







# Thank you for your attention!

Dr. Sen Wang s.wang@hw.ac.uk