

Lecture 9: Machine Learning and Modern Visual Recognition Techniques

9.1 Neural Network Basics

Work on Artificial Neural Networks, generally known as just 'Neural Networks', has been inspired by the fact that the human brain processes information in a completely different way than computers. Due to this, human brains are much more efficient at certain tasks, such as vision, language processing, etc. than computers. Artificial neural networks aim to mimic this biological model by employing a large number of simple interconnected processing units or 'neurons'. We may thus offer the following definition of a neural network viewed as an adaptive machine[1]:

A neural network is a massively parallel distributed processor made up of simple processing units that has a natural propensity for storing experiential knowledge and making it available for use. It resembles the brain in two respects:

1. Knowledge is acquired by the network from its environment through a learning process.
2. Inter-neuron connection strengths, known as synaptic weights, are used to store the acquired knowledge.

9.1.1 Perceptron - Analogy to a Biological Neuron

The figure below shows a simplified diagram of a biological neuron:

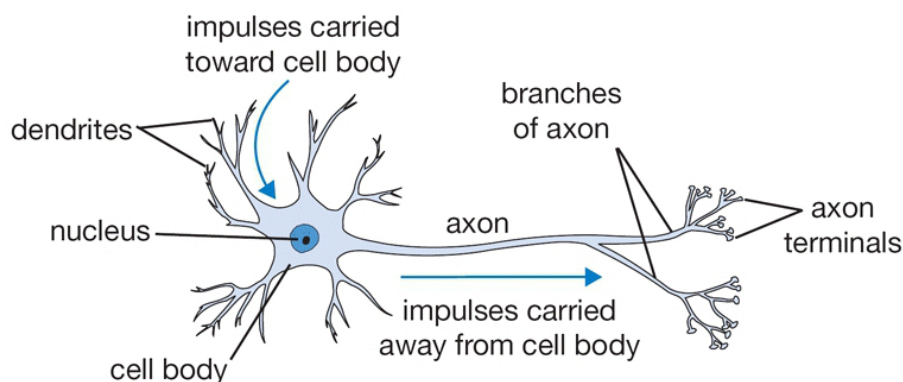


Figure 9.1: Simplified model of a biological neuron[2]

There are approximately 86 billion neurons in the human nervous system and they are connected with approximately 10^{14} - 10^{15} synapses[2].

The figure below shows a commonly used mathematical model of a neuron:

Signals travel along the axons (e.g. x_0) and interact multiplicatively (e.g. w_0x_0) with the dendrites of the other neuron based on the synaptic strength at that synapse (e.g. w_0). The idea is that the synaptic

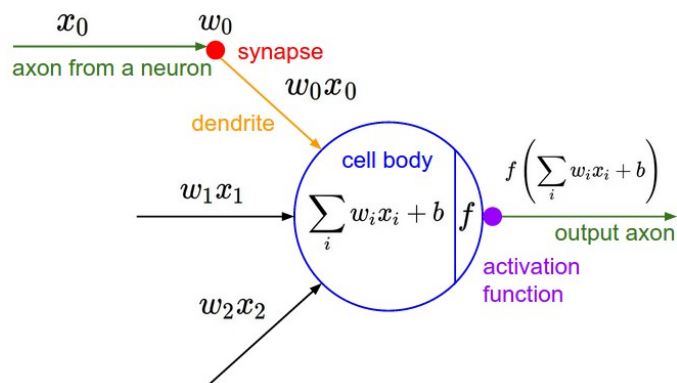


Figure 9.2: Mathematical model of a neuron[2]

strengths (the weights w) are learnable and control the strength of influence (and its direction: excitatory (positive weight) or inhibitory (negative weight)) of one neuron on another. We can model the firing rate of the neuron with an activation function f , which represents the frequency of the spikes along the axon. Historically, a common choice of activation function was the sigmoid function, but the ReLU function is more commonly used nowadays. This is because the gradient of the sigmoid function becomes small as its value increases which reduces the training speed of the model.

9.1.2 Single Layer Network / Logistic Regression

Inspired by this model of the neuron, a 1 layer neural net can be built which takes in binary inputs and provides binary output by taking a weighted sum and passing it through an activation function (as seen below). The theory is that if these weights are tuned perfectly, then it should be able to classify the inputs correctly.

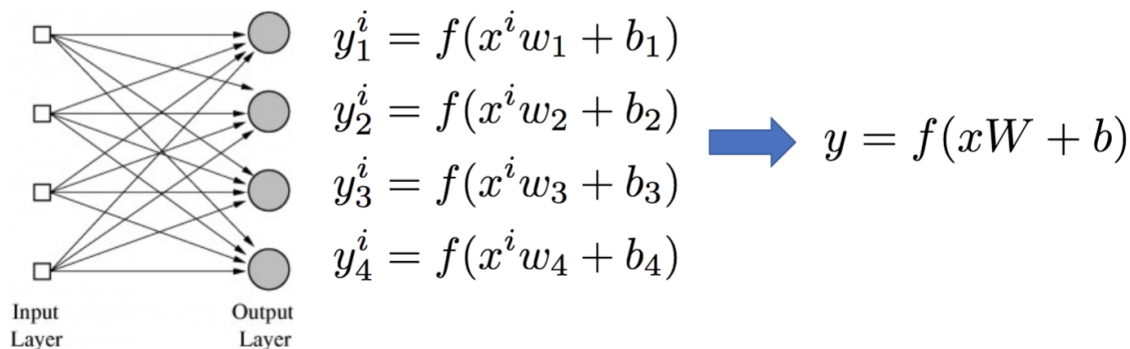


Figure 9.3: Single layer neural net[?]

9.1.3 Multi-layer Neural Network / Deep Learning

We can get more resolution on our classifications by putting multiple of these layers together. In such regard, each layer can "learn" to classify a different part of the input. Take for an example the problem

of classifying handwritten digits, in particular classifying the number 3. The first layer can be responsible for classifying a curved top, the second a curved bottom, etc.

A sample is shown below (from : <http://www.cs.cmu.edu/~aharley/vis/conv/flat.html>)

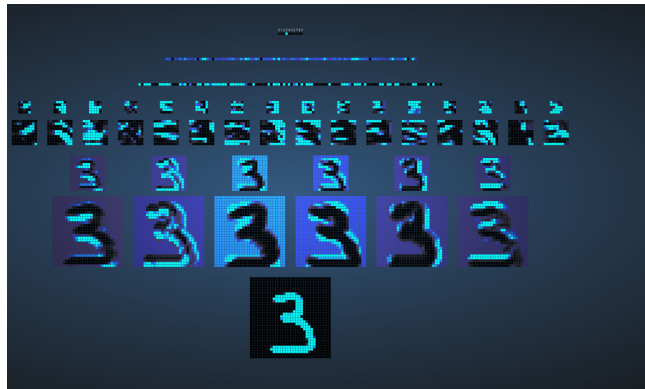
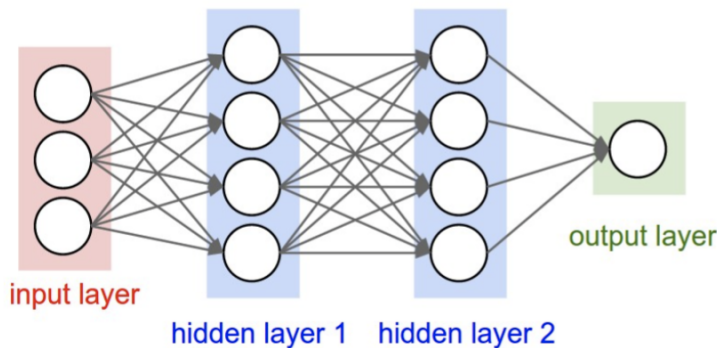


Figure 9.4: Deep Neural Net for Classifying Handwritten Numbers[]

Representing this is simply a bunch of single layer neural networks linked together where the output of 1 layer is the input of the second layer.



$$h_1 = f_1(xW_1 + b_1)$$

$$h_2 = f_2(h_1W_2 + b_2)$$

$$y = f_3(h_2W_3 + b_3)$$

Figure 9.5: Deep Learning[]

9.1.4 Activation Functions

If we didn't use non-linear activation functions, the output of the neural network would just be a linear function of the input. Such a network is basically just a Linear regression Model and no matter how many layer and nodes we add to such a network, it is equivalent to having a network with a single node.

Hence, we use Activation functions to make the network more powerful give it the ability to represent non-linear complex functional mappings between inputs and outputs.

Another important feature of an Activation function is that it should be differentiable. This is because we perform back-propagation to calculate the gradients of the Loss function with respect to the weights

and then accordingly optimize weights using gradient descend or any other Optimization technique to decrease the Loss Function.

The following are the most commonly used activation functions:

1. **Sigmoid Function:** The Sigmoid Function takes a real-valued number and maps it to a range between 0 and 1.

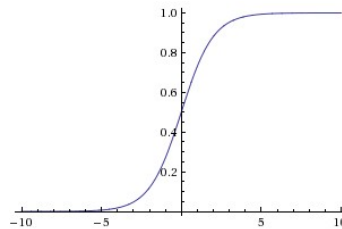


Figure 9.6: The sigmoid function[2]

$$f(x) = \frac{1}{1 + e^{-x}}$$

Sigmoid Activation functions have historically been used extensively but have now fallen out of favor due to the following drawbacks:

- (a) They are prone to saturate as the gradients are very small away from the center. When this occurs, almost no signal will flow through the neuron to its weights and recursively to its data.
 - (b) They are not zero-centered.
2. **Tanh function:** The Tanh function looks quite similar to a sigmoid function. It maps a real-valued number to a range between -1 and 1.

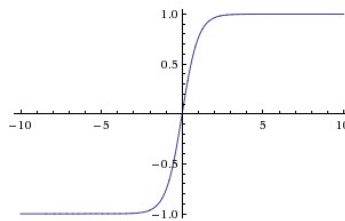


Figure 9.7: The tanh function[2]

$$f(x) = \tanh(x)$$

It saturates in the same way as a sigmoid does but it has the advantage of being zero-centered. Hence, in practice, it is always preferred to a sigmoid function.

3. **ReLU Function:** The Rectified Linear Unit is simply a threshold at zero and has become very popular in the last few years.

$$f(x) = \max(0, x)$$

It has the following advantages:

- (a) It greatly accelerates the convergence of stochastic gradient descent compared to tanh/sigmoid.

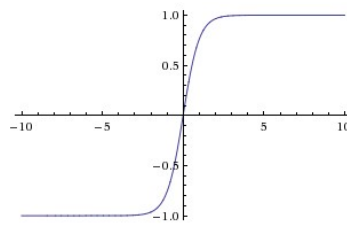


Figure 9.8: The ReLU function[2]

- (b) It is computationally cheaper to implement than tanh/sigmoid as it can be implemented by simply thresholding a matrix of activations at zero.

The main disadvantage of the ReLU is that ReLU units can irreversibly die during training since they can get knocked off the data manifold.

4. **Leaky ReLU Function:** The Leaky ReLU attempts to fix the problem of dying neurons by providing a small negative slope in the $x < 0$ region.

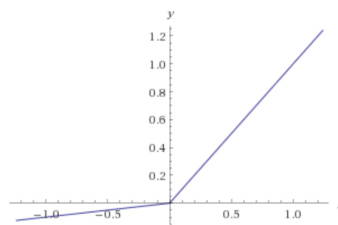


Figure 9.9: The Leaky ReLU function[2]

$$f(x) = \max(0.1x, x)$$

9.1.5 Training Neural Networks

We now need a way to teach these neural networks weights to make correct predictions. This is done by running the network on known data points and updating the weights of the neurons based on the correctness of the model in classifying the input.

Running this however on each input one at a time can take a long period of time for training. Instead these models are run on a sample batch of data. In particular the steps are the following.

1. **Sample a batch of data.** We can use a smaller batch of data at each iteration to train the model to decrease training time.
2. **Run the input it through the graph and compute the loss.** The loss function provides serves as a quality metric and provides us an idea on the error of our model. That is, how correct it is.
3. **Backpropagate** Activation functions are chosen to be differentiable to make this step easier. Backpropagation is done by taking the gradient with respect to the weights of each neuron in each layer.

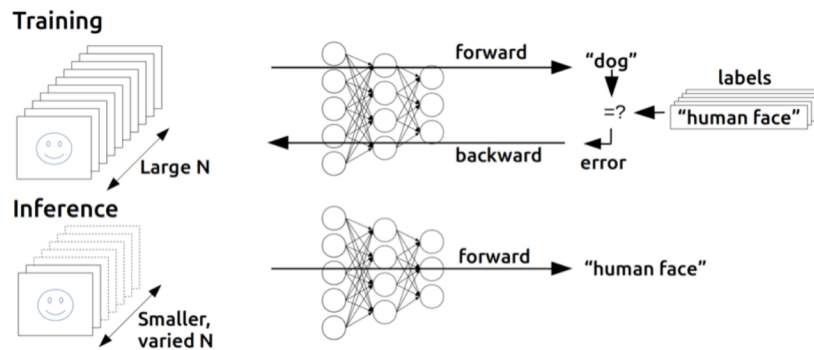


Figure 9.10: Training a neural net[]

4. **Update these parameters using SGD (Stochastic Gradient Descent)** Update each parameter using the gradient calculated in the previous step

5. **Repeat many times**

9.1.6 Overfitting

Creating complex models sometimes have the disadvantage of being extremely accurate only on the test data. That is, the model is able to classify its entire training data with high accuracy however is unable to generalize and fails with test / real data.

On the other hand too simple models are unable to learn the more subtle trends in the data.

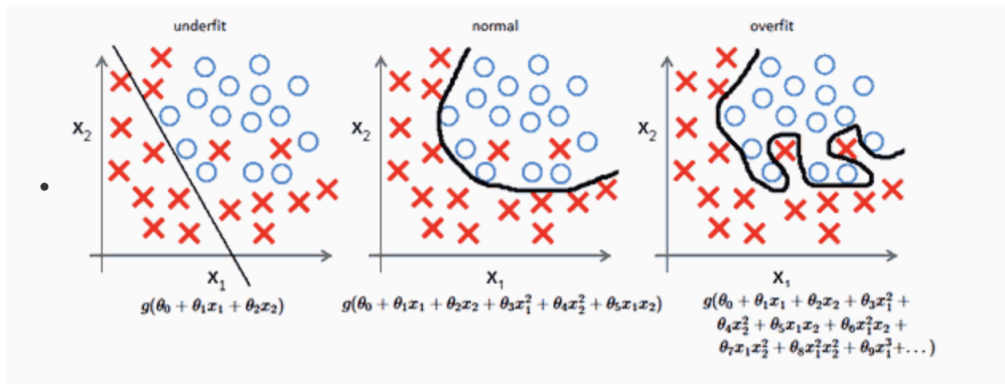


Figure 9.11: Overfitting[]

<http://mlwiki.org/index.php/Overfitting>

In the example above the left model was too simple and was unable to classify some of the values correctly. On the other the right model despite perfectly classifying on this training dataset is unlikely to perform as well in real-world situations. The middle offers a good classification.

To combat overfitting, regularization is usually introduced in models. Regularization is a technique that penalizes complex models or prevents them from being made. Some of these techniques include

Dropout Randomly select neurons that will not be activated at each pass during training.

L2 Regularization Penalize the square magnitude of all parameters which forces neuron size to be smaller.

Max Norm Constraints Limit the maximum weight of a neuron.

9.2 Convolutional Neural Networks

A simple Convolutional Neural Networks is a sequence of layers, and every layer of a Convolutional Neural Networks transforms one volume of activations to another through a differentiable function. There are four main types of layers to build ConvNet architectures: **Convolutional Layer**, **Nonlinearity Layer**, **Pooling Layer**, and **Fully-Connected Layer**. These layers are stacked to form a full ConvNet architecture. In this way, Convolutional Neural Networks transform the original image layer by layer from the original pixel values to the final class scores.

In this section, we mainly discuss the three types of layer in Convolutional Neural Networks, including Convolutional layer, Pooling layer and Fully-connected layer.

9.2.1 Convolutional Layer

Parameter Sharing Parameter sharing scheme is used in Convolutional Layers to control the number of parameters. It is hard to make perceptrons scalable when the input image size is large given the fact the number of parameters grow quickly with the input size. However, we can dramatically reduce the number of parameters by making one reasonable assumption. If we know the input is image data, we can assume some spatial symmetries. In other word, if one feature is useful to compute at some spatial position (x, y) , then it should also be useful to compute at a different position (x_2, y_2) , so same parameters could be used.

Convolution Details The convolution operation essentially performs dot products between the filters and local regions of the input. In other words, each element is computed by element-wise multiplying the local regions of the input (blue part in Figure 9.12) with the filter (red part in Figure 9.12), summing it up, and then offsetting the result by the bias. The output is the green part in Figure 9.12.

Noticing we give the demo in 2D format. In real cases, the input would be a 3D tensor while the depth of the output volume is a hyperparameter. It corresponds to the number of filters we would like to use, each learning to look for something different in the input. Take Figure 9.13 as an example, if we have $6 \times 5 \times 5$ filters, we will get 6 separate activation maps as output.

Spatial arrangement Besides depth, there are two hyperparameters to control the size of the output volume: the stride and zero-padding

- **stride** Stride specifies how we slide the filter. When the stride is 1 then we move the filters one pixel at a time. When the stride is 2 (or uncommonly 3 or more, though this is rare in practice) then the filters jump 2 pixels at a time as we slide them around.
- **zero-padding** In some cases, the size of filters and stride don't fit neatly and symmetrically across the input. For example, if the input size $H = W = 10$ while the filter size is $F = 3$ and it takes stride $S = 3$, it would be impossible to use stride $S = 2$, since

$$\begin{aligned}(H - F)/S + 1 &= (10 - 3)/3 + 1 = 3.33 \text{ (Not an integer)} \\ (W - F)/S + 1 &= (10 - 3)/3 + 1 = 3.33 \text{ (Not an integer)}\end{aligned}\tag{9.1}$$

Therefore, we need to use zero-padding. We add additional zero-paddings in the contour of the image. If we take zero-padding $P = 1$, we would have

$$\begin{aligned} (H - F + 2 * P) / S + 1 &= (10 - 3 + 2 * 1) / 3 + 1 = 4 \text{ (An integer)} \\ (W - F + 2 * P) / S + 1 &= (10 - 3 + 2 * 1) / 3 + 1 = 4 \text{ (An integer)} \end{aligned} \quad (9.2)$$

In general, setting zero padding to be $P = (F-1)/2$ when the stride is $S = 1$ ensures that the input volume and output volume will have the same size spatially. It is very common to use zero-padding in this way.

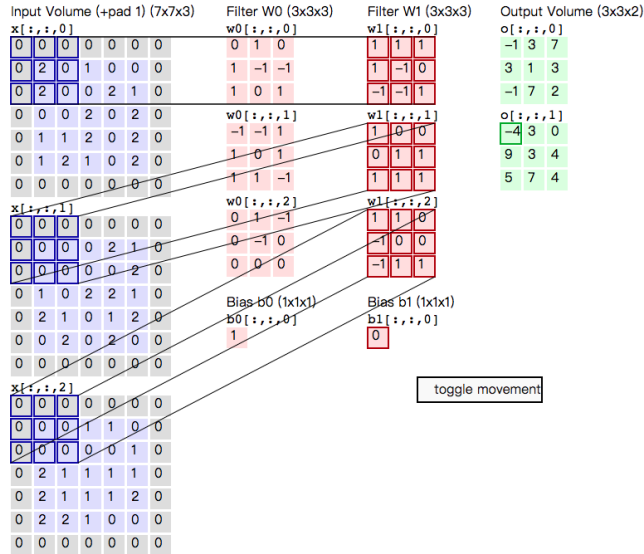


Figure 9.12: Convolution Demo[2]

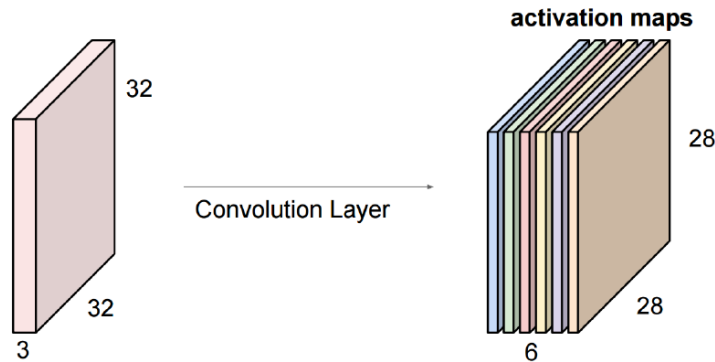


Figure 9.13: Number of filter and depth of output[2]

9.2.2 Pooling Layer

In practice, it is common to periodically insert a Pooling Layer in-between successive Conv layers in a ConvNet architecture. It is mainly used to progressively reduce the spatial size of the representation to

reduce the amount of parameters and computation load, and hence to also control overfitting. The Pooling Layer operates independently on every depth slice of the input and resizes it spatially, using the MAX operation if max pooling, AVG if average pooling, etc. The most common form is a pooling layer with filters of size 2×2 applied with a stride of 2 downsamples every depth slice in the input by 2 along both width and height, discarding 75% of the activations. In this case, every pooling operation would be taking a max or average over 4 numbers (little 2×2 region in some depth slice). The depth dimension remains unchanged. In a nutshell, *pooling layer downsamples the volume spatially, independently in each depth slice of the input volume.*

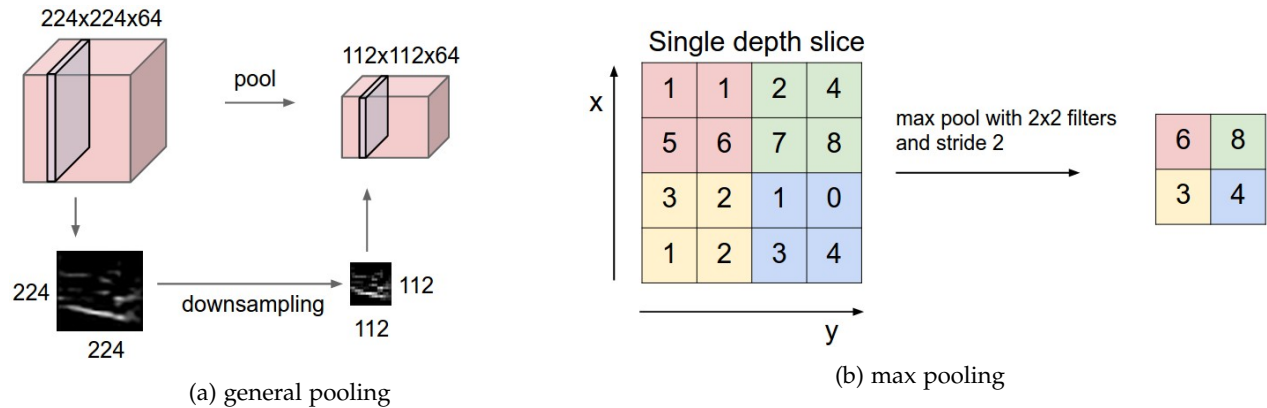


Figure 9.14: Pooling Layer[2]

The Fig.9.14a shows how pooling works generally. A 2×2 pooling filter with stride 2 is applied on the $224 \times 224 \times 64$ input volume, yielding the $112 \times 112 \times 64$ output volume. Note that the width and height of the input volume shrinks to half while the depth is preserved. The most common downsampling operation is MAX. The Fig.9.14b here shows max pooling with a stride of 2. That is, each max is taken over 4 numbers in a 2×2 square. For instance, the top-left red square squashes to a single value, 6, which is $\max\{1, 1, 5, 6\}$.

9.2.3 Fully-Connected Layer

The last few layers of a ConvNet architecture are typically Fully-Connected (FC) Layers. As seen in regular neural networks, FC Layers serve as a linear classifier for classification or regression. Their activations can be computed with a matrix multiplication followed by a bias offset, *i.e.* $f = Wx + b$. See the Neural Network Basics section of the note for more information.

9.2.4 Mainstream Model Architectures

Empowered with impressive ability of feature extraction and optimization, deep convolutional neural networks has become the cornerstone of data-driven visual recognition techniques nowadays. Since AlexNet's overwhelming success in ILSVRC 2012, a number of milestone model architectures have been proposed. The most common are:

- **AlexNet[3]:** Developed by Alex Krizhevsky, Ilya Sutskever and Geoff Hinton. The AlexNet significantly outperformed the second runner-up (top 5 error of 16% compared to runner-up with 26%

error) ILSVRC challenge in 2012. It is the first work that popularized Convolutional Neural Networks in Computer Vision.

- **VGGNet**[4]: Proposed by Karen Simonyan and Andrew Zisserman. The VGGNet was the runner-up in ILSVRC 2014. Its main contribution was in showing that the depth of the network is a critical component for good performance. Researchers and engineers tend to design deeper rather than shallower models ever since.
- **GoogLeNet**[5]: As its name shows, it is authored by Szegedy *et al.* from Google. GoogLeNet was the ILSVRC 2014 winner with main contribution in developing Inception Module that dramatically reduced the number of parameters in the network.
- **ResNet**[6]: Developed by Kaiming He *et al.* ResNet was the winner of ILSVRC 2015. It features special skip connections and a heavy use of batch normalization, which is now known as Residual Module. ResNets are currently by far state of the art Convolutional Neural Network models and are the default choice for using ConvNets in practice (as of January 2018).

There also exist more model architectures proposed for specific tasks, such as YOLO[7] and Faster R-CNN[8] in the field of object detection and localization.

References

- [1] Simon S. Haykin. *Neural Networks and Learning Machines*. 2009.
- [2] Andrej Karpathy. Stanford university cs231n course notes, May 2016. Accessed: 2018-02-07.
- [3] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran Associates, Inc., 2012.
- [4] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014.
- [5] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015.
- [7] Joseph Redmon, Santosh Kumar Divvala, Ross B. Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. *CoRR*, abs/1506.02640, 2015.
- [8] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 91–99. Curran Associates, Inc., 2015.

Contributors

Chi Zhang, Honghao Wei, Matthew Tan, Taylor Howell, Brian Jackson, Saifan Rafiq