

BASIC INTRODUCTION TO NEURAL NETWORKS

Muhammed Shifas PV



University of Crete, Dept of Computer Science
shifaspv@csd.uoc.gr

CS-HY578: Speech Signal Processing, 9 May 2022

ABOUT ME

Presently: Apple Inc., Cambridge, UK.

2017-2021: PhD in Speech Processing, University of Crete, Greece.

2014-2016: M.Tech in Signal Processing, NIT Calicut, India.

2009-2013: B.Tech in Communication Engineering, Calicut University, India.

OUTLINE

- 1 BASICS OF NEURAL NETWORKS
- 2 FULLY CONNECTED NEURAL NETWORK
- 3 CONVOLUTIONAL NEURAL NETWORK
- 4 RECURRENT NEURAL NETWORK
- 5 THANKS

OUTLINE

- 1 BASICS OF NEURAL NETWORKS
- 2 FULLY CONNECTED NEURAL NETWORK
- 3 CONVOLUTIONAL NEURAL NETWORK
- 4 RECURRENT NEURAL NETWORK
- 5 THANKS

MODELLING BIOLOGICAL NEURON

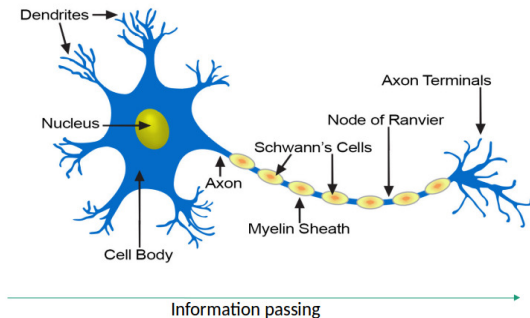


FIGURE: The biological structure of a neuron

MATHEMATICAL EQUIVALENT OF A NEURON

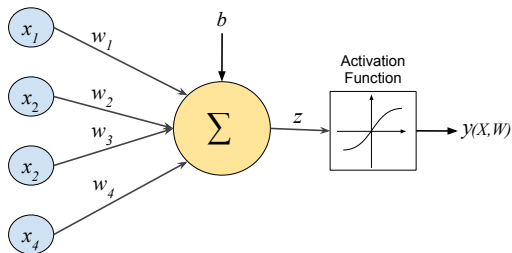


FIGURE: Mathematical equivalent of a neuron

$$y(\mathbf{x}^{(k)}) = \phi(z_i) = \phi(\mathbf{w}_i^T \mathbf{x}^{(k)} + b_i) = \phi\left(\sum_{j=1}^n w_{ij} x_j^{(k)} + b_i\right) \quad (1)$$

- Output is a function of the input (data) and the weights.
- **Training:** optimize the weights such that to reach to the desired output.

MATHEMATICAL EQUIVALENT OF A NEURON

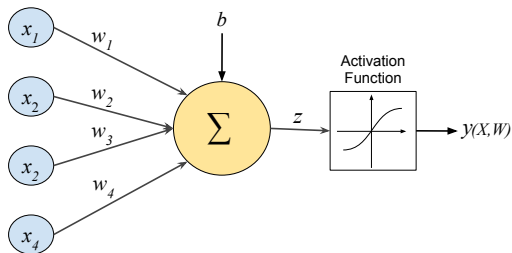


FIGURE: Mathematical equivalent of a neuron

$$y(\mathbf{x}^{(k)}) = \phi(z_i) = \phi(\mathbf{w}_i^T \mathbf{x}^{(k)} + b_i) = \phi\left(\sum_{j=1}^n w_{ij} x_j^{(k)} + b_i\right) \quad (1)$$

- Output is a function of the input (data) and the weights.
- **Training:** optimize the weights such that to reach to the desired output.

NEURAL NETWORK: NETWORK OF NEURONS

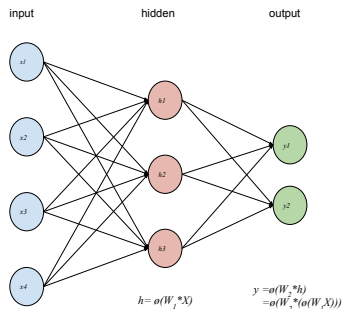


FIGURE: Neural network

HIGHLIGHT

- Intuition: any complex function can be approximated as a series of simple non-linear functions.

NEURAL NETWORK: NETWORK OF NEURONS

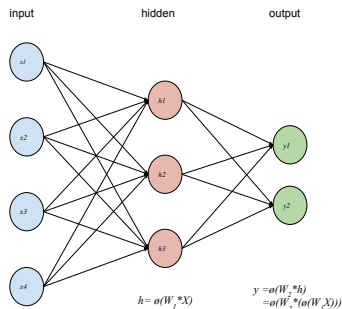


FIGURE: Neural network

HIGHLIGHT

- Intuition: any complex function can be approximated as a series of simple non-linear functions.

MATHAMATICAL INSIGHT

$$h = \begin{pmatrix} h_1 \\ h_2 \\ h_3 \end{pmatrix} = \Phi \begin{pmatrix} w_{10} & w_{11} & w_{12} & w_{13} \\ w_{20} & w_{21} & w_{22} & w_{23} \\ w_{30} & w_{31} & w_{32} & w_{33} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} \quad (2)$$

$$y = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \Phi \begin{pmatrix} w_{10} & w_{11} & w_{12} \\ w_{20} & w_{21} & w_{22} \end{pmatrix} \begin{pmatrix} h_1 \\ h_2 \\ h_3 \end{pmatrix} \quad (3)$$

$$y = \Phi \begin{pmatrix} w_{10} & w_{11} & w_{12} \\ w_{20} & w_{21} & w_{22} \end{pmatrix} \begin{pmatrix} w_{10} & w_{11} & w_{12} & w_{13} \\ w_{20} & w_{21} & w_{22} & w_{23} \\ w_{30} & w_{31} & w_{32} & w_{33} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} \quad (4)$$

NON-LINEAR ACTIVATION FUNCTIONS: Φ

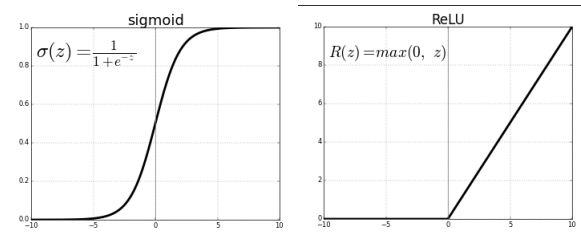


FIGURE: Commonly used activation functions

WHICH FUNCTION ARE WE LOOKING FOR ?

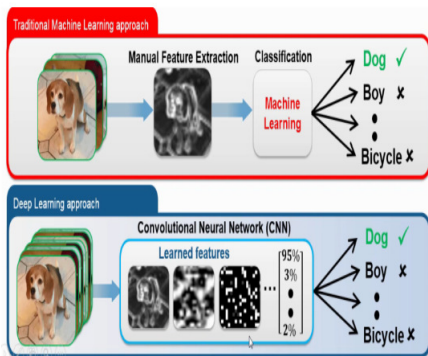
- **Assumption:** there is a statistics, hidden in our data
- The statistics to model depends on the task: Speaker identification, emotion detection, enhancement

MODEL INPUT DIMENSION

- Traditional approach: manually extracts the features and feed into the network.
- Advance models: feed the raw samples as it is into the models, letting the network to extracts the task relevant features.

TRAINING THE MODEL: TEACHING

- **Training:** Teaching the model by exploring to the already know data pair
 - Supervised: data (input, output)
 - Unsupervised: data (input,)



TRAINING THE MODEL: WEIGHT TUNING

$$\hat{y} = \begin{pmatrix} \hat{y}_1 \\ \hat{y}_2 \end{pmatrix} = \Phi \begin{pmatrix} w_{10} & w_{11} & w_{12} \\ w_{20} & w_{21} & w_{22} \end{pmatrix} \begin{pmatrix} h_1 \\ h_2 \\ h_3 \end{pmatrix} \quad (5)$$

- Loss: the measure of deviation of network prediction \hat{y} from the true training set label $y = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}$
- Penalize each wrong prediction (tune the weight) so that it to be a good predictor at the end.

PARAMETER OPTIMIZATION: GRADIENT DESCENT

Given a set of training input and output pairs

$$\{\{\mathbf{x}_1, \mathbf{t}_1\} \dots, \{\mathbf{x}_n, \mathbf{t}_n\}\} \quad (6)$$

compute the loss function for each prediction

$$E = \frac{1}{2} \sum_{p=1}^n \sum_{i=1}^K (y_i(\mathbf{x}_p) - t_{pi})^2 \quad (7)$$

$$\frac{\partial E}{\partial w_i} = \left(\frac{\partial E}{\partial z} \right) \left(\frac{\partial z}{\partial w_i} \right) \quad (8)$$

with sigmoid as the activation function ϕ , $y(\mathbf{z}) = \frac{1}{1 + \exp(-z)}$

$$\frac{\partial E}{\partial w_i} = (y(\mathbf{x}) - t)y(\mathbf{x})(1 - y(\mathbf{x}))x_i \quad (9)$$

Walk on the direction to which gradient descend

$$w_{i+1} = w_i - \gamma \frac{\partial E}{\partial w_i} \quad (10)$$

THE LOGISTIC LOSS FUNCTION (CLASSIFIER)

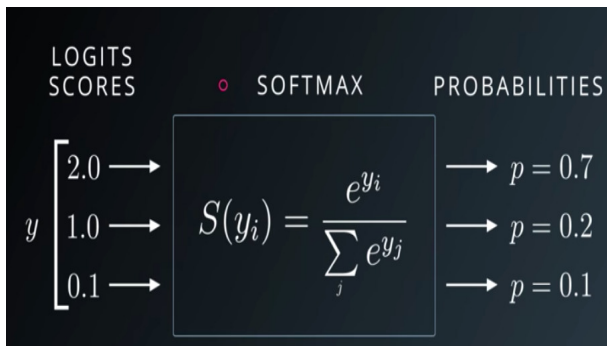


	Image #1	Image #2	Image #3
Dog	-0.39	-4.61	1.03
Cat	1.49	3.28	-2.37
Horse	4.21	1.46	-2.27

SVM loss: Minimize the objective

$$L(y, \hat{y}) = \sum_{i \neq c} \max(0, \hat{y}_i - \hat{y}_c + \Delta) \quad (11)$$

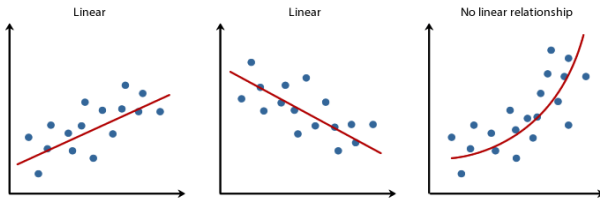
PROBABILISTIC LOSS FUNCTION



The cross entropy loss:

$$L(y, \hat{y}) = - \sum_i (y_i \log(\hat{p}_i) + (1 - y_i) \log(1 - \hat{p}_i)) \quad (12)$$

THE REGRESSION LOSS



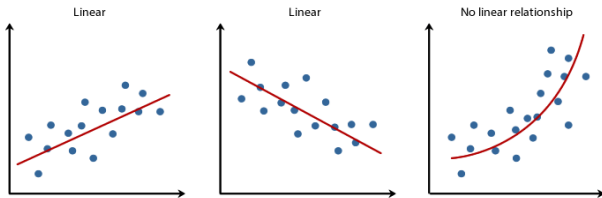
Copyright 2014. Laerd Statistics.

The target to be predicted is a continuous value function: eg. speech enhancement.

The final layer of regression model:

$$\hat{y} = \begin{pmatrix} \hat{y}_1 \end{pmatrix} = \Phi \begin{pmatrix} w_{10} & w_{11} & w_{12} \end{pmatrix} \begin{pmatrix} h_1 \\ h_2 \\ h_3 \end{pmatrix} \quad (13)$$

THE REGRESSION LOSS



Copyright 2014. Laerd Statistics.

The target to be predicted is a continuous value function: eg. speech enhancement.

The final layer of regression model:

$$\hat{y} = \begin{pmatrix} \hat{y}_1 \end{pmatrix} = \Phi \begin{pmatrix} w_{10} & w_{11} & w_{12} \end{pmatrix} \begin{pmatrix} h_1 \\ h_2 \\ h_3 \end{pmatrix} \quad (13)$$

REGRESSION LOSS

Mean square error:

$$MSE = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n} \quad (14)$$

Mean absolute error:

$$MAE = \frac{\sum_{i=1}^n |y_i - \hat{y}_i|}{n} \quad (15)$$

n = Total data points

OUTLINE

- 1 BASICS OF NEURAL NETWORKS
- 2 FULLY CONNECTED NEURAL NETWORK**
- 3 CONVOLUTIONAL NEURAL NETWORK
- 4 RECURRENT NEURAL NETWORK
- 5 THANKS

FULLY CONNECTED NEURAL NETWORK

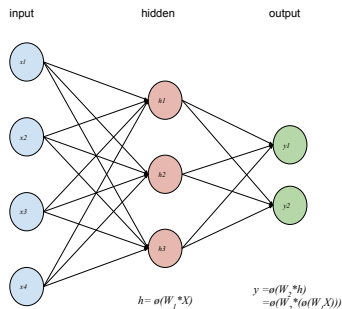


FIGURE: Fully Connected Network

NETWORK IDENTITY

- All nodes from the previous layer is connected to the next layer.

MATHEMATICAL INSIGHT

$$\hat{\mathbf{h}} = \begin{pmatrix} \hat{h}_1 \\ \hat{h}_2 \\ \hat{h}_3 \end{pmatrix} = \begin{pmatrix} w_{10} & w_{11} & w_{12} & w_{13} \\ w_{20} & w_{21} & w_{22} & w_{23} \\ w_{30} & w_{31} & w_{32} & w_{33} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} \quad (16)$$

$$\hat{\mathbf{y}} = \begin{pmatrix} \hat{y}_1 \\ \hat{y}_2 \end{pmatrix} = \begin{pmatrix} w_{10} & w_{11} & w_{12} \\ w_{20} & w_{21} & w_{22} \end{pmatrix} \begin{pmatrix} \hat{h}_1 \\ \hat{h}_2 \\ \hat{h}_3 \end{pmatrix} \quad (17)$$

The number of parameters are linear with input/ hidden layer size

OUTLINE

- 1 BASICS OF NEURAL NETWORKS
- 2 FULLY CONNECTED NEURAL NETWORK
- 3 CONVOLUTIONAL NEURAL NETWORK**
- 4 RECURRENT NEURAL NETWORK
- 5 THANKS

WHY DOES A NEW NETWORK?

The fully connected network has some draw backs:

- It always gives a merged representation of the previous input/hidden layer
- Failed to capture the local information in the input signal.
- The complexity of the model increases rapidly as we build deeper networks.

CONVOLUTIONAL NEURAL NETWORK

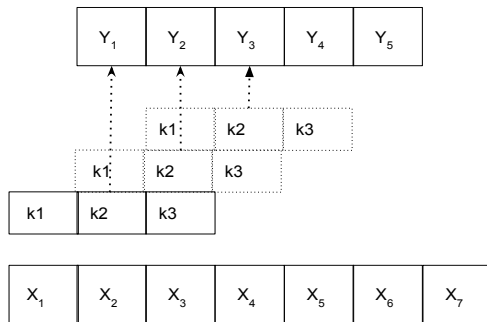


FIGURE: Convolution Network

$$Y[n] = (X * k)[n] = \sum_{i=1}^n X(n)k(n - m) \quad (18)$$

CONVOLUTION WITH MULTIPLE KERNELS

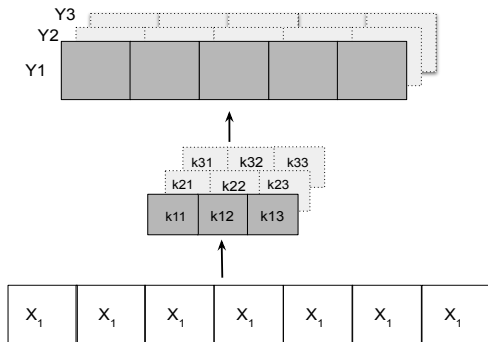
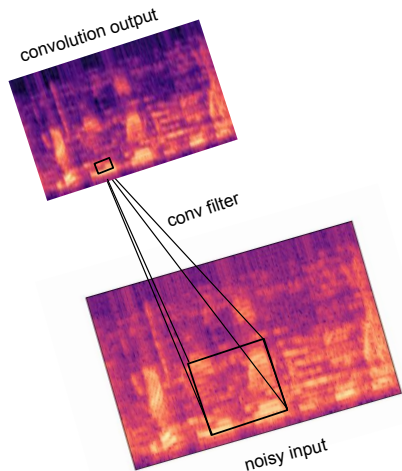


FIGURE: Convolution Network

$$Y_i[n] = (X * k_i)[n] = \sum_{m=1}^n X(m)k_i(n-m) \quad (19)$$

Number of kernels = number of channels

CONVOLUTION CHANNELS: 2D



- Number of parameters are independent of the input size.
- Network parameter is independent of the input dimension.
- The kernel size is customizable.

FIGURE: 2D Convolution

DILATED CONVOLUTION

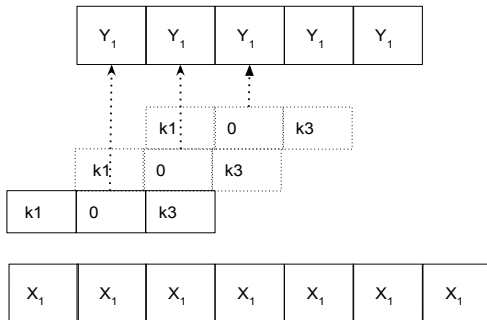


FIGURE: Dilated convolution Network

OUTLINE

- ① BASICS OF NEURAL NETWORKS
- ② FULLY CONNECTED NEURAL NETWORK
- ③ CONVOLUTIONAL NEURAL NETWORK
- ④ RECURRENT NEURAL NETWORK
- ⑤ THANKS

WHY RECURRENCY?

- speech articulations are a highly correlated over time.
- Estimation of the current input phoneme can tell something about the phonemes follows.

-eg. We will meet

- We must store and pass the information at the current instant to be consulted in the future predictions.
- Adding the Markove structure into the neural network

RECURRENCE IN FULLY CONNECTED NETWORK

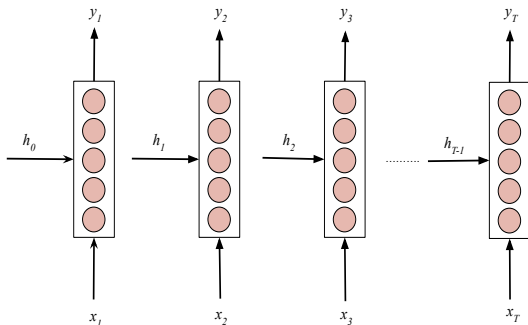


FIGURE: Fully connected Recurrent Network

FULLY CONNECTED LONG SHORT-TERM MEMORY (FC-LSTM)

$$\begin{aligned}
 i_t &= \Phi(W_{xi}X_t + W_{hi}h_t + W_{ci}oc_{t-1} + b_i) \\
 f_t &= \Phi(W_{xf}X_t + W_{hf}h_t + W_{cf}oc_{t-1} + b_f) \\
 o_t &= f_t oc_{t-1} + i_t otanh(W_{xc}X_t + W_{hc}h_{t-1} + b_0) \\
 y_t &= \Phi(W_{xo}X_t + W_{ho}h_t + W_{co}oc_{t-1} + b_y)
 \end{aligned} \tag{20}$$

OUTLINE

- 1 BASICS OF NEURAL NETWORKS
- 2 FULLY CONNECTED NEURAL NETWORK
- 3 CONVOLUTIONAL NEURAL NETWORK
- 4 RECURRENT NEURAL NETWORK
- 5 THANKS

ACKNOWLEDGMENTS



ENRICH EUROPEAN
TRAINING
NETWORK

This project has received funding from the EU's H2020 research and innovation programme under the MSCA GA 675324

Thank for your attention