

Training One Neural Layer

Deep Learning (ENGG*6600*07)

School of Engineering,
University of Guelph, ON, Canada

Course Instructor: Benyamin Ghogh
Fall 2023

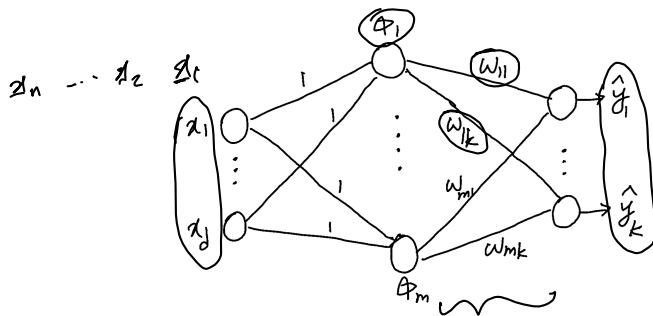
Introduction

- In the lecture of “Training one neuron”, we were introduced to the neural networks with only one layer and only one neuron.
- In this lecture, we are introduced to the neural networks with only one [learnable] layer but possibly multiple neurons.
- Two of these networks are Radial Basis Function (RBF) network and Self-Organizing Map (SOM).

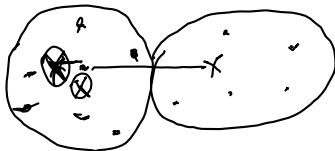
Radial Basis Function Network

Radial Basis Function Network

- A **Radial Basis Function (RBF)** network was first proposed in 1988 [1, 2].
- It has two layers but the first layer has fixed weights equal to one. The second layer has learnable weights.
- The first layer connects the data $\mathbf{x} \in \mathbb{R}^d$ to m basis kernel functions $\{\phi_i(\mathbf{x})\}_{i=1}^m$.
- The second layer connects the basis functions $\{\phi_i(\mathbf{x})\}_{i=1}^m$ to the output neurons $\{\hat{y}_j\}_{j=1}^k$.
- The weight w_{ij} denotes the weight connecting $\phi_i(\mathbf{x})$ to \hat{y}_j .



Radial Basis Function Network



- The basis functions can be various kernel functions such as:

Gaussian distribution:

$$\phi_i(\mathbf{x}) = e^{-\frac{\|\mathbf{x}_i - \mu_i\|_2^2}{2\sigma_i^2}}, \quad (1)$$

→ RBF function (kernel)

Logistic (sigmoid) function:

$$\phi_i(\mathbf{x}) = \frac{1}{1 + e^{-\frac{\|\mathbf{x}_i - \mu_i\|_2^2}{2\sigma_i^2}}}, \quad (2)$$

where $\mu_i \in \mathbb{R}^d$ and $\sigma_i^2 \in \mathbb{R}$ are the mean and variance for $\phi_i(\mathbf{x})$.

- At the first step, the means $\{\mu_i\}_{i=1}^m$ are found by applying a clustering method, such as K-means, on the training data with m clusters. The variances of clusters determine the variances $\{\sigma_i\}_{i=1}^m$.

$$\sigma_i^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_i)^2$$

Radial Basis Function Network

- RBF networks can be considered as an **additive model**. An additive model, first proposed in [3], maps data to a space with several basis functions and then tries to learn a weighted average of those bases.
- In RBF, the output is obtained as:

$$\hat{y}_j = \underbrace{w_{1j} \phi_1(\mathbf{x}) + \dots + w_{mj} \phi_m(\mathbf{x})}_{\text{additive model}} = \sum_{i=1}^m \underbrace{w_{ij} \phi_i(\mathbf{x})}_{\text{additive model}}. \quad (3)$$

- Consider n data points together in a matrix $\mathbf{X} \in \mathbb{R}^{d \times n}$. In matrix form:

$$\hat{\mathbf{Y}} = \underbrace{\mathbf{W}^T}_{k \times m} \underbrace{\boldsymbol{\Phi}}_{m \times n} \rightarrow k \times n \quad (4)$$

where:

$$\underbrace{\mathbb{R}^{k \times n}}_{\hat{\mathbf{Y}}} = \begin{bmatrix} \hat{y}_{11} & \dots & \hat{y}_{1n} \\ \vdots & \ddots & \vdots \\ \hat{y}_{k1} & \dots & \hat{y}_{kn} \end{bmatrix}, \quad \underbrace{\mathbb{R}^{m \times k}}_{\mathbf{W}} = \begin{bmatrix} w_{11} & \dots & w_{1k} \\ \vdots & \ddots & \vdots \\ w_{m1} & \dots & w_{mk} \end{bmatrix},$$

$$\underbrace{\mathbb{R}^{m \times n}}_{\boldsymbol{\Phi}} = \begin{bmatrix} \phi_1(x_1) & \dots & \phi_1(x_n) \\ \vdots & \ddots & \vdots \\ \phi_m(x_1) & \dots & \phi_m(x_n) \end{bmatrix}.$$

Radial Basis Function Network

$$\Phi^T W W^T \Phi \rightarrow \Phi \Phi^T W W^T \rightarrow W \Phi \Phi^T W$$

- Least squares error minimization between the label of data $\mathbf{Y} \in \mathbb{R}^{k \times n}$ and the output of network $\hat{\mathbf{Y}} \in \mathbb{R}^{k \times n}$:

$$\underset{\mathbf{W}}{\text{minimize}} \quad \|\mathbf{Y} - \hat{\mathbf{Y}}\|_F^2 = \|\mathbf{Y} - \mathbf{W}^T \Phi\|_F^2. \quad (5)$$

- Simplification of the cost function:

$$\begin{aligned} \star \|\mathbf{Y} - \mathbf{W}^T \Phi\|_F^2 &= \text{tr}((\mathbf{Y} - \mathbf{W}^T \Phi)^T (\mathbf{Y} - \mathbf{W}^T \Phi)) = \text{tr}((\mathbf{Y}^T - \Phi^T \mathbf{W})(\mathbf{Y} - \mathbf{W}^T \Phi)) \\ &= \text{tr}(\mathbf{Y}^T \mathbf{Y} - \mathbf{Y}^T \mathbf{W}^T \Phi - \Phi^T \mathbf{W} \mathbf{Y} + \Phi^T \mathbf{W} \mathbf{W}^T \Phi) \\ &\stackrel{(a)}{=} \text{tr}(\mathbf{Y}^T \mathbf{Y}) - \text{tr}(\mathbf{W}^T \Phi \mathbf{Y}^T) - \text{tr}(\mathbf{W} \mathbf{Y} \Phi^T) + \text{tr}(\mathbf{W}^T \Phi \Phi^T \mathbf{W}), \end{aligned}$$

where (a) is because of the linearity and cyclic property of the trace operator.

- Solving this optimization problem:

$$\begin{aligned} \frac{\partial}{\partial \mathbf{W}} \|\mathbf{Y} - \mathbf{W}^T \Phi\|_F^2 &= -\Phi \mathbf{Y}^T - \Phi \mathbf{Y}^T + 2 \Phi \Phi^T \mathbf{W} \stackrel{\text{set}}{=} \mathbf{0} \Rightarrow \Phi \Phi^T \mathbf{W} = \Phi \mathbf{Y}^T \\ \Rightarrow \mathbf{W} &= (\Phi \Phi^T)^{-1} \Phi \mathbf{Y}^T. \end{aligned} \quad (6)$$

$$\text{tr}(ABC) = \text{tr}(CAB) = \text{tr}(B \overset{CA}{\cancel{AC}}) = \text{tr}(ABC)$$

$$\text{tr}(\omega^T (\Phi Y^T)) = \text{tr}(\omega^T A)$$

$$\frac{\partial \text{tr}(\omega^T \Phi Y^T)}{\partial \omega} = \begin{cases} A \rightarrow 3 \times 5 \checkmark \\ A^T \quad \times \end{cases}$$

\downarrow
 3×5

$$\left\{ \begin{array}{l} \Phi Y^T \\ Y^T \Phi \quad \times \\ \Phi^T Y \quad \times \\ \vdots \end{array} \right.$$

Radial Basis Function Network

$$[] \quad \underbrace{AX = B} \rightarrow X = A^{-1}B \rightarrow \underbrace{A^T A X = A^T B}_{\star} \rightarrow X = \underbrace{(A^T A)^{-1} A^T}_{A^\dagger} B$$

- If $\Phi\Phi^T$ is a singular matrix, the pseudo-inverse can be used:

$$W = (\Phi\Phi^T)^\dagger \Phi Y^T, \quad (7)$$

where † denotes the pseudo-inverse of matrix.

- The output of the RBF network, for either the training or test data, is:

$$\hat{Y} = W^T \Phi = ((\Phi\Phi^T)^{-1} \Phi Y^T)^T \Phi = Y \Phi^T (\Phi\Phi^T)^{-1} \Phi = Y \Phi^T (\Phi\Phi^T)^{-1} \Phi. \quad (8)$$

$$[] \quad A A^T \quad A^\dagger = A^T (A A^T)^{-1} \quad X = \underbrace{A^T (A A^T)^{-1}}_{A^\dagger} B$$

$$W = (\Phi \Phi^T)^{-1} \Phi Y^T$$

$$\Sigma (\Phi \Phi^T + \lambda I)^{-1} \Phi Y^T$$

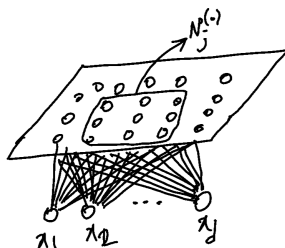
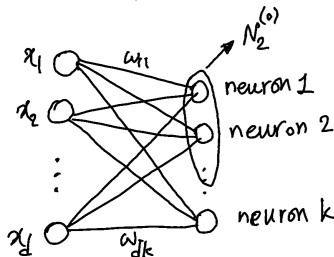
$$(AB)^T = B^T A^T$$

$$(AB)^{-1} = B^{-1} A^{-1}$$

Self-Organizing Map

Self-Organizing Map

- Self-Organizing Map (SOM) is a neural network with one layer. It is used for unsupervised clustering, where the name “self-organizing” comes from.
- It was proposed by Teuvo Kohonen in 1982 [4]; therefore, it is also called the Kohonen network [5].
- It is one layer connecting $\mathbf{x} = [x_1, \dots, x_d]^T \in \mathbb{R}^d$ to k neurons. Let w_{ij} denote the weight connecting x_i to the j -th neuron. Each neuron represents a cluster. SOM trains the weights to cluster the input data \mathbf{x} to one of the k clusters.
- The neurons can be put in 1D or 2D structure.
- Every neuron has a neighborhood around it in the 1D or 2D structure of neurons. This neighborhood is decreased gradually during the training phase. Let $\mathcal{N}_j^{(\tau)}$ denote the neighborhood of the j -th neuron at iteration τ .



Self-Organizing Map

- Step 1 of training: Initialize all weights to small random values. Set all neighborhoods $\{\mathcal{N}_j^{(0)}\}_{j=1}^k$ to half of the neuron structure grid. Set the initial learning rate $\eta^{(0)}$ to a number in range (0, 1].
- Step 2 of training: Select one the input data points, $\mathbf{x} = [x_1, \dots, x_d]^T$, and feed it to the network. Select the winning neuron z by:

$$\longrightarrow \textcircled{z} := \arg \min_j \sum_{i=1}^d \underbrace{\|x_i - w_{ij}\|_2}_{\textcircled{w_{ij}}} \cdot (x_i - w_{ij})^2 \quad (9)$$

- Step 3 of training: Update the weights of the neurons in the neighborhood of the winning neuron z :

$$w_{ij}^{(\tau+1)} := \begin{cases} \textcircled{w_{ij}^{(\tau)} + \eta^{(\tau)}(x_i - w_{ij}^{(\tau)})} & \text{if } j \in \mathcal{N}_z^{(\tau)} \\ \textcircled{w_{ij}^{(\tau)}} & \text{Otherwise,} \end{cases} \quad (10)$$

for all $i \in \{1, \dots, d\}$.

- Step 4 of training: Decrease the learning rate and the neighborhoods:

$$\star \eta^{(\tau+1)} := \eta^{(0)} \left(1 - \frac{\tau}{t}\right), \quad (11)$$

$$\mathcal{N}_j^{(\tau+1)} := \mathcal{N}_j^{(\tau+1)}/2, \quad \forall j \in \{1, \dots, k\}, \quad (12)$$

where t denotes the total number of training iterations.

- Step 5 of training: Increase τ and go to step 2.

Acknowledgment

- Some slides of this slide deck were inspired by teachings of Prof. Ali Ghodsi (at University of Waterloo, Department of Statistics), Prof. Fakhri Karray (at University of Waterloo, Department of Electrical and Computer Engineering), and Prof. Saeed Bagheri Shouraki (at Sharif University of Technology, Department of Electrical Engineering).

References

- [1] D. S. Broomhead and D. Lowe, "Radial basis functions, multi-variable functional interpolation and adaptive networks," tech. rep., Royal Signals and Radar Establishment Malvern (United Kingdom), 1988.
- [2] D. Broomhead and D. Lowe, "Multivariable functional interpolation and adaptive networks, complex systems, vol. 2," *Complex Systems*, vol. 2, pp. 321–355, 1988.
- [3] J. H. Friedman and W. Stuetzle, "Projection pursuit regression," *Journal of the American statistical Association*, vol. 76, no. 376, pp. 817–823, 1981.
- [4] T. Kohonen, "Self-organized formation of topologically correct feature maps," *Biological cybernetics*, vol. 43, no. 1, pp. 59–69, 1982.
- [5] T. Kohonen and T. Honkela, "Kohonen network," *Scholarpedia*, vol. 2, no. 1, p. 1568, 2007.