

Lecture 9: Neural Network Revisit

*Lecturer: Bo Dai**Scribes: Shuojiang Liu, Donaven Lobo*

Note: *LaTeX template courtesy of UC Berkeley EECS Department.*

Disclaimer: *These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor.*

9.1 Recap

The present lesson marks the conclusion of Module 1 of the course. The module 1 is comprised of three principal topics:

- Convex Optimization
- Density Parametrization
- Sampling Strategies

9.2 New Content

In today's lecture, we will demonstrate that a neural network is, in essence, just a function approximator or function parametrization. The manner in which it is learned is contingent upon the specific learning problem at hand.

For instance, polynomial is also a simple neural network.

9.2.1 Linear Model

The basic function of linear model is:

$$f_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x} \quad (9.1)$$

where we have data pairs $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$, where $\mathbf{x} \in \mathcal{X}$, $\mathbf{y} \in \mathcal{Y}$.

If $\mathcal{X} = \mathbb{R}^D$, $\mathcal{Y} = \mathbb{R}$, where D is the dimension of the number of input features, the model is **regression**.

For regression, we use quadratic loss for the loss function. The empirical risk associated with the use of quadratic loss is equivalent to the mean squared error (MSE):

$$loss(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \|y_i - f_{\boldsymbol{\theta}}(\mathbf{x}_i)\|^2 \quad (9.2)$$

From a different perspective, we can model the uncertainty in the prediction as Gaussian:

$$\mathcal{N}(y|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(y-\mu)^2}, \mu = f_{\boldsymbol{\theta}}(\mathbf{x}_n) \quad (9.3)$$

We can get the conditional probability distribution:

$$p(y_n|\mathbf{x}_n; \boldsymbol{\theta}) = \mathcal{N}(y_n|f_{\boldsymbol{\theta}}(\mathbf{x}_n), \sigma^2) \quad (9.4)$$

Similarly, if $\mathcal{X} = \mathbb{R}^D$, $\mathcal{Y} = \{1, 2, \dots, C\}$, where D is the dimension of the number of input features and C is the number of classes, the model is **classification**.

From the frequentist's point of view, the negative log likelihood of the training set is defined as:

$$loss(\boldsymbol{\theta}) = -\frac{1}{N} \sum_{i=1}^N \log p(y_i|f_{\boldsymbol{\theta}}(\mathbf{x}_i)) \quad (9.5)$$

By minimizing this quantity, we can compute the maximum likelihood estimate (MLE):

$$\hat{\boldsymbol{\theta}}_{MLE} = \arg \min_{\boldsymbol{\theta}} loss(\boldsymbol{\theta}) \quad (9.6)$$

9.2.2 Basic Neural Network Modeling

A neural network $f(\mathbf{x}; \boldsymbol{\theta})$, can be modeled as:

$$f(\mathbf{x}; \boldsymbol{\theta}) = f_M(f_{M-1}(\dots(f_1(\mathbf{x})\dots))) \quad (9.7)$$

where:

$$f_l(\mathbf{x}) = f(\mathbf{x}; \boldsymbol{\theta}_l) = \varphi_l(\mathbf{W}_l \mathbf{x} + \mathbf{b}_l) \quad (9.8)$$

The function φ_l is the activation function such as Sigmoid: $\sigma(\mathbf{z}) = \frac{1}{1+e^{-\mathbf{z}}}$, and ReLU: $\text{ReLU}(\mathbf{z}) = \max(\mathbf{z}, \mathbf{0})$.

A multilayer perceptron is a type of feedforward neural network comprising fully connected neurons with a nonlinear activation function, organized in at least three layers.

9.2.3 Kernel Methods

Kernel method is a kind of non-parametric method that involves using linear classifiers to solve nonlinear problems.

Kernel methods rely on mapping input data into a higher-dimensional space. Let $\mathbf{x} \in \mathbb{R}^n$ be an input feature vector. A feature mapping function ϕ transforms \mathbf{x} into a higher-dimensional space:

$$\phi: \mathbf{x} \mapsto \phi(\mathbf{x}) \in \mathbb{R}^m, m \gg n \quad (9.9)$$

In this space, data points that are not linearly separable in the original space can become linearly separable. Instead of computing the mapping $\phi(\mathbf{x})$ explicitly, kernel methods compute the inner product between two transformed points directly using a kernel function \mathcal{K} :

$$\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle \quad (9.10)$$

In SVMs, the decision function for a classification problem is expressed as:

$$f_{\boldsymbol{\theta}}(\mathbf{x}) = \sum_{i=1}^N \alpha_i y_i \mathcal{K}(\mathbf{x}_i, \mathbf{x}) + b \quad (9.11)$$

The objective is to find α_i that maximize the margin between classes in the transformed space.

Deep neural networks can be conceptualized as a form of learned feature mapping, wherein the feature space is learned from data during the training phase. In contrast, kernel methods employ predefined functions (e.g., Gaussian, polynomial) to transform the data into a higher-dimensional space.

9.2.4 Special Structred Neural Networks (Mentioned in the Slides)

9.2.4.1 Convolutional Neural Network (CNN)

CNNs are designed for grid-like data such as images. The key operations include convolution and pooling.

Given an input $\mathbf{X} \in \mathbb{R}^{H \times W \times C}$ and filter $\mathbf{W} \in \mathbb{R}^{k_H \times k_W \times C}$, the convolution output \mathbf{Y} is:

$$\mathbf{Y}(i, j) = \sum_{m=1}^{k_H} \sum_{n=1}^{k_W} \sum_{c=1}^C \mathbf{X}(i+m, j+n, c) \cdot \mathbf{W}(m, n, c) \quad (9.12)$$

where:

- \mathbf{X} : Input tensor (image)
- \mathbf{W} : Filter weights
- \mathbf{Y} : Feature map

Typically, ReLU is used as activation function:

$$f(x) = \max(0, x) \quad (9.13)$$

Pooling (max pooling) reduces the spatial dimensions by selecting the maximum value in a window.

$$\mathbf{P}(i, j) = \max_{m,n} \mathbf{Y}(i+m, j+n) \quad (9.14)$$

9.2.4.2 Residual Neural Network (ResNet)

ResNet introduces shortcut connections to solve the problem of vanishing gradients in deep networks. The key idea is to add a shortcut (identity) connection:

$$\mathbf{y}_l = f(\mathbf{x}_l, \mathbf{W}_l) + \mathbf{x}_l \quad (9.15)$$

where:

- \mathbf{x}_l : Input to the layer
- \mathbf{y}_l : Output from the layer
- $f(\mathbf{x}_l, \mathbf{W}_l)$: Transformation (convolution, activation, etc.) with learnable weights \mathbf{W}_l

9.2.4.3 Recurrent Neural Network (RNN)

RNNs handle sequential data by maintaining a hidden state that evolves over time steps. For an input sequence \mathbf{x}_t at time step t , the hidden state \mathbf{h}_t is updated as:

$$\mathbf{h}_t = f(\mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{W}_x \mathbf{x}_t + \mathbf{b}_h) \quad (9.16)$$

where:

- \mathbf{x}_t : Input at time t
- \mathbf{h}_t : Hidden state at time t
- $\mathbf{W}_h, \mathbf{W}_x$: Weight matrices
- \mathbf{b}_h : Bias term
- $f(\cdot)$: Non-linear activation (e.g., tanh)

The output \mathbf{y}_t at time t is typically given by:

$$\mathbf{y}_t = \varphi(\mathbf{W}_y \mathbf{h}_t + \mathbf{b}_y) \quad (9.17)$$

where $\varphi(\cdot)$ is an activation function (e.g., softmax for classification).

9.2.4.4 Self-Attention (Building Blocks for Transformer)

Given queries $\mathbf{Q} \in \mathbb{R}^{n \times d_k}$, keys $\mathbf{K} \in \mathbb{R}^{n \times d_k}$, and values $\mathbf{V} \in \mathbb{R}^{n \times d_v}$, the attention output \mathbf{Z} is:

$$\mathbf{Z} = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}} \right) \mathbf{V} \quad (9.18)$$

where:

- $\mathbf{Q}, \mathbf{K}, \mathbf{V}$: Query, key, and value matrices
- d_k : Dimensionality of queries/keys (scaling factor)

9.2.4.5 Transformer

In transformer, multiple attention heads are used to capture different aspects of the input:

$$\text{MultiHead}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Concat}(\mathbf{Z}_1, \dots, \mathbf{Z}_h) \mathbf{W}_O \quad (9.19)$$

where \mathbf{Z}_i is the output of the i -th attention head, and \mathbf{W}_O is a learned projection matrix.

After self-attention, a feed-forward network is applied to each position independently:

$$\text{FFN}(\mathbf{X}) = \max(0, \mathbf{X} \mathbf{W}_1 + \mathbf{b}_1) \mathbf{W}_2 + \mathbf{b}_2 \quad (9.20)$$

where:

- $\mathbf{W}_1, \mathbf{W}_2$: Weight matrices
- $\mathbf{b}_1, \mathbf{b}_2$: Bias terms

The Transformer combines multi-head self-attention and FFN with residual connections and layer normalization:

$$\mathbf{Y} = \text{LayerNorm}(\mathbf{X} + \text{MultiHead}(\mathbf{X}, \mathbf{X}, \mathbf{X})) \quad (9.21)$$

$$\mathbf{Z} = \text{LayerNorm}(\mathbf{Y} + \text{FFN}(\mathbf{Y})) \quad (9.22)$$

where:

- \mathbf{X} : Input to the layer
- \mathbf{Y}, \mathbf{Z} : Intermediate outputs

Layer normalization is applied after both the self-attention and feed-forward sub-layers:

$$\text{LayerNorm}(\mathbf{X}) = \frac{\mathbf{X} - \mu}{\sigma + \epsilon} \cdot \gamma + \beta \quad (9.23)$$

where:

- μ : Mean of the elements in \mathbf{X}
- σ : Standard deviation of elements in \mathbf{X}
- γ and β : Learnable scaling and shifting parameters

- ϵ : Small constant to avoid division by zero

Positional encoding is used to inject sequence order:

$$\text{PE}(i, 2k) = \sin\left(\frac{i}{10000^{2k/d_{\text{model}}}}\right) \quad (9.24)$$

$$\text{PE}(i, 2k + 1) = \cos\left(\frac{i}{10000^{2k/d_{\text{model}}}}\right) \quad (9.25)$$

where:

- i : Position index
- k : Dimension index
- d_{model} : Dimensionality of the model