





HW4 TBPosted

Lecture Notes III - Neural Networks

Marina Meilă mmp@stat.washington.edu

> Department of Statistics University of Washington

October, 2023



Multi-layer neural networks

A zoo of multilayer networks

Reading HTF Ch.: 11.3 Neural networks, Murphy Ch.: (16.5 neural nets), Bach Ch.: –, Deep Learning Book (Goodfellow, Bengio, Courville) 6.1-4, ResNet 7.6, ConvNet 9., Autoencoders 14.1, Dive Into Deep Learning 4.1-4.3.

Two-layer Neural Networks

► The activation function (a term borrowed from neuroscience) is any continuous, bounded and strictly increasing function on ℝ. Almost universally, the activation function is the logistic (or sigmoid)

$$\phi(u) = \frac{1}{1 + e^{-u}}$$
(1)

because of its nice additional computational and statistical properties.

► We build a two-layer neural network in the following way: Inputs Bottom layer¹ $z_j = \phi(w_j^T x)$ $j = 1: m, w_j \in \mathbb{R}^d \longrightarrow [ayer^4 1 (hidden)]$ Top layer $f = \#\beta^T z$ $\beta \in \mathbb{R}^m \longrightarrow [ayer^4 2 (output)]$ Output $f \in TP$ In other words, the neural network implements the function $w_j = 1: \frac{1}{2} = \frac{1}{2} \frac{\beta \in \mathbb{R}^m}{\beta_j z_j} = \sum_{j=1}^m \beta_j \phi(\sum_{k=1}^d w_{kj} x_k) \in (-\infty, \infty)$ (2)

Note that this is just a linear combination of logistic functions.

$$X \in \mathbb{R}^d$$
 $\mathcal{Z}_{im}^= \varphi(W_{im}^T X) = lopistic representation
 $W \in \mathbb{R}^d$ $\mathcal{B} \in \mathbb{R}^m$$

¹In neural net terminology, each variable z_j is a unit, the bottom layer is hidden, while top one is visible, and the units in this layer are called hidden/visible units as well. Sometimes the inputs are called input units; imagine neurons or individual circuits in place of each x, y, z variable.

Variables = units
hayors

$$x_{j} = \varphi(W_{j}^{T}x) \leftarrow logisfic
 $y = \beta^{T} 2$
parameters =
= weights $(W_{1}\beta)$
hidden units = represent
features
 $W = \begin{bmatrix} W_{j}^{T} \end{bmatrix} \in \mathbb{R}^{m \times d}$
 $z \in \mathbb{R}^{m}$ such that $W = \beta^{T}z$
 $with = \frac{1}{2} \sum_{units} \sum$$$

Output layer options

• linear layer as in (2) $f = \sum_{j} \beta_{j} z_{j} = \beta^{T} z \in \mathbb{R}$ regretion

▶ logistic layer: in classification $f(x) \in [0,1]$ is interpreted as the probability of the + class.

$$f(x) = \phi\left(\sum_{j=1}^{m} \beta_j z_j\right) = \phi\left(\sum_{j=1}^{m} \beta_j \phi(\sum_k w_{kj} x_k)\right) \in \left(0, 1, 1, \dots, 1, 2, \dots, 2, 2, \dots, 2, 2, \dots, 2, 2, \dots, 2, \dots,$$

softmax layer in multiway classification

The softmax function $\phi(z) : \mathbb{R}^r \to (0,1)^r$

$\phi_k(u)$	=	$\frac{e^{u_k}}{\sum_{j=1}^m e^{u_j}}$	-
	6	= 1:r	

Properties

$$\sum_{j=1}^{m} \phi_{j}(u) = 1 \text{ for all } u$$

$$for u_{k} \gg u_{j}, j \neq k \phi_{k}(u) \rightarrow 1.$$

$$derivatives \frac{\partial \phi_{j}}{\partial u_{k}} = \phi_{k} \delta_{jk} - \phi_{j} \phi_{k}$$

$$figmoid \qquad \phi'(u) = \phi(u) \qquad \phi^{2}(u)$$

$$u \in \mathbb{R}$$

Softmax function
$$\varphi \leftarrow ARJSIVELY$$

 $u \in \mathbb{R}^{r}$
 $\varphi : \mathbb{R}^{r} \rightarrow (0,1)^{r}$
 $\varphi_{k} = \underbrace{\mathbb{Q}^{u_{k}}}_{k'=1}$, $k = 1:n$
 $\varphi_{k}(2(x)) = \mathbb{P}[\gamma = k | x]$
 $confidence in $g = k$
 $g_{k}(2(x)) = 0$
 $g_$$

Generalized Linear Models (GLM)

A GLM is a regression where the "noise" distribution is in the exponential fami ly.

 \triangleright $y \in \mathbb{R}, y \sim P_{\theta}$ with

$$P_{\theta}(y) = e^{\theta y - \ln \psi(\theta)}$$
⁽⁵⁾

• the parameter θ is a linear function of $x \in \mathbb{R}^d$

$$\theta = \beta^T x \tag{6}$$

• We denote $E_{\theta}[y] = \mu$. The function $g(\mu) = \theta$ that relates the mean parameter to the natural parameter is called the **link function**.

The log-likelihood (w.r.t. β) is

$$I(\beta) = \ln P_{\theta}(y|x) = \theta y - \ln \psi(\theta) \quad \text{where } \theta = \beta' x \tag{7}$$

and the gradient w.r.t. β is therefore

$$\nabla_{\beta}I = \nabla_{\theta}I\nabla_{\beta}(\beta^{T}x) = (y-\mu)x$$
(8)

This simple expression for the gradient is the generalization of the gradient expression you obtained for the two layer neural network in the homework. [Exercise: This means that the sigmoid function is the *inverse link function* defined above. Find what is the link function that corresponds to the neural network.]



Marina Meila: Lecture I







$$\begin{aligned}
\varphi &= \frac{1}{(u)} \quad 1 + e^{-wu} \\
& W \in \mathbb{R} \\
& W > 0 \\
& W &= 0 \quad \Rightarrow \quad \varphi = \frac{1}{2} \\
& W &= 1000 \\
& U &= +\epsilon \quad e^{-u} = 0 \quad \varphi = 1 \\
& U &= 1000 \\
& U &= -\epsilon \quad e^{-u} = 0 \quad \varphi = 1 \\
& U &= 0 \quad \varphi = 1 \\
& U &= 0 \quad \varphi = 1 \\
& U &= 0 \quad \varphi = 1 \\
& U &= 0 \quad \varphi = 1 \\
& U &= 0 \quad \varphi = 1 \\
& U &= 0 \quad \varphi = 1 \\
& U &= 0 \quad \varphi = 1 \\
& U &= 0 \quad \varphi = 1 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
& U &= 0 \quad \varphi = 0 \\
&$$

Multi-layer/Deep neural networks

The construction can be generalized recursively to arbitrary numbers of layers. Each layer is a linear combination of the outputs from a previous layer (a multivariate operation), followed by a non-linear transformation via the logistic function ϕ . Let $x \equiv x^{(0)}, y \equiv x^{(L)}, n_0 = n, n_L = 1$ and define the recursion:

$$x_{j}^{(l)} = \phi\left((w_{j}^{(l)})^{T} x^{(l-l)}\right), \text{ for } j = 1: n_{l}$$
(9)

The vector variable $x^{(l)} \in \mathbb{R}^{n_l}$ is the ouput of layer *l* of the network. As before, the sigmoid of the last layer may be omitted.

