

MSCI 546: Advanced Machine Learning

Sirisha Rambhatla

University of Waterloo

Lecture
Feed-Forward Neural Networks

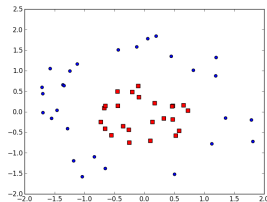
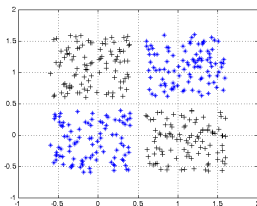
Outline

- 1 Neural Networks: Motivation
- 2 Neural Networks: Notation
- 3 Neural Networks: Learning and Backpropagation

Outline

- 1 Neural Networks: Motivation
- 2 Neural Networks: Notation
- 3 Neural Networks: Learning and Backpropagation

Linear models are not always adequate



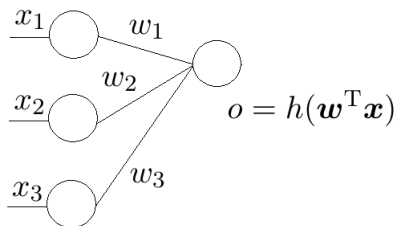
We can use a nonlinear mapping:

$$\phi(x) : x \in \mathbb{R}^D \rightarrow z \in \mathbb{R}^M$$

But what kind of nonlinear mapping ϕ should be used? Can we actually learn this nonlinear mapping?

THE most popular nonlinear models nowadays: **neural nets**

Linear model as a one-layer neural net



$$h(a) = a \text{ for linear model}$$

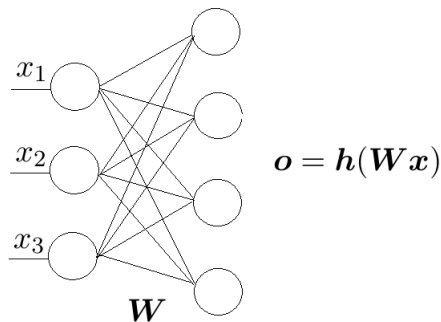
To create non-linearity, can use

- Rectified Linear Unit (**ReLU**): $h(a) = \max\{0, a\}$
- sigmoid function: $h(a) = \frac{1}{1+e^{-a}}$
- TanH: $h(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$
- many more

Outline

- 1 Neural Networks: Motivation
- 2 Neural Networks: Notation**
- 3 Neural Networks: Learning and Backpropagation

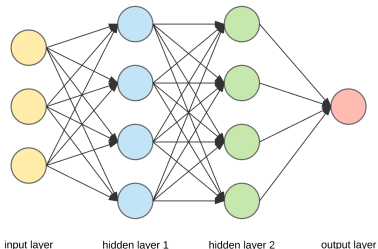
More output nodes



$W \in \mathbb{R}^{4 \times 3}$, $h : \mathbb{R}^4 \rightarrow \mathbb{R}^4$ so $h(a) = (h_1(a_1), h_2(a_2), h_3(a_3), h_4(a_4))$

Can think of this as a nonlinear basis: $\Phi(x) = h(Wx)$

More layers



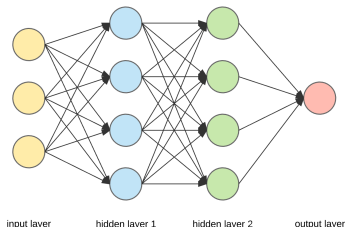
Becomes a network:

- each node is called a **neuron**
- h is called the **activation function**
 - can use $h(a) = 1$ for one neuron in each layer to *incorporate bias term*
 - output neuron can use $h(a) = a$
- #layers refers to #hidden_layers (plus 1 or 2 for input/output layers)
- **deep** neural nets can have many layers and *millions* of parameters
- this is a **feedforward, fully connected** neural net, there are many variants

Math formulation

An L-layer neural net can be written as

$$f(x) = h_L(W_L h_{L-1}(W_{L-1} \cdots h_1(W_1 x)))$$



To ease notation, for a given input x , define recursively

$$o_0 = x, \quad a_\ell = W_\ell o_{\ell-1}, \quad o_\ell = h_\ell(a_\ell) \quad (\ell = 1, \dots, L)$$

where

- $W_\ell \in \mathbb{R}^{D_\ell \times D_{\ell-1}}$ is the weights between layer $\ell - 1$ and ℓ
- $D_0 = D, D_1, \dots, D_L$ are numbers of neurons at each layer
- $a_\ell \in \mathbb{R}^{D_\ell}$ is input to layer ℓ
- $o_\ell \in \mathbb{R}^{D_\ell}$ is output to layer ℓ
- $h_\ell : \mathbb{R}^{D_\ell} \rightarrow \mathbb{R}^{D_\ell}$ is activation functions at layer ℓ

Outline

- 1 Neural Networks: Motivation
- 2 Neural Networks: Notation
- 3 Neural Networks: Learning and Backpropagation**

Learning the model

No matter how complicated the model is, our goal is the same: minimize

$$\mathcal{E}(\mathbf{W}_1, \dots, \mathbf{W}_L) = \frac{1}{N} \sum_{n=1}^N \mathcal{E}_n(\mathbf{W}_1, \dots, \mathbf{W}_L)$$

where

$$\mathcal{E}_n(\mathbf{W}_1, \dots, \mathbf{W}_L) = \begin{cases} \|\mathbf{f}(\mathbf{x}_n) - \mathbf{y}_n\|_2^2 & \text{for regression} \\ \ln \left(1 + \sum_{k \neq y_n} e^{f(\mathbf{x}_n)_k - f(\mathbf{x}_n)_{y_n}} \right) & \text{for classification} \end{cases}$$

Learning the model

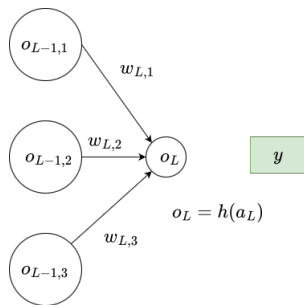
Let's take the regression task:

$$\mathcal{E}(\mathbf{W}_1, \dots, \mathbf{W}_L) = \frac{1}{N} \sum_{n=1}^N \|\mathbf{f}(\mathbf{x}_n) - \mathbf{y}_n\|_2^2$$

where the goal is to learn the function $\mathbf{f}(\cdot)$

Learning the model: How to optimize such a complicated function?

The function $f(\cdot)$ corresponds to the output of the neural network o_L . Due to the non-linear activations $f(\cdot)$ is in general *nonconvex*, and we use **SGD**!



Note that the output is our regression task is a scalar, therefore o_L is a scalar, but in general the output can be a vector. There is also a bias term, but we ignore it for brevity.

Analyzing the last layer for gradient descent

$$\mathbf{o}_0 = \mathbf{x}, \quad \mathbf{a}_\ell = \mathbf{W}_\ell \mathbf{o}_{\ell-1}, \quad \mathbf{o}_\ell = \mathbf{h}_\ell(\mathbf{a}_\ell) \quad (\ell = 1, \dots, L)$$

For the last layer, for square loss-based regression case we are considering, we can evaluate the gradient w.r.t. $w_{L,1}$ as follows

$$\frac{\partial \mathcal{E}_n}{\partial w_{L,1}} = \frac{\partial \mathcal{E}_n}{\partial o_L} \frac{\partial o_L}{\partial a_L} \frac{\partial a_L}{\partial w_{L,1}} \quad (\text{Chain Rule})$$

$$\frac{\partial \mathcal{E}_n}{\partial w_{L,1}} = \frac{\partial (o_L - y_n)^2}{\partial o_L} \frac{\partial o_L}{\partial a_L} \frac{\partial a_L}{\partial w_{L,1}} = 2(o_L - y_n)h'_L(a_L)o_{L-1,1}$$

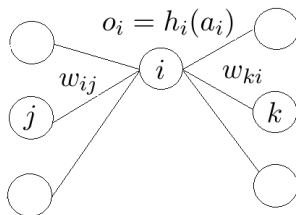
We will follow a similar procedure for remaining weights $w_{2,L}$ and $w_{3,L}$, which gives us gradient w.r.t. all weights.

Notice how the weights in the last layer here constitute a vector. Is this always the case? **No!** If the output layer is a vector instead of a scalar, we will have a matrix!

Computing the gradients for inner layers

We essentially do a similar application of chain rule

- We move from output layers to inner layers, and use the evaluated gradients as we move inward

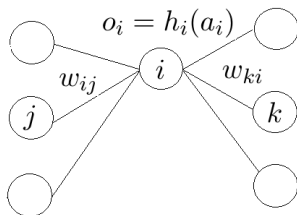


Computing the derivative

$$\mathbf{o}_0 = \mathbf{x}, \quad \mathbf{a}_\ell = \mathbf{W}_\ell \mathbf{o}_{\ell-1}, \quad \mathbf{o}_\ell = \mathbf{h}_\ell(\mathbf{a}_\ell) \quad (\ell = 1, \dots, L)$$

Drop the subscript ℓ for layer for simplicity.

Find the **derivative of \mathcal{E}_n w.r.t. to w_{ij}**



$$\frac{\partial \mathcal{E}_n}{\partial w_{ij}} = \underbrace{\frac{\partial \mathcal{E}_n}{\partial a_i}}_{\frac{\partial \mathcal{E}_n}{\partial o_i} \frac{\partial o_i}{\partial a_i}} \frac{\partial a_i}{\partial w_{ij}} = \frac{\partial \mathcal{E}_n}{\partial o_i} \frac{\partial o_i}{\partial a_i} \frac{\partial (w_{ij} o_j)}{\partial w_{ij}} = \underbrace{\frac{\partial \mathcal{E}_n}{\partial o_i}}_{\frac{\partial \mathcal{E}_n}{\partial a_i}} h'_i(a_i) o_j$$

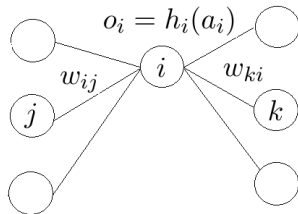
$$\frac{\partial \mathcal{E}_n}{\partial o_i} = \left(\sum_k \frac{\partial \mathcal{E}_n}{\partial a_k} \frac{\partial a_k}{\partial o_i} \right) = \left(\sum_k \frac{\partial \mathcal{E}_n}{\partial a_k} w_{ki} \right)$$

Computing the derivative

Adding the subscript for layer:

$$\frac{\partial \mathcal{E}_n}{\partial w_{\ell,ij}} = \frac{\partial \mathcal{E}_n}{\partial a_{\ell,i}} o_{\ell-1,j}$$

$$\frac{\partial \mathcal{E}_n}{\partial a_{\ell,i}} = \left(\sum_k \frac{\partial \mathcal{E}_n}{\partial a_{\ell+1,k}} w_{\ell+1,ki} \right) h'_{\ell,i}(a_{\ell,i})$$



For the last layer, for square loss

$$\frac{\partial \mathcal{E}_n}{\partial a_{L,i}} = \frac{\partial (h_{L,i}(a_{L,i}) - y_{n,i})^2}{\partial a_{L,i}} = 2(h_{L,i}(a_{L,i}) - y_{n,i}) h'_{L,i}(a_{L,i})$$

Computing the derivative

Using **matrix notation** greatly simplifies presentation and implementation:

$$\frac{\partial \mathcal{E}_n}{\partial \mathbf{W}_\ell} = \frac{\partial \mathcal{E}_n}{\partial \mathbf{a}_\ell} \mathbf{o}_{\ell-1}^\top$$

$$\frac{\partial \mathcal{E}_n}{\partial \mathbf{a}_\ell} = \begin{cases} \left(\mathbf{W}_{\ell+1}^\top \frac{\partial \mathcal{E}_n}{\partial \mathbf{a}_{\ell+1}} \right) \circ \mathbf{h}'_\ell(\mathbf{a}_\ell) & \text{if } \ell < L \\ 2(\mathbf{h}_L(\mathbf{a}_L) - \mathbf{y}_n) \circ \mathbf{h}'_L(\mathbf{a}_L) & \text{else} \end{cases}$$

where $\mathbf{v}_1 \circ \mathbf{v}_2 = (v_{11}v_{21}, \dots, v_{1D}v_{2D})$ is the element-wise product (a.k.a. Hadamard product)

Putting everything into SGD

The **backpropagation** algorithm (**Backprop**)

Initialize $\mathbf{W}_1, \dots, \mathbf{W}_L$. Repeat:

- ① randomly pick one data point $n \in [\mathbf{N}]$
- ② **forward propagation**: for each layer $\ell = 1, \dots, L$
 - compute $\mathbf{a}_\ell = \mathbf{W}_\ell \mathbf{o}_{\ell-1}$ and $\mathbf{o}_\ell = \mathbf{h}_\ell(\mathbf{a}_\ell)$ ($\mathbf{o}_0 = \mathbf{x}_n$)
- ③ **backward propagation**: for each $\ell = L, \dots, 1$
 - compute

$$\frac{\partial \mathcal{E}_n}{\partial \mathbf{a}_\ell} = \begin{cases} \left(\mathbf{W}_{\ell+1}^\top \frac{\partial \mathcal{E}_n}{\partial \mathbf{a}_{\ell+1}} \right) \circ \mathbf{h}'_\ell(\mathbf{a}_\ell) & \text{if } \ell < L \\ 2(\mathbf{h}_L(\mathbf{a}_L) - \mathbf{y}_n) \circ \mathbf{h}'_L(\mathbf{a}_L) & \text{else} \end{cases}$$

- update weights

$$\mathbf{W}_\ell \leftarrow \mathbf{W}_\ell - \eta \frac{\partial \mathcal{E}_n}{\partial \mathbf{W}_\ell} = \mathbf{W}_\ell - \eta \frac{\partial \mathcal{E}_n}{\partial \mathbf{a}_\ell} \mathbf{o}_{\ell-1}^\top$$

More tricks to optimize neural nets

Many variants based on backprop

- SGD with **minibatch**: randomly sample a batch of examples to form a stochastic gradient
- SGD with **momentum**
- ...

References

- Relevant sections of Chapter 13 (I built my notes from various sources for simplicity)