

Neural Networks

Sandro Cumani

sandro.cumani@polito.it

Politecnico di Torino

Neural networks

Neural Networks (NN) provide a method to approximate a non-linear function ϕ

A Neural Network can be interpreted as a non-linear parametric function $\phi(\mathbf{x}, \mathbf{\Pi})$

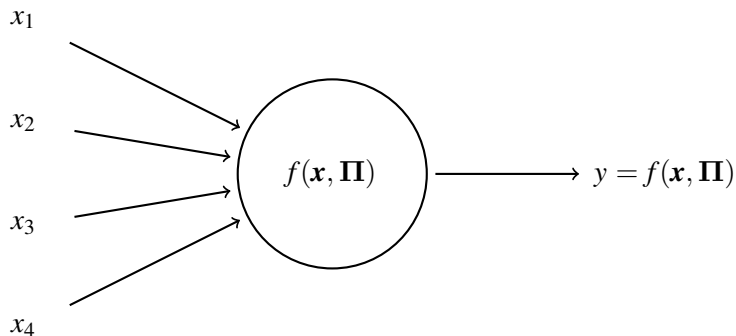
The parameters of the non-linear transformation are learned from the data

The function is represented by means of a directed graph

Each node is associated to a function that operates on the input nodes and provides the node output

Neural networks

The basic unit of a neural network is a computation node



The node computes a parametric function of its inputs $f(\mathbf{x}, \mathbf{\Pi})$

Neural networks

Typically, the non-linear function is expressed in terms of a parametric affine combination of the node inputs and a scalar, non-linear, non-parametric transformation

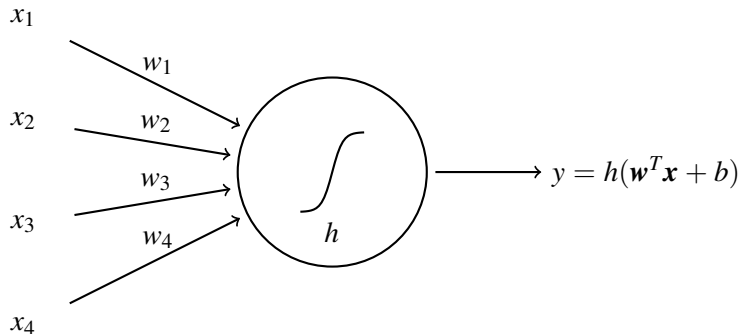
$$f(\mathbf{x}, \mathbf{\Pi}) = h(\mathbf{w}^T \mathbf{x} + b)$$

$\mathbf{\Pi} = (\mathbf{w}, b)$ are the function parameters: \mathbf{w} is a vector containing the weights of the affine combination and b is a bias (scalar)

In this case, we can graphically associate the weights parameters with the node input arcs

Neural networks

Let $\mathbf{w} = [w_1 \dots w_d]$, where d is the dimensionality of the node input. The node representation becomes



The weights can be interpreted as the “strength” of the connection — $w_i = 0$ implies that the corresponding input x_i does not influence the final result

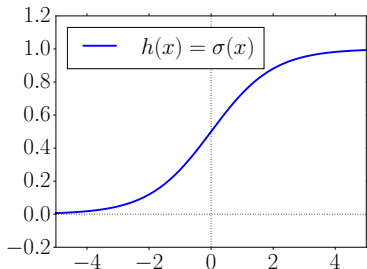
Neural networks

Several possible functions have been proposed for the non-linearity function h

- Sigmoid function $h(x) = \frac{1}{1+e^{-x}}$
- Hyperbolic tangent $h(x) = \tanh(x)$
- Rectified linear $h(x) = \max(0, x)$

Neural networks

Sigmoid function $h(x) = \frac{1}{1+e^{-x}}$

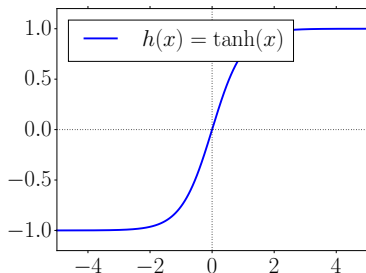


A sigmoid function-activated node can be interpreted as a binary logistic regression model for the input data x

$$h(\mathbf{w}^T \mathbf{x} + b) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x} + b}}$$

Neural networks

Hyperbolic tangent $h(x) = \tanh(x)$



It can be interpreted as a symmetric version of the sigmoid

$$\sigma(x) = \frac{\tanh\left(\frac{x}{2}\right) + 1}{2}, \quad \tanh(x) = 2\sigma(2x) - 1$$

Neural networks

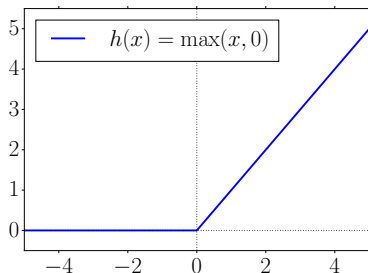
For small inputs both sigmoid and hyperbolic tangent behave almost linearly

For large inputs the functions saturate (non-linear behavior)

Computationally expensive, may pose issues during training

Neural networks

Rectified linear $h(x) = \max(0, x)$



More computationally efficient, less prone to training issues due to vanishing gradients, but not differentiable everywhere

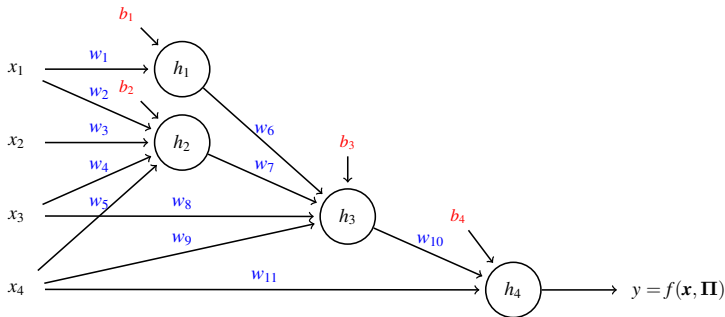
Widely adopted in practice

Neural networks

A neural network is a directed graph of nodes

We will consider only graphs without loops (acyclic)

The graph defines a sequence of operations



$$f(\mathbf{x}, \mathbf{\Pi}) = h_4(w_{10}h_3(w_6h_1(w_1x_1+b_1)+w_7h_2(w_2x_1+w_3x_2+w_4x_3+w_5x_4+b_2) \\ + w_8x_3 + w_9x_4 + b_3) + w_{11}x_4 + b_4)$$

Neural networks

An acyclic graph is also known as *feed forward network*

Information “flows” from the input to the output nodes

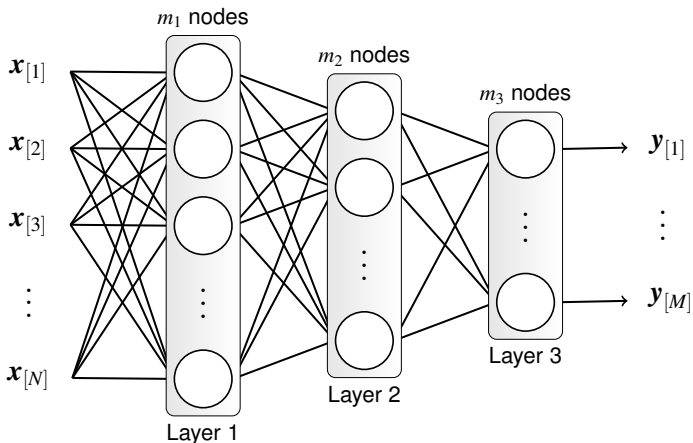
It's useful to organize nodes in *layers*

Layer: group of nodes that share the same inputs

Typically, a layer receives as input the output of a previous layer

Neural networks

Multi-layer network: units are organized in layers



Connections are defined between layers (no loops)

Neural networks

The layers are identified with a progressive index

We can represent the input \mathbf{n}_j of the layer as the vector of the inputs of the nodes associated to that layer

We can represent the output \mathbf{x}_j of the layer as the vector of outputs of its nodes

Each node j in layer i has an associated bias $b_{i,j}$ and a vector of weights

$$\mathbf{w}_{i,j} = \begin{bmatrix} \mathbf{w}_{i,j,1} \\ \vdots \\ \mathbf{w}_{i,j,m_{i-1}} \end{bmatrix}$$

where m_i denoted the number of nodes in a layer

Neural networks

We can arrange the weight vectors in a layer matrix, and the biases in an array

$$\mathbf{W}_j = [\mathbf{w}_{j,1} \dots \mathbf{w}_{j,m_j}] \quad , \quad \mathbf{b}_j = \begin{bmatrix} \mathbf{b}_{j,1} \\ \dots \\ \mathbf{b}_{j,m_j} \end{bmatrix}$$

The input \mathbf{n}_j of layer j can be expressed in terms of the output \mathbf{x}_{j-1} of layer $j - 1$ as

$$\mathbf{n}_j = \mathbf{W}_j^T \mathbf{x}_{j-1} + \mathbf{b}_j = \begin{bmatrix} \mathbf{w}_{j,1}^T \mathbf{x}_{j-1} + \mathbf{b}_{j,1} \\ \vdots \\ \mathbf{w}_{j,m_j}^T \mathbf{x}_{j-1} + \mathbf{b}_{j,m_j} \end{bmatrix}$$

with $\mathbf{x}_0 = \mathbf{x}$ representing the network input

Neural networks

The layer output is computed from the layer input by applying element-wise the nodes non-linearity

Assuming that all nodes share the same non-linear activation function h :

$$\mathbf{x}_j = \begin{bmatrix} h(\mathbf{n}_{j,1}) \\ \vdots \\ h(\mathbf{n}_{j,m_j}) \end{bmatrix} = \mathbf{h}(\mathbf{n}_j)$$

Note that we use the same letter, but in bold \mathbf{h} , to indicate the *vector-valued* function that applies the *scalar* function h to all its inputs

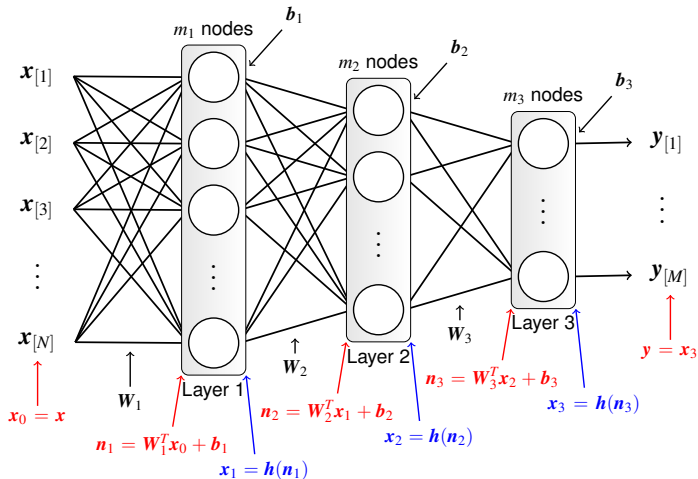
Each layer of the network computes a parametric linear transformation, followed by an element-wise non-parametric non-linear transformation¹

$$\mathbf{f}_j(\mathbf{x}_{j-1}, \mathbf{W}_j, \mathbf{b}_j) = \mathbf{h}(\mathbf{W}_j^T \mathbf{x}_{j-1} + \mathbf{b}_j)$$

¹These networks can be extended to include parametric non-linearities

Neural networks

Multi-layer feed-forward network:



Neural networks

To compute the network function we “forward” the input vector through the different layers. For a network with m layers:

$$\mathbf{x}_0 = \mathbf{x}$$

$$\mathbf{x}_1 = \mathbf{f}_1(\mathbf{x}_0, \mathbf{W}_1, \mathbf{b}_1) = \mathbf{h}(\mathbf{W}_1^T \mathbf{x}_0 + \mathbf{b}_1)$$

$$\mathbf{x}_2 = \mathbf{f}_2(\mathbf{x}_1, \mathbf{W}_2, \mathbf{b}_2) = \mathbf{h}(\mathbf{W}_2^T \mathbf{x}_1 + \mathbf{b}_2)$$

$$\vdots$$

$$\mathbf{x}_{m-1} = \mathbf{f}_{m-1}(\mathbf{x}_{m-2}, \mathbf{W}_{m-1}, \mathbf{b}_{m-1}) = \mathbf{h}(\mathbf{W}_{m-1}^T \mathbf{x}_{m-2} + \mathbf{b}_{m-1})$$

$$\mathbf{x}_m = \mathbf{f}_m(\mathbf{x}_{m-1}, \mathbf{W}_m, \mathbf{b}_m) = \mathbf{h}(\mathbf{W}_m^T \mathbf{x}_{m-1} + \mathbf{b}_m)$$

$$\mathbf{y} = \mathbf{x}_m = \mathbf{f}(\mathbf{x}, \mathbf{\Pi})$$

$\mathbf{f}(\mathbf{x}, \mathbf{\Pi})$ is the function that corresponds to the network, and $\mathbf{\Pi}$ is the set of layer weights \mathbf{W}_i and bias vectors \mathbf{b}_i

$\mathbf{f}_j(\mathbf{x}_{j-1}, \mathbf{W}_j, \mathbf{b}_j)$ is the function of the j -th network layer

Neural networks

The layers between the input and output layer are called *hidden* layers

A feed-forward neural network with hidden layers is also called multilayer perceptron (MLP)

It can be shown that any continuous function can be approximated up to a desired degree by a MLP of sufficient size

Neural networks

We can employ MLPs to represent a non-linear, parametric transformation of our input data

We can then classify our samples with a linear model in the transformed feature space induced by the MLP function

In a similar way, we can employ neural networks to compute non linear mappings to and from a lower dimensional space to achieve a form of non-linear dimensionality reduction

Neural networks

Binary problem: combine a binary logistic regression model with the MLP transformation

$$\log \frac{P(C = 1|X = \mathbf{x})}{P(C = 0|X = \mathbf{x})} = \mathbf{w}^T \mathbf{f}(\mathbf{x}, \mathbf{\Pi}) + b$$

The class posterior probability can be expressed as

$$P(C = 1|X = \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{f}(\mathbf{x}, \mathbf{\Pi}) + b)$$

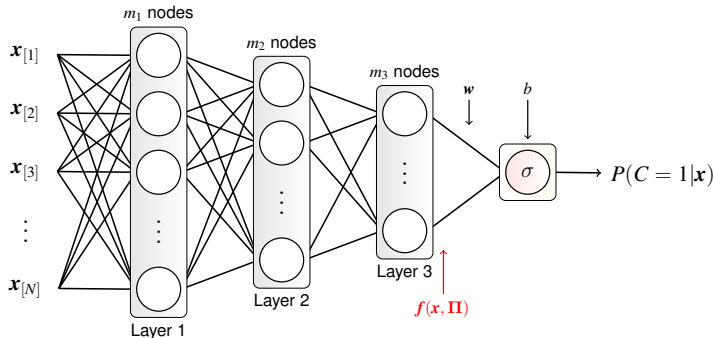
The network $f(\cdot, \mathbf{\Pi})$ represents the non-linear transformation of our data

(\mathbf{w}, b) are the parameters of a linear logistic regression classifier

Neural networks

The map from the network output $f(\mathbf{x}, \mathbf{\Pi})$ to the class posterior probability $P(C = 1|\mathbf{x})$ can be represented as a sigmoid-activated network layer with a single node:

$$P(C = 1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{f}(\mathbf{x}, \mathbf{\Pi}) + b) = \hat{f}(\mathbf{x}, \hat{\mathbf{\Pi}}), \quad \hat{\mathbf{\Pi}} = (\mathbf{\Pi}, \mathbf{w}, b)$$



Neural networks

To estimate the network parameters we can optimize the logistic loss (also referred to as binary cross-entropy loss)

$$\begin{aligned}\mathbf{\Pi}^*, \mathbf{w}^*, b^* &= \arg \min_{\mathbf{\Pi}, \mathbf{w}, b} \frac{1}{N} \sum_{i=1}^N \log \left(1 + e^{-z_i(\mathbf{w}^T \mathbf{f}(\mathbf{x}_i, \mathbf{\Pi}) + b)} \right) \\ &= \arg \min_{\mathbf{\Pi}, \mathbf{w}, b} -\frac{1}{N} \sum_{i=1}^N \left[c_i \log \sigma(\mathbf{w}^T \mathbf{f}(\mathbf{x}_i, \mathbf{\Pi}) + b) \right. \\ &\quad \left. + (1 - c_i) \log(1 - \sigma(\mathbf{w}^T \mathbf{f}(\mathbf{x}_i, \mathbf{\Pi}) + b)) \right]\end{aligned}$$

where \mathbf{x}_i is the i -th sample in the dataset (pay attention not to confuse \mathbf{x}_i in the loss expression with the output of the i -th network layer for sample \mathbf{x})

Neural networks

In terms of the network \hat{f} we can express the optimization as

$$\hat{\Pi}^* = \arg \min_{\hat{\Pi}} -\frac{1}{N} \sum_{i=1}^N \left[c_i \log \hat{f}(\mathbf{x}_i, \hat{\Pi}) + (1 - c_i) \log(1 - \hat{f}(\mathbf{x}_i, \hat{\Pi})) \right]$$

As for linear logistic regression models, we cannot directly optimize the objective function

Again, we can rely on numerical optimization

Typically, numerical solvers require that we are able to compute both the objective function value (forward run), and its gradient

Neural networks

We thus turn our attention to the computation of the gradient of a loss function that depends on the outputs of a neural network

Let $\mathcal{L}(\mathbf{W}_1, \mathbf{b}_1, \dots, \mathbf{W}_m, \mathbf{b}_m)$ be the objective function that we want to minimize

For binary classification, \mathcal{L} is the average logistic loss computed from the network \hat{f} (i.e., the extended network that includes $m - 1$ feature transformation layers and, as last layer, the single-node, sigmoidal linear classification layer)

$$\begin{aligned}\mathcal{L}(\mathbf{W}_1, \mathbf{b}_1, \dots, \mathbf{W}_m, \mathbf{b}_m) = \\ \arg \min_{\mathbf{W}_1, \mathbf{b}_1, \dots, \mathbf{W}_m, \mathbf{b}_m} -\frac{1}{N} \sum_{i=1}^N \left[c_i \log \hat{f}(\mathbf{x}_i, \mathbf{W}_1, \mathbf{b}_1, \dots, \mathbf{W}_m, \mathbf{b}_m) \right. \\ \left. + (1 - c_i) \log(1 - \hat{f}(\mathbf{x}_i, \mathbf{W}_1, \mathbf{b}_1, \dots, \mathbf{W}_m, \mathbf{b}_m)) \right]\end{aligned}$$

Neural networks

In terms of the network layer, the loss can be represented as

$$\mathcal{L}(\mathbf{W}_1, \mathbf{b}_1, \dots, \mathbf{W}_m, \mathbf{b}_m) = \frac{1}{N} \sum_{i=1}^N \ell(\hat{f}(\mathbf{x}_i, \mathbf{W}_1, \mathbf{b}_1, \dots, \mathbf{W}_m, \mathbf{b}_m), c_i)$$

where ℓ is a loss function, e.g., the logistic loss

$$\ell(y, c) = -[c \log y + (1 - c) \log(1 - y)]$$

The gradient of \mathcal{L} requires computing the partial derivatives of \mathcal{L} with respect to the network parameters $\mathbf{W}_{i,jk}$ and $\mathbf{b}_{i,j}$

The derivatives can be computed using the standard chain rule

Neural networks

Let's consider the derivatives with respect to the parameters of the last layer first

The parameters are $\mathbf{W}_m, \mathbf{b}_m$

We want to compute the partial derivatives $\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{m,ij}}$ and $\frac{\partial \mathcal{L}}{\partial \mathbf{b}_{m,i}}$

We have

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{m,ij}} = \frac{1}{N} \sum_{i=1} \frac{\partial \ell}{\partial \mathbf{W}_{m,jk}}$$

We thus need to compute the term $\frac{\partial \ell}{\partial \mathbf{W}_{m,jk}}$, i.e., the derivative of the loss function ℓ

Neural networks

We recall that, as a function of \mathbf{W}_m , the loss can be expressed as

$$\ell(f(\mathbf{x}, \mathbf{W}_1, \mathbf{b}_1, \dots, \mathbf{W}_m, \mathbf{b}_m), c) = \ell(\mathbf{x}_m, c) = \ell(\mathbf{h}(\mathbf{W}_m^T \mathbf{x}_{m-1} + \mathbf{b}_m), c_i)$$

where $\mathbf{x}_m = \mathbf{h}(\mathbf{W}_m^T \mathbf{x}_{m-1} + \mathbf{b}_m)$ is the output of the last layer of the network, which depends directly on \mathbf{W}_m and the output of the $(m - 1)$ -th layer \mathbf{x}_{m-1}

Pay attention that \mathbf{x}_m is the *output of the last layer* for input \mathbf{x} , *it's not* the m -th sample of the dataset

In the following we will consider a single sample \mathbf{x} to avoid having too many indices For the same reason, we will not consider single elements of $\mathbf{W}_{m,jk}$ but derive an expression for the *gradient* of ℓ

Neural networks

The gradient of ℓ is defined as the (row) vector of partial derivatives of ℓ with respect to all the parameters

In the following, we consider the gradient components that correspond to the terms of \mathbf{W}_m (and \mathbf{b}_m):

$$\nabla_{\mathbf{W}_m} \ell = \left[\frac{\partial \ell}{\partial \mathbf{W}_{m,11}} \cdots \frac{\partial \ell}{\partial \mathbf{W}_{m,1d_m}} , \frac{\partial \ell}{\partial \mathbf{W}_{m,21}} \cdots \frac{\partial \ell}{\partial \mathbf{W}_{m,2d_m}} \cdots \frac{\partial \ell}{\partial \mathbf{W}_{m,d_{m-1}1}} \cdots \frac{\partial \ell}{\partial \mathbf{W}_{m,d_{m-1}d_m}} \right]$$
$$\nabla_{\mathbf{b}_m} \ell = \left[\frac{\partial \ell}{\partial \mathbf{b}_{m,1}} \cdots \frac{\partial \ell}{\partial \mathbf{b}_{m,d_m}} \right]$$

where d_m and d_{m-1} are the size of layers m and $m - 1$

Neural networks

Let's consider a scalar function $g : \mathbb{R}^q \rightarrow \mathbb{R}$, a vector function $\mathbf{h} : \mathbb{R}^p \rightarrow \mathbb{R}^q$,

with components $\mathbf{h}(\mathbf{z}) = \begin{bmatrix} h_1(\mathbf{z}) \\ \vdots \\ h_p(\mathbf{z}) \end{bmatrix}$

Let $f(\mathbf{z}) = g(\mathbf{h}(\mathbf{z}))$. The chain rule allows computing the gradient of f w.r.t. \mathbf{z} from the gradient of $g(\mathbf{h})$ w.r.t. its input \mathbf{h}

$$\nabla_{\mathbf{h}} g = \begin{bmatrix} \frac{\partial g}{\partial h_1} & \cdots & \frac{\partial g}{\partial h_q} \end{bmatrix}$$

and the matrix of partial derivatives (Jacobian matrix) of $\mathbf{h}(\mathbf{z})$

$$\frac{\partial \mathbf{h}}{\partial \mathbf{z}} = \begin{bmatrix} \frac{\partial h_1}{\partial z_1} & \cdots & \frac{\partial h_1}{\partial z_p} \\ \vdots & \ddots & \vdots \\ \frac{\partial h_q}{\partial z_1} & \cdots & \frac{\partial h_q}{\partial z_p} \end{bmatrix}$$

The gradient of f is obtained as

$$\nabla_{\mathbf{z}} f = \nabla_{\mathbf{h}} g \cdot \frac{\partial \mathbf{h}}{\partial \mathbf{z}}$$

Neural networks

We can apply the chain rule to the derivation of the partial gradient of the loss function w.r.t. the terms \mathbf{W}_m . We recall that

$$\mathbf{x}_m = \mathbf{h}(\mathbf{n}_m) , \quad \mathbf{n}_m = \mathbf{W}_m^T \mathbf{x}_{m-1} + \mathbf{b}_m$$

Thus

$$\nabla_{\mathbf{W}_m} \ell = \nabla_{\mathbf{x}_m} \ell \cdot \frac{\partial \mathbf{x}_m}{\partial \mathbf{n}_m} \cdot \frac{\partial \mathbf{n}_m}{\partial \mathbf{W}_m}$$

The left-most term is the loss function gradient. For example, for binary logistic regression \mathbf{x}_m is a scalar, and the gradient is thus a 1-element vector:

$$\begin{aligned} \ell(\mathbf{x}_m, c) &= c \log \mathbf{x}_m + (1 - c) \log(1 - \mathbf{x}_m) \\ \nabla_{\mathbf{x}_m} \ell &= \left[\frac{c}{\mathbf{x}_m} - \frac{1-c}{1-\mathbf{x}_m} \right] \end{aligned}$$

Neural networks

The term $\frac{\partial \mathbf{x}_m}{\partial \mathbf{n}_m}$ is the matrix of partial derivatives of the nodes non-linearity \mathbf{h}

Since function \mathbf{h} consists of element-wise non-linearities, we can compute the matrix of partial derivatives as

$$\frac{\partial \mathbf{x}_m}{\partial \mathbf{n}_m} = \begin{bmatrix} \left. \frac{dh(x)}{dx} \right|_{\mathbf{n}_{m,1}} & 0 & \dots & 0 \\ 0 & \left. \frac{dh(x)}{dx} \right|_{\mathbf{n}_{m,2}} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \left. \frac{dh(x)}{dx} \right|_{\mathbf{n}_{m,d_m}} \end{bmatrix}$$

$\left. \frac{dh(x)}{dx} \right|_y$ is the derivative of the non-linear function $h(x)$ evaluated at y .

For example: $h(x) = \sigma(x) \rightarrow \frac{dh}{dx}(x) = \sigma(x)(1 - \sigma(x))$
 $h(x) = \tanh(x) \rightarrow \frac{dh}{dx}(x) = 1 - \tanh^2(x)$

Neural networks

We observe that, given the gradient of ℓ with respect to the network outputs, to compute the gradient w.r.t. \mathbf{W}_m we must first compute the product

$$\nabla_{\mathbf{x}_m} \ell \cdot \frac{\partial \mathbf{x}_m}{\partial \mathbf{n}_m}$$

Since $\frac{\partial \mathbf{x}_m}{\partial \mathbf{n}_m}$ is diagonal, we can compute efficiently the result by multiplying each element of $\nabla_{\mathbf{x}_m} \ell$ by the corresponding element of the diagonal of $\frac{\partial \mathbf{x}_m}{\partial \mathbf{n}_m}$

Let

$$\mathbf{v}_m = \nabla_{\mathbf{x}_m} \ell \cdot \frac{\partial \mathbf{x}_m}{\partial \mathbf{n}_m}$$

We further need to compute

$$\nabla_{\mathbf{W}_m} \ell = \mathbf{v}_m \cdot \frac{\partial \mathbf{n}_m}{\partial \mathbf{W}_m}$$

i.e. the product of the vector \mathbf{v}_m and the matrix of partial derivatives of *input of the network last layer* with respect to *the layer weights*

Neural networks

Since \mathbf{W}_m is a matrix, it's useful to represent in matrix form the subset of the gradient that corresponds to derivatives w.r.t. $\mathbf{W}_{m,jk}$:

$$\nabla_{\mathbf{W}_m} \ell = \begin{bmatrix} \nabla_{\mathbf{w}_1} \ell \\ \vdots \\ \nabla_{\mathbf{w}_{d_m}} \ell \end{bmatrix} = \begin{bmatrix} \frac{\partial \ell}{\partial \mathbf{W}_{m,11}} & \cdots & \frac{\partial \ell}{\partial \mathbf{W}_{m,d_{m-1}1}} \\ \vdots & \ddots & \vdots \\ \frac{\partial \ell}{\partial \mathbf{W}_{m,1d_m}} & \cdots & \frac{\partial \ell}{\partial \mathbf{W}_{m,d_{m-1}d_m}} \end{bmatrix}$$

Note: a *row* of $\nabla_{\mathbf{W}_m} \ell$ contains the derivatives with respect to the elements of a *column* of $\mathbf{W}_m = [\mathbf{w}_1 \ \dots \ \mathbf{w}_{d_m}]$

The gradient $\nabla_{\mathbf{W}_m} \ell$ can be computed in matrix form as

$$\nabla_{\mathbf{W}_m} \ell = \mathbf{v}_m \cdot \frac{\partial \mathbf{n}_m}{\partial \mathbf{W}_m} = \mathbf{x}_{m-1} \mathbf{v}_m$$

Similarly, we can show that

$$\nabla_{\mathbf{b}_m} \ell = \mathbf{v}_m$$

Neural networks

Let's now consider the set of derivatives with respect to the previous-from-last layer parameters $\mathbf{W}_{m-1}, \mathbf{b}_{m-1}$

As before, we can express the dependency of the loss on these parameters through \mathbf{x}_m :

$$\mathbf{x}_m = \mathbf{h}(\mathbf{n}_m)$$

$$\mathbf{n}_m = \mathbf{W}_m^T \mathbf{x}_{m-1} + \mathbf{b}_m$$

$$\mathbf{x}_{m-1} = \mathbf{h}(\mathbf{n}_{m-1})$$

$$\mathbf{n}_{m-1} = \mathbf{W}_{m-1}^T \mathbf{x}_{m-2} + \mathbf{b}_{m-1}$$

Applying the chain rule, we can write

$$\nabla_{\mathbf{W}_m} \ell = \nabla_{\mathbf{x}_m} \ell \cdot \frac{\partial \mathbf{x}_m}{\partial \mathbf{n}_m} \cdot \frac{\partial \mathbf{n}_m}{\partial \mathbf{x}_{m-1}} \cdot \frac{\partial \mathbf{x}_{m-1}}{\partial \mathbf{n}_{m-1}} \cdot \frac{\partial \mathbf{n}_{m-1}}{\partial \mathbf{W}_{m-1}}$$

Neural networks

We can observe that the product of the first three terms corresponds to the set of partial derivatives of the loss w.r.t. the $(m - 1)$ -th layer outputs

$$\begin{aligned}\nabla_{\mathbf{x}_{m-1}} \ell &= \nabla_{\mathbf{x}_m} \ell \cdot \frac{\partial \mathbf{x}_m}{\partial \mathbf{n}_m} \cdot \frac{\partial \mathbf{n}_m}{\partial \mathbf{x}_{m-1}} \\ \nabla_{\mathbf{W}_{m-1}} \ell &= \nabla_{\mathbf{x}_{m-1}} \ell \cdot \frac{\partial \mathbf{x}_{m-1}}{\partial \mathbf{n}_{m-1}} \cdot \frac{\partial \mathbf{n}_{m-1}}{\partial \mathbf{W}_{m-1}}\end{aligned}$$

We can compare with the derivatives w.r.t \mathbf{W}_m :

$$\nabla_{\mathbf{W}_m} \ell = \nabla_{\mathbf{x}_m} \ell \cdot \frac{\partial \mathbf{x}_m}{\partial \mathbf{n}_m} \cdot \frac{\partial \mathbf{n}_m}{\partial \mathbf{W}_m}$$

Neural networks

We observe that in both cases we can express the gradient as a vector-matrix multiplication

The vector term represents the *gradient of the loss w.r.t. the layer outputs*

The matrix terms are *partial derivatives of the layer function with respect to the layer parameters, computed at the previous layer outputs*

Iterating this process, we can verify that

$$\begin{aligned}\nabla_{\mathbf{W}_k} \ell &= \mathbf{v}_k \cdot \frac{\partial \mathbf{n}_k}{\partial \mathbf{W}_k} \\ \nabla_{\mathbf{b}_k} \ell &= \mathbf{v}_k \cdot \frac{\partial \mathbf{n}_k}{\partial \mathbf{b}_k}\end{aligned}$$

with

$$\mathbf{v}_k = \nabla_{\mathbf{x}_k} \ell \cdot \frac{\partial \mathbf{x}_k}{\partial \mathbf{n}_k}$$

In matrix form we can express the gradients as

$$\nabla_{\mathbf{W}_k} \ell(\mathbf{x}_m, c) = \mathbf{x}_{k-1} \mathbf{v}_k, \quad \nabla_{\mathbf{b}_k} \ell(\mathbf{x}_m, c) = \mathbf{v}_k$$

We can thus compute the gradient by computing the terms \mathbf{v}_k for each layer

Neural networks

In general, we can verify that, for layer k , we have

$$\begin{aligned}\nabla_{\mathbf{x}_k} \ell &= \nabla_{\mathbf{x}_{k+1}} \ell \cdot \frac{\partial \mathbf{x}_{k+1}}{\partial \mathbf{n}_{k+1}} \cdot \frac{\partial \mathbf{n}_{k+1}}{\partial \mathbf{x}_k} \\ \mathbf{v}_k &= \mathbf{v}_{k+1} \cdot \frac{\partial \mathbf{n}_{k+1}}{\partial \mathbf{x}_k} \cdot \frac{\partial \mathbf{x}_k}{\partial \mathbf{n}_k}\end{aligned}$$

Furthermore, we can show that

$$\begin{aligned}\nabla_{\mathbf{x}_k} \ell &= \mathbf{v}_{k+1} \cdot \frac{\partial \mathbf{n}_{k+1}}{\partial \mathbf{x}_k} \\ \mathbf{v}_k &= \nabla_{\mathbf{x}_k} \ell \cdot \frac{\partial \mathbf{x}_k}{\partial \mathbf{n}_k}\end{aligned}$$

Neural networks

Each term \mathbf{v}_k corresponds to the gradient of the loss with respect to the k -th layer *inputs*

$$\mathbf{v}_k = \nabla_{\mathbf{n}_k} \ell$$

Starting from the last layer, we can efficiently compute $\nabla_{\mathbf{x}_k} \ell$ and $\mathbf{v}_k = \nabla_{\mathbf{n}_k} \ell$ through an iterative procedure that starts from $\nabla_{\mathbf{x}_m} \ell$:

The diagram illustrates the backpropagation process through three layers. It shows how the gradient of the loss with respect to the inputs of one layer is used to compute the gradient with respect to the inputs of the previous layer. Vertical arrows indicate the calculation of \mathbf{v} from $\nabla_{\mathbf{x}}$, while diagonal arrows show the propagation of $\nabla_{\mathbf{x}}$ to the next layer.

$$\begin{array}{ccccc} \nabla_{\mathbf{x}_m} \ell & & \nabla_{\mathbf{x}_{m-1}} \ell = \mathbf{v}_m \cdot \frac{\partial \mathbf{n}_m}{\partial \mathbf{x}_{m-1}} & & \nabla_{\mathbf{x}_{m-2}} \ell = \mathbf{v}_m \cdot \frac{\partial \mathbf{n}_{m-1}}{\partial \mathbf{x}_{m-2}} \\ \downarrow & \nearrow & \downarrow & \nearrow & \downarrow \\ \mathbf{v}_m = \nabla_{\mathbf{x}_m} \ell \cdot \frac{\partial \mathbf{x}_m}{\partial \mathbf{n}_m} & & \mathbf{v}_{m-1} = \nabla_{\mathbf{x}_{m-1}} \ell \cdot \frac{\partial \mathbf{x}_{m-1}}{\partial \mathbf{n}_{m-1}} & & \mathbf{v}_{m-2} = \nabla_{\mathbf{x}_{m-2}} \ell \cdot \frac{\partial \mathbf{x}_{m-2}}{\partial \mathbf{n}_{m-2}} \end{array}$$

The procedure is called *back-propagation*

At each step, we are propagating the loss gradient backwards through the layers

Neural networks

We already showed how to compute the terms $\frac{\partial \mathbf{x}_k}{\partial \mathbf{n}_k}$

It remains to see how to compute the terms $\frac{\partial \mathbf{n}_k}{\partial \mathbf{x}_{k-1}}$

Since $\mathbf{n}_k = \mathbf{W}_k^T \mathbf{x}_{k-1} + \mathbf{b}_k$ we can show that

$$\frac{\partial \mathbf{n}_k}{\partial \mathbf{x}_{k-1}} = \mathbf{W}_k$$

and the iterative procedure becomes

$$\mathbf{v}_k = \nabla_{\mathbf{x}_k} \ell \cdot \frac{\partial \mathbf{x}_k}{\partial \mathbf{n}_k}$$

$$\nabla_{\mathbf{x}_{k-1}} \ell = \mathbf{v}_k \mathbf{W}_k$$

Neural networks

Forward run (compute the network outputs):

$$\mathbf{n}_1 = \mathbf{W}_1^T \mathbf{x}_0 + \mathbf{b}_1$$

$$\mathbf{x}_1 = \mathbf{h}(\mathbf{n}_1)$$

$$\mathbf{n}_2 = \mathbf{W}_2^T \mathbf{x}_1 + \mathbf{b}_2$$

$$\mathbf{x}_2 = \mathbf{h}(\mathbf{n}_2)$$

...

$$\mathbf{n}_m = \mathbf{W}_m^T \mathbf{x}_{m-1}$$

$$\mathbf{x}_m = \mathbf{h}(\mathbf{n}_m)$$

(Column vectors)

Backward run (compute the terms required for gradient computation):

$$\mathbf{v}_m = \nabla_{\mathbf{x}_m} \ell \cdot \frac{\partial \mathbf{x}_m}{\partial \mathbf{n}_m}$$

$$\nabla_{\mathbf{x}_{m-1}} \ell = \mathbf{v}_m \mathbf{W}_m$$

$$\mathbf{v}_{m-1} = \nabla_{\mathbf{x}_{m-1}} \ell \cdot \frac{\partial \mathbf{x}_{m-1}}{\partial \mathbf{n}_{m-1}}$$

$$\nabla_{\mathbf{x}_{m-2}} \ell = \mathbf{v}_{m-1} \mathbf{W}_{m-1}$$

...

$$\mathbf{v}_1 = \nabla_{\mathbf{x}_2} \ell \cdot \frac{\partial \mathbf{x}_2}{\partial \mathbf{n}_1}$$

$$\nabla_{\mathbf{x}_0} \ell = \mathbf{v}_1 \mathbf{W}_1$$

(Row vectors)

Neural networks

Once we are able to compute gradients, we can apply a numerical solver to find a local minimum of our objective

While for small datasets methods like L-BFGS may provide fast and robust results, for larger datasets these methods may become expensive

Faster approaches based on Gradient Descent (GD) are typically employed

GD: starting from an initial value for the network weights $\mathbf{W}_1^0, \mathbf{b}_1^0 \dots \mathbf{W}_m^0, \mathbf{b}_m^0$ the weights are iteratively updated according to

$$\mathbf{W}_k^t = \mathbf{W}_k^{t-1} - \alpha_t \nabla_{\mathbf{W}_k^{t-1}}^T \mathcal{L}, \quad \mathbf{b}_k^t = \mathbf{b}_k^{t-1} - \alpha_t \nabla_{\mathbf{b}_k^{t-1}}^T \mathcal{L}$$

Neural networks

The coefficient α_t is called *learning rate*, and controls the strength of the weights update

GD convergence is guaranteed if

$$\sum_t \alpha_t = \infty \quad \sum_t \alpha_t^2 < \infty$$

However, in practice the number of iterations can be heavily influenced by the learning rate schedule

Since neural networks typically require large training sets, the standard GD approach is not practical, as it has the same drawback of L-BFGS: it requires a full iteration over the whole dataset to compute the loss gradient

$$\nabla \mathcal{L} = \frac{1}{N} \sum_{i=1}^N \nabla \ell(\mathbf{x}_i, c_i)$$

Neural networks

To address this issue, training is usually performed using *Stochastic Gradient Descent* (SGD) over *batches*

A batch is a set of randomly selected samples

We approximate the gradient

$$\nabla \mathcal{L} = \frac{1}{N} \sum_{i=1}^N \nabla \ell(\mathbf{x}_i, c_i)$$

with

$$\nabla \mathcal{L} \approx \frac{1}{K} \sum_{\mathbf{x}_i \in B} \nabla \ell(\mathbf{x}_i, c_i)$$

where B is a set of K samples

Typical batch sizes range in the tens to few hundreds

Neural networks

Training is also usually organized in *epochs*. At each epoch:

1. Randomly sample a batch using samples that have been not employed during the current epoch yet
2. Compute the batch gradient and update the weights using the approximated gradient
3. Update the learning rate
4. Repeat from 1. until the whole dataset has been used

A limitation of SGD is that its performance relies heavily on the selection of a good learning rate schedule

Large values of α may overshoot the local minimum, but small values may make slow progress

Neural networks

We can extend SGD by incorporating a *momentum* term

The momentum term performs an exponential smoothing of the gradient

At iteration t we compute the update

$$\Delta_{\mathbf{W}_k}^t = \eta \Delta_{\mathbf{W}_k}^{t-1} - \alpha_t \nabla_{\mathbf{W}_k}^T \mathcal{L}$$

$$\Delta_{\mathbf{b}_k}^t = \eta \Delta_{\mathbf{b}_k}^{t-1} - \alpha_t \nabla_{\mathbf{b}_k}^T \mathcal{L}$$

$$\mathbf{W}_k^t = \mathbf{W}_k^{t-1} + \Delta_{\mathbf{W}_k}^t$$

$$\mathbf{b}_k^t = \mathbf{b}_k^{t-1} + \Delta_{\mathbf{b}_k}^t$$

where η is a constant factor

More sophisticated approaches have been recently introduced (RMSProp, Adam) to improve the convergence rate of SGD

Neural networks

We can extend neural networks to multiclass classification

As for the binary case, we can assume that a neural network $f(\mathbf{x}, \mathbf{\Pi})$ computes a non-linear feature transformation

We can pair the network output with a multiclass logistic regression model

The model defines the class posterior probabilities as

$$P(C = k | \mathbf{W}, \mathbf{b}, \mathbf{\Pi}, \mathbf{x}) = \frac{e^{\mathbf{w}_k^T f(\mathbf{x}, \mathbf{\Pi}) + b_k}}{\sum_{j=1}^K e^{\mathbf{w}_j^T f(\mathbf{x}, \mathbf{\Pi}) + b_j}}$$

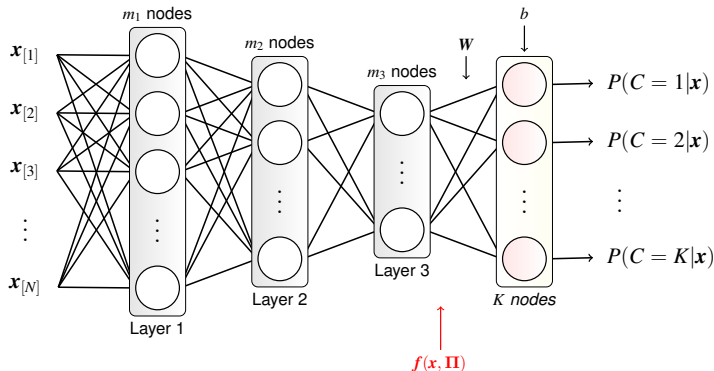
where $\mathbf{W} = [\mathbf{w}_1 \ \dots \ \mathbf{w}_K]$ and K is the number of classes

We can estimate the model parameters by minimizing the cross-entropy

$$\arg \min_{\mathbf{W}, \mathbf{b}, \mathbf{\Pi}} - \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K z_{ik} \log \frac{e^{\mathbf{w}_k^T f(\mathbf{x}, \mathbf{\Pi}) + b_k}}{\sum_{j=1}^K e^{\mathbf{w}_j^T f(\mathbf{x}, \mathbf{\Pi}) + b_j}}$$

Neural networks

As for binary logistic regression, we can cast the whole model as a single neural network \hat{f} that incorporates also the logistic regression parameters \mathbf{W}, \mathbf{b}



Neural networks

The last layer activation function is called softmax activation

The output of each node is computed by applying the softmax function to its inputs

$$\mathbf{x}_m = s(\mathbf{n}_m) = \begin{bmatrix} \frac{e^{n_{m,1}}}{\sum_{i=1}^K e^{n_{m,i}}} \\ \vdots \\ \frac{e^{n_{m,K}}}{\sum_{i=1}^K e^{n_{m,i}}} \end{bmatrix}$$

Note that the layer does not strictly follow the feed-forward network topology we defined earlier, as the output values of the nodes of the last layer depend on the inputs of the other nodes of the layer

However, if we abstract the network at the layers level this does not introduce practical differences

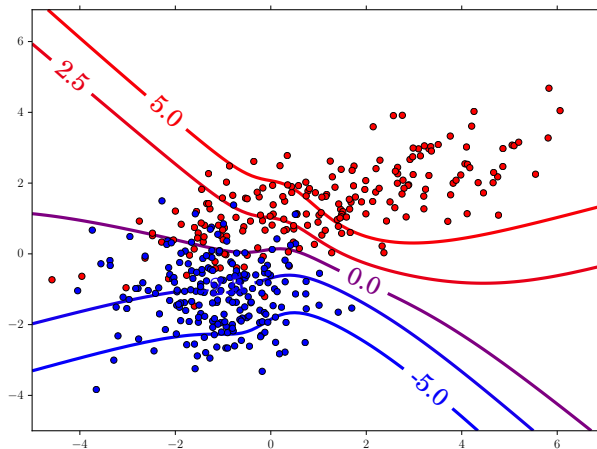
Neural networks

As in the binary case, we can employ back-propagation to compute the gradient of the loss with respect to the model parameters

Stochastic gradient descent or similar optimizers can then be employed to iteratively train the network weight and bias terms

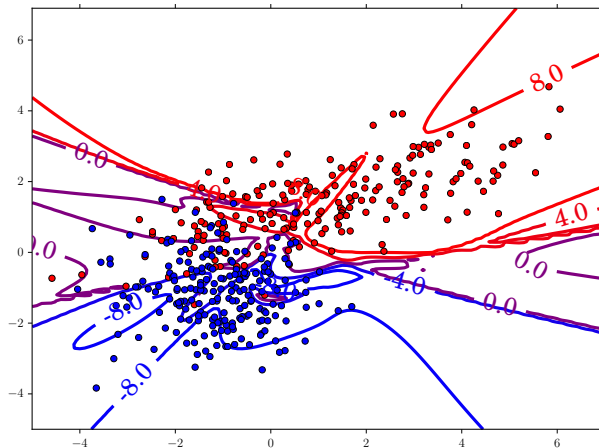
Neural networks

A binary 2D example



Neural networks

Overfitting can be much more dramatic than for linear logistic regression



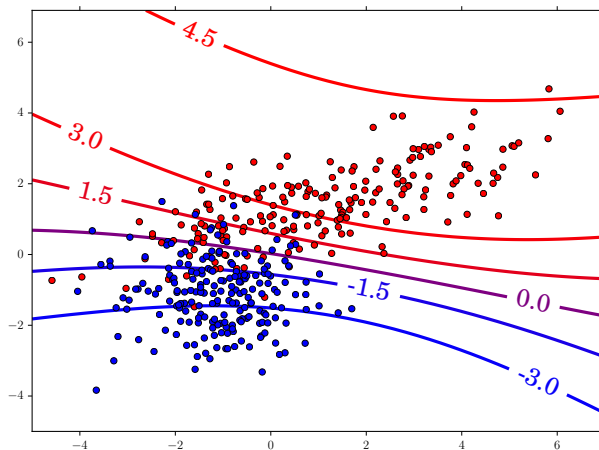
Different regularization strategies can be adopted

- L2 weights regularization
- Dropout
- Early stopping (computing error on validation set)

Neural networks

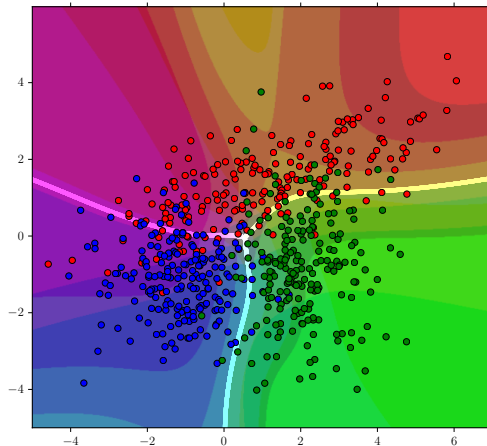
L2 regularization: penalize the squared norm of the weights

$$\arg \min_{W_1, b_1 \dots W_m, b_m} \mathcal{L}(W_1, b_1 \dots W_m, b_m) + \frac{\lambda}{2} \sum_{i=1}^m \|W_m\|^2$$



Neural networks

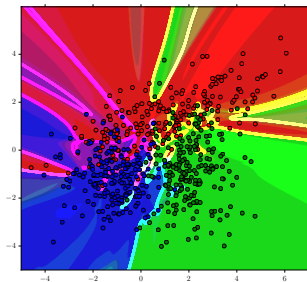
Multiclass (simple network)



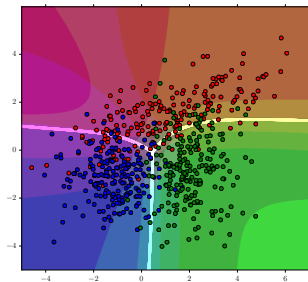
Neural networks

Multiclass (deep network)

No regularization



L2 regularization



Neural networks

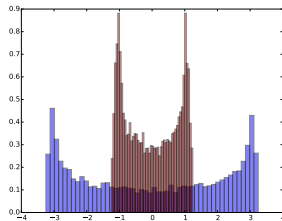
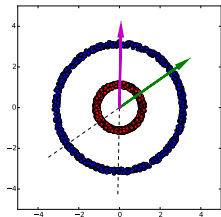
MNIST — Error rates for Neural Networks²

| | No Reg. | L2 ($\lambda = 1e^{-5}$) | Dropout ($p = 0.5$) |
|-----------------------------------|-------------|----------------------------|-----------------------|
| MLP (Tanh) 512–512–512 | 1.9% [1.8%] | 2.0% [1.7%] | 1.5% [1.5%] |
| MLP (ReLU) 512–512–512 | 1.6% [1.5%] | 1.7% [1.6%] | 1.7% [1.4%] |
| MLP (ReLU) 1024–1024–1024–1024 | 1.6% [1.5%] | 1.6% [1.4%] | 1.4% [1.4%] |

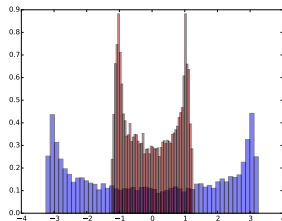
²Training set was split into development (90% of the data) and validation (10% of the data) sets to select the best performing model. The performance of the model with lowest error rate on the test set is shown in brackets.

Neural networks

Linear transformations are not always suited for our data



PCA



LDA

Neural networks

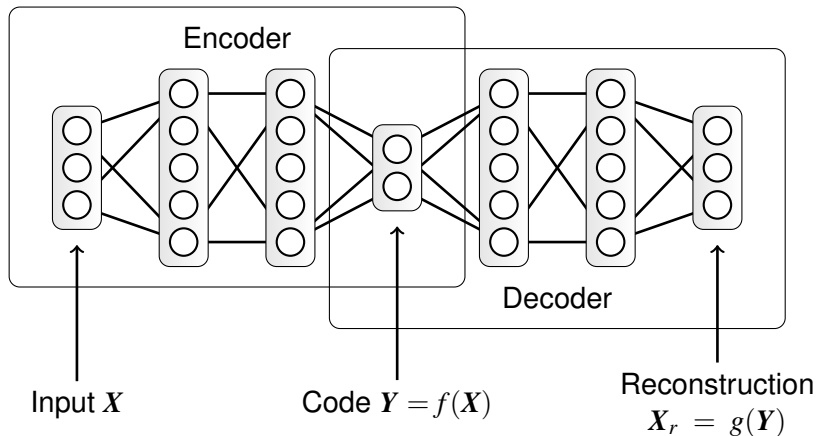
Neural networks can be used to discover the underlying data structure

Autoencoders are similar to PCA

- Unsupervised training (no guarantee that the resulting low-dimensional embedding is useful for classification)
- Minimize reconstruction error
- Used mainly to remove noise from samples (denoising autoencoders)
- Can be used also for dimensionality reduction

Neural networks

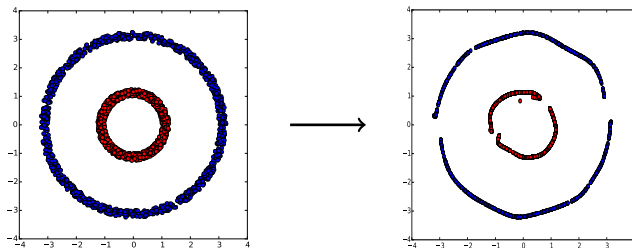
Autoencoder structure:



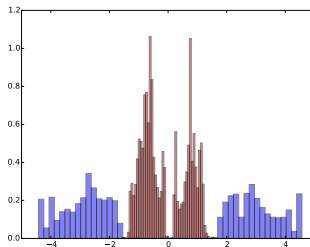
Neural networks

- The input layer has n nodes
- The central hidden layer has $m \ll n$ nodes (bottleneck)
- The central hidden layer acts as a compact representation of the input
- We optimize the network in order to minimize the reconstruction error $\|X - X_r\|^2$
- Denoising autoencoder: provide as input a noisy version of the samples and minimize the reconstruction error with respect to the clean sample

Neural networks



BN layer:



Neural networks

Autoencoder on MNIST — 50 bottleneck nodes

0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9

Train Set



0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9

0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9

Test Set



0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9

Neural networks

Autoencoder on MNIST — 50 bottleneck nodes — sub-sampled

0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9

Train Set



0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9

0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9

Test Set



0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9

Neural networks

Overfitting: A low reconstruction error on the training set *does not guarantee* a low reconstruction error on a different set

The problem is more evident with complex models.

Some strategies can be employed to reduce overfitting issues:

- Early stopping
 - Monitor loss over an held-out validation set
 - Stop the optimization when the error over the validation set starts increasing
- Model regularization (L2, dropout, ...)
- Choose the simplest model that is suitable for the task (Occam's razor)

Neural networks

Ad-hoc architectures have been proposed for different tasks

For example, for image processing dense layers are not too effective, since image characteristics are typically local and may appear in different places in the image

Fully connected layers would require too many parameters

For image processing convolutional networks are typically used

Neural networks

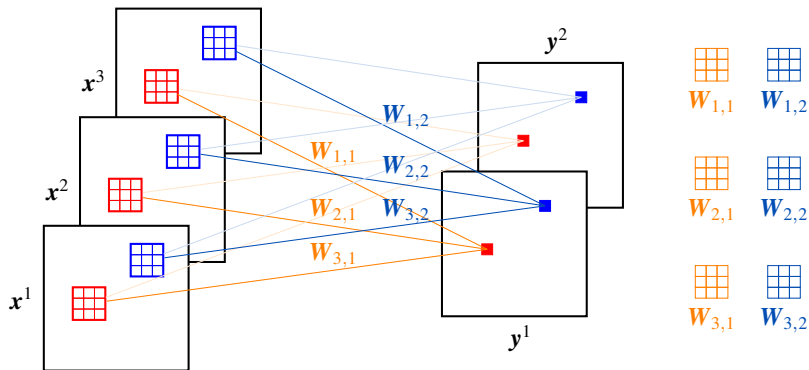
Convolutional networks behave like learned image filters

A convolutional layer receives as input a set of m 2-D channels (e.g., the image color channels) and outputs a set of n 2-D channels

The output is obtained by applying learnable convolutional filters on the different channels

Neural networks

A convolutional layer with 3x3 filters:



Neural networks

The input channels can be interpreted as 2-D maps

The output channels can also be interpreted as 2-D maps

For each combination input/output channel we have a convolutional filter (or, more precisely, a cross-correlation filter) - for example, a 3x3 filter can be

$$\mathbf{W}_{i,j} = \begin{bmatrix} W_{i,j1,1} & W_{i,j1,2} & W_{i,j1,3} \\ W_{i,j2,1} & W_{i,j2,2} & W_{i,j2,3} \\ W_{i,j3,1} & W_{i,j3,2} & W_{i,j3,3} \end{bmatrix}$$

We define the filter width as half the size of the filter matrix minus one, i.e. for a 3x3 filter the width is 1, for a 5x5 filter the width is 2 and so on

Neural networks

The output channel y^c is computed as

$$y_{i,j}^c = h \left(\sum_{k=1}^m \left\langle \mathbf{x}_{[i-l,i+l],[j-l,j+l]}^k, \mathbf{W}_{k,c} \right\rangle + b_c \right), \quad c = 1 \dots n$$

- l is the filter width
- $\mathbf{x}_{[i-l,i+l],[j-l,j+l]}^k$ is the sub-matrix of size $(2l+1) \times (2l+1)$ of the input channel \mathbf{x}^k centered at (i,j)
- b_c is a bias term
- h is a non-linear function (e.g. ReLU or sigmoid)
- $\langle \cdot, \cdot \rangle$ is the inner product $\langle \mathbf{A}, \mathbf{B} \rangle = \sum_{i,j} \mathbf{A}_{ij} \mathbf{B}_{ij}$
- m is the number of input channels, n is the number of output channels

Neural networks

The filters are also called kernels

Cross-correlation is not well defined for elements that are on the boundary of the input channels

To address this, either we compute cross-correlations only for valid inputs, or we add a *padding*, i.e., we employ default values for out-of-boundary elements (e.g. zeros)

We may also avoid computing cross-correlations for all input elements by computing correlations only for positions $(s \cdot i, s \cdot j)$, where s is a constant called stride. In practice, we compute correlations only for 1 in s input rows and columns

Neural networks

Convolutional neural networks typically also employ pooling layers, which aggregate information corresponding to different positions in the input channels

Pooling layers reduce the dimensionality of the input channels, thus increasing the *receptive field* (i.e. the components of the original input feature that may affect the output of the neuron) of a single neuron

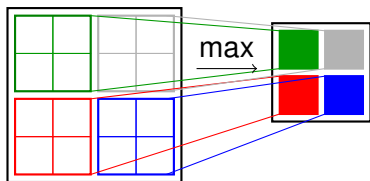
Neural networks

A commonly used pooling layer is the max-pooling layer

The max-pooling layer divides each input channel in $k \times k$ blocks and computes, as output, the maximum of the values of each block

For example, 2×2 pooling layer computes its outputs as

$$y_{i,j}^c = \max(x_{2i,2j}^c, x_{2i+1,2j}^c, x_{2i,2j+1}^c, x_{2i+1,2j+1}^c)$$



Note that the output channels in this case have half the number of rows and columns of the input channels

Neural networks

MNIST — Error rates for Neural Networks³

| | No Reg. | L2 ($\lambda = 1e^{-5}$) | Dropout ($p = 0.5$) |
|-----------------------------------|-------------|----------------------------|-----------------------|
| MLP (Tanh) 512–512–512 | 1.9% [1.8%] | 2.0% [1.7%] | 1.5% [1.5%] |
| MLP (ReLU) 512–512–512 | 1.6% [1.5%] | 1.7% [1.6%] | 1.7% [1.4%] |
| MLP (ReLU) 1024–1024–1024–1024 | 1.6% [1.5%] | 1.6% [1.4%] | 1.4% [1.4%] |
| ConvNet (ReLU) | 1.1% [1.0%] | 1.1% [1.0%] | 0.9% [0.8%] |

³Training set was split into development (90% of the data) and validation (10% of the data) sets to select the best performing model. The performance of the model with lowest error rate on the test set is shown in brackets.

Neural networks

Recent trends have seen networks becoming deeper rather than larger

When training deep networks gradient methods may incur in problems

Typically, back-propagation results in very small values (numerically zeros) for the gradient of the initial layers for networks that employ traditional sigmoid or hyperbolic tangent non linearities

Gradient methods may not be able to progress

Neural networks

Alongside using different activation functions such as ReLU, a typical approach is to introduce residual connection

Rather than computing the input-output transformation function, residual networks compute a *residual* function

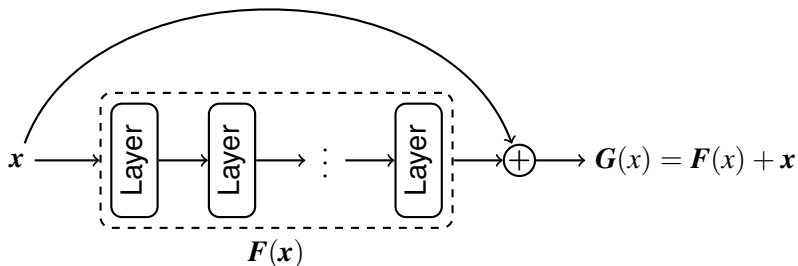
Let $G(x)$ represent the function that we want to compute, with $G : \mathbb{R}^d \rightarrow \mathbb{R}^d$ (note that the function input and output have the same size, as typically happens for convolutional networks)

The residual blocks provide a model for $F(x) = G(x) - x$

Function $G(x)$ can then be computed as $G(x) = F(x) + x$

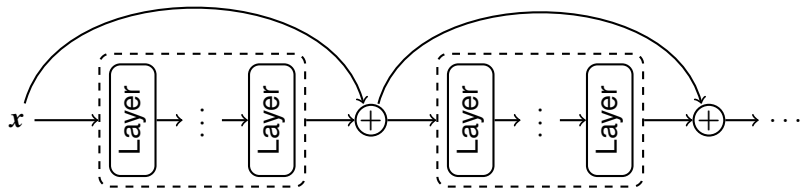
Neural networks

In practice, the effect can be obtained by introducing *skip* connections



Neural networks

We can combine several residual blocks to obtain a deep residual network



Residual blocks allow for effective training of network with hundreds of layers