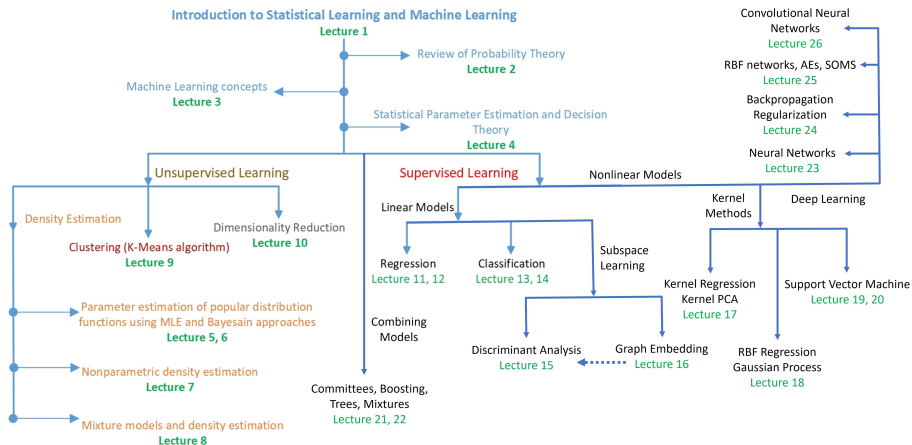


Statistical Learning and Machine Learning

Lecture 23 - Neural Networks 1

October 13, 2021

Course overview and where do we stand



Linear and kernel methods

Linear Basis Function models:

- Useful analytical properties (convex optimization)
- Usually fast to compute
- Manual design of the basis functions

Kernel methods:

- Useful analytical properties (convex optimization)
- Efficient for large D and small N
- Inefficient (even intractable) for large N
- A-priori design of the kernel functions

Why Neural Networks?

Neural networks:

- Adaptive (parametric) basis functions
- Hierarchical data transformations
- Usually more compact (and efficient) than kernel methods of similar performance
- Non-convex optimization
- Neural networks comprising of many hidden layers (*Deep Learning*) have proven to be effective in many machine learning problems

The artificial neuron

The basic building block of a neural network is called *neuron*:

- It receives as input a vector, e.g. $\mathbf{x} \in \mathbb{R}^D$ and applies the following transformation:

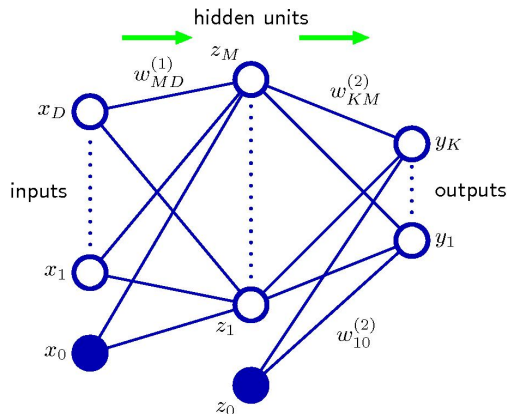
$$\alpha = \sum_{i=1}^D w_i x_i + w_0 = \mathbf{w}^T \mathbf{x} + w_0 \quad (1)$$

$$\mathbf{z} = h(\alpha) \quad (2)$$

where:

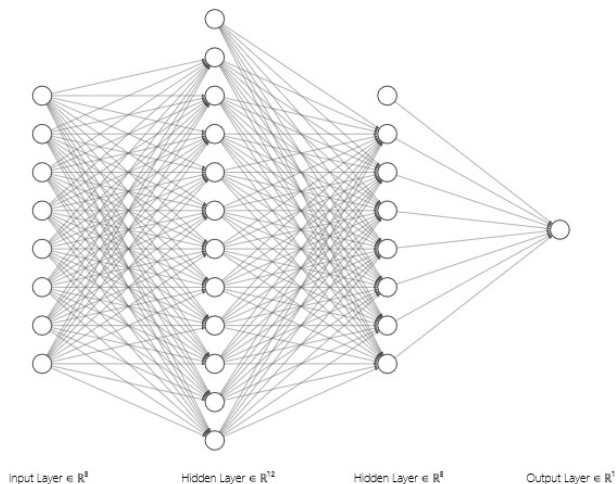
- $\{w, w_0\}$ are the parameters of the neuron
- w is called *weight* and w_0 is called *bias*
- α is known as the *activation*
- $h(\cdot)$ is a nonlinear *activation function*
- $h(\cdot)$ can take many forms, depending on the position of the neuron in the neural network and the problem at hand

A two-layer feed-forward neural network

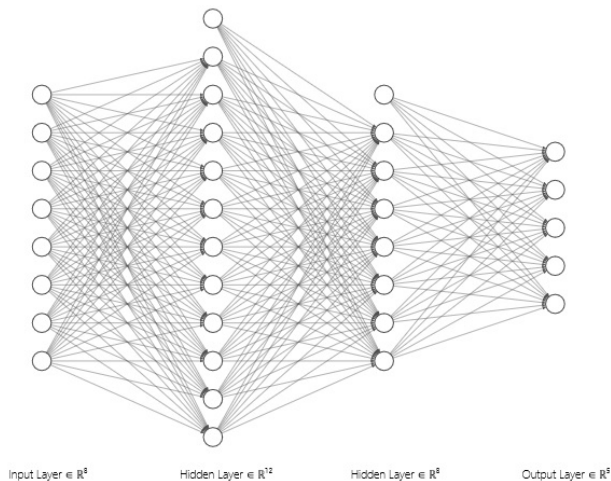


$$y_k(\mathbf{x}, \mathbf{w}) = \sigma \left(\sum_{j=1}^M W_{kj}^{(2)} h \left(\sum_{i=1}^D W_{ji}^{(1)} x_i + w_{j0}^{(1)} \right) + w_{k0}^{(2)} \right) \quad (3)$$

Multi-layer feed-forward neural network



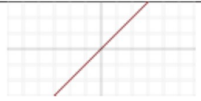
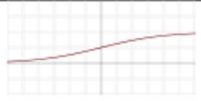
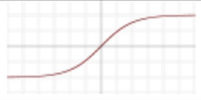


Multi-layer feed-forward neural network



Activation functions

Name	Equation	Derivative
Identity	$f(x) = x$	$f'(x) = 1$
Logistic	$f(x) = \frac{1}{1+e^{-x}}$	$f'(x) = f(x)(1 - f(x))$
HyperbTan	$f(x) = \tanh(x) = \frac{2}{1+e^{-2x}} - 1$	$f'(x) = 1 - f(x)^2$
ArcTan	$f(x) = \tanh^{-1}(x)$	$f'(x) = \frac{1}{x^2+1}$
ReLU	$f(x) = \begin{cases} 0, & \text{if } x \leq 0 \\ x, & \text{if } x > 0 \end{cases}$	$f'(x) = \begin{cases} 0, & \text{if } x \leq 0 \\ 1, & \text{if } x > 0 \end{cases}$
Softmax	$f_i(\mathbf{x}) = \frac{e^{x_k}}{\sum_{l=1}^K e^{x_l}}, \quad k = 1, \dots, K$	$\frac{\partial f_i(\mathbf{x})}{\partial x_j} = f_i(\mathbf{x})(\delta_{ij} - f_j(\mathbf{x}))$

Activation functions

Name	Plot
Identity	
Logistic	
HypTan	
ArcTan	
ReLU	

A two-layer feed-forward neural network

If we 'absorb' the bias parameters in the corresponding weight vectors using additional dimensions (having a value equal to 1) on the input and hidden-layer output vectors:

$$y_k(x, w) = \sigma \left(\sum_{j=0}^M W_{kj}^{(2)} h \left(\sum_{i=0}^D W_{ji}^{(1)} x_i \right) \right) = \sigma \left(W^{(2)T} \underbrace{h(W^{(1)T} \tilde{x})}_{\tilde{z}} \right) \quad (4)$$

where $\tilde{z} \in \mathbb{R}^{M+1}$ and:

$$W^{(1)} = [w_1^{(1)}, \dots, w_M^{(1)}] \in \mathbb{R}^{(D+1) \times M}$$

$$W^{(2)} = [w_1^{(2)}, \dots, w_K^{(2)}] \in \mathbb{R}^{(M+1) \times K}$$

Note that the activation functions are applied *element-wise* (on each dimension)

Multilayer Perceptron

The above neural network is called *Multilayer Perceptron* (or MLP):

- an important property is that the activation functions of all neurons are differentiable w.r.t. their parameters

The use of nonlinear activation functions is crucial:

- If we use linear activation functions in the above two-layer neural network:

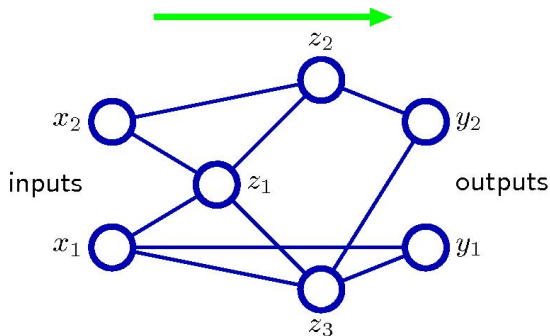
$$y_k(\mathbf{x}, \mathbf{w}) = \mathbf{W}^{(2)T} \mathbf{W}^{(1)T} \tilde{\mathbf{x}} = \mathcal{W}^T \tilde{\mathbf{x}} \quad (5)$$

where $\mathcal{W} = \mathbf{W}^{(2)T} \mathbf{W}^{(1)T} \in \mathbb{R}^{(D+1) \times K}$.

- Thus, any neural network with more than one layers and linear activation functions correspond to an one-layer neural network.

Skip connections

The connections in the network do not need to be restricted between neurons of successive layers (however they need to be feed-forward):

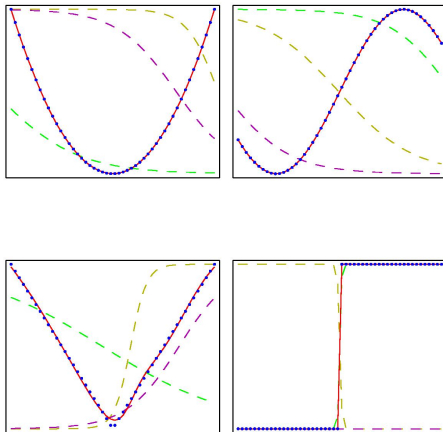


Universal approximators

The approximation properties of feed-forward networks have been widely studied:

- Neural networks are said to be *universal approximators*.
- A two-layer network with linear outputs can uniformly approximate any continuous function on a compact input domain to arbitrary accuracy provided the network has a sufficiently large number of hidden units
- This holds for a wide range of hidden unit activation functions (excluding polynomials)

Universal approximators



(left to right) $f(x) = x^2$, $f(x) = \sin(x)$, $f(x) = |x|$ and $f(x) = H(x)$
Red: approximation of $f(x)$ using a two-layer network with 3 hidden neurons (\tanh act. function - color lines) and linear output

Weight-space symmetries

Let us consider the above two-layer network (having \tanh activation functions):

- If we use $\tilde{W}^{(1)} = -W^{(1)}$, then the output of the hidden layer is $\tilde{\zeta} = -\tilde{z}$ (because $\tanh(-z_i) = -\tanh(z_i)$)
- Using $\tilde{W}^{(2)} = -W^{(2)}$ leads to the same output as before.

There exist multiple combinations of different weight values which can give the same result!

Neural Network Training: Regression

Given a set of data points x_n , $n = 1, \dots, N$, the corresponding target values t_n and including all the weights of the neural network in a parameter w :

- we assume that t follows a Gaussian distribution:

$$p(t|x, w) = \mathcal{N}(t|y(x, w), \beta^{-1}) \quad (6)$$

where β is the precision (inverse variance) of the distribution.

- The likelihood function is (we use linear output neurons):

$$p(t|X, w, \beta) = \prod_{n=1}^N p(t_n|x_n, w, \beta) \quad (7)$$

- The error function is the negative log-likelihood (discarding the terms not depending on w and scaling factors):

$$E(w) = \frac{1}{2} \sum_{n=1}^N \left(y(x_n, w) - t_n \right)^2 \quad (8)$$

Neural Network Training: Regression

Minimizing $E(w)$ will lead to w_{ML} :

- Due to the use of nonlinear activation functions for the hidden neurons, w_{ML} will be a local minimum of $E(w)$ (non-convex optimization)
- Using w_{ML} we can optimize for β by minimizing the negative log-likelihood function:

$$\frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^N \left(y(x_n, w_{ML}) - t_n \right)^2 \quad (9)$$

The optimizations w.r.t. w and β are performed by using an iterative optimization process.

Neural Network Training: Regression

When the targets are vectors \mathbf{t}_n , $n = 1, \dots, N$:

- We obtain \mathbf{w}_{ML} by minimizing:

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \|\mathbf{y}(\mathbf{x}_n, \mathbf{w}) - \mathbf{t}_n\|^2 \quad (10)$$

- We use \mathbf{w}_{ML} to optimize for β_{ML} :

$$\frac{1}{\beta_{ML}} = \frac{1}{NK} \|\mathbf{y}(\mathbf{x}_n, \mathbf{w}_{ML}) - \mathbf{t}_n\|^2 \quad (11)$$

Neural Network Training: Binary classification

Consider that the targets are $t_n \in \{0, 1\}$, where:

- $t_n = 1$ means that $x_n \in \mathcal{C}_1$
- $t_n = 0$ means that $x_n \in \mathcal{C}_2$

The neural network will have one output neuron with logistic sigmoid activation function:

$$y = \sigma(\alpha) = \frac{1}{1 + \exp(-\alpha)} \quad (12)$$

ensuring that $0 \leq y(x, w) \leq 1$.

Neural Network Training: Binary classification

The conditional distribution of t w.r.t. x and w is then a Bernoulli distribution:

$$p(t|x, w) = y(x, w)^t (1 - y(x, w))^{1-t} \quad (13)$$

The negative log-likelihood function becomes the *cross-entropy* error function:

$$E(w) = - \sum_{n=1}^N \left(t_n \ln y(x_n, w) + (1 - t_n) \ln(1 - y(x_n, w)) \right) \quad (14)$$

Neural Network Training: Multi-class classification

The target vectors $\mathbf{t}_n \in \mathbb{R}^K$ follow the 1-of- K coding scheme.

The error function is:

$$E(\mathbf{w}) = - \sum_{n=1}^N \sum_{k=1}^K t_{kn} \ln y_k(\mathbf{x}_n, \mathbf{w}) \quad (15)$$

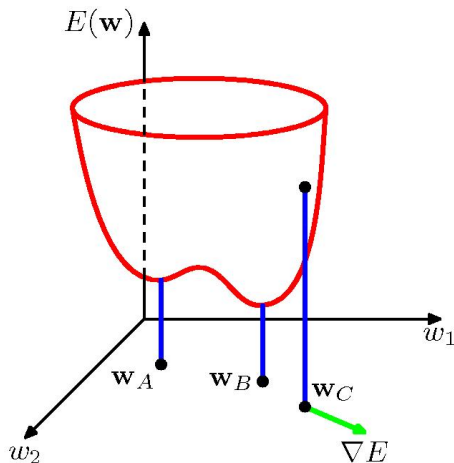
The output neuron activation function is the *softmax* function:

$$y_k(\mathbf{x}, \mathbf{w}) = \frac{\exp(\alpha_k(\mathbf{x}, \mathbf{w}))}{\sum_j \exp(\alpha_j(\mathbf{x}, \mathbf{w}))} \quad (16)$$

ensuring that $0 \leq y_k(\mathbf{x}, \mathbf{w}) \leq 1$ and $\sum_k y_k(\mathbf{x}, \mathbf{w}) = 1$.

Parameter optimization

Local minima correspond to $\nabla E(\mathbf{w}) = 0$.



Parameter optimization

Process:

- 1 Choose an initial set of parameter values $w^{(0)}$
- 2 Update the parameters until reaching a local minimum using:

$$w^{(\tau+1)} = w^\tau + \Delta w^{(\tau)} \quad (17)$$

Several optimization techniques which use different choices for $\Delta w^{(\tau)}$:

- Local quadratic approximation
- Gradient descent optimization
- Error Backpropagation
- Recent schemes: RMSprop, Adagrad, Adadelata, Adam, Adamax, Nadam

Parameter optimization: Local quadratic approximation

Taylor expansion of $E(w)$ around some point \hat{w} :

$$E(w) \cong E(\hat{w}) + (w - \hat{w})^T b + \frac{1}{2} (w - \hat{w})^T H (w - \hat{w}) + \text{higher order terms} \quad (18)$$

where:

$$b \equiv \nabla E|_{w=\hat{w}} \quad \text{and} \quad H = \nabla \nabla E \quad (19)$$

Thus:

$$\nabla E \cong b + H(w - \hat{w}) \quad (20)$$

Parameter optimization: Gradient descent

Gradient descent updates w using:

$$w^{(\tau+1)} = w^\tau - \eta \nabla E(w^{(\tau)}) \quad (21)$$

To find a sufficiently good minimum, it may be necessary to run a gradient-based algorithm multiple times, each time using a different randomly chosen starting point.

When

$$E(w) = \sum_{n=1}^N E_n(w) \quad (22)$$

stochastic gradient descent (SGD) can be used:

$$w^{(\tau+1)} = w^\tau - \eta \nabla E_n(w^{(\tau)}) \quad (23)$$

Mini-batch gradient descent uses small chunks (e.g. 64) data points for each update.