

Neural Network Training

Sargur Srihari

Topics

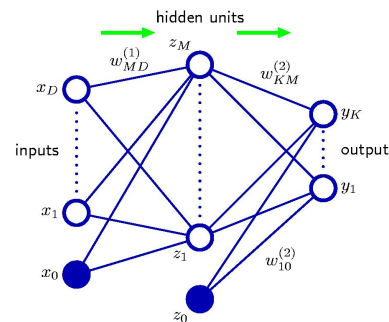
- Neural network parameters
- Probabilistic problem formulation
- Determining the error function
 - Regression
 - Binary classification
 - Multi-class classification
- Parameter optimization
- Local quadratic approximation
- Use of gradient optimization
- Gradient descent optimization

Neural Network parameters

- Linear models for regression and classification can be represented as

$$y(\mathbf{x}, \mathbf{w}) = f\left(\sum_{j=1}^M w_j \phi_j(\mathbf{x})\right)$$

- which are linear combinations of basis functions $\phi_j(\mathbf{x})$
- In a neural network the basis functions $\phi_j(\mathbf{x})$ depend on parameters
 - During training allow these parameters to be adjusted along with the coefficients w_j

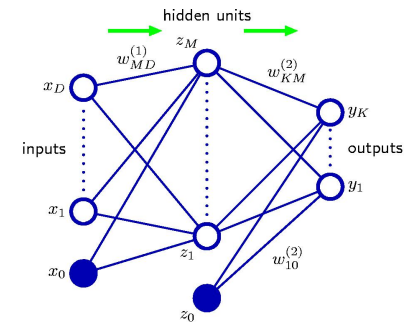


Network Training: Sum of squared errors

- Neural networks perform a transformation
 - vector \mathbf{x} of input variables to vector \mathbf{y} of output variables
 - For sigmoid activation function

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

D input variables
 M hidden units



- Where vector \mathbf{w} consists of all weight and bias parameters
- To determine \mathbf{w} , simple analogy with curve fitting
 - minimize sum-of-squared errors function
 - Given set of input vectors $\{\mathbf{x}_n\}$, $n=1, \dots, N$ and target vectors $\{\mathbf{t}_n\}$ minimize the error function

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

N training vectors

- Consider a more general probabilistic interpretation

Probabilistic View: From activation function f determine Error Function E (as defined by likelihood function)

1. Regression

- f : activation function is identity $y(x, w) = w^T \phi(x) = \sum_{j=1}^M w_j \phi_j(x)$
- E : Sum-of-squares error/Maximum Likelihood $E(w) = \frac{1}{2} \sum_{n=1}^N \|y(x_n, w) - t_n\|^2$

2. (Multiple Independent) Binary Classifications

- f : activation function is Logistic sigmoid $y(x, w) = \sigma(w^T \phi(x)) = \frac{1}{1 + \exp(-w^T \phi(x))}$
- E : Cross-entropy error function $E(w) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$

3. Multiclass Classification

- f : Softmax outputs $y_k(x, w) = \frac{\exp(w_k \phi(x))}{\sum_j \exp(w_j \phi(x))}$
- E : Cross-entropy error function $E(w) = -\sum_{n=1}^N \sum_{k=1}^K t_{kn} \ln y_k(x_n, w)$

1. Probabilistic View: Regression

- Output is a single target variable t that can take any real value
- Assuming t is Gaussian distributed with an \mathbf{x} -dependent mean

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

- Likelihood function

$$p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^N N(t_n | y(\mathbf{x}_n, \mathbf{w}), \beta^{-1})$$

- Taking negative logarithm, we get the negative log-likelihood

$$\frac{\beta}{2} \sum_{n=1}^N \{y(\mathbf{x}_n, \mathbf{w}) - t_n\}^2 - \frac{N}{2} \ln \beta + \frac{N}{2} \ln(2\pi)$$

- which can be used to learn parameters \mathbf{w} and β

Regression Error Function

- Likelihood Function could be used to learn parameters w and β
 - Usually done in a Bayesian treatment
- In neural network literature minimizing error is used
 - They are equivalent here. Sum of squared errors is

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

- Its smallest value occurs when $\nabla E(w) = 0$
- Since $E(w)$ is non-convex:
 - Solution w_{ML} found using iterative optimization $w^{t+1} = w^t + \underline{\Delta} w^t$
 - Gradient descent (discussed later in this lecture))
 - Another solution is back-propagation
 - Since regression output is same as activation $y_k = a_k$, so

$$\frac{\partial E}{\partial a_k} = y_k - t_k$$

$$a_k = \sum_{i=1}^M w_{ki}^{(2)} x_i + w_{k0}^{(2)} \text{ where } k = 1, \dots, K$$

- Having found w_{ML} the value of β_{ML} can also be found using

$$\frac{\ln \beta / 2\pi}{\beta} = \frac{1}{N} \sum_{n=1}^N \{y(x_n, w_{ML}) - t_n\}^2$$

2. Binary Classification

- Single target variable t where $t=1$ denotes C_1 and $t=0$ denotes C_2
- Consider network with single output whose activation function is logistic sigmoid

$$y = \sigma(a) = \frac{1}{1 + \exp(-a)}$$

- so that $0 \leq y(\mathbf{x}, \mathbf{w}) \leq 1$
- Interpret $y(\mathbf{x}, \mathbf{w})$ as conditional probability $p(C_1 | \mathbf{x})$
- Conditional distribution of targets given inputs

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

Binary Classification Error Function

- Error function is negative log-likelihood which in this case is a Cross-Entropy error function

$$E(\mathbf{w}) = -\sum_{n=1}^N \left\{ t_n \ln y_n + (1 - t_n) \ln(1 - y_n) \right\}$$

- where y_n denotes $y(\mathbf{x}_n, \mathbf{w})$
- Using cross-entropy error function instead of sum of squares leads to faster training and improved generalization

2. K Separate Binary Classifications

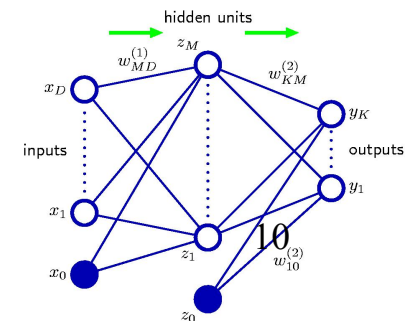
- Network has K outputs each with a logistic sigmoid activation function
- Associated with each output is a binary class label t_k

$$p(\mathbf{t} \mid \mathbf{x}, \mathbf{w}) = \prod_{k=1}^K y_k(\mathbf{x}, \mathbf{w})^{t_k} [1 - y_k(\mathbf{x}, \mathbf{w})]^{1-t_k}$$

- Taking negative logarithm of likelihood function

$$E(\mathbf{w}) = -\sum_{n=1}^N \sum_{k=1}^K \{t_{nk} \ln y_{nk} + (1 - t_{nk}) \ln(1 - y_{nk})\}$$

where y_{nk} denotes $y_k(\mathbf{x}_n, \mathbf{w})$



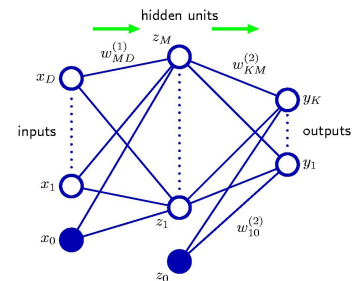
3. Multiclass Classification

- Each input assigned to one of K classes
- Binary target variables have 1-of- K coding scheme

$$t_k \in \{0,1\}$$

- Network outputs are interpreted as $y_k(\mathbf{x}, \mathbf{w}) = p(t_k = 1 | \mathbf{x})$
- Leads to following error function

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



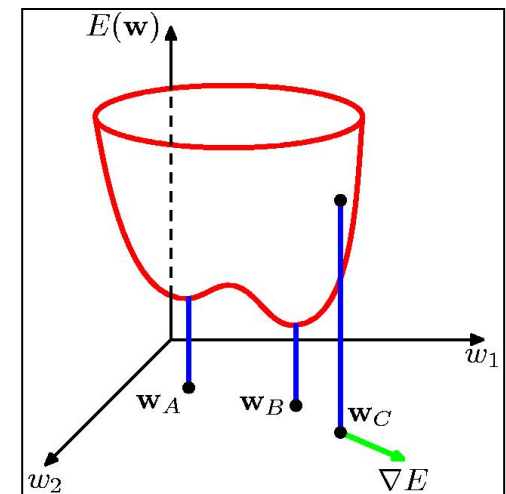
- Output unit activation function is given by softmax

$$y_k(\mathbf{x}, \mathbf{w}) = \frac{\exp(a_k(\mathbf{x}, \mathbf{w}))}{\sum_j \exp(a_j(\mathbf{x}, \mathbf{w}))}$$

$$y_k(\mathbf{x}, \mathbf{w}) = \frac{\exp(w_k \phi(\mathbf{x}))}{\sum_j \exp(w_j \phi(\mathbf{x}))}$$

Parameter Optimization

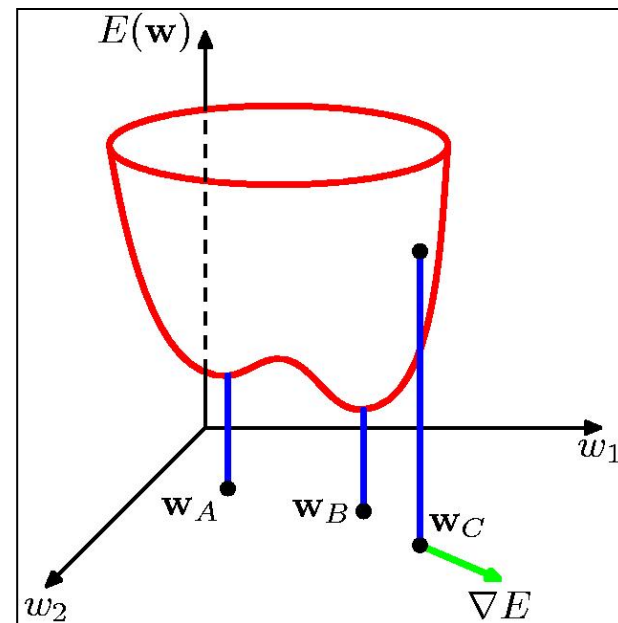
- Task: Find weight vector \mathbf{w} which minimizes the chosen function $E(\mathbf{w})$
- Geometrical picture of error function
- Error function has a highly nonlinear dependence



Parameter Optimization: Geometrical View

$E(\mathbf{w})$: surface sitting over weight space

- \mathbf{w}_A : a local minimum
- \mathbf{w}_B global minimum
- Need to find minimum
- At point \mathbf{w}_C local gradient
 - is given by vector $\nabla E(\mathbf{w})$
 - points in direction of greatest rate of increase of $E(\mathbf{w})$
 - Negative gradient points to rate of greatest decrease



Finding \mathbf{w} where $E(\mathbf{w})$ is smallest

Small step from \mathbf{w} to $\mathbf{w} + \delta\mathbf{w}$ leads to

- change in error function

$$\delta E \approx \delta\mathbf{w}^T \nabla E(\mathbf{w})$$

- Minimum of $E(\mathbf{w})$ will occur when

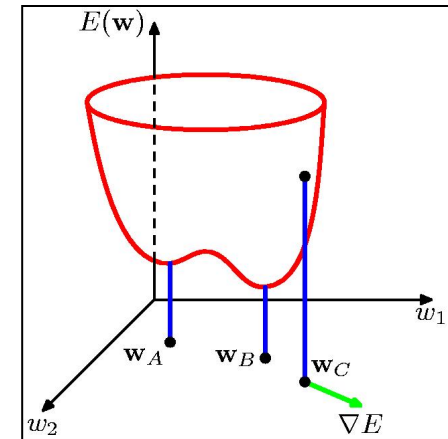
$$\nabla E(\mathbf{w}) = 0$$

- Points at which gradient vanishes are stationary points: minima, maxima, saddle

Complex surface

No hope of finding analytical solution to equation

$$\nabla E(\mathbf{w}) = 0$$



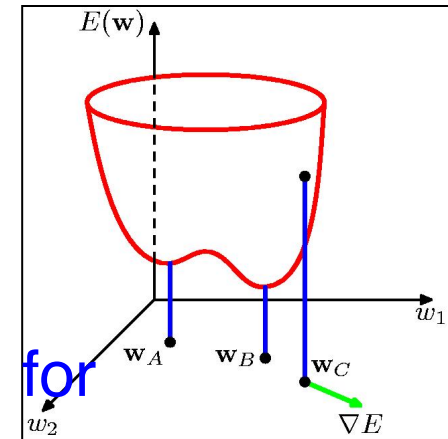
Iterative Numerical Procedure for Minima

- Since there is no analytical solution choose initial $\mathbf{w}^{(0)}$ and update it using

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

where τ is the iteration step

- Different algorithms involve different choices for weight vector update $\Delta \mathbf{w}^{(\tau)}$



- Weight vector update $\Delta \mathbf{w}^{(\tau)}$ is usually based on gradient $\nabla E(\mathbf{w})$ evaluated at weight vector $\mathbf{w}^{(\tau+1)}$
- To understand importance of gradient information consider Taylor's series expansion of error function

Leads to local quadratic approximation

Discussion Overview

Preliminary concepts for Backpropagation

1 Local quadratic approximation

- Provides insight into optimization problem
- $O(W^3)$ where W is dimensionality of w
- Based on Taylor's series: $f(x)$ is approximated

$$f(a) + \frac{f'(a)}{1!}(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \frac{f^{(3)}(a)}{3!}(x-a)^3 + \dots$$

2 Use of gradient information

- Leads to significant improvements in speed of locating minima of error function
- Backpropagation is $O(W^2)$

3 Gradient descent optimization

- Simplest approach of using gradient information

Definitions of Gradient and Hessian

- First derivative of a scalar function $E(\mathbf{w})$ with respect to a vector $\mathbf{w}=[w_1, w_2]^T$ is a vector called the *Gradient* of $E(\mathbf{w})$

$$\nabla E(\mathbf{w}) = \frac{d}{d\mathbf{w}} E(\mathbf{w}) = \begin{bmatrix} \frac{\partial E}{\partial w_1} \\ \frac{\partial E}{\partial w_2} \end{bmatrix}$$

If there are M elements in the vector
then Gradient is a $M \times 1$ vector

- Second derivative of $E(\mathbf{w})$ is a matrix called the *Hessian*

$$H = \nabla \nabla E(\mathbf{w}) = \frac{d^2}{d\mathbf{w}^2} E(\mathbf{w}) = \begin{bmatrix} \frac{\partial^2 E}{\partial w_1^2} & \frac{\partial^2 E}{\partial w_1 \partial w_2} \\ \frac{\partial^2 E}{\partial w_2 \partial w_1} & \frac{\partial^2 E}{\partial w_2^2} \end{bmatrix}$$

Hessian is a matrix with
 M^2 elements

1. Local Quadratic Optimization

- Taylor's Series Expansion of $E(\mathbf{w})$ around some point $\hat{\mathbf{w}}$ in weight space (with cubic and higher terms omitted)

$$E(\mathbf{w}) \cong E(\hat{\mathbf{w}}) + (\mathbf{w} - \hat{\mathbf{w}})^T \mathbf{b} + \frac{1}{2} (\mathbf{w} - \hat{\mathbf{w}})^T \mathbf{H} (\mathbf{w} - \hat{\mathbf{w}})$$

where \mathbf{b} is the gradient of E evaluated at $\hat{\mathbf{w}}$ $\mathbf{b} \equiv \nabla E|_{\mathbf{w}=\hat{\mathbf{w}}}$

- \mathbf{b} is a vector of W elements

\mathbf{H} is the Hessian matrix $\mathbf{H} = \nabla \nabla E$ with elements $(\mathbf{H})_{ij} = \left. \frac{\partial^2 E}{\partial w_i \partial w_j} \right|_{\mathbf{w}=\hat{\mathbf{w}}}$

- \mathbf{H} is a $W \times W$ matrix

- Consider local quadratic approximation around \mathbf{w}^* , a minimum of the error function

$$E(\mathbf{w}) \cong E(\mathbf{w}^*) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^*)^T \mathbf{H} (\mathbf{w} - \mathbf{w}^*)$$

- where \mathbf{H} is evaluated at \mathbf{w}^* and the linear term vanishes

- Let us interpret this geometrically

- Consider eigen value equation for the Hessian matrix $\mathbf{H} \mathbf{u}_i = \lambda_i \mathbf{u}_i$

- where the eigen vectors \mathbf{u}_i are orthonormal $\mathbf{u}_i^T \mathbf{u}_j = \delta_{ij}$

- Expand $(\mathbf{w} - \mathbf{w}^*)$ as a linear combination the eigenvectors $\mathbf{w} - \mathbf{w}^* = \sum_i \alpha_i \mathbf{u}_i$

Error Function Approximation by a Quadratic

Error Functions

- Linear Regression $y(x, w) = w^T \phi(x) = \sum_{j=1}^M w_j \phi_j(x)$

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

- Binary Classification $y(x, w) = \sigma(w^T \phi(x)) = \frac{1}{1 + \exp(-w^T \phi(x))}$

$$E(w) = - \sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$$

- Multiclass Classification $y_k(x, w) = \frac{\exp(w_k \phi(x))}{\sum_j \exp(w_j \phi(x))}$

$$E(w) = - \sum_{n=1}^N \sum_{k=1}^K t_{kn} \ln y_k(x_n, w)$$

Approximated by quadratic

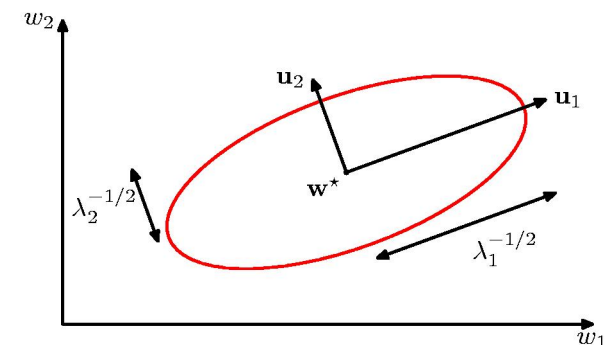
$$E(w) \cong E(w^*) + \frac{1}{2} (w - w^*)^T H (w - w^*)$$

Where $H = \nabla \nabla E$ is a $W \times W$ matrix

Whose contours of constant error are ellipses

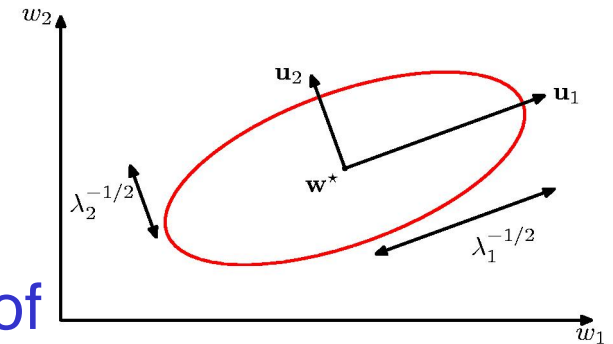
with axes aligned with eigen vectors u_i of H whose

lengths are inversely proportional to sq roots of eigenvalues



Neighborhood of a minimum \mathbf{w}^*

- $\mathbf{w} - \mathbf{w}^*$ is a coordinate transformation
 - Origin is translated to \mathbf{w}^*
 - Axes rotated to align with eigenvectors of Hessian



- Error function can be written as

$$E(\mathbf{w}) = E(\mathbf{w}^*) + \frac{1}{2} \sum_i \lambda_i \alpha_i^2$$

- Matrix $\mathbf{H} = \nabla \nabla E$ is positive definite iff

$$\mathbf{v}^T \mathbf{H} \mathbf{v} > 0 \text{ for all } \mathbf{v}$$

- Since eigenvectors form a complete set $\mathbf{v} = \sum_i c_i \mathbf{u}_i$

- Then an arbitrary vector \mathbf{v} can be written as $\mathbf{v}^T \mathbf{H} \mathbf{v} = \sum_i c_i^2 \lambda_i$

- The stationary point \mathbf{w}^* will be a minimum if the Hessian matrix is positive definite (or all its eigenvalues are positive)

Condition for a point \mathbf{w}^* to be a minimum

- For a one-dimensional weight space, a stationary point \mathbf{w}^* will be minimum if

$$\left. \frac{\partial^2 E}{\partial w^2} \right|_{\mathbf{w}^*} > 0$$

- Corresponding result in D dimensions is that the Hessian matrix evaluated at \mathbf{w}^* is positive definite
 - A matrix \mathbf{H} is positive definite iff $\mathbf{v}^T \mathbf{H} \mathbf{v} > 0$ for all \mathbf{v}

Complexity of Quadratic Approximation

$$E(\mathbf{w}) \cong E(\hat{\mathbf{w}}) + (\mathbf{w} - \hat{\mathbf{w}})^T \mathbf{b} + \frac{1}{2} (\mathbf{w} - \hat{\mathbf{w}})^T \mathbf{H} (\mathbf{w} - \hat{\mathbf{w}})$$

where \mathbf{b} is the gradient of E evaluated at $\hat{\mathbf{w}}$

- \mathbf{b} is a vector of W elements $\mathbf{b} \equiv \nabla E|_{\mathbf{w}=\hat{\mathbf{w}}}$

\mathbf{H} is the Hessian matrix with elements

- \mathbf{H} is a $W \times W$ matrix $\mathbf{H} = \nabla \nabla E$

- Error surface is specified by \mathbf{b} and \mathbf{H}
- They contain total of $W(W+3)/2$ independent elements
 - W is total number of adaptive parameters in network
- Minimum depends on $O(W^2)$ parameters
- Need to perform $O(W^2)$ function evaluations, each requiring $O(W)$ steps.
- Computational effort needed is $O(W^3)$
- W in a 10 x 10 x 10 network needs 100+100=200 weights which means 8 million steps

2. Use of Gradient Information

- Gradient of error function can be evaluated efficiently using back-propagation
 - Use of gradient information can lead to significant improvements in speed with which minimum of error function can be located
- In quadratic approximation to error function
 - Computational effort needed is $O(W^3)$
- By using gradient information minimum can be found in $O(W^2)$ steps
 - 4,000 steps for 10 x 10 x 10 network with 200 weights

3. Gradient Descent Optimization

- Simplest approach to using gradient information
- Take a small step in the direction of the negative gradient

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

- η is the learning rate
- There are batch and on-line versions

Summary

- Neural network parameters have many parameters
 - can be determined analogous to linear regression parameters
- Probabilistic formulation leads to appropriate error functions for linear regression, binary and multi-class classification
- Parameter optimization can be viewed as minimizing error function in weight space
- At the minimum Hessian is positive definite
- Local quadratic optimization needs $O(W^3)$ steps
- Using gradient information more efficient $O(W^2)$ algorithm can be designed

