

Statistical Machine Learning

Christian Walder

Machine Learning Research Group
CSIRO Data61

and

College of Engineering and Computer Science
The Australian National University

Canberra
Semester One, 2020.

(Many figures from C. M. Bishop, "Pattern Recognition and Machine Learning")



Outlines

Overview
Introduction
Linear Algebra
Probability
Linear Regression 1
Linear Regression 2
Linear Classification 1
Linear Classification 2
Kernel Methods
Sparse Kernel Methods
Mixture Models and EM 1
Mixture Models and EM 2
Neural Networks 1
Neural Networks 2
Principal Component Analysis
Autoencoders
Graphical Models 1
Graphical Models 2
Graphical Models 3
Sampling
Sequential Data 1
Sequential Data 2



Part VIII

Neural Networks 2

Review

Error Backpropagation

*Regularisation in Neural
Networks*

*Bayesian Neural
Networks*

Neural Networks 2: Overview



- Recall: we would like gradients w.r.t. parameters so that we can optimise.
- Today: gradients of neural network parameters via the **backpropagation** of gradients algorithm.
- Regularisation/model selection.
- Incorporating invariances/domain knowledge.
- Bayesian neural network (Laplace's method).

Review

Error Backpropagation

*Regularisation in Neural
Networks*

*Bayesian Neural
Networks*

Good News

We study back propagation for pedagogical reasons: in practice one uses automatic differentiation which is far more general and efficient (see *e.g.* the especially easy to use **PyTorch**).



- The composition of two functions is given by

$$f \circ g(x) = f(g(x))$$

- Let f and g be differentiable functions with derivatives f' and g' respectively

- Chain rule

$$\frac{d}{dx}f(g(x)) = f'(g(x))g'(x)$$

- If we write $u = g(x)$ and $y = f(u)$,

$$\frac{dy}{dx} = \frac{dy}{du} \cdot \frac{du}{dx}$$

- Multivariate case we also need is the **total derivative**, e.g.

$$\frac{d}{dt}f(x(t), y(t)) = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt},$$

Review

Error Backpropagation

Regularisation in Neural
Networks

Bayesian Neural
Networks



- Goal: Efficiently update the weights in order to find a local minimum of some error function $E(\mathbf{w})$ utilizing the gradient of the error function.
- Core ideas :
 - 1 Propagate the errors backwards through the network to efficiently calculate the gradient.
 - 2 Update the weights using the calculated gradient.
- Sequential procedure : Calculate gradient and update weights for each data/target pair.
- Batch procedure : Collect gradient information for all data/target pairs for the same weight setting. Then adjust the weights.
- Main question in both cases: How to calculate the gradient of $E(\mathbf{w})$ given one data/target pair?

Review

Error Backpropagation

Regularisation in Neural
Networks

Bayesian Neural
Networks



- Assume the error is a sum over errors for each data/target pair

$$E(\mathbf{w}) = \sum_{n=1}^N E_n(\mathbf{w}).$$

- After applying input \mathbf{x}_n to the network, we get the output \mathbf{y}_n and calculate the error $E_n(\mathbf{w})$.
- What is the gradient for one such term $E_n(\mathbf{w})$?
- Note : In the following, we will drop the subscript n in order to unclutter the equations.
- Notation: Input pattern is \mathbf{x} .
Scalar x_i is the i^{th} component of the input pattern \mathbf{x} .

Review

Error Backpropagation

Regularisation in Neural
Networks

Bayesian Neural
Networks



Review

Error Backpropagation

Regularisation in Neural
NetworksBayesian Neural
Networks

Backprop - One Layer - Scalar View

- Simple linear model **without** hidden layers
- One layer only, identity function as activation function!

$$y_k = \sum_l w_{kl} x_l,$$

and error after applying input \mathbf{x}_n

$$E_n(\mathbf{w}) = \frac{1}{2} \sum_k (y_k - t_k)^2.$$

- The gradient with respect to w_{ji} is now

$$\begin{aligned} \frac{\partial E_n(\mathbf{w})}{\partial w_{ji}} &= \sum_k (y_k - t_k) \frac{\partial}{\partial w_{ji}} y_k = \sum_k (y_k - t_k) \frac{\partial}{\partial w_{ji}} \sum_l w_{kl} x_l \\ &= \sum_k (y_k - t_k) \sum_l x_l \delta_{jk} \delta_{il} \\ &= (y_j - t_j) x_i. \end{aligned}$$

Backprop - One Layer - Vector Calculus



- Vector setup:

$$\begin{aligned}\mathbf{y} &= \mathbf{W} \mathbf{x} \\ \mathbf{W} &\in \mathbb{R}^{D_2 \times D_1} \\ \mathbf{x} &\in \mathbb{R}^{D_1} \\ \mathbf{y} &\in \mathbb{R}^{D_2}\end{aligned}$$

- Error after applying input training pair (\mathbf{x}, \mathbf{t})

$$E_n(\mathbf{W}) = \frac{1}{2} \|\mathbf{y} - \mathbf{t}\|^2.$$

- Using the vector calculus rules gives

$$\begin{aligned}\nabla_{\mathbf{W}} E_n(\mathbf{W}) &= \nabla_{\mathbf{W}} \frac{1}{2} \|\mathbf{y} - \mathbf{t}\|^2 \\ &= (\mathbf{y} - \mathbf{t}) \nabla_{\mathbf{W}} \mathbf{y} \\ &= (\mathbf{y} - \mathbf{t}) \mathbf{x}^T.\end{aligned}$$

Review

Error Backpropagation

Regularisation in Neural
Networks

Bayesian Neural
Networks



Review

Error Backpropagation

Regularisation in Neural Networks

Bayesian Neural Networks

*FYI Only:**Backprop - One Layer - Directional Derivative*

- Do the same using the directional derivative:

$$\mathbf{y} = \mathbf{W} \mathbf{x} \quad \mathbf{W} \in \mathbb{R}^{D_2 \times D_1},$$

and error after applying input training pair (\mathbf{x}, \mathbf{t})

$$E_n(\mathbf{W}) = \frac{1}{2}(\mathbf{y} - \mathbf{t})^\top (\mathbf{y} - \mathbf{t}) = \frac{1}{2}(\mathbf{W}\mathbf{x} - \mathbf{t})^\top (\mathbf{W}\mathbf{x} - \mathbf{t}).$$

- Define $\nabla_d f(\mathbf{x}) = \frac{d}{dh} f(\mathbf{x} + h\mathbf{d})$.
- Relate this to the gradient by $\nabla_d f(\mathbf{x}) = \langle \nabla f(\mathbf{x}), \mathbf{d} \rangle$.
- The directional derivative with respect to \mathbf{W} is now

$$\nabla_\xi E_n(\mathbf{W}) = \frac{1}{2} ((\xi \mathbf{x})^\top (\mathbf{y} - \mathbf{t}) + (\mathbf{y} - \mathbf{t})^\top \xi \mathbf{x}) = \mathbf{x}^\top \xi^\top (\mathbf{y} - \mathbf{t})$$

- With canonical inner product $\langle A, B \rangle = \text{tr} \{A^\top B\}$ the gradient of $E_n(\mathbf{W})(\xi)$ is

$$DE_n(\mathbf{W})(\xi) = \text{tr} \left\{ \underbrace{\mathbf{x}^\top \xi^\top (\mathbf{y} - \mathbf{t})}_{\text{scalar}} \right\} = \text{tr} \left\{ \xi^\top \underbrace{(\mathbf{y} - \mathbf{t}) \mathbf{x}^\top}_{\text{gradient}} \right\}$$



- The gradient

$$\nabla_{\mathbf{W}} E_n(\mathbf{W}) = (\mathbf{y} - \mathbf{t}) \mathbf{x}^\top$$

or in components

$$\frac{\partial E_n(\mathbf{w})}{\partial w_{ji}} = (y_j - t_j) x_i.$$

looks like the product of the output error $(y_j - t_j)$ with the input x_i associated with an edge for w_{ji} in the network diagram.

- Can we generalise this idea to nonlinear activation functions?

Review

Error Backpropagation

*Regularisation in Neural
Networks*

*Bayesian Neural
Networks*



- Now consider a network with **nonlinear** activation functions $h(\cdot)$ composed with the sum over the inputs z_i in one layer and z_j in the next layer connected by edges with weights w_{ji}

$$a_j = \sum_i w_{ji} z_i$$
$$z_j = h(a_j).$$

- Use the **chain rule** to calculate the gradient

$$\frac{\partial E_n(\mathbf{w})}{\partial w_{ji}} = \frac{\partial E_n(\mathbf{w})}{\partial a_j} \frac{\partial a_j}{\partial w_{ji}} = \delta_j z_i,$$

where we defined the **error** (a slight misnomer hailing from the derivative of the squared error) $\delta_j = \frac{\partial E_n(\mathbf{w})}{\partial a_j}$

- Same intuition as before: gradient is output error times the input associated with the edge for w_{ji} .

Review

Error Backpropagation

Regularisation in Neural
Networks

Bayesian Neural
Networks



- Need to calculate the errors δ in **every** layer.

$$\frac{\partial E_n(\mathbf{w})}{\partial w_{ji}} = \delta_j z_i \qquad \delta_j = \frac{\partial E_n(\mathbf{w})}{\partial a_j}$$

- Start the recursion; for output units with squared error:

$$\delta_k = y_k - t_k.$$

- For the hidden units we use the **total derivative**, e.g.

$$\frac{d}{dt} f(x(t), y(t)) = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt},$$

to calculate

$$\delta_j = \frac{\partial E_n(\mathbf{w})}{\partial a_j} = \sum_k \frac{\partial E_n(\mathbf{w})}{\partial a_k} \frac{\partial a_k}{\partial a_j} = \sum_k \delta_k \frac{\partial a_k}{\partial a_j},$$

using the definition of δ_k .

Review

Error Backpropagation

Regularisation in Neural
Networks

Bayesian Neural
Networks



- Express a_k as a function of the incoming a_j

$$a_k = \sum_j w_{kj} z_j = \sum_j w_{kj} h(a_j),$$

- and differentiate

$$\frac{\partial a_k}{\partial a_j} = w_{kj} \frac{\partial h(a_j)}{\partial a_j} = w_{kj} \frac{\partial h(s)}{\partial s} \Bigg|_{s=a_j} = w_{kj} h'(a_j).$$

- Finally, we get for the error in the previous layer

$$\delta_j = h'(a_j) \sum_k w_{kj} \delta_k.$$

Review

Error Backpropagation

Regularisation in Neural
Networks

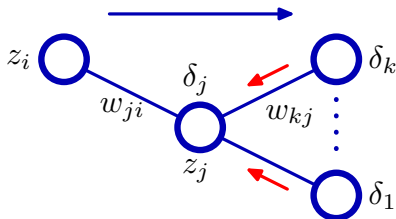
Bayesian Neural
Networks



- The **backpropagation formula**

$$\delta_j = h'(a_j) \sum_k w_{kj} \delta_k.$$

- Form of $h'(\cdot)$ is available, because we choose the activation function $h(\cdot)$ to be a fixed function such as a \tanh *etc.*.



Review

Error Backpropagation

Regularisation in Neural
Networks

Bayesian Neural
Networks

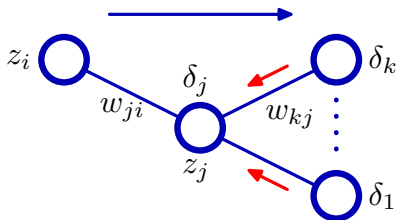
Error Backpropagation Algorithms



- 1 Apply the input vector \mathbf{x} to the network and forward propagate through the network to calculate all activations and outputs of each unit.
- 2 Compute the gradients of the error at the output.
- 3 Backpropagate the gradients backwards through the network using the **backpropagation formula**.
- 4 Calculate all components of ∇E_n by

$$\frac{\partial E_n(\mathbf{w})}{\partial w_{ji}} = \delta_j z_i$$

- 5 Update the weights \mathbf{w} using $\frac{\partial E_n(\mathbf{w})}{\partial w_{ji}}$.





- For batch processing, we repeat backpropagation for each pattern in the training set and then sum over all patterns

$$\frac{\partial E(\mathbf{w})}{\partial w_{ji}} = \sum_{n=1}^N \frac{\partial E_n(\mathbf{w})}{\partial w_{ji}}$$

- Backpropagation can be generalised by assuming that each node has a different activation function $h(\cdot)$.

Review

Error Backpropagation

Regularisation in Neural
Networks

Bayesian Neural
Networks

Easy Backprop

Let $\mathbf{z}^{(0)} = \mathbf{x} = \text{input}$

$$\mathbf{a}^{(l)} = \mathbf{W}^{(l)} \mathbf{z}^{(l-1)}$$

$$\mathbf{z}^{(l)} = h(\mathbf{a}^{(l)})$$

$$E = \mathcal{L}(\mathbf{a}^{(L)}) = \mathcal{L}(\mathbf{y}) \stackrel{\text{e.g.}}{=} \frac{1}{2} \|\mathbf{y} - \mathbf{t}\|^2.$$



The gradients of E w.r.t. the parameters are (total derivative)

$$\frac{\partial E}{\partial \mathbf{W}^{(l)}} = \frac{\partial E}{\partial \mathbf{a}^{(l)}} \cdot \frac{\partial \mathbf{a}^{(l)}}{\partial \mathbf{W}^{(l)}} = \boldsymbol{\delta}^{(l)} \mathbf{z}^{(l-1)\top},$$

where (neglecting transposes — assume conformant shapes)

$$\boldsymbol{\delta}^{(l)} \equiv \frac{\partial E}{\partial \mathbf{a}^{(l)}} = \frac{\partial \mathbf{a}^{(l+1)}}{\partial \mathbf{a}^{(l)}} \frac{\partial \mathbf{a}^{(l+2)}}{\partial \mathbf{a}^{(l+1)}} \cdots \frac{\partial \mathbf{a}^{(L)}}{\partial \mathbf{a}^{(L-1)}} \frac{\partial \mathcal{L}(\mathbf{a}^{(L)})}{\partial \mathbf{a}^{(L)}}$$

has the recursion $\boldsymbol{\delta}^{(L)} = \frac{\partial \mathcal{L}(\mathbf{a}^{(L)})}{\partial \mathbf{a}^{(L)}}$ along with

$$\boldsymbol{\delta}^{(l-1)} = \frac{\partial \mathbf{a}^{(l)}}{\partial \mathbf{a}^{(l-1)}} \boldsymbol{\delta}^{(l)}$$

$$\frac{\partial \mathbf{a}^{(l)}}{\partial \mathbf{a}^{(l-1)}} = \frac{\partial \mathbf{W}^{(l)} h(\mathbf{a}^{(l-1)})}{\partial \mathbf{a}^{(l-1)}} = \text{diag}\{h'(\mathbf{a}^{(l-1)})\} \mathbf{W}^{(l)\top}.$$





- For dense weight matrices, the complexity of calculating the gradient $\frac{\partial E_n(\mathbf{w})}{\partial w_{ji}}$ via backpropagation is of $O(W)$ where W is the number of weights.
- Compare this to **numerical differentiation** using e.g.

$$\frac{\partial E_n(\mathbf{w})}{\partial w_{ji}} = \frac{E_n(w_{ji} + \epsilon) - E_n(w_{ji} - \epsilon)}{2\epsilon} + O(\epsilon^2)$$

which needs $O(W^2)$ operations, and is less accurate.

FYI only — as in the previous lecture: In general we have the “cheap gradient principle”. See (Griewank, A., 2000. *Evaluating Derivatives: Principles and Techniques of Algorithmic Differentiation*, Section 5.1).

Review

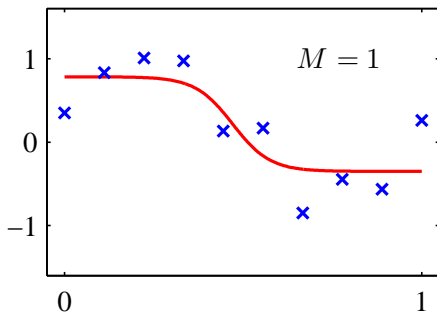
Error Backpropagation

Regularisation in Neural
Networks

Bayesian Neural
Networks



- Number of input and output nodes determined by the application.
- Number of hidden nodes is a free parameter.



Training a two-layer network with 1 hidden node.

Review

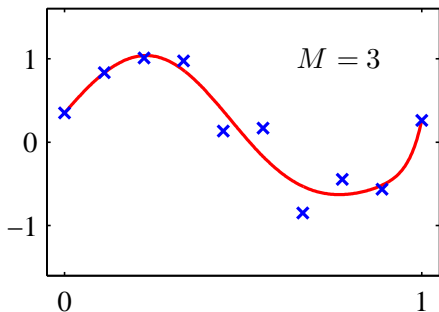
Error Backpropagation

Regularisation in Neural Networks

Bayesian Neural Networks



- Number of input and output nodes determined by the application.
- Number of hidden nodes is a free parameter.



Training a two-layer network with 3 hidden nodes.

Review

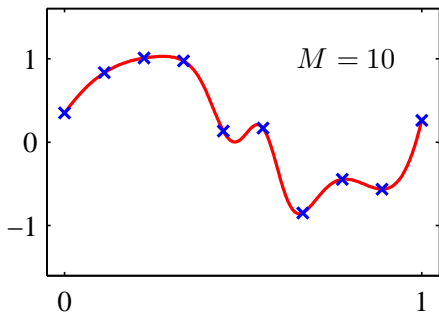
Error Backpropagation

Regularisation in Neural Networks

Bayesian Neural Networks



- Number of input and output nodes determined by the application.
- Number of hidden nodes is a free parameter.



Training a two-layer network with 10 hidden nodes.

Review

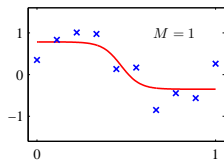
Error Backpropagation

Regularisation in Neural Networks

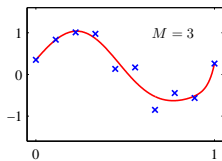
Bayesian Neural Networks



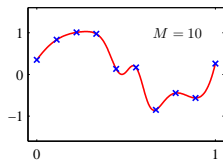
- Model complexity matters again.



$M = 1$



$M = 3$



$M = 10$

- As before, we can use the **regularised error**

$$\tilde{E}(\mathbf{w}) = E(\mathbf{w}) + \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}$$

Review

Error Backpropagation

Regularisation in Neural
Networks

Bayesian Neural
Networks

Regularisation via Early Stopping



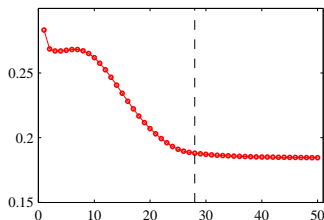
- Stop training at the minimum of the validation set error.

Review

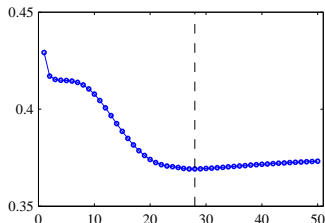
Error Backpropagation

Regularisation in Neural Networks

Bayesian Neural Networks



Training set error.



Validation set error.



- If input data should be invariant with respect to some transformations, we can utilise this for training.
- Use training patterns including these transformations (e.g. handwritten digits translated in the input space).
- Or create extra artificial input data by applying several transformations to the original input data.
- Alternatively, preprocess the input data to remove the transformation.
- Or use **convolutional neural networks** (e.g. in image processing where close pixels are more correlated than far away pixels; therefore extract local features first and later feed into a network extracting higher-order features).

Review

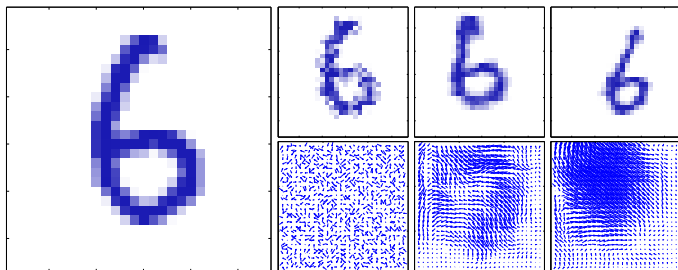
Error Backpropagation

*Regularisation in Neural
Networks*

*Bayesian Neural
Networks*



- Create synthetic data by warping handwritten digits.



Left: Original digitised image. Right : Examples of warped images (above) and their corresponding displacement fields (below).

Review

Error Backpropagation

*Regularisation in Neural
Networks*

*Bayesian Neural
Networks*



- Predict a single target t from a vector of inputs \mathbf{x}
- Assume conditional distribution to be Gaussian with precision β

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

- Prior distribution over weights \mathbf{w} is also assumed to be Gaussian

$$p(\mathbf{w} | \alpha) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I})$$

- For an i.i.d training data set $\{\mathbf{x}_n, t_n\}_{n=1}^N$, the likelihood of the targets $\mathcal{D} = \{t_1, \dots, t_N\}$ is

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

- Posterior distribution

$$p(\mathbf{w} | \mathcal{D}, \alpha, \beta) \propto p(\mathbf{w} | \alpha) p(\mathcal{D} | \mathbf{w}, \beta)$$

Review

Error Backpropagation

Regularisation in Neural
Networks

Bayesian Neural
Networks



- But $y(\mathbf{x}, \mathbf{w})$ is **nonlinear**, and therefore we can no longer calculate the posterior in closed form.
- Use Laplace approximation
 - 1 Find a (local) maximum \mathbf{w}_{MAP} of the posterior via **numerical optimisation**.
 - 2 Evaluate the matrix of second derivatives of the negative log posterior distribution.
- Find a (local) maximum using the log-posterior

$$\ln p(\mathbf{w} | \mathcal{D}, \alpha, \beta) = -\frac{\alpha}{2} \mathbf{w}^\top \mathbf{w} - \frac{\beta}{2} \sum_{n=1}^N (y(\mathbf{x}_n, \mathbf{w}) - t_n)^2 + \text{const}$$

- Find the matrix of second derivatives of the negative log posterior distribution

$$\mathbf{A} = -\nabla \nabla \ln p(\mathbf{w} | \mathcal{D}, \alpha, \beta) = \alpha \mathbf{I} + \beta \mathbf{H}$$

where \mathbf{H} is the Hessian matrix of the sum-of-squares error function with respect to the components of \mathbf{w} .

Review

Error Backpropagation

Regularisation in Neural Networks

Bayesian Neural Networks



- Having \mathbf{w}_{MAP} , and \mathbf{A} , we can approximate the posterior by a Gaussian

$$q(\mathbf{w} | \mathcal{D}, \alpha, \beta) = \mathcal{N}(\mathbf{w} | \mathbf{w}_{MAP}, \mathbf{A}^{-1})$$

- For the predictive distribution further linearly approximate

$$y(\mathbf{x}, \mathbf{w}) \approx y(\mathbf{x}, \mathbf{w}_{MAP}) + \mathbf{g}^\top (\mathbf{w} - \mathbf{w}_{MAP})$$

where $\mathbf{g} = \nabla_{\mathbf{w}} y(\mathbf{x}, \mathbf{w})|_{\mathbf{w}=\mathbf{w}_{MAP}}$.

Then

$$p(t | \mathbf{x}, \mathcal{D}, \alpha, \beta) = \mathcal{N}(t | y(\mathbf{x}, \mathbf{w}_{MAP}), \sigma^2(\mathbf{x}))$$

where

$$\sigma^2(\mathbf{x}) = \beta^{-1} + \mathbf{g}^\top \mathbf{A}^{-1} \mathbf{g}.$$

(Recall the multivariate normal conditionals.)

- variance due to the intrinsic noise on the target: β^{-1}
- variance due to the model parameter \mathbf{w} : $\mathbf{g}^\top \mathbf{A}^{-1} \mathbf{g}$

Review

Error Backpropagation

Regularisation in Neural
Networks

Bayesian Neural
Networks