Pattern Information Processing^{1,61} Neural Networks

Masashi Sugiyama (Department of Computer Science)

Contact: W8E-505 <u>sugi@cs.titech.ac.jp</u> http://sugiyama-www.cs.titech.ac.jp/~sugi/

Linear/Non-Linear Models ¹⁶²

Linear model: $f_{\alpha}(x)$ is linear with respect to α (Note: not necessarily linear with respect to x)

$$f_{oldsymbol{lpha}}(oldsymbol{x}) = \sum_{i=1}^b lpha_i arphi_i(oldsymbol{x})$$

Non-linear model: Otherwise

Today's Plan

163

Neural networks

Least-squares in neural networks: Error back-propagation algorithm

Non-Linear Models

A popular choice: A hierarchical model

$$f_{\boldsymbol{w}}(\boldsymbol{x}) = \sum_{i=1}^{b} \alpha_i \varphi(\boldsymbol{x}; \boldsymbol{\beta}_i)$$

$$oldsymbol{w} = (oldsymbol{lpha}^ op, oldsymbol{eta}_1^ op, oldsymbol{eta}_2^ op, \dots, oldsymbol{eta}_b^ op)^ op$$

Basis functions are parameterized by β_i .

Three-Layer Networks

Such a hierarchical model can be represented as a 3-layer network.



Sigmoid Function

166

A typical basis function: Sigmoid function

$$arphi(oldsymbol{x};oldsymbol{eta},h) = rac{1}{1 + \exp\left(-\langle oldsymbol{x},oldsymbol{eta}
ight
angle - h
ight)}$$



Perceptrons

- The behavior of the sigmoid functions is similar to neurons in our brain.
- For this reason, hierarchical models with sigmoid functions are called artificial neural networks or perceptrons.
- Mathematically, a 3-layer neural network can approximate any continuous functions with arbitrary small error ("a universal approximator").

Gaussian Radial Basis Function¹⁶⁸

Another popular basis function: Gaussian radial basis function

$$arphi(oldsymbol{x};oldsymbol{eta},c) = \exp\left(-rac{\|oldsymbol{x}-oldsymbol{eta}\|^2}{2c^2}
ight)$$



Least-Squares Learning
$$f_{w}(x) = \sum_{i=1}^{b} \alpha_{i} \varphi(x; \beta_{i})$$

$$oldsymbol{w} = (oldsymbol{lpha}^ op, oldsymbol{eta}_1^ op, oldsymbol{eta}_2^ op, \dots, oldsymbol{eta}_b^ op)^ op$$

Least-squares learning is often used for training hierarchical models.

$$\min_{\boldsymbol{w}} J_{LS}(\boldsymbol{w})$$

 $J_{LS}(\boldsymbol{w}) = \sum_{i=1}^{n} (f_{\boldsymbol{w}}(\boldsymbol{x}_i) - y_i)^2$

How to Obtain A Solution ¹⁷⁰

No analytic solution is known.

Simple gradient search is usually used.



One of the local minima can be found.

Error Back-Propagation

171

Efficient calculation of gradient for sigmoid basis functions:



Error Back-Propagation (cont.)¹⁷²

- When the output values of the network are calculated, the input points are propagated following the forward path.
- On the other hand, when the gradients are calculated, the output error $(f_w(x_i) y_i)$ is propagated backward.
- For this reason, this algorithm is called the error back-propagation.
- However, it is gradient descent so global convergence is not guaranteed.

Stochastic Gradient Descent ¹⁷³

- In the usual gradient method, all training examples are used at the same time.
- In practice, the following stochastic method would be computationally advantageous.
 - Randomly choose one of the training examples (say, (\boldsymbol{x}_i, y_i))
 - Update the parameter vector by

$$oldsymbol{w}^{new} \longleftarrow oldsymbol{w}^{old} - arepsilon
abla J_i(oldsymbol{w}) = (f_{oldsymbol{w}}(oldsymbol{x}_i) - y_i)^2$$

Repeat this procedure until convergence.

Avoiding Overfitting

$$f_{\boldsymbol{w}}(\boldsymbol{x}) = \sum_{i=1}^{b} \alpha_i \varphi(\boldsymbol{x}; \boldsymbol{\beta}_i)$$

$$oldsymbol{w} = (oldsymbol{lpha}^ op, oldsymbol{eta}_1^ op, oldsymbol{eta}_2^ op, \dots, oldsymbol{eta}_b^ op)^ op.$$

LS overfits to noisy samples.

• Regularization:

$$\min_{oldsymbol{w}}[J_{LS}(oldsymbol{w})+\lambda||oldsymbol{w}||^2]$$

$$J_{LS}(\boldsymbol{w}) = \sum_{i=1}^{n} \left(f_{\boldsymbol{w}}(\boldsymbol{x}_i) - y_i \right)^2 \qquad \lambda > 0$$

 Early stopping: Stop gradient descent before it converges

Special Structure of Neural Networks

175

Parameters and functions are not one-to-one:

$$f_{oldsymbol{w}}(oldsymbol{x}) = \sum_{i=1}^{b} lpha_i arphi(oldsymbol{x};oldsymbol{eta}_i)$$

• If $\alpha_i = 0$, any β_i gives the same function.

• Any permutation of units in the hidden layer gives the same function.

This non-identifiability causes great difficulty in the mathematical analysis of neural networks.

S. Watanabe, Algebraic Geometry and Statistical Learning Theory, Cambridge University Press, 2009

Special Structure of Neural Networks (cont.)

176

- Theory: Neural networks were shown to go together well with Bayesian learning:
 - Given a prior distribution of parameters, compute the posterior distribution of parameters.

 $p(\boldsymbol{w}|\text{Data}) = p(\text{Data}|\boldsymbol{w})p(\boldsymbol{w})$

Practice: Training neural networks from many different initial values and taking their ensemble sometimes work well.

Pre-Training of Hidden Layers¹⁷⁷

- Neural networks are prone to a poor local optimum solution.
- Recent idea ("a deep belief network"):
 - First, pre-training hidden layers one by one in an unsupervised manner.
 - Then all layers are fine-tuned by supervised training.
- The pre-training idea was shown to work well in experiments.

Hinton and Salakhutdinov. Reducing the dimensionality of data with neural networks. Science, 2006.

Homework

178

 Prove that the gradients for sigmoid basis functions are given as



Homework (cont.)

179

- 2. For your own toy 1-dimensional data, perform simulations using
 - 3-layer neural networks with sigmoids/Gaussians
 - Error back-propagation

and analyze the results, e.g., by changing

- Target functions
- Number of training samples
- Noise level
- Number of neurons in the second layer (b)
- Initialization heuristics
- Number of ensembles

Notification of Final Assignment

180

- 1. Apply supervised learning techniques to your data set and analyze it.
- 2. Write your opinion about this course.

 Final report deadline: Aug 3rd (Fri.)
 E-mail submission is also accepted! sugi@cs.titech.ac.jp

Mini-Workshop on Data Mining⁸¹

- On July 10th and 24th, we will have a miniworkshop on data mining.
- Several students present their own data mining results.
- Those who give a talk at the workshop will have very good grades!

Mini-Workshop on Data Mining⁸²

- Application (just to declare that you want to give a presentation) deadline: June 19th.
- Presentation: 10-15 minutes (?).
 - Specification of your dataset
 - Methods used
 - Outcome
- Slides should be in English.
- Better to speak in English, but Japanese is also allowed.