# Alternatives to BackProp

# BackPropagation Basic Principle

#### Faster Alternatives

The Problem The Solution

# Measuring Performance

Generalization

#### Regularization

Limiting the Complexity Punishing Large Weights Optimal Pruning Örjan Ekeberg BackPropagation Faster Alternatives Measuring Performance

gularizatio

ANN fk

## BackPropagation Basic Principle

Faster Alternatives The Problem The Solution

Measuring Performance Generalization

## Regularization

Limiting the Complexity Punishing Large Weights Optimal Pruning Örjan Ekeberg

ANN fk

# BackPropagation

Measuring Performance

# ANN fk

Örjan Ekeberg

BackPropagation Basic Principle

Faster Alternatives

Measuring Performance

Regularizatio

# ANN fk Örjan Ekeberg

BackPropagation Basic Principle Faster Alternatives

> suring formance



- Multilayer feedforward network
- Arbitrary decision boundaries/functions





- Classification
- Function approximation
- Trained with prescribed outputs
- ► Batch or Incremental learning

- ► All computation can be performed locally
- Slow convergence
- Requires short step lengths



ANN fk Örjan Ekeberg

Basic Principle

Faster Alternatives



# Learning:

Minimize the error (*E*) as a function of all weights  $(\vec{w})$ 

- 1. Compute the direction in weight space where the error increases most:  $\operatorname{grad}_{\vec{w}}(E)$
- 2. Change the weights in the opposite direction

 $w_i \leftarrow w_i - \eta \frac{\partial E}{\partial w_i}$ 

BackPropagation Basic Principle

Faster Alternatives The Problem The Solution

Aeasuring Performance Generalization

#### Regularizatio

Limiting the Complexity Punishing Large Weights Optimal Pruning

#### ANN fk

Örjan Ekeberg

#### BackPropagation Faster Alternatives The Problem The Solution Measuring Performance

Regularization

#### ANN fk

Örjan Ekeberg

BackPropagation Basic Principle Faster Alternative:

Measuring Performance

Regularizatio

Does the gradient point in the right direction?



- Incremental learning
- ► Large steps
- High-dimensional space

Idéa: Make use of the second derivative too

Ordinary gradient following

$$\Delta w = -\eta \frac{\partial E}{\partial w}$$

Newtons method

$$\Delta w = \left(\frac{\partial^2 E}{\partial w^2}\right)^{-1} \frac{\partial E}{\partial w}$$

Makes it necessary to invert a very large matrix!

ANN fk Örjan Ekeberg BackPropagation Faster Alternatives The Problem The Solution Measuring Performance Regularization

ANN fk

Örjan Ekeberg

aster Alternative

The Solution

- Normalization
- De-correlation

These techniques only help on the global scale

Works for "toy-problems" but not when there is a lot of structure in the task

## Conjugate Gradient Method

- Established numerical method
- Incremental updates are made in directions where they do not counteract each other
- Does not require explicit computation of the second derivative

#### ANN fk

Örjan Ekeberg

BackPropagation Faster Alternative The Problem The Solution Measuring Performance

#### ANN fk

Örjan Ekeberg

BackPropagation Faster Alternativ The Problem **The Solution** Measuring Performance Conjugate Gradient Method:

- Initiate:
  - $\vec{r} \leftarrow -\frac{\partial E}{\partial \vec{w}}$  $\vec{s} \leftarrow \vec{r}$
- Repeat:
  - Find  $\eta$  which minimizes  $E(\vec{w} + \eta \vec{s})$
  - $\vec{w} \leftarrow \vec{w} + \eta \vec{s}$  $\vec{r} \leftarrow -\frac{\partial E}{\partial \vec{w}}$  $\beta = \max\left[\frac{\vec{r}^T \cdot (\vec{r} - \vec{r}_{old})}{\vec{r}_{old}^T \cdot \vec{r}_{old}}, 0\right]$  $\vec{s} \leftarrow \vec{r} + \beta \vec{s}$

BackPropagation Basic Principle

Faster Alternatives The Problem The Solution

Measuring Performance Generalization

#### Regularization

data set

Limiting the Complexity Punishing Large Weights Optimal Pruning

Separation of training and testing data

Training

Performance should always be measured on a separate test

Testing

Örjan Ekeberg

ANN fk

BackPropagation

Measuring Performance Generalization

#### ANN fk

Örjan Ekeberg

ackPropagation aster Alternatives

Measuring Performance Generalization Regularization

#### ANN fk Örjan Ekeberg

#### BackPropagation Faster Alternatives Measuring Performance Generalization

How can one measure the performance of a neural network?

- Evaluation of a classifier
- Positive and negative errors
- The error for evaluation does not have to be the same as the error minimized during learning

How large should the test data set be?



The Problem The Solution

## Regularization

Limiting the Complexity Punishing Large Weights **Optimal Pruning** 

# ANN fk Örjan Ekeberg Faster Alternatives Generalization



ANN fk

Örjan Ekeberg

Faster Alternatives

Performance

Regularization

Limiting the Complexity Punishing Large Weights Optimal Pruning

## Maximal utilization of available data



Average over different partitionings

Many weights  $\Rightarrow$  Bad generalization

Risk of making errors

- $\triangleright \mathcal{E}_s$  Empirical risk (measurable)
- $\triangleright \mathcal{E}_c$  Structural risk

$$R(\vec{w}) = \mathcal{E}_{s}(\vec{w}) + \lambda \mathcal{E}_{c}(\vec{w})$$

# $\lambda$ — regularization parameter

#### ANN fk

Örjan Ekeberg

Measuring Performance Generalization

#### ANN fk

Örjan Ekeberg

Measuring Performance

Limiting the Complexity

Augment the cost function with a complexity term Weight Decay (Hinton, 1989):

 $\mathcal{E}_c = \sum_i w_i^2$ 

Weight Elimination (Weigend et al., 1991):

$$\mathcal{E}_{c} = \sum_{i} \frac{(w_{i}/w_{0})^{2}}{1 + (w_{i}/w_{0})^{2}}$$

# **Optimal Brain Surgeon**

(Hassibi et al., 1992)

- Improved version of Optimal Brain Damage
- ▶ Takes into account that other weights may need readjustment
- Requires an estimate of mixed second derivatives (Hessian-matrix)
- ► Can be efficiently estimated using the errors of the individual patterns



ANN fk

Örjan Ekeberg

aster Alternative

Punishing Large Weights

Optimal Pruning

# Alternative technique

Remove "unnecessary" weights after training

Optimal Brain Damage (LeCun et al., 1990)

Idéa: Remove the least important weight

- Estimate how much the error increases when a weight is set to zero
- $\blacktriangleright \Delta w_i = -w_i$
- Error increase  $-w_i \cdot \frac{\partial \mathcal{E}}{\partial w_i} = 0$  (since BP has converged) Oops
- Second derivative  $\frac{w_i^2}{2} \cdot \frac{\partial^2 \mathcal{E}}{\partial w_i^2}$

#### ANN fk

Örjan Ekeberg

Faster Alternative

**Optimal Pruning**