

EFFICIENT LEARNING AND SECOND-ORDER METHODS

Yann Le Cun

Adaptive Systems Research Dept

AT&T Bell Laboratories

Holmdel, NJ

USA

OVERVIEW

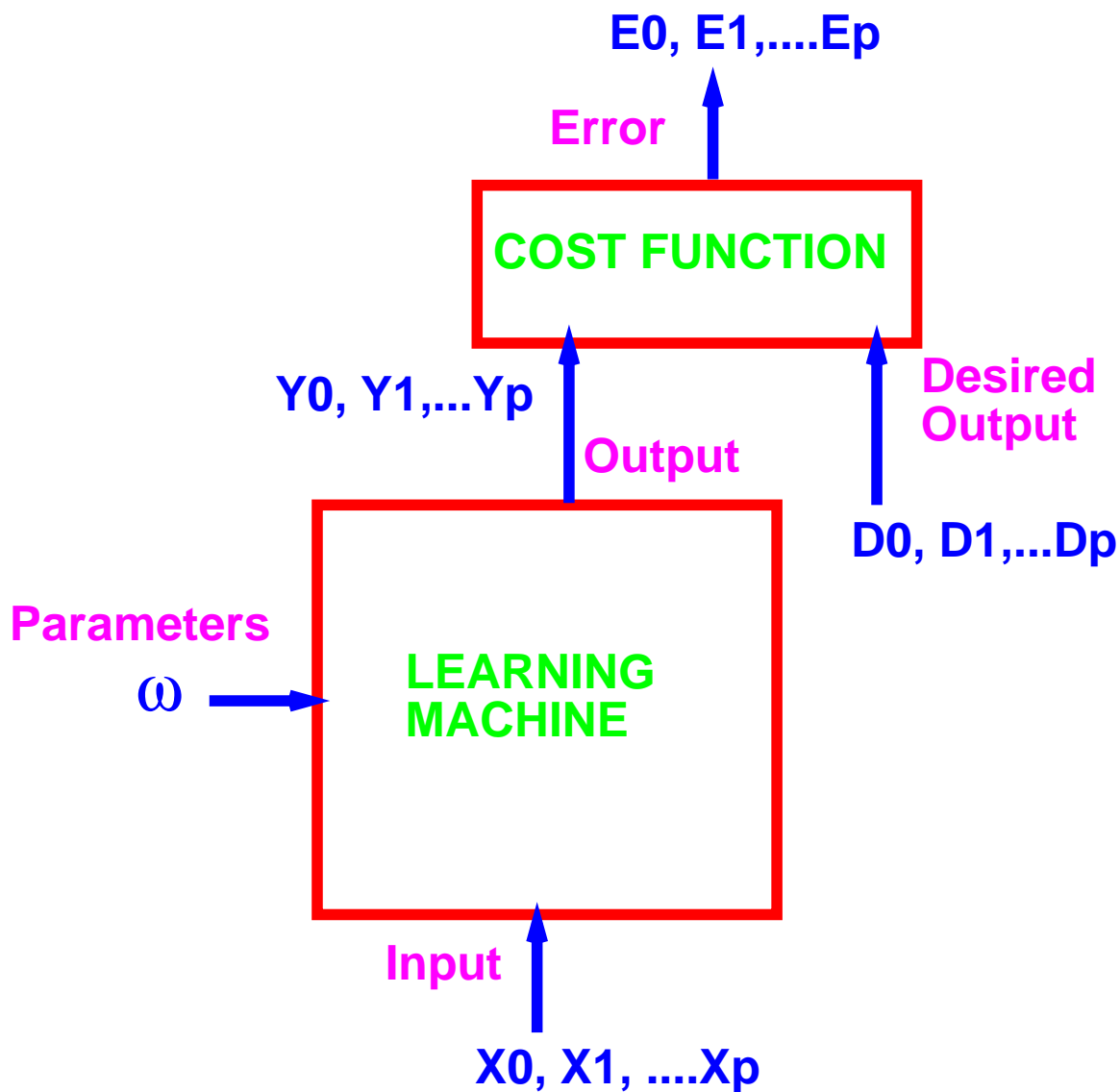
- 1 – Plain Backprop: how to make it work
 - Basic concepts, terminology and notation
 - Intuitive analysis and practical tricks
- 2– The convergence of gradient descent
 - A little theory
 - Quadratic forms, Hessians, and Eigenvalues
 - maximum learning rate, minimum learning time
 - How GD works in simple cases
 - a single 2–input neuron
 - stochastic vs batch update
 - Transformation laws
 - shifting, scaling, and rotating the input
 - the non–invariance of GD
 - The minimal multilayer network
- 3 – Second order methods
 - Newton’s algorithm, and why it does not work.
 - parameter space transformations
 - computing the second derivative information
 - diagonal terms, quasi–linear Hessians, partial Hessians
 - analysis of multilayer net Hessians
 - Classical optimization methods
 - Conjugate Gradient methods
 - Quasi Newton methods: BFGS et al.
 - Levenberg–Marquardt methods
- 4 – Applying second order methods to multilayer nets
 - (non)applicability of 2nd order techniques to backprop
 - Mini batch methods
 - An on–line diagonal Levenberg–Marquardt method
 - computing the maximum learning rate and the principal eigenvalues

1

PLAIN BACKPROP:

HOW TO MAKE IT WORK

BASIC CONCEPTS, TERMINOLOGY, NOTATIONS

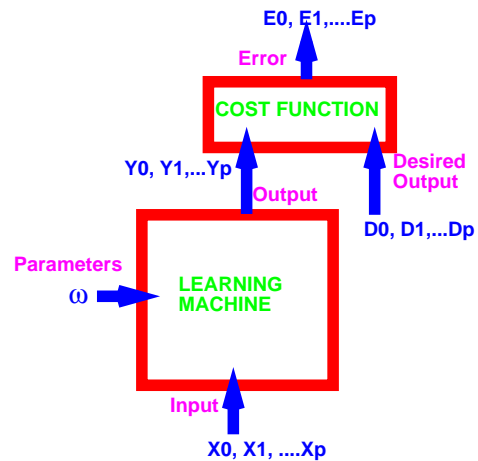


Average Error:
$$E = \frac{1}{p} \sum E_k$$

GRADIENT DESCENT LEARNING

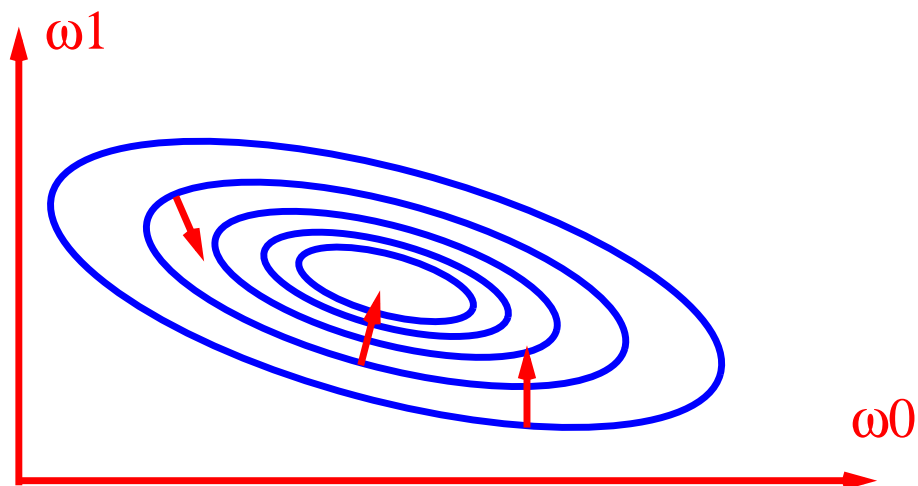
Average Error:

$$E(\omega) = \frac{1}{p} \sum E_k(\omega)$$

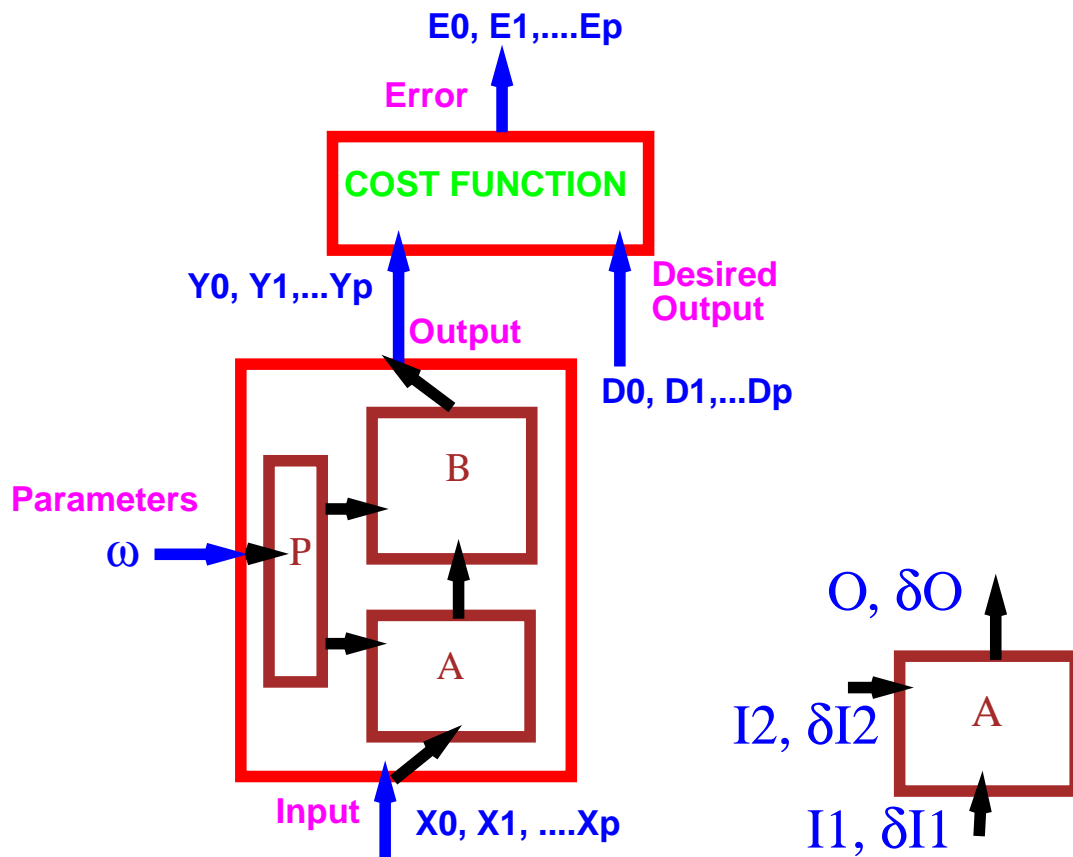


Gradient Descent:

$$\omega(\tau+1) = \omega(\tau) - \eta \frac{\partial E}{\partial \omega}$$



COMPUTING THE GRADIENT WITH BACKPROPAGATION



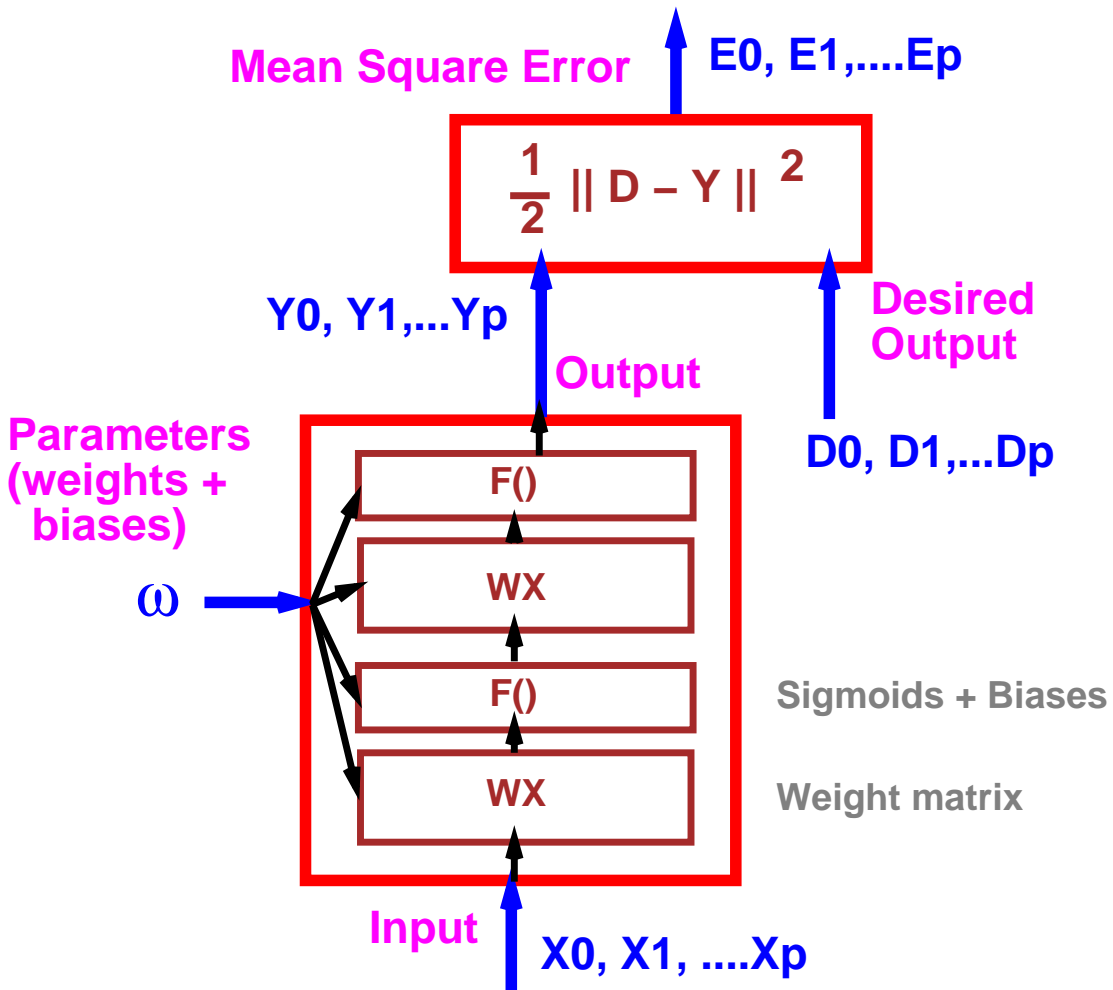
- The learning machine is composed of modules (e.g. layers)
- Each module can do two things:
 - 1- compute its outputs from its inputs (FPROP)

$$O = A(I_1, I_2)$$

- 2- compute gradient vectors at its inputs from gradient vectors at its outputs (BPROP)

$$\delta I_1 = \frac{\partial A}{\partial I_1} \delta O \quad \delta I_2 = \frac{\partial A}{\partial I_2} \delta O$$

AN INTERESTING SPECIAL CASE: MULTILAYER NETWORKS



Weight matrices:

FPROP $O = \omega I$

BPROP $\delta I = \omega' \delta O ; \delta \omega = \delta O' I$

Sigmoids + Bias:

FPROP $O = f(I+B)$

BPROP $\delta I = f'(I+B) \delta O ; \delta B = \delta I$

STOCHASTIC UPDATE BATCH UPDATE

STOCHASTIC GRADIENT

$$\omega(\tau+1) = \omega(\tau) - \eta \frac{\partial E\tau}{\partial \omega}$$

```
Repeat {  
  for all examples in training set {  
    forward prop      // compute output  
    backward prop     // compute gradients  
    update parameters }  
}
```

The parameters are updated after each presentation of an example

FULL GRADIENT

$$\omega(\rho+1) = \omega(\rho) - \eta \frac{\partial E}{\partial \omega}$$

```
Repeat {  
  for all examples in training set {  
    forward prop      // compute output  
    backward prop     // compute gradients  
    accumulate gradient }  
  update parameters }
```

The gradients are accumulated over the whole training set before a parameter update is performed

A FEW PRACTICAL TRICKS

BackProp is a simple algorithm, but convergence can take ages if it is not used properly.

The error surface of a multilayer network is non quadratic and non-convex, and has often many dimensions.

THERE IS NO MIRACLE TECHNIQUE for finding the minimum. Heuristics (tricks) must be used.

Depending on the details of the implementation, the convergence time can vary by orders of magnitude, especially on small problems.

On large, real-world problems, the convergence time is usually much better than one would expect from extrapolating the results on small problems.

Here is a list of some common traps, and some ideas about how to avoid them.

The theoretical justifications for many of these tricks will be given later in the talk.

STOCHASTIC vs BATCH UPDATE

Stochastic update is usually MUCH faster than Batch update. Especially on large, redundant data sets.

Here is why:

Imagine you are given a training set with 1000 examples.

Imagine this training set is in fact composed of 10 copies of a set of 100 patterns

- **BATCH:** the computation for one update will be 10 times larger than necessary.
- **STOCHASTIC:** the redundancy in the training set will be taken advantage of. One epoch on the large set will be like 10 epochs on the smaller set.

Batch will be AT LEAST 10 times slower than Stochastic

In real life, repetitions rarely occur, but very often the training examples are highly redundant (many patterns are similar to one another), which has the same effect.

In practice speed differences of orders of magnitude between Batch and Stochastic are not uncommon

small batches can be used without penalty, provided the patterns in a minibatch are not too similar.

STOCHASTIC vs BATCH UPDATE (continued)

STOCHASTIC:

ADVANTAGES:

- much faster convergence on large redundant datasets
- stochastic trajectory allows escaping from local minima

DISADVANTAGES:

- keeps bouncing around unless the learning rate is reduced
- theoretical conditions for convergence are not as clear as for batch
- convergence proofs are probabilistic
- most nice acceleration tricks or second-order methods do not work with stochastic gradient
- it is harder to parallelize than batch

BATCH:

ADVANTAGES:

- guaranteed convergence to a local minimum under simple conditions.
- lots of tricks and second order methods to accelerate it
- easy convergence proofs

DISADVANTAGES:

- painfully slow on large problems

Despite the long list of disadvantages for STOCHASTIC, that is what most people use (and rightfully so, at least on large problems).

SHUFFLING THE EXAMPLES

RULE: at any time, chose the training example with the maximum information content.

For example:

- the one with the largest error
- the one that is maximally different from its predecessors

A SIMPLE TRICK:

(applicable to stochastic gradient on classification tasks)

Shuffle the training set so that successive examples never (or rarely) belong to the same class.

A MORE REFINED TRICK:

use an "emphasizing" scheme:
show difficult patterns more often than easy patterns

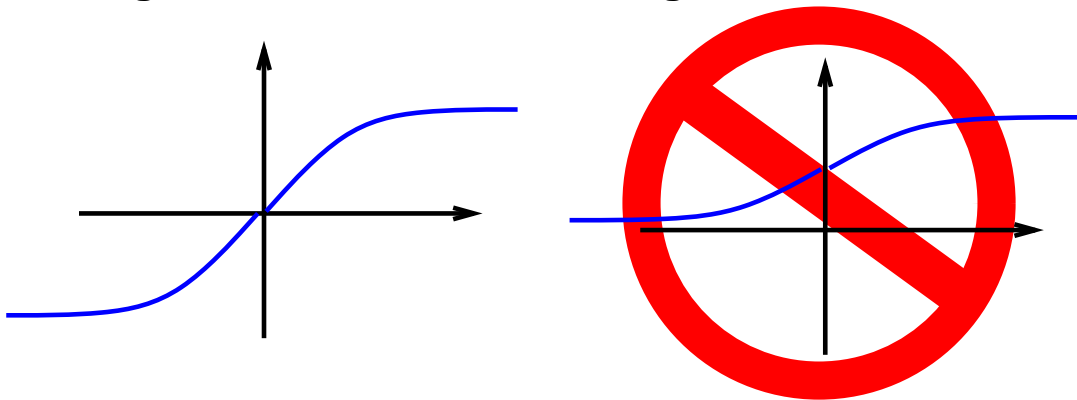
[whether a pattern is easy or hard can be determined with the error it produced during the previous iterations]

Problem with emphasizing techniques:

- they perturb the distribution of inputs
- the presence of outliers or of mislabeled examples can be catastrophic.

THE SIGMOID

Symmetric sigmoids (like tanh) often yield faster convergence than the standard logistic function.



MORE GENERALLY: the mean of each input to a neuron should be small compared to its standard deviation [more on this later]

Symmetric sigmoids are more likely to produce "small mean" signals than are positive sigmoids.

Sigmoids (and their derivatives) can be efficiently computed as ratios of polynomials

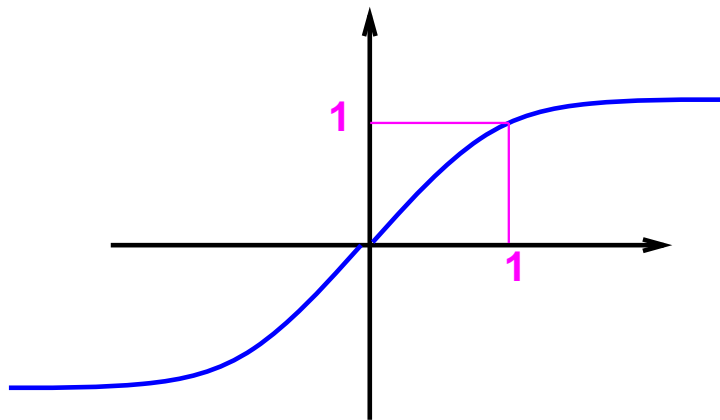
Problems with symmetric sigmoids:

- The error surface is **VERY FLAT** around the origin. (the origin is a saddle point which is attractive in almost all directions)
- **Avoid small weights**

THE SIGMOID (continued)

The one I use: a rational approximation to

$$f(x) = 1.7159 \tanh\left(\frac{2}{3}x\right)$$



Properties:

- $f(1) = 1$, $f(-1) = -1$
- 2nd derivative is maximum at $x=1$
- the effective gain is close to 1

The precise choice of the sigmoid is almost irrelevant, but some choices are more convenient than others

It is sometimes helpful to add a small linear term to avoid flat spots, e.g.

$$f(x) = \tanh(x) + ax$$

NORMALIZING THE INPUTS

Each input variable should be sifted so that its mean (averaged over the training set) is close to 0 (or is small compared to its standard deviation).

Here is why:

Consider the extreme case where the input variables are always positive.

The weights of neuron in the first hidden layer can only increase together or decrease together (for a given input pattern the gradients all have the same sign).

This means that if the weight vector has to change its direction, it will have to do it by zigzaging (read: SLOW).

Shifts of the input variables to a neuron introduce a preferred direction for weight changes, which slows down the learning.

This is also why we prefer symmetric sigmoids: what is true for input units is also true for other units in the network.

NORMALIZING THE INPUTS (continued)

The speed at which the output of a particular weight varies with gradient descent is proportional to the **COVARIANCE** of its input.
[more on this later]

$$\text{covariance: } \frac{1}{P} \sum_K \chi_K^2$$

To equalize the learning speeds of input weights, the input variables should be scaled to have approximately equal covariances.

To equalize the learning speeds of these weights with that of the weights in the next layers, this covariance should be comparable to the expected covariances of the hidden units states (around 1 with the type of sigmoid proposed earlier).

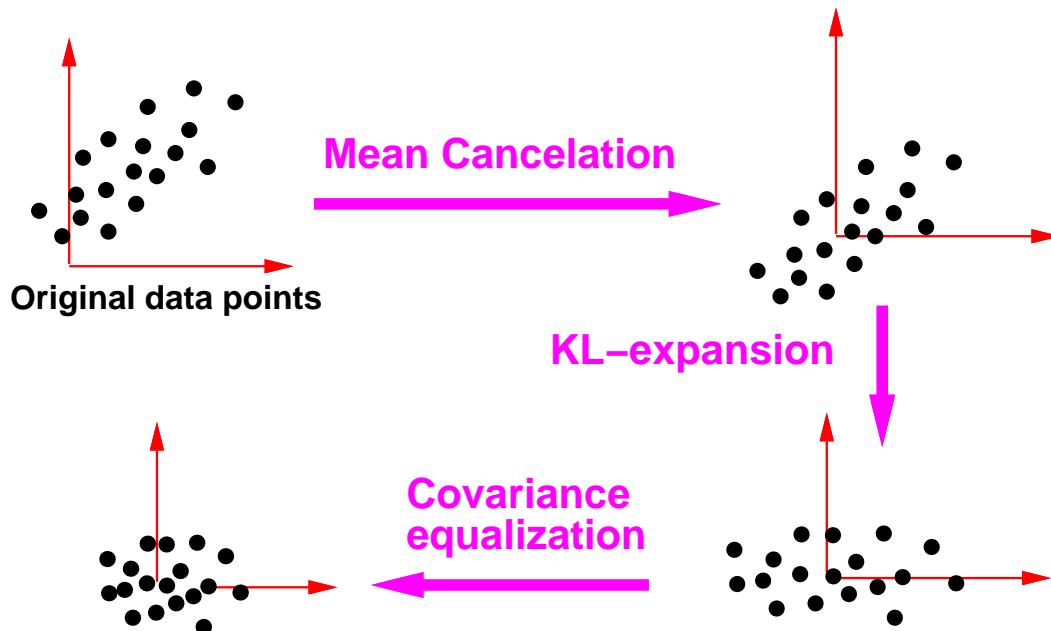
An exception to this rule is when some inputs are known to be of lesser significance than others. Scaling them down makes them less "visible" to the learning process.

NORMALIZING THE INPUTS (continued)

Input variables should be **UNCORRELATED**
if possible

Correlations between input variables also
introduce "preferred directions for weight changes"

Decorrelation can be performed by a Principal
Component Analysis (Karhunen–Loeve expansion).



Sometimes the input have a particular meaning
that would be destroyed by the KL-expansion
(e.g.: if the input variables are the pixels of an image
and the network uses local connections)

CHOOSING THE TARGET VALUES

Avoid saturating the output units. Choose target values within the range of the sigmoid

In classification applications, the desired outputs are often binary.

Common sense would suggest to set the target values on the asymptotes of the sigmoid.

However this has several adverse effects:

- 1 – this tends to drive the output weights to infinity (and to saturate the hidden units as well). When a training example happens not to saturate the outputs (say an outlier), it will produce ENORMOUS gradients (due to the large weights).**
- 2 – outputs will tend to be binary EVEN WHEN THEY ARE WRONG. This means that a mistake will be difficult to correct, and that the output levels cannot be used as reliable confidence factors.**

Saturating the units erases the differences between typical and non-typical examples.

Setting the target values at the point of maximum second derivative on the sigmoid (-1 and +1 for the sigmoid proposed earlier) is the best way to take advantage of the non linearity without saturating the units

INITIALIZING THE WEIGHTS

Initialize the weights so that the expected standard deviation of the weighted sums is at the transition between the linear part and the saturated part of the sigmoid function

Large initial weights saturate the units, leading to small gradients and slow learning.
Small weights correspond to a very flat area of the error surface (especially with symmetric sigmoids)

Assuming the sigmoid proposed earlier is used, the expected standard deviation of the inputs to a unit is around 1, and the desired standard deviation of its weighted sum is also around 1.

Assuming the inputs are independent, the expected standard deviation of the weighted sum is

$$\sigma = \left(\sum_i \omega_i^2 \right)^{1/2} = \phi^{1/2} \bar{\omega}$$

where ϕ is the number of input to the unit, and $\bar{\omega}$ is the standard deviation of its incoming weights.

To ensure that σ is close to 1, the weights to a unit can be drawn from a distribution with standard deviation

$$\bar{\omega} = \phi^{-1/2}$$

CHOOSING LEARNING RATES

Equalize the learning speeds.

Each weight (or parameter) should have its own learning rate. η

Some weights may require a small learning rate to avoid divergence, while others may require a large learning rate to converge at reasonable speed.

Because of possible correlations between input variables, the learning rate of a unit should be inversely proportional to the square root of the number of inputs to the unit.

If shared weights are used (as in TDNNs and convolutional networks), the learning rate of a weight should be inversely proportional to the square root of the number of connection sharing that weight.

Learning rates in the lower layers should generally be larger than that in the higher layers.

The rationale for many of these rules of thumb will become clearer later.
Several techniques are available to reduce "learning rate fiddling".

NEURONS AND WEIGHTS

Although most systems use neurons based on dot products and sigmoids, many other types of units (or layers) can be used.

A particularly interesting example is when the dot product of the input by the weight vector is replaced by a Euclidean distance, and the sigmoid by an exponential (Gaussian units or RBF).

These units can replace (or coexist with) standard units, and they can be trained with gradient descent:

exponential + Bias:

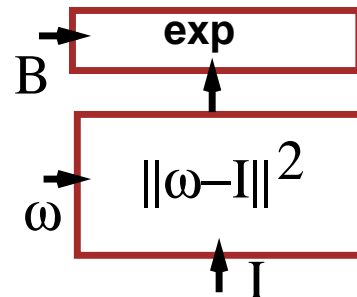
$$\text{FPROP } O = \exp(I+B)$$

$$\text{BPROP } \delta I = O \delta O ; \delta B = \delta I$$

weight vector:

$$\text{FPROP } O = (\omega - I)'(\omega - I)$$

$$\text{BPROP } \delta I = 2 \delta O' (I - \omega) ; \\ \delta \omega = -\delta I$$



PROS and CONS

- Locality: each unit is only affected by a small part of the input space. This can be good (for learning speed) and bad (a lot of RBF are needed to cover high dimensional spaces)
- gradient descent learning may fail if the RBFs are not properly initialized (using clustering techniques e.g. K-means). There are LOTS of local minima.
- RBFs are more appropriate in higher layers, and sigmoids in lower layers (higher dimension).

MORE STANDARD TRICKS

- Momentum**
 - Increases speed in batch mode.**
seems marginally useful but not indispensable in stochastic mode.

- Adaptive learning rates:**
 - a separate learning rate for each weight is increased if the gradient is steady, decreased if the gradient changes sign often [Jacobs 88]. This only works with BATCH.**

 - a global learning rate is adjusted using line searches. Again, this only works for BATCH.**

2

**THE CONVERGENCE OF
GRADIENT DESCENT**

A LITTLE THEORY

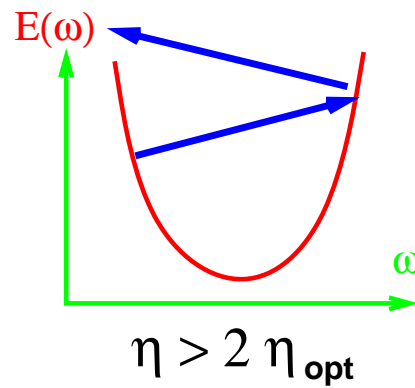
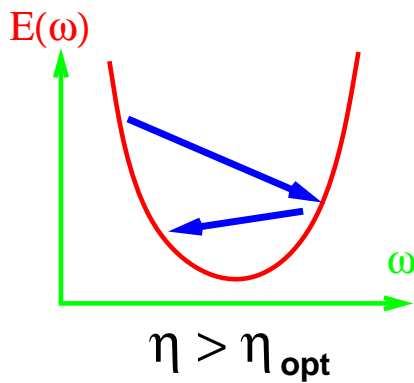
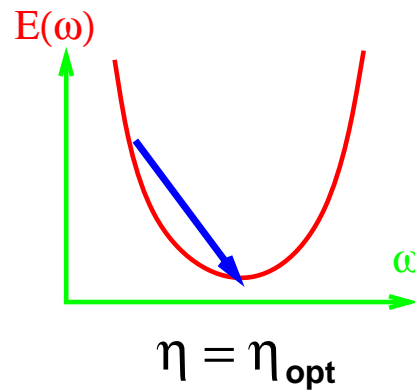
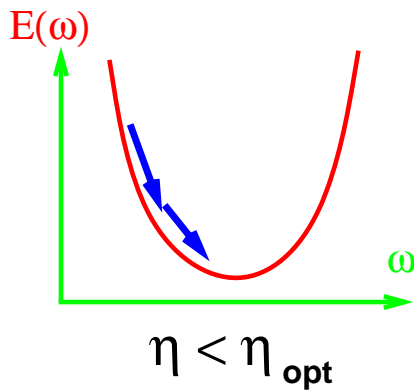
Gradient Descent in one dimension

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

weight vector

learning rate

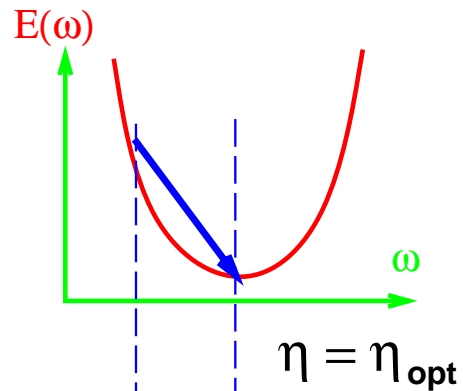
gradient of objective function



OPTIMAL LEARNING RATE IN 1D

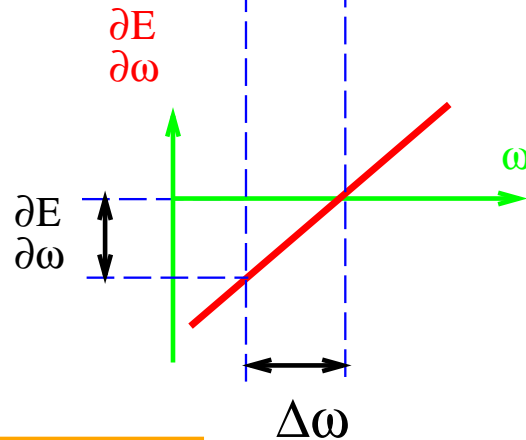
Weight change:

$$\Delta\omega = \eta \frac{\partial E}{\partial \omega}$$



Assuming E is quadratic:

$$\frac{\partial^2 E}{\partial \omega^2} \Delta\omega = \frac{\partial E}{\partial \omega}$$



Optimal Learning Rate

$$\eta_{\text{opt}} = \left(\frac{\partial^2 E}{\partial \omega^2} \right)^{-1}$$

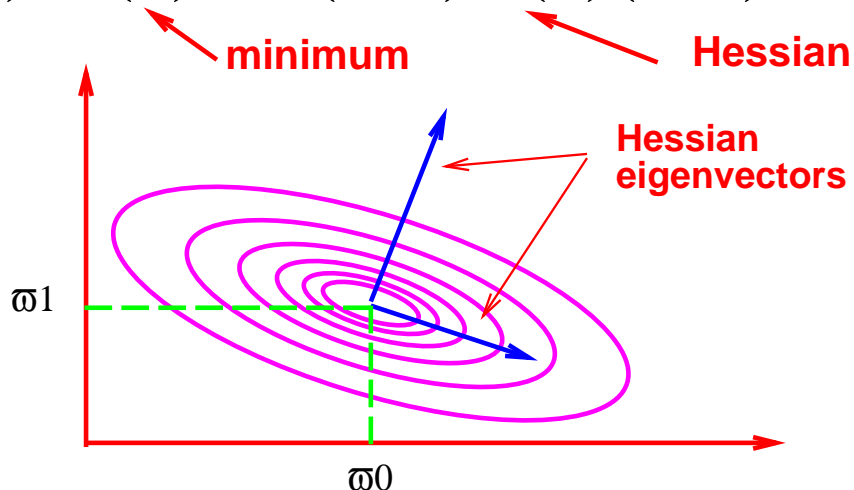
Maximum Learning Rate

$$\eta_{\text{max}} = 2 \eta_{\text{opt}}$$

CONVERGENCE OF GRADIENT DESCENT

Local quadratic approximation of the cost function around a minimum:

$$E(\omega) \approx E(\bar{\omega}) + 1/2(\omega - \bar{\omega})' H(\bar{\omega}) (\omega - \bar{\omega})$$



HESSIAN

Second derivative matrix

$$H_{ij} = \frac{\partial^2 E}{\partial \omega_i \partial \omega_j}$$

The Hessian is a symmetric NxN matrix

Gradient Descent weight update:

$$\omega_{k+1} = \omega_k - \eta \frac{\partial E}{\partial \omega} = \omega_k - \eta H(\omega_k) (\omega_k - \bar{\omega})$$

$$(\omega_{k+1} - \bar{\omega}) = (I - \eta H(\omega_k)) (\omega_k - \bar{\omega})$$

Convergence \iff if the prefactor of the right handside shrinks any vector

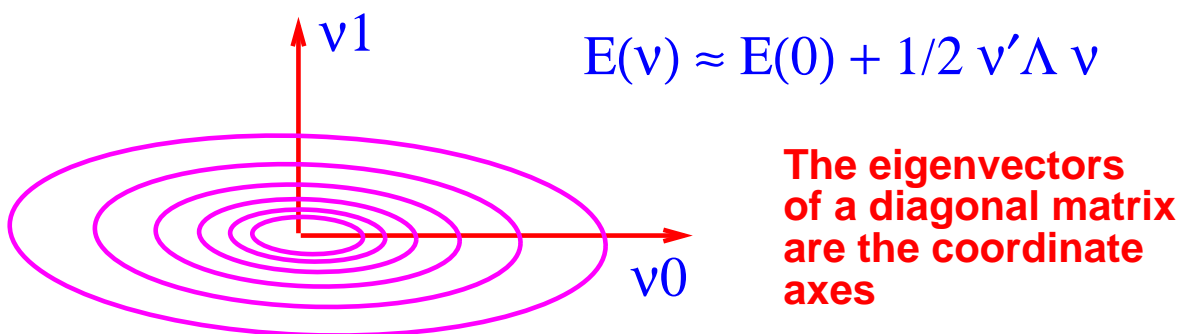
CONVERGENCE OF GRADIENT DESCENT (continued)

Let Θ be the rotation matrix that make H diagonal:

$$\Theta H \Theta' = \Lambda ; \quad \Theta' \Theta = I$$

$$E(\omega) \approx E(\bar{\omega}) + 1/2[(\omega - \bar{\omega})' \Theta'] [\Theta H(\bar{\omega}) \Theta'] [\Theta(\omega - \bar{\omega})]$$

Now denote: $v = \Theta (\omega - \bar{\omega})$



Gradient update in the transformed space: $v_{k+1} = (I - \eta \Lambda) v_k$

Gradient Descent in N dimensions can be viewed as N independent unidimensional Gradient Descents along the eigenvectors of the Hessian.

Convergence is obtained for $\eta < 2 / \lambda_{\max}$
where λ_{\max} is the largest eigenvalue of the Hessian

CONVERGENCE SPEED OF GRADIENT DESCENT

The maximum learning rate to ensure convergence is

$$\eta_{\max} = 2 / \lambda_{\max}$$

The one that yields the fastest convergence in the direction of highest curvature is

$$\eta_{\text{opt}} = 1 / \lambda_{\max}$$

With this choice, the convergence time will be determined by the directions of **SMALL** eigenvalues (they will be the slowest to converge).

The convergence time is proportional to:

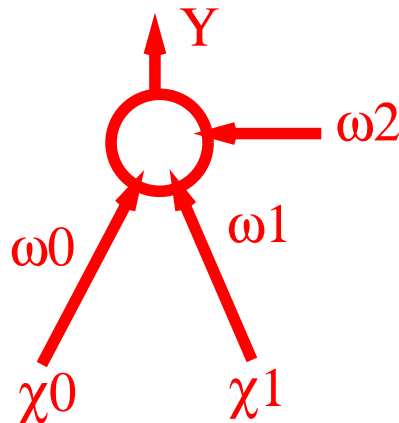
$$\frac{1}{\eta \lambda_{\min}} > \frac{\lambda_{\max}}{2\lambda_{\min}}$$

where λ_{\min} is the smallest "non-negligible" eigenvalue

The convergence time is proportional to the ratio of the largest eigenvalue to smallest "non-negligible" eigenvalue of the Hessian

CONVERGENCE OF GRADIENT DESCENT A SIMPLE EXAMPLE

A single 2-input
linear neuron:



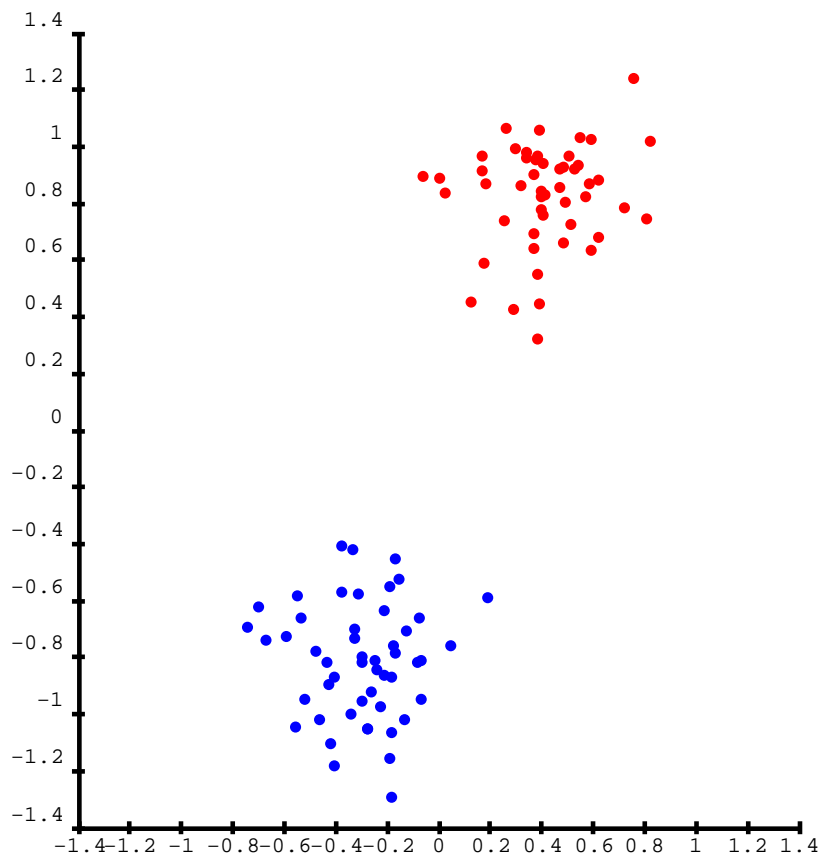
$$E(\omega) = \frac{1}{2P} \sum_p \| d_p - y_p(\omega) \|^2 = \frac{1}{2P} \sum_p \| d_p - \omega' \chi_p \|^2$$

$$E(\omega) = \frac{1}{2P} \left[\sum_p d_p^2 - 2 \left(\sum_p d_p \chi_p \right)' \omega + \omega' \left(\sum_p \chi_p \chi_p' \right) \omega \right]$$

$$H = \frac{1}{P} \sum_p \chi_p \chi_p'$$

**The Hessian of a single linear neuron is
the covariance matrix of the inputs**

Dataset #1



Examples of each class are drawn from a Gaussian distribution centered at $(-0.4, -0.8)$, and $(0.4, 0.8)$.

Eigenvalues of covariance matrix: 0.83 and 0.036

Batch gradient descent

data set: set-1 (100 examples, 2 gaussians)
network: 1 linear unit, 2 inputs, 1 output.
2 weights, 1 bias.

Learning rate:

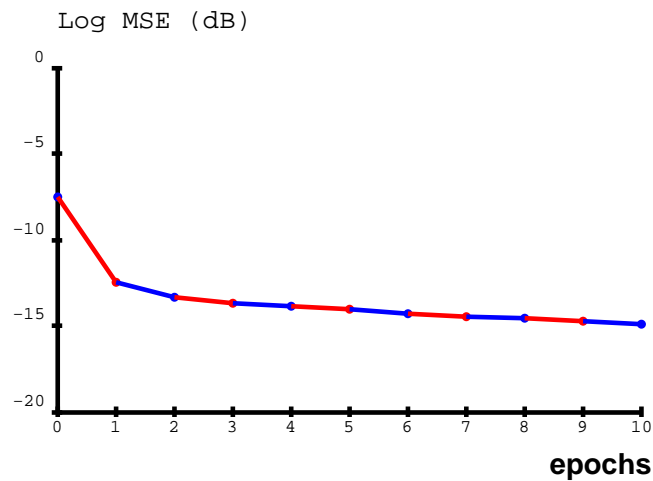
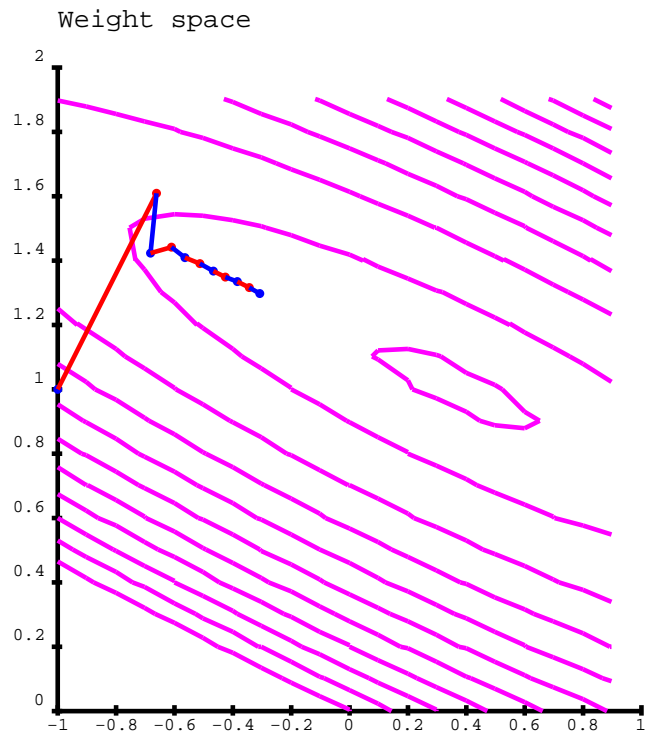
$$\eta = 1.5$$

Hessian largest eigenvalue:

$$\lambda_{\max} = 0.84$$

Maximum admissible Learning rate:

$$\eta_{\max} = 2.38$$



Batch gradient descent

data set: set-1 (100 examples, 2 gaussians)
network: 1 linear unit, 2 inputs, 1 output.
2 weights, 1 bias.

Learning rate:

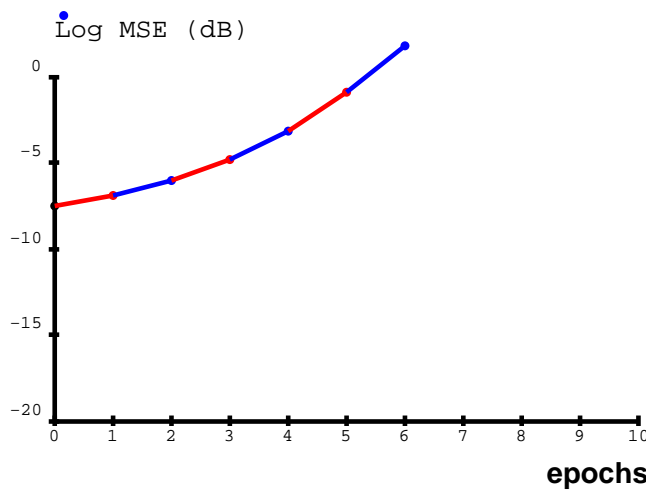
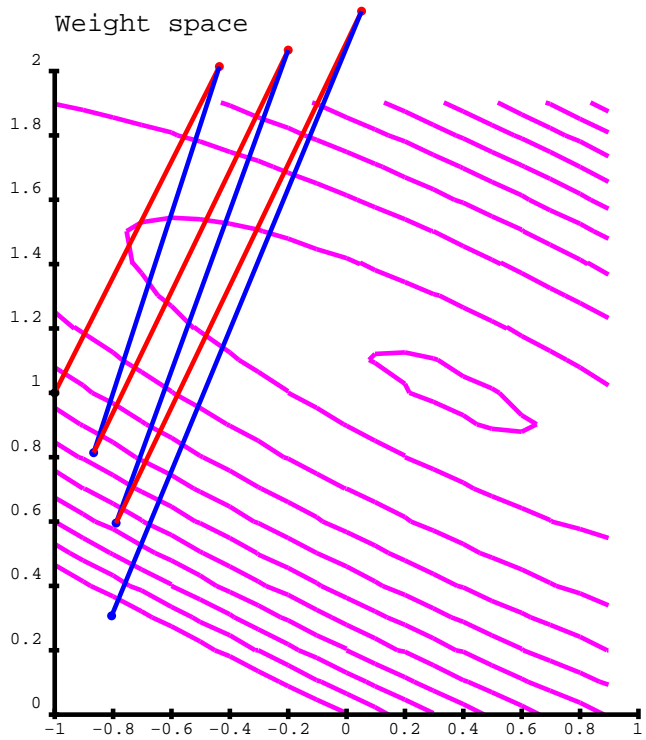
$$\eta = 2.5$$

Hessian largest eigenvalue:

$$\lambda_{\max} = 0.84$$

Maximum admissible Learning rate:

$$\eta_{\max} = 2.38$$



Stochastic gradient descent

data set: set-1 (100 examples, 2 gaussians)
network: 1 linear unit, 2 inputs, 1 output.
2 weights, 1 bias.

Learning rate:

$$\eta = 0.2$$

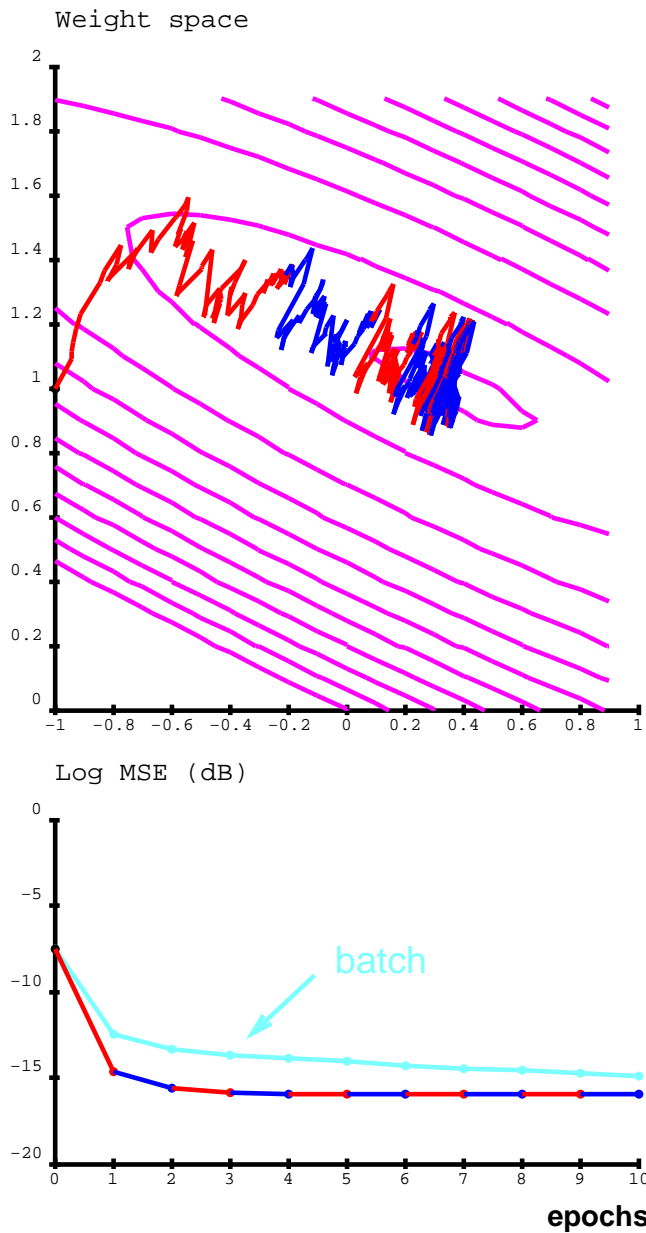
(equivalent to a batch learning rate of 20)

Hessian largest eigenvalue:

$$\lambda_{\max} = 0.84$$

Maximum admissible Learning rate (for batch):

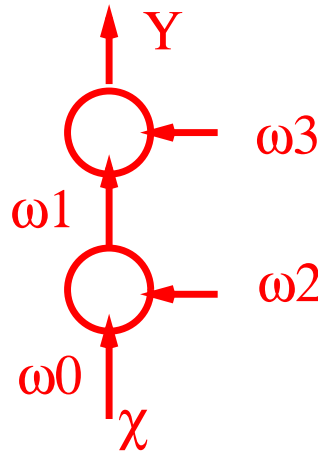
$$\eta_{\max} = 2.38$$



CONVERGENCE OF GRADIENT DESCENT: MINIMAL MULTILAYER NETWORK

1 input
1 hidden unit
1 output

2 weights
2 biases



Sigmoid: $1.71 \tanh(2/3 x)$

Targets: -1 for class 1, $+1$ for class 2

TRAINING SET: 20 examples.

Class1:

10 examples drawn from a Gaussian distribution with mean -1 , and standard deviation 0.4

Class2:

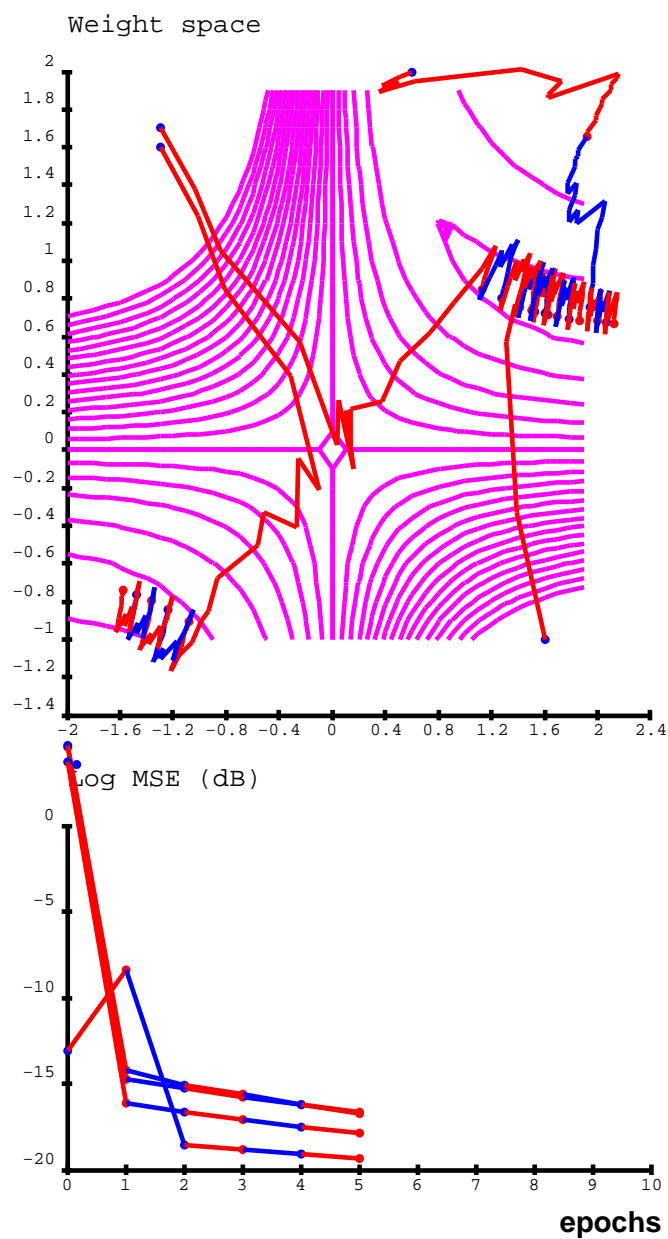
10 examples drawn from a Gaussian distribution with mean $+1$, and standard deviation 0.4

Stochastic gradient: 1-1-1 network

data set: 20 examples, 2 1D-gaussians)
network: 1 input, 1 hidden, 1 output
2 weights 2 biases

Learning rate:

$$\eta = 0.4$$



INPUT TRANSFORMATIONS

ERROR SURFACE TRANSFORMATIONS

A shift (non zero mean) in the input variables creates a VERY LARGE eigenvalue which leads to eccentric paraboloids, and to slow convergence

Subtract the means from the input variables

For a single linear neuron, If the input variables have zero means, the eigenvectors of the Hessian are the principal axes of the cloud of training vectors

widely spread variances for the input variables lead to widely spread Hessian eigenvalues

Normalize the variances of the input variables

Correlations between input variables lead to eccentric paraboloid with ROTATED axes

Decorrelate the input variables

The gradient is NOT the best descent direction

**Use a separate learning rate for each weight.
No use spending time and effort to compute an accurate gradient estimate**

3

SECOND ORDER METHODS

NEWTON ALGORITHM

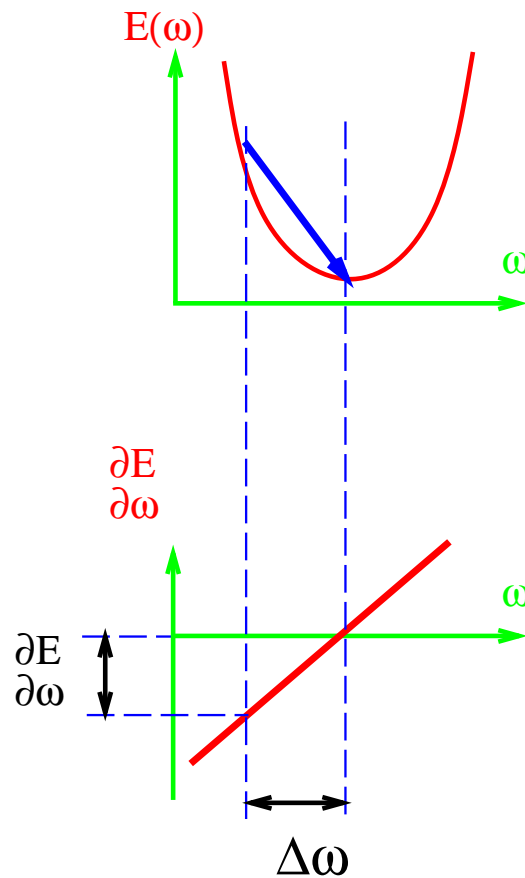
Newton Algorithm in one dimension

Optimal weight change:

Assuming E is quadratic:

$$\frac{\partial^2 E}{\partial \omega^2} \Delta \omega = \frac{\partial E}{\partial \omega}$$

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



If E is not perfectly quadratic:

$$\Delta \omega = \eta \left(\frac{\partial^2 E}{\partial \omega^2} \right)^{-1} \frac{\partial E}{\partial \omega} \quad 0 < \eta < 1$$

NEWTON ALGORITHM

Local quadratic approximation of the cost function around the current point:

$$E(\omega + \Delta\omega) \approx E(\omega) + \nabla E(\omega) \Delta\omega + 1/2 \Delta\omega' H(\omega) \Delta\omega$$

\nearrow
weight change

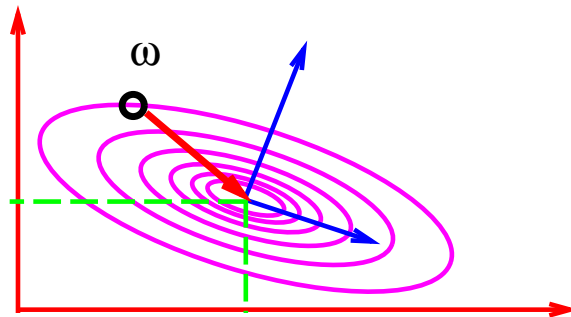
\nwarrow
gradient

\nearrow
Hessian

IDEA: find the weight change that minimizes the above.
(i.e.: find the weight change for which the gradient is 0)

Solve: $\nabla E(\omega) + H(\omega) \Delta\omega = 0$ for $\Delta\omega$

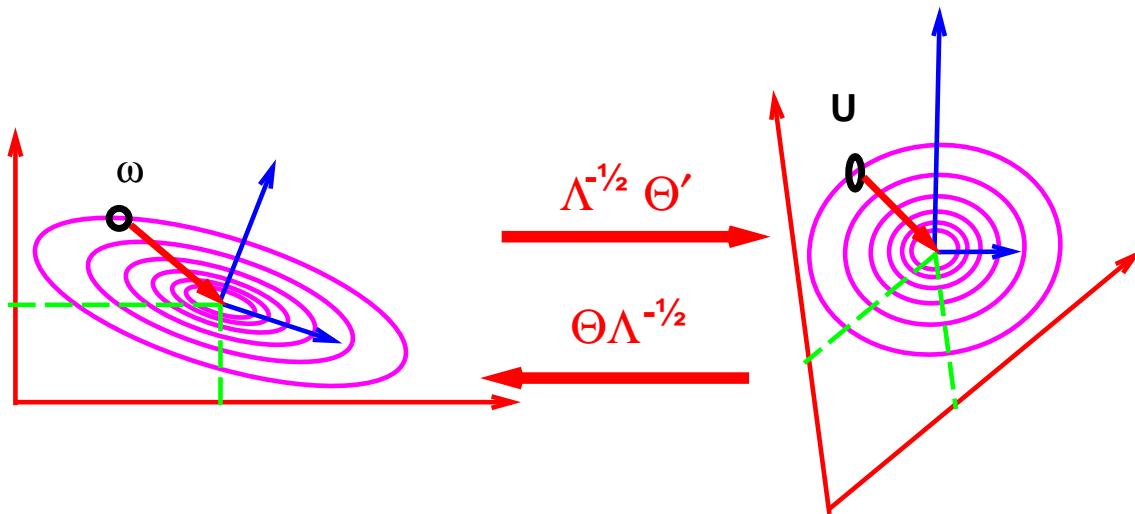
$$\Delta\omega = H(\omega)^{-1} \nabla E(\omega)$$



NEWTON ALGORITHM AND PARAMETER SPACE TRANSFORMS

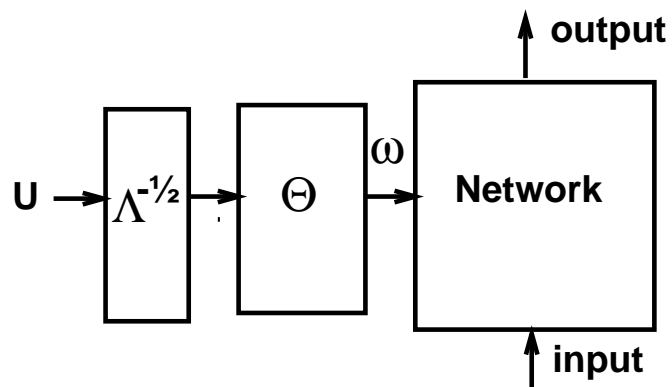
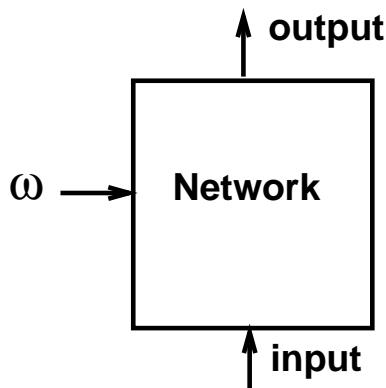
Diagonalized Hessian

$$H(\omega) = \Theta' \Lambda^{1/2} \Lambda^{1/2} \Theta \quad H^{-1}(\omega) = \Theta \Lambda^{-1/2} \Lambda^{-1/2} \Theta'$$



Newton Algorithm here

....is like Gradient Descent
there



NEWTON ALGORITHM

- it converges in 1 iteration if the error is quadratic
- unlike Gradient Descent, it is invariant with respect to linear transformations of the input vectors, i.e. the convergence time is not affected by shifts, scaling and rotations of the input vectors.

BUT:

- it requires computing, storing and inverting the $N \times N$ Hessian (or solving an $N \times N$ linear system). The complexity is $O(N^3)$, which is impractical with more than a few variables)
- there is **ABSOLUTELY NO GUARANTEE** of convergence when the error is non-quadratic
- in fact, it diverges if the Hessian has some null or negative eigenvalues (the error surface is flat or curved downward in some directions). The Hessian **MUST** be Positive Definite. (this is obviously not the case in multilayer nets)

The Newton Algorithm in its original form is unusable for Neural Net learning.

but its basic idea is useful for understanding more sophisticated algorithms

COMPUTING THE HESSIAN INFORMATION IN MULTILAYER NETWORKS

There are many techniques to compute the full Hessian, or parts of it, or approximations to it, in multilayer networks.

We will review the following simple methods:

- finite difference**
- square Jacobian approximation
(for the Gauss–Newton and
Levenberg–Marquardt algorithms)**
- computing the diagonal term (or block
diagonal terms) by backpropagation**
- computing the product of the Hessian by a
vector without computing the Hessian**

There exist more complex techniques to compute semi-analytically the full Hessian [Bishop 92, Buntine&Weigend 93, others.....] but they are REALLY complicated, and require many forwardprop/backprop passes.

FINITE DIFFERENCE

The k-th line of the Hessian is the derivative of the GRADIENT with respect to the k-th parameter

$$(\text{Line } k \text{ of } H) = \frac{\partial (\nabla E(\omega))}{\partial \omega_k}$$

Finite difference approximation:

$$(\text{Line } k \text{ of } H) = \frac{\nabla E(\omega + \delta \phi_k) - \nabla E(\omega)}{\delta}$$

$$\phi_k = (0, 0, 0, \dots, 1, \dots, 0)$$

RECIPE for computing the k-th line of the Hessian

- 1- compute total gradient (multiple fprop/bprop)
- 2- add Delta to k-th parameter
- 3- compute total gradient
- 4- subtract result of line 1 from line 3, divide by Delta.

due to numerical errors, the resulting Hessian may not be perfectly symmetric. It should be symmetrized.

SQUARE JACOBIAN APPROXIMATION FOR GAUSS-NEWTON AND LEVENBERG-MARQUARDT ALGOS.

Assume the cost function is the Mean Square Error:

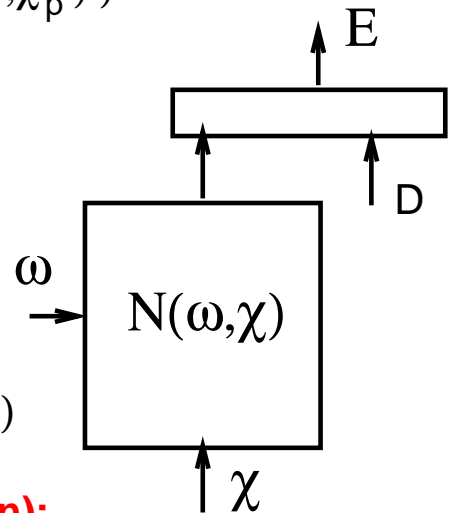
$$E(\omega) = 1/2 \sum_p (D_p - N(\omega, \chi_p))' (D_p - N(\omega, \chi_p))$$

Gradient:

$$\frac{\partial E(\omega)}{\partial \omega} = -\sum_p (D_p - N(\omega, \chi_p))' \frac{\partial N(\omega, \chi_p)}{\partial \omega}$$

Hessian:

$$H(\omega) = \sum_p \frac{\partial N(\omega, \chi_p)}{\partial \omega}' \frac{\partial N(\omega, \chi_p)}{\partial \omega} + \sum_p (D_p - N(\omega, \chi_p))' \frac{\partial^2 N(\omega, \chi_p)}{\partial \omega \partial \omega}$$



Simplified Hessian (square of the Jacobian):

$$H(\omega) = \sum_p \frac{\partial N(\omega, \chi_p)}{\partial \omega}' \frac{\partial N(\omega, \chi_p)}{\partial \omega}$$

**Jacobian: $N \times O$ matrix
(O : number of outputs)**

- the resulting approximate Hessian is positive semi-definite
- dropping the second term is equivalent to assuming that the network is a linear function of the parameters

RECIPE for computing the k-th column of the Jacobian:

for all training patterns {

forward prop

set gradients of output units to 0;

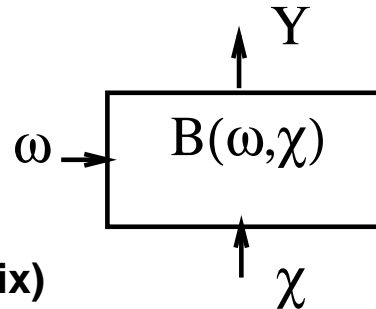
set gradient of k-th output unit to 1;

back propagate; accumulate gradient;

}

BACKPROPAGATING SECOND DERIVATIVES

A multilayer system composed of functional blocs. Consider one of the blocs with I inputs, O outputs, and N parameters



Assuming we know $\frac{\partial^2 E}{\partial Y^2}$ (OxO matrix)

what are $\frac{\partial^2 E}{\partial \omega^2}$ (NxN matrix) and $\frac{\partial^2 E}{\partial \chi^2}$ (IxI matrix)

Chain rule for 2nd derivatives:

$$\frac{\partial^2 E}{\partial \omega^2} = \frac{\partial Y'}{\partial \omega} \frac{\partial^2 E}{\partial Y^2} \frac{\partial Y}{\partial \omega} + \frac{\partial E}{\partial Y} \frac{\partial^2 Y}{\partial \omega^2}$$

\uparrow \uparrow \uparrow \uparrow \uparrow \uparrow
 NxN NxO OxO OxN 1xO OxNxN

ignore this!

The above can be used to compute a bloc diagonal subset of the Hessian

If the term in the red square is dropped, the resulting Hessian estimate will be positive semi-definite

If we are only interested in the diagonal terms, it reduces to:

$$\frac{\partial^2 E}{\partial \omega_{ii}^2} = \sum_k \frac{\partial^2 E}{\partial Y_{kk}^2} \left(\frac{\partial Y_{kk}}{\partial \omega_{ii}} \right)^2 \text{ (and same with } \chi \text{ instead of } \omega)$$

BACKPROPAGATING THE DIAGONAL HESSIAN IN NEURAL NETS

(with the square Jacobian approximation)

[LeCun 87, Becker&LeCun 88, LeCun 89]

Sigmoids (and other scalar functions)

$$\frac{\partial^2 E}{\partial Y_k^2} = \frac{\partial^2 E}{\partial Z_k^2} (f'(Y_k))^2$$

Weighted sums

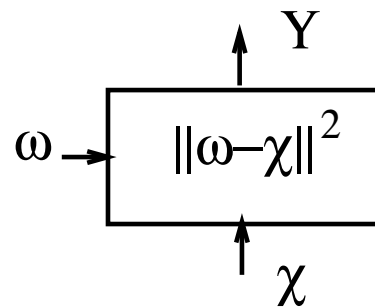
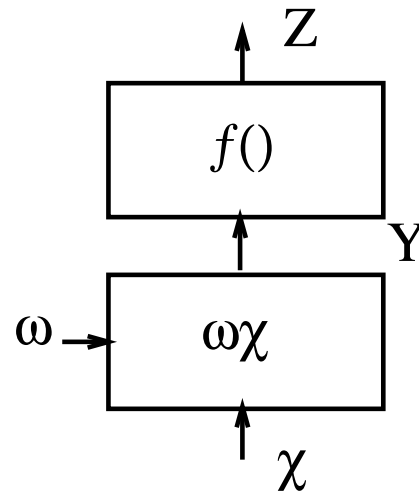
$$\frac{\partial^2 E}{\partial \omega_{ki}^2} = \frac{\partial^2 E}{\partial Y_k^2} \chi_i^2$$

$$\frac{\partial^2 E}{\partial \chi_i^2} = \sum_k \frac{\partial^2 E}{\partial Y_k^2} \omega_{ki}^2$$

RBFs

$$\frac{\partial^2 E}{\partial \omega_{ki}^2} = \frac{\partial^2 E}{\partial Y_k^2} (\chi_i - \omega_{ki})^2$$

$$\frac{\partial^2 E}{\partial \chi_i^2} = \sum_k \frac{\partial^2 E}{\partial Y_k^2} (\chi_i - \omega_{ki})^2$$



(the 2nd derivatives with respect to the weights should be averaged over the training set)

SAME COST AS REGULAR BACKPROP

the "OBD" network pruning techniques uses this procedure [LeCun,Denker&Solla 90]

COMPUTING THE PRODUCT OF THE HESSIAN BY A VECTOR

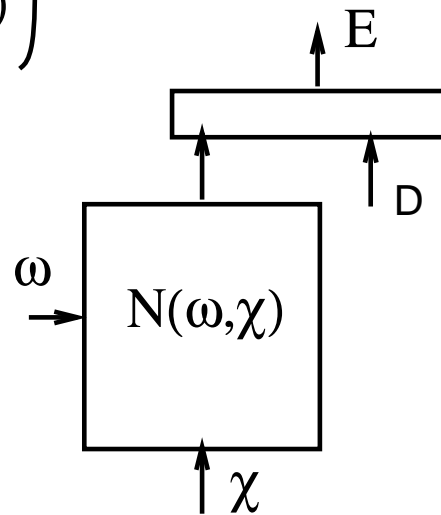
(without computing the Hessian itself)

Finite difference:

$$H\Psi \approx \frac{1}{\alpha} \left(\frac{\partial E}{\partial \omega} (\omega + \alpha\Psi) - \frac{\partial E}{\partial \omega} (\omega) \right)$$

RECIPE for computing the product of a vector Ψ by the Hessian:

- 1– compute gradient
- 2– add $\alpha\Psi$ to the parameter vector
- 3– compute gradient with perturbed parameters
- 4– subtract result of 1 from 3, divide by α



This method can be used to compute the principal eigenvector and eigenvalue of H by the power method.

By iterating $\Psi \leftarrow H\Psi / \|\Psi\|$ Ψ

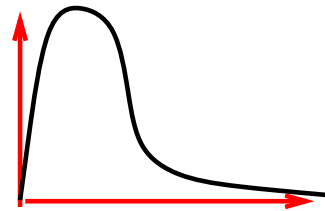
will converge to the principal eigenvector of H and $\|\Psi\|$ to the corresponding eigenvalue [LeCun, Simard&Pearlmutter 93]

A more accurate method which does not use finite differences (and has the same complexity) has recently been proposed [Pearlmutter 93]

ANALYSIS OF THE HESSIAN IN MULTILAYER NETWORKS

- What does the Hessian of a multilayer network look like?
- How does it change with the architecture and the details of the implementation?
- Typically, the distribution of eigenvalues of a multilayer network looks like this:

a few small eigenvalues, a large number of medium ones, and a small number of very large ones



These large ones are the killers

They come from:

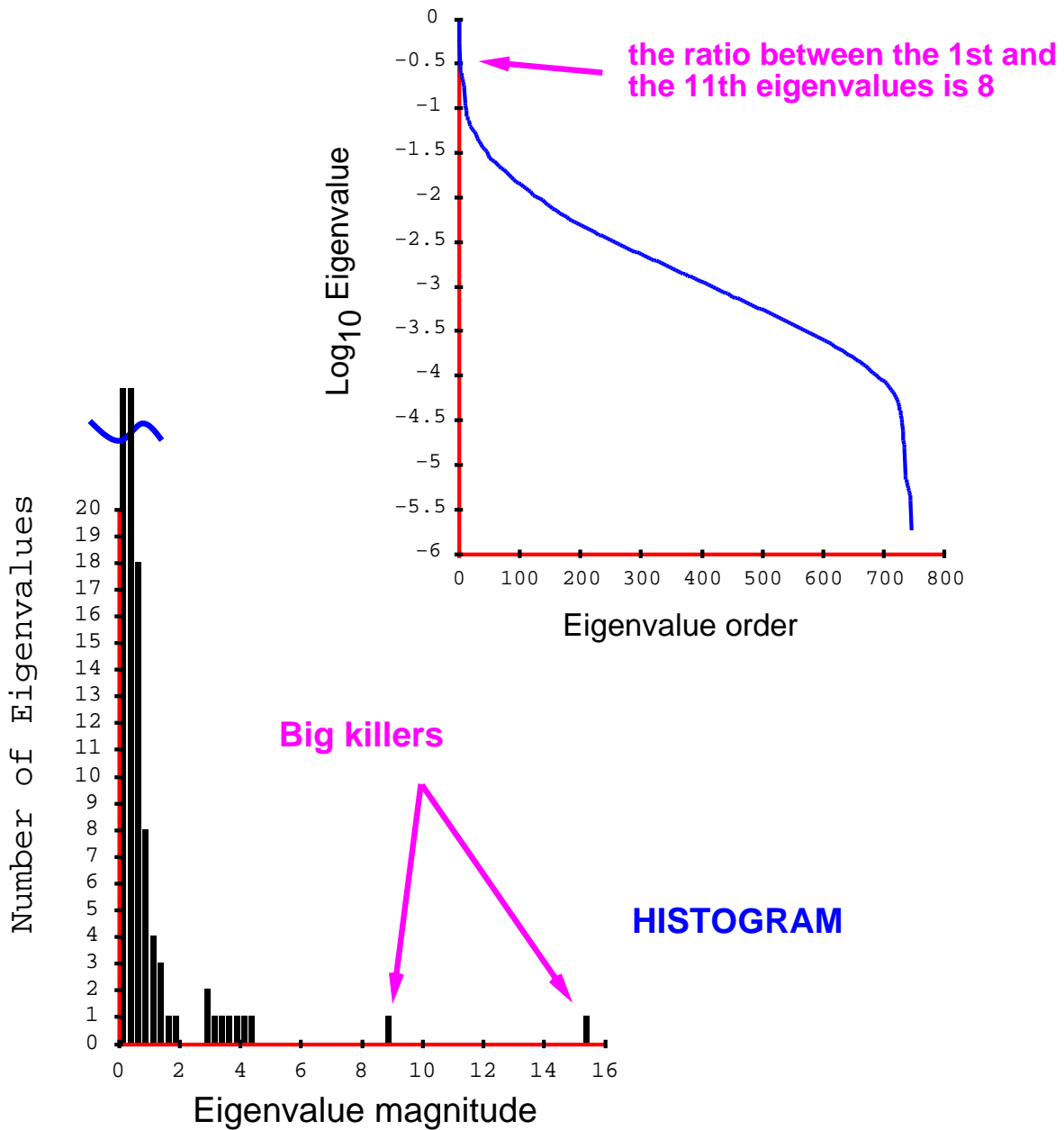
- non-zero mean inputs or neuron states
- wide variations second derivatives from layer to layer
- correlations between state variables

for more details see [LeCun, Simard&Pearlmutter 93]
[LeCun, Kanter&Solla 91]

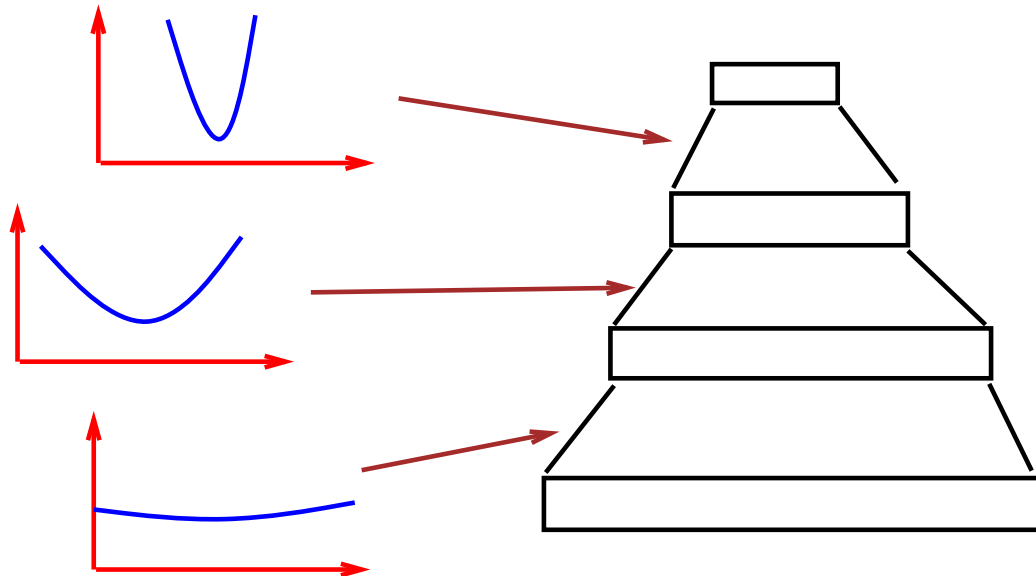
EIGENVALUE SPECTRUM

Network: 256–128–64–10 with local connections and shared weights (around 750 parameters)

Data set: 320 handwritten digits



MULTILAYER NETWORKS HESSIAN



The second derivative is often smaller in lower layers. The first layer weights learn very slowly, while the last layer weights change very quickly.

This can be compensated for using the diagonal 2nd derivatives (more on this later)

CLASSICAL 2ND ORDER OPTIMIZATION METHODS

- **Conjugate Gradient methods**
 - **$O(N)$ methods**
 - **do not use the Hessian directly**
 - **attempts to find descent directions that minimally disturb the result of the previous iterations**
 - **uses line search. Works only in BATCH.**
- **Gauss–Newton and Levenberg–Marquardt methods**
 - **use the square Jacobian approximation**
 - **works only for mean–square error**
 - **mainly designed for BATCH**
 - **$O(N^3)$**
- **Quasi–Newton methods (BFGS)**
 - **iteratively computes an estimate of the inverse Hessian.**
 - **requires a line search. Works only in BATCH.**
 - **$O(N^2)$**

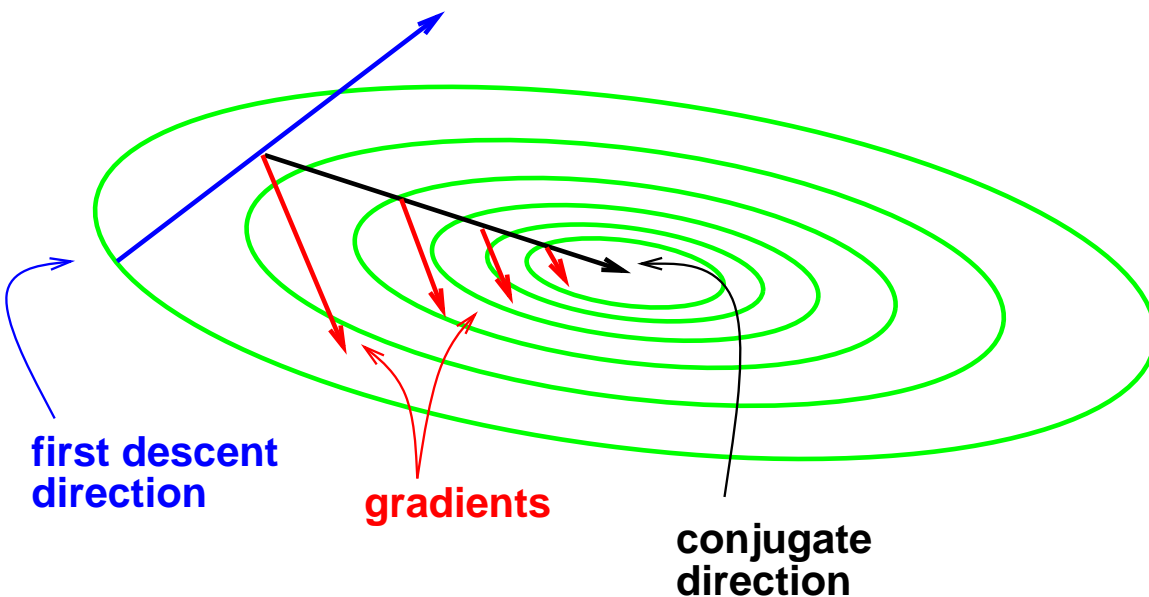
**[Dennis&Schnabel 83], [Fletcher 87],
[Press et al. 88] [Battiti 92]**

CONJUGATE GRADIENT

[Hestenes&Stiefel 52], [Fletcher&Reeves 64]
[Polak 71] (see [Fletcher 87])

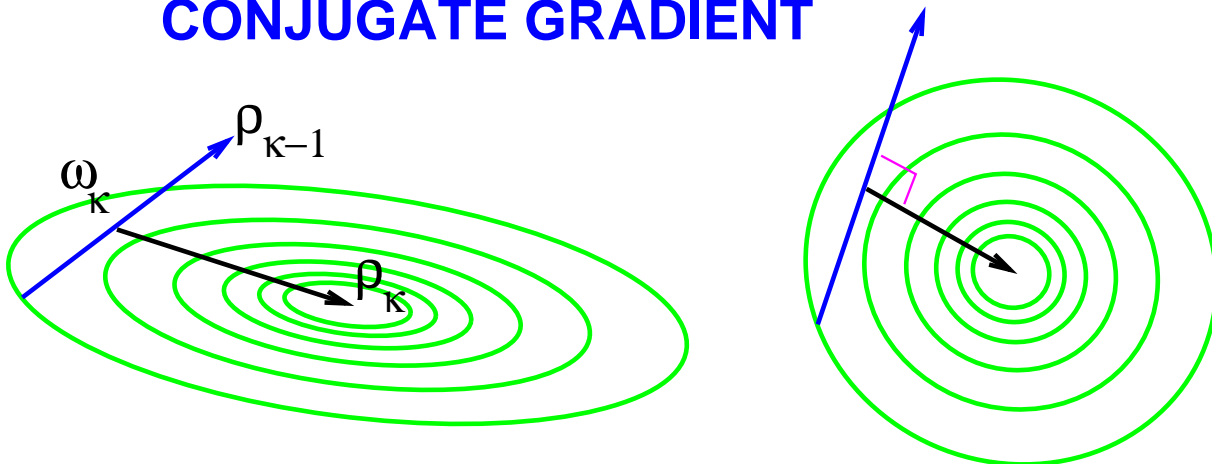
MAIN IDEA: to find a descent direction which does not spoil the result of the previous iterations

Pick a descent direction (say the gradient), find the minimum along that direction (line search). Now, find a direction along which the gradient does not change its direction, but merely its length (conjugate direction). Moving along that direction will not spoil the result of the previous iteration



There are two slightly different formulae:
Fletcher-Reeves & Polak-Ribiere

CONJUGATE GRADIENT



Conjugate directions are like **ORTHOGONAL** directions in the space where the Hessian is the identity matrix

p and q are conjugate $\iff p'Hq = 0$

ρ_k : descent direction at iteration k

$$\rho_k = -\nabla E(\omega_k) + \beta_k \rho_{k-1}$$

$$\beta_k = \frac{\nabla E(\omega_k)' \nabla E(\omega_k)}{\nabla E(\omega_{k-1})' \nabla E(\omega_{k-1})} \quad \text{(Fletcher-Reeves)}$$

$$\beta_k = \frac{(\nabla E(\omega_k) - \nabla E(\omega_{k-1}))' \nabla E(\omega_k)}{\nabla E(\omega_{k-1})' \nabla E(\omega_{k-1})} \quad \text{(Polak-Ribiere)}$$

- A good line search must be done along each descent direction (works only in batch)
- Convergence in N iterations is guaranteed for a quadratic function with N variables

CONJUGATE GRADIENT

- Conjugate gradient is simple and effective. the Polak–Ribiere formula seems to be more robust for non–quadratic functions.
- The conjugate gradient formulae can be viewed as "smart ways" to choose the momentum.
- Conjugate gradient has been used with success in the context of multilayer network training [Kramer&Sangiovani–Vincentelli 88, Barnard&Cole 88, Bengio&Moore 89, Møller 92, Hinton's group in Toronto.....]
- It seems particularly appropriate for moderate size problems with relatively low redundancy in the data.
Typical applications include function approximation, robotic control, times–series prediction, and other real–valued problems (especially if a high accuracy solution is sought).
On large classification problems, stochastic backprop is faster.
- The main drawback of CG is that it is a BATCH method, partly due to its requirement for an accurate line search.
There have been attempts to solve that problem using "mini–batches" [Møller 92].

BFGS and Quasi-Newton methods

[Fletcher 87],[Dennis&Schnabel 83],
[Watrous 88],[Battiti 92].

Quasi-Newton (or secant) methods attempt to keep a positive definite estimate of the INVERSE HESSIAN directly, without requiring matrix inversion, and by only resorting to the gradient information.

They work as follows:

- 1- pick a positive definite matrix M (say $M=I$)
- 2- set search direction $\rho = M \nabla E(\omega)$
- 3- line search along ρ giving $\omega \leftarrow \omega - \eta\rho$
- 4- update estimate of inverse Hessian M

Compared to Newton method, Quasi-Newton methods only require the first derivative, they use positive definite approximations to the inverse Hessian (which means they go downhill), and they require $O(N^2)$ operations per iteration.

All of the quasi-Newton methods are BATCH methods

There exist several Quasi-Newton methods, but the most successful is the Broyden-Fletcher-Goldfarb-Shanno (BFGS) method.

BFGS

past parameter vector:	ω_{k-1}
present parameter vector:	ω_k
past parameter increment:	$\delta = \omega_k - \omega_{k-1}$
past gradient:	$\nabla E(\omega_{k-1})$
present gradient:	$\nabla E(\omega_k)$
past gradient increment:	$\phi = \nabla E(\omega_k) - \nabla E(\omega_{k-1})$
present inverse Hessian:	M_{k-1}
future inverse Hessian	M_k

update inverse Hessian estimate

$$M_k = M_{k-1} + \left(1 + \frac{\phi' M \phi}{\delta' \phi}\right) \frac{\delta \delta'}{\delta' \phi} - \left(\frac{\delta \phi' M + M \phi \delta}{\delta \phi}\right)$$

compute descent direction $\rho_{k+1} = M_k \nabla E(\omega_k)$

line search $\omega_{k+1} = \omega_k - \eta_{k+1} \rho_{k+1}$

- it is an $O(N^2)$ algorithm BUT
- it requires storing an $N \times N$ matrix
- it is a BATCH algorithm (requires a line search)

Only practical for VERY SMALL networks with non-redundant training sets.

Several variations exist that attempt to reduce the storage requirements:

- limited storage BFGS [Nocedal 80]
- memoryless BFGS, OSS [Battiti 92]

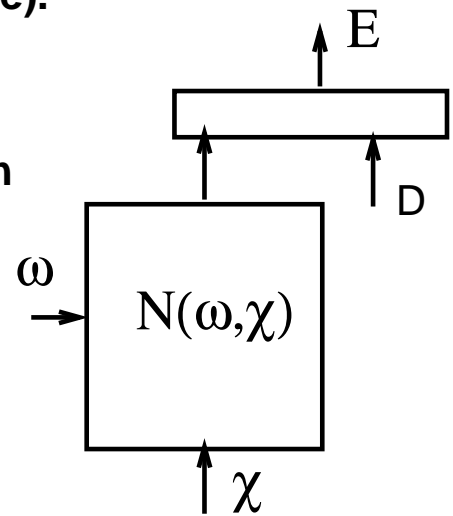
GAUSS-NEWTON AND LEVENBERG-MARQUARDT METHODS

These methods only apply to Mean-Square Error objective functions (non-linear least square).

Gauss-Newton algorithm:

like Newton but the Hessian is approximated by the square of the jacobian (which is always positive semidefinite)

$$\Delta\omega = \left(\sum_p \frac{\partial N(\omega, \chi_p)}{\partial \omega} \frac{\partial N(\omega, \chi_p)}{\partial \omega} \right)^{-1} \nabla E(\omega)$$



Levenberg-Marquardt algorithm:

like Gauss-Newton, but has a safeguard parameter to prevent it from blowing up if some eigenvalues are small

$$\Delta\omega = \left(\sum_p \frac{\partial N(\omega, \chi_p)}{\partial \omega} \frac{\partial N(\omega, \chi_p)}{\partial \omega} + \mu I \right)^{-1} \nabla E(\omega)$$

- Both are $O(N^3)$ algorithms
- they are widely used in statistics for regression
- they are only practical for small numbers of parameters.
- they do not require a line search, so in principle they can be used in stochastic mode (although that has not been tested)

4

**APPLYING SECOND ORDER
METHODS TO MULTILAYER NETS**

(NON)APPLICABILITY OF 2nd ORDER METHODS TO NEURAL-NET LEARNING

BAD NEWS:

- Full Hessian techniques (GN, LM, BFGS) can only apply to small networks. But small networks are not the ones we need to speed up most.**
- Most 2nd order techniques (CG, BFGS....) require a line search, and therefore are not directly usable in stochastic mode**
- Many heuristic tricks (adaptive learning rates....) also apply to batch only.**

On large classification problems, a carefully tuned stochastic gradient is hard to beat.

On smaller problems requiring accurate real-valued outputs (function approximation, control...), conjugate gradient (with Polak-Ribiere) offers the best combination of speed, reliability and simplicity.

This section is devoted to 2nd order techniques specifically designed for large neural-net training

MINI BATCH METHODS

Attempts at applying Conjugate Gradient to large and redundant problems have been made [Kramer&Sangiovani-Vincentelli 88], [Møller 92]

They use "mini batches": subsets of increasing sizes are used.

Møller proposes a systematic way of choosing the size of the mini batch.

He also uses a variant of CG which he calls "scaled CG". Essentially, the line search is replaced by a 1D Levenberg-Marquardt-like algorithm.

A STOCHASTIC DIAGONAL LEVENBERG–MARQUARDT METHOD

[LeCun 87, Becker&LeCun 88, LeCun 89]

THE MAIN IDEAS:

- use formulae for the backpropagation of the diagonal Hessian (shown earlier) to keep a running estimate of the second derivative of the error with respect to each parameter.
- use these term in a "Levenberg–Marquardt" formula to scale each parameter's learning rate

Each parameter (weight) ω_{ki} has its own learning rate η_{ki} computed as:

$$\eta_{ki} = \frac{\epsilon}{\frac{\partial^2 E}{\partial \omega_{ki}^2} + \mu}$$

ϵ is a global "learning rate"

$\frac{\partial^2 E}{\partial \omega_{ki}^2}$ is an estimate of the diagonal second derivative with respect to weight (ki)

μ is a "Levenberg–Marquardt" parameter to prevent η_{ki} from blowing up if the 2nd derivative is small

A STOCHASTIC DIAGONAL LEVENBERG-MARQUARDT METHOD

The second derivatives $\frac{\partial^2 E}{\partial \omega_{ki}^2}$ can be computed using a running average formula over a subset of the training set prior to training:

$$\frac{\partial^2 E}{\partial \omega_{ki}^2} \leftarrow (1-\gamma) \frac{\partial^2 E}{\partial \omega_{ki}^2} + \gamma \frac{\partial^2 E^p}{\partial \omega_{ki}^2}$$

new estimate of 2nd der. previous estimate small constant instantaneous 2nd der. for pattern p

The instantaneous second derivatives are computed using the formula in the slide entitled: "BACKPROPAGATING THE DIAGONAL HESSIAN IN NEURAL NETS"

Since the second derivatives evolve slowly, there is no need to reestimate them often.

They can be estimated once at the beginning by sweeping over a few hundred patterns.

Then, they can be reestimated every few epochs.

The additional cost over regular backprop is negligible.

Is usually about 3 times faster than carefully tuned stochastic gradient.

Stochastic Diagonal Levenberg–Marquardt

data set: set-1 (100 examples, 2 gaussians)
network: 1 linear unit, 2 inputs, 1 output.
2 weights, 1 bias.

Learning rates:

$$\eta_0 = 0.12$$

$$\eta_1 = 0.03$$

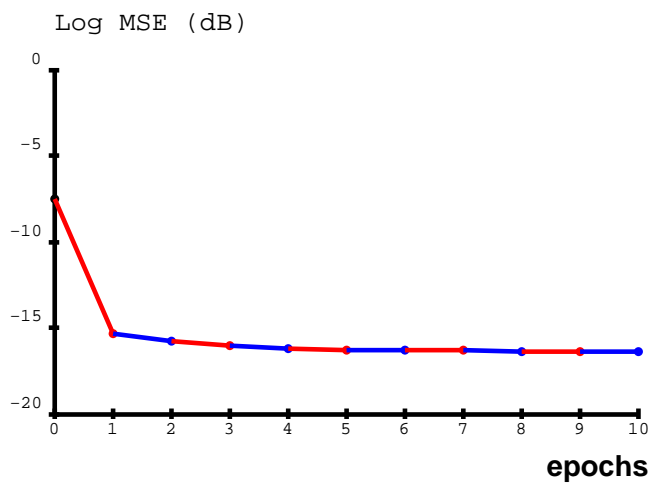
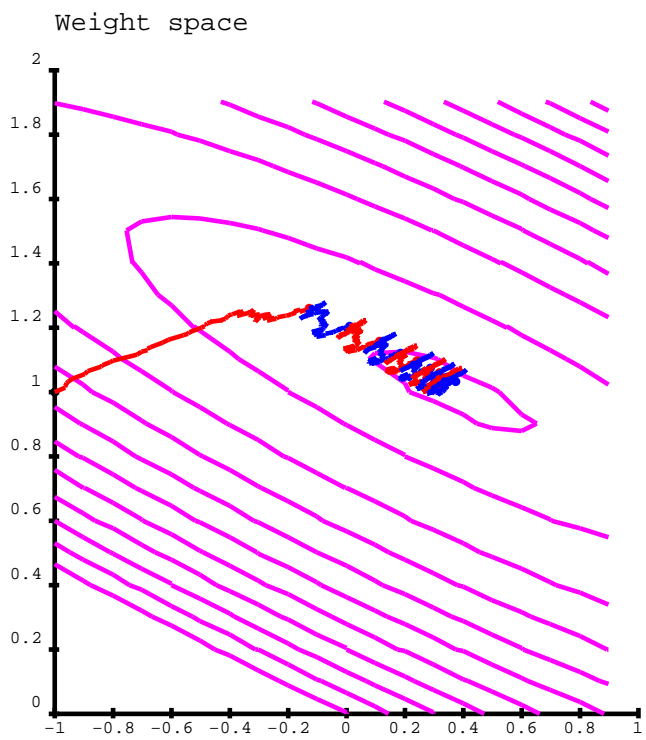
$$\eta_2 = 0.02$$

Hessian largest eigenvalue:

$$\lambda_{\max} = 0.84$$

Maximum admissible Learning rate (batch):

$$\eta_{\max} = 2.38$$



Stochastic Diagonal Levenberg–Marquardt

data set: set-1 (100 examples, 2 gaussians)
network: 1 linear unit, 2 inputs, 1 output.
2 weights, 1 bias.

Learning rates:

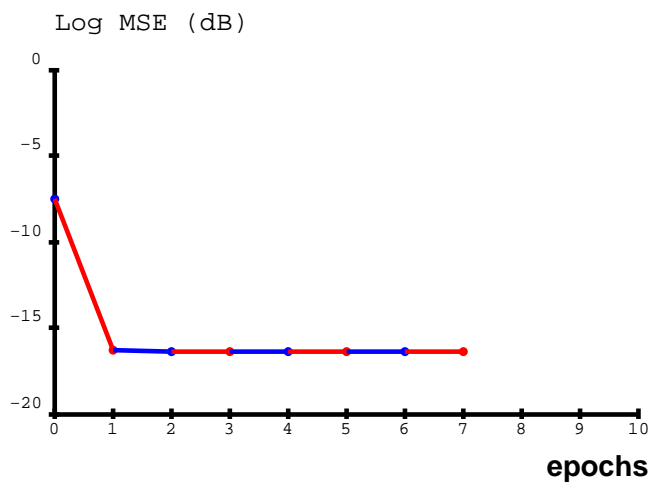
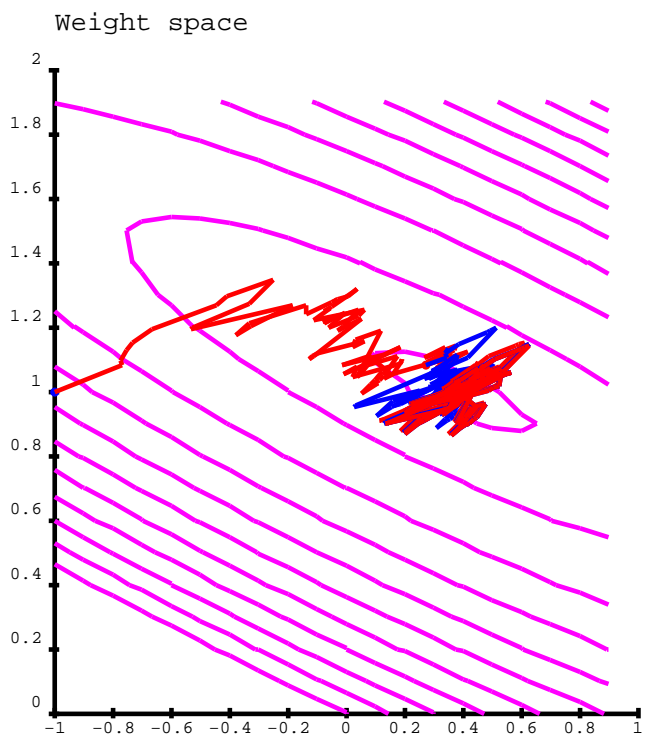
$$\eta_0 = 0.76$$
$$\eta_1 = 0.18$$
$$\eta_2 = 0.12$$

Hessian largest eigenvalue:

$$\lambda_{\max} = 0.84$$

Maximum admissible Learning rate (batch):

$$\eta_{\max} = 2.38$$



COMPUTING THE PRINCIPAL EIGENVALUE/VECTOR OF THE HESSIAN

without computing the Hessian

IDEA #1 (the power method):

1 – Choose a vector Ψ at random

2 – iterate: $\Psi \leftarrow H \frac{\Psi}{\|\Psi\|}$

Annotations:

- OLD ESTIMATE OF EIGENVECTOR (points to Ψ in the fraction)
- HESSIAN (points to H)
- ESTIMATE OF EIGENVALUE (points to $\|\Psi\|$)
- NEW ESTIMATE OF EIGENVECTOR (points to the new Ψ)

Ψ will converge to the principal eigenvector (or a vector in the principal eigenspace)

$\|\Psi\|$ will converge to the corresponding eigenvalue

COMPUTING THE PRODUCT $H\Psi$

IDEA #2 (Taylor expansion):

$$\Psi \leftarrow \frac{1}{\alpha} \left(\frac{\partial E}{\partial \omega} \left(\omega + \alpha \frac{\Psi}{\|\Psi\|} \right) - \frac{\partial E}{\partial \omega} (\omega) \right)$$

NEW ESTIMATE OF EIGENVECTOR

OLD ESTIMATE OF EIGENVECTOR

"SMALL" CONSTANT

PERTURBED GRADIENT

GRADIENT

One iteration of this procedure requires 2 forward props and 2 backward props for each pattern in the training set.

This converges very quickly to a good estimate of the largest eigenvalue of H

ON-LINE COMPUTATION OF Ψ

IDEA #3 (running average):

$$\Psi \leftarrow (1-\gamma)\Psi + \gamma \frac{1}{\alpha} \left(\frac{\partial E^p}{\partial \omega} \left(\omega + \alpha \frac{\Psi}{\|\Psi\|} \right) - \frac{\partial E^p}{\partial \omega} (\omega) \right)$$

NEW ESTIMATE OF EIGENVECTOR

OLD ESTIMATE OF EIGENVECTOR

"SMALL" CONSTANTS

PERTURBED GRADIENT FOR CURRENT PATTERN

GRADIENT FOR CURRENT PATTERN

This procedure converges VERY quickly to the largest eigenvalue of the AVERAGE Hessian.

The properties of the average Hessian determine the behavior of ON-LINE gradient descent (stochastic, or per-sample update).

EXPERIMENT: A shared-weight network with 5 layers of weights, 64638 connections and 1278 free parameters. Training set: 1000 handwritten digits.

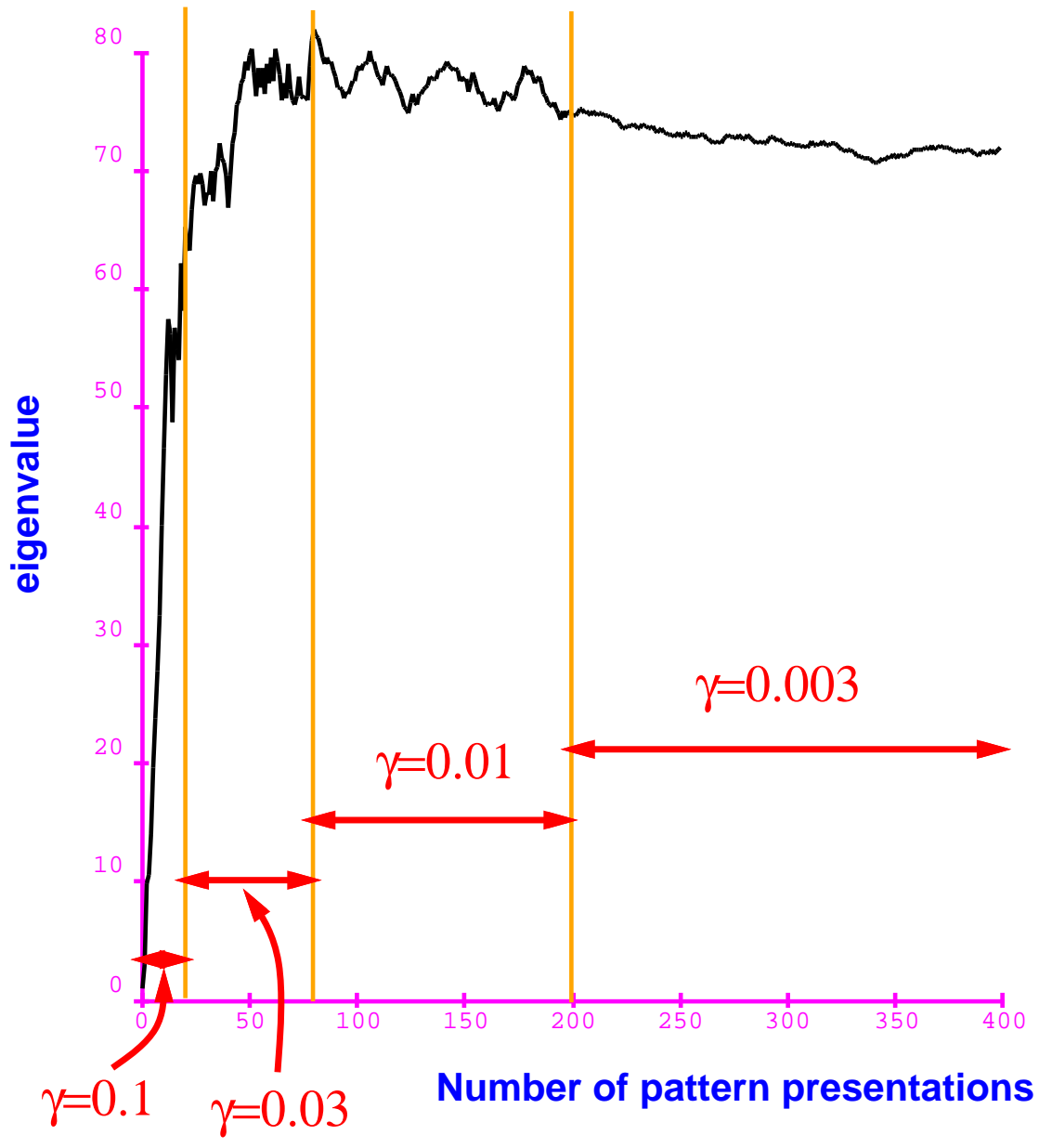
Correct order of magnitude is obtained in less than 100 pattern presentations (10% of training set size)

The fluctuations of the average Hessian over the training set are small.

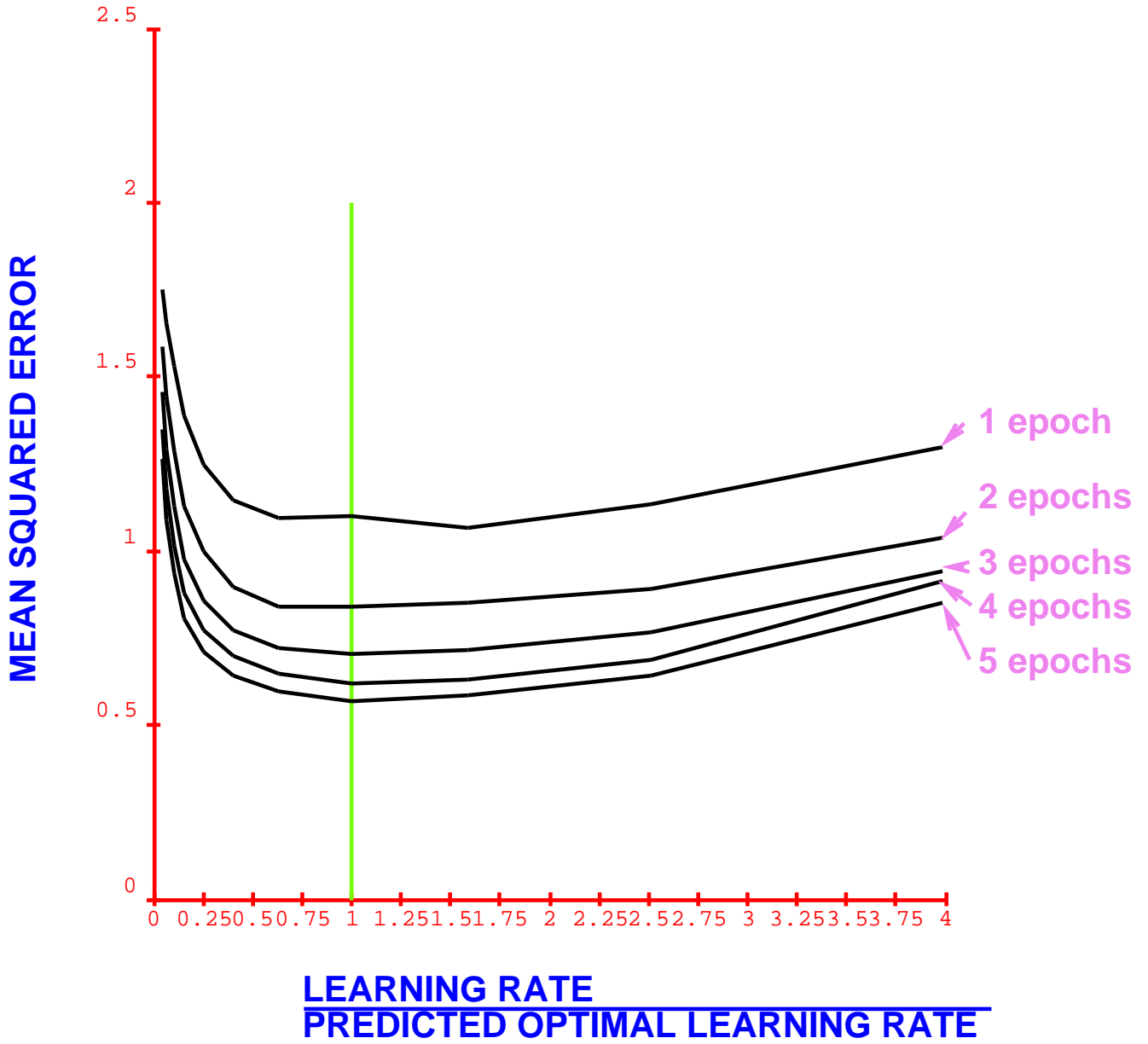
RECIPE

$$\Psi \leftarrow (1-\gamma)\Psi + \gamma \frac{1}{\alpha} \left(\frac{\partial E^p}{\partial \omega} \left(\omega + \alpha \frac{\Psi}{\|\Psi\|} \right) - \frac{\partial E^p}{\partial \omega} (\omega) \right)$$

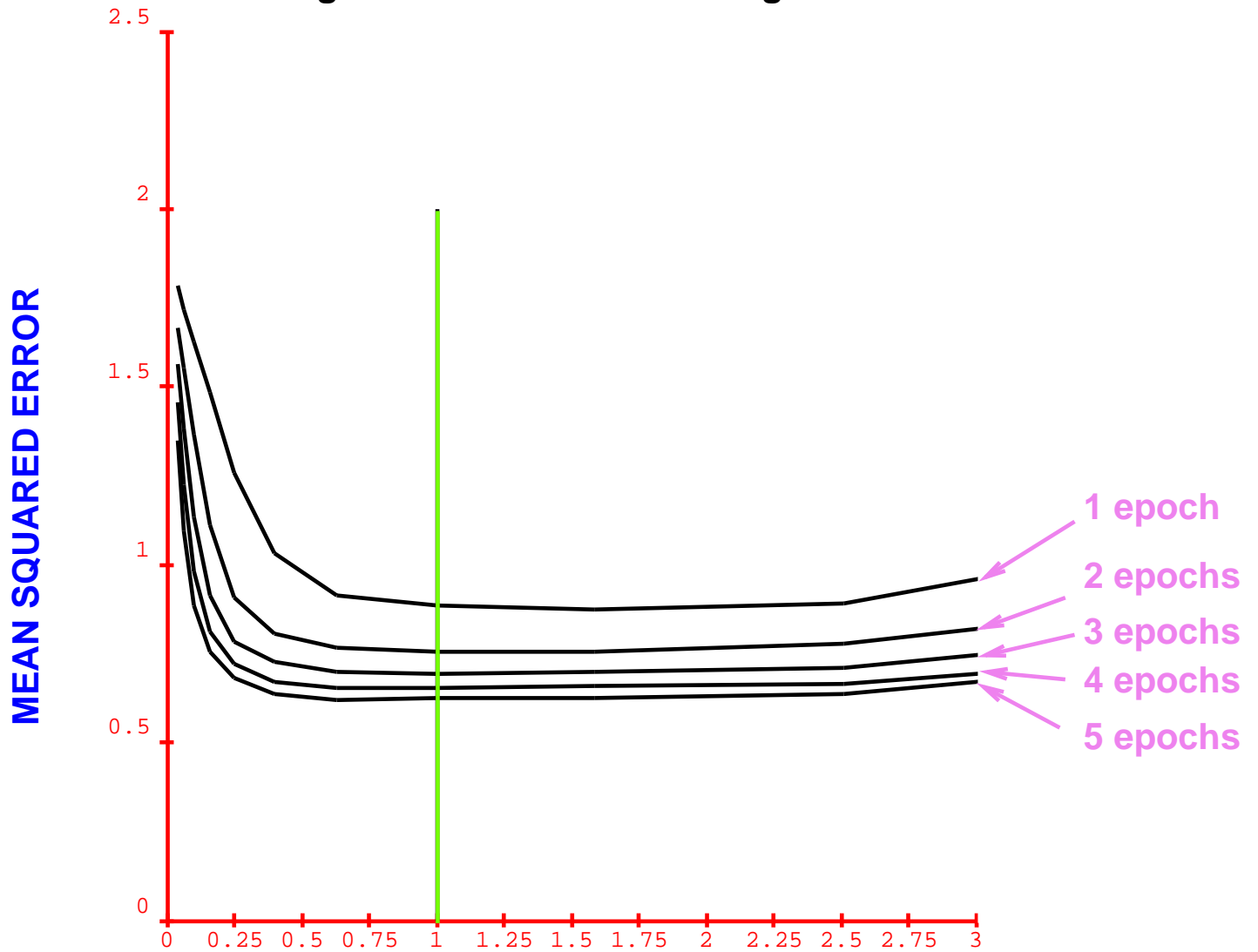
- 1 – Pick initial eigenvector estimate at random
- 2 – present input pattern, and desired output. perform forward prop and backward prop. Save gradient vector $G(w)$
- 3 – add $\alpha \frac{\Psi}{\|\Psi\|}$ to current weight vector
- 4 – perform forward prop and backward prop with perturbed weight vector. Save gradient vector $G'(w)$
- 5 – compute difference $G'(w) - G(w)$. and divide by α update running average of eigenvector with the result
- 6 – goto 2 unless a reasonably stable result is obtained
- 7 – the optimal learning rate is $\frac{1}{\|\Psi\|}$



Network: 784x30x10 fully connected
Training set: 300 handwritten digits



Network: 1024x1568x392x400x100x10
with 64638 (local) connections
and 1278 shared weights
Training set: 1000 handwritten digits



LEARNING RATE
PREDICTED OPTIMAL LEARNING RATE