# 4. Gradient based Learning

- Networks of **continuous units**
- **Regression** problems
- **Gradient descent,** *backpropagation of error*
- The role of the **learning rate**
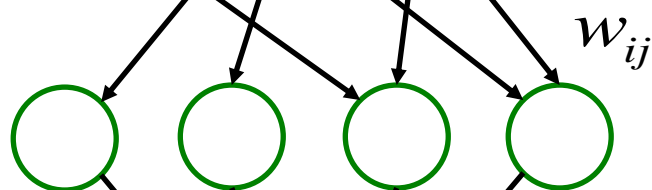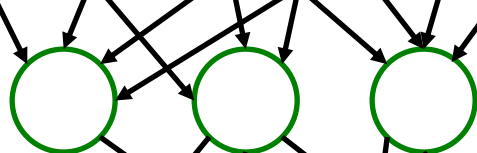- **Online learning,** *stochastic approximation*

**input layer** (external stimulus)

**layered architecture**
(here: 6–3–4–1)
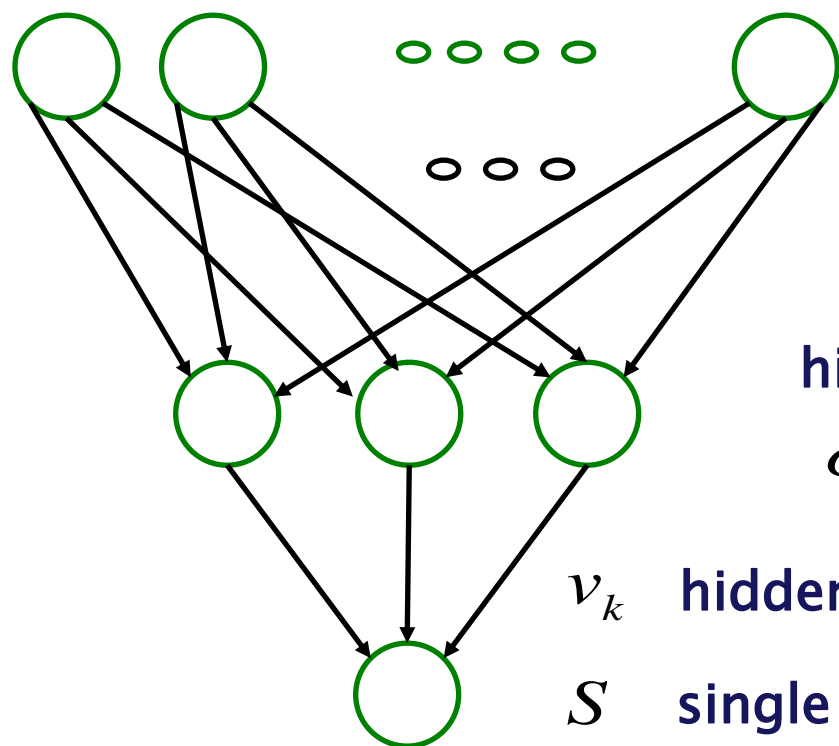
**directed connections**
(here: only to next layer)

**hidden units**
(internal representation)

$w_{ij}$

$$S_i = g\left(\sum w_{ij} S_j\right)$$

↑ previous layer only

**output unit(s)**
(function of input vector)

input units $\quad \xi_j \in \mathbb{R}, \boldsymbol{\xi} \in \mathbb{R}^N$

$w_j^{(k)}$    input to hidden weights

hidden layer units

$$\sigma_k = g\left(\sum w_j^{(k)} \xi_j\right)$$

$v_k$    hidden to output weights

$S$    single output unit

output = **non-linear function** of input variables:

$$S = g\left(\sum_{k=1}^{K} v_k \ \sigma_k\right) = g\left(\sum_{k=1}^{K} v_k g\left(\sum w_j^{(k)} \xi_j\right)\right)$$

parameterized by set of all weights (and threshold)

continuous activation functions, e.g. $\qquad g(x) = \tanh(\gamma \mathbf{x})$
for all nodes in the network

given a network architecture, the weights (and thresholds)
parameterize a continuous function:

$$\xi \in \mathbb{R}^N \quad \rightarrow \quad \sigma(\xi) \in \mathbb{R} \qquad \text{(here: single output unit)}$$

Learning as **regression problem**

set of examples $\quad \left\{ \xi^{\mu}, \tau(\xi^{\mu}) \right\}_{\mu=1}^{P} \quad$ with cont. labels $\quad \tau \in \mathbb{R}$

**training**:
(approximately) implement $\quad \sigma(\xi^{\mu}) = \tau(\xi^{\mu}) \text{ for all } \mu$

**generalization**:
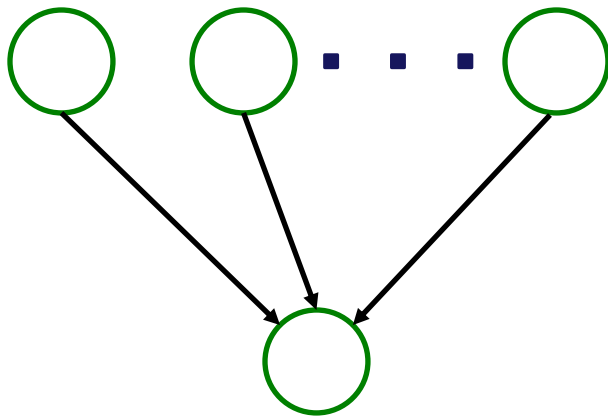application to novel data $\qquad \sigma(\xi) \approx \tau(\xi)$

training strategy:  employ an **error measure**
for comparison of student/teacher outputs

just <u>one</u> very popular and plausible choice:

**quadratic deviation:** $\quad e(\sigma, \tau) = \dfrac{1}{2}(\sigma - \tau)^2$

**cost function:** $\quad E = \dfrac{1}{P}\sum\limits_{\mu=1}^{P} e^{\mu} = \dfrac{1}{P}\sum\limits_{\mu=1}^{P}\dfrac{1}{2}\big(\sigma(\xi^{\mu}) - \tau(\xi^{\mu})\big)^2$

– defined for a given set of example data

– guides the training process

– is a **differentiable function** of weights and thresholds

– training by **gradient descent** minimization of $E$

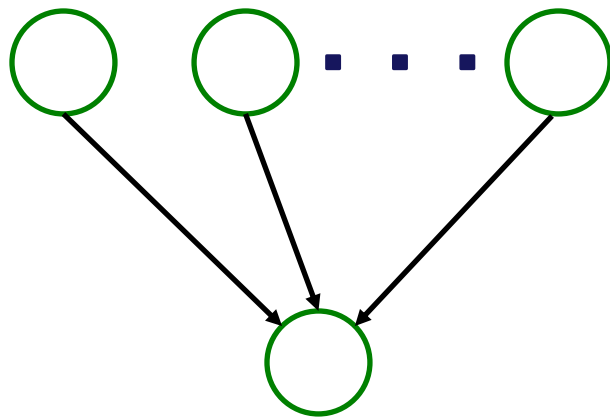$$\xi_j \in \mathbb{R}, \boldsymbol{\xi} \in \mathbb{R}^N$$

$$\mathbf{w} \in \mathbb{R}^N$$

$$\sigma = g\left(\sum_{j=1}^{N} w_j \xi_j\right)$$

$$E(\mathbf{w}) = \frac{1}{P} \sum_{\mu=1}^{P} \frac{1}{2}\left(g(\mathbf{w}\cdot\boldsymbol{\xi}^\mu) - \tau(\boldsymbol{\xi}^\mu)\right)^2$$

$$\frac{\partial E}{\partial w_k} = \frac{1}{P} \sum_{\mu=1}^{P} \left(g(\mathbf{w}\cdot\boldsymbol{\xi}^\mu) - \tau(\boldsymbol{\xi}^\mu)\right) g'(\mathbf{w}\cdot\boldsymbol{\xi}^\mu) \xi_k^\mu$$

$$\nabla_w E = \frac{1}{P} \sum_{\mu=1}^{P} \left(g(\mathbf{w}\cdot\boldsymbol{\xi}^\mu) - \tau(\boldsymbol{\xi}^\mu)\right) g'(\mathbf{w}\cdot\boldsymbol{\xi}^\mu) \boldsymbol{\xi}^\mu$$

$$\xi_j \in \mathbb{R}, \boldsymbol{\xi} \in \mathbb{R}^N$$

$$\mathbf{w} \in \mathbb{R}^N$$

$$\sigma = g\left(\sum_{j=1}^{N} w_j \xi_j\right)$$

frequent choice:     $\mathrm{g(x)=tanh(x)}$     $\mathrm{g'(x)=1\text{-}tanh^2(x)}$

$$\nabla_w E = \frac{1}{P} \sum_{\mu=1}^{P} \left(\tanh(\mathbf{w}\cdot\boldsymbol{\xi}^\mu) - \tau(\boldsymbol{\xi}^\mu)\right) \left(1\text{-}\tanh^2(\mathbf{w}\cdot\boldsymbol{\xi}^\mu)\right) \boldsymbol{\xi}^\mu$$
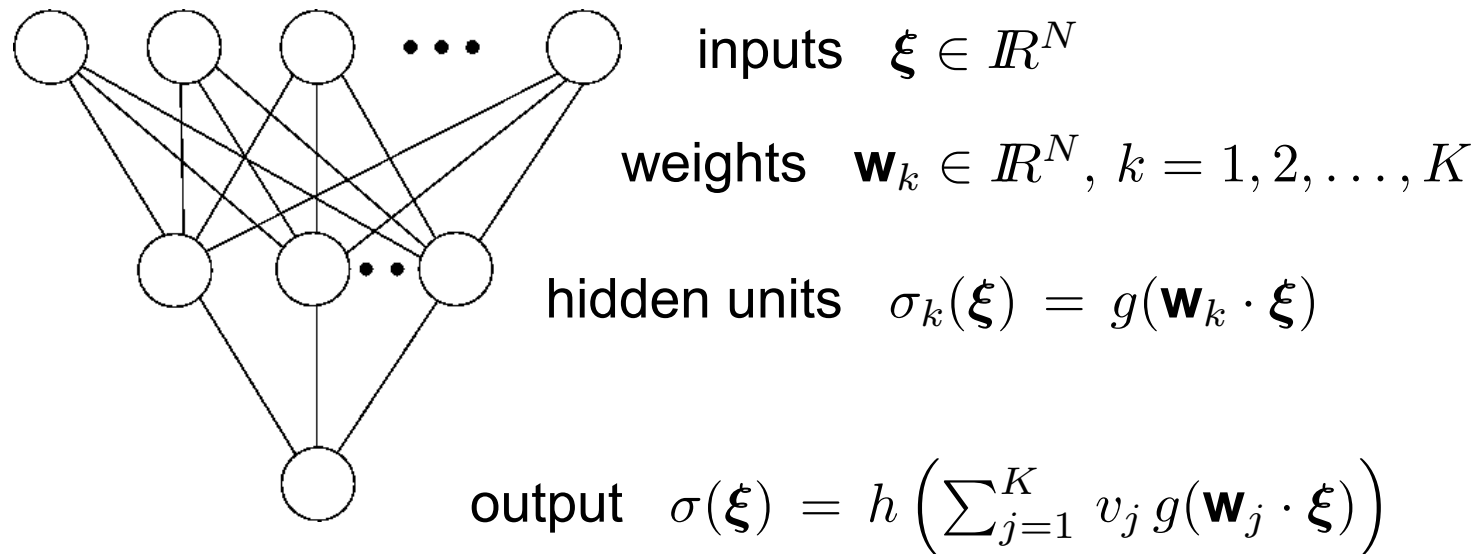
single example contribution:

$$\nabla_w e^\mu = \left(\tanh(\mathbf{w}\cdot\boldsymbol{\xi}^\mu) - \tau(\boldsymbol{\xi}^\mu)\right) \left(1\text{-}\tanh^2(\mathbf{w}\cdot\boldsymbol{\xi}^\mu)\right) \boldsymbol{\xi}^\mu$$

**Backpropagation of Error**

convenient calculation of the gradient in multilayer networks ($\leftarrow$ chain rule)

example:   continuous two-layer network with $K$ hidden units



inputs   $\boldsymbol{\xi} \in I\!\!R^N$

weights   $\mathbf{w}_k \in I\!\!R^N,\ k = 1, 2, \ldots, K$

hidden units   $\sigma_k(\boldsymbol{\xi}) = g(\mathbf{w}_k \cdot \boldsymbol{\xi})$

output   $\sigma(\boldsymbol{\xi}) = h\left(\sum_{j=1}^{K} v_j\, g(\mathbf{w}_j \cdot \boldsymbol{\xi})\right)$

**Exercise:**   derive   $\nabla_{\mathrm{w}_k} E$   and $\frac{\partial E}{\partial v_k}$

the weigths   $\mathbf{w}_k$   and   $v_k$   are used ...

       – *downward*  for the calculation of hidden states and output
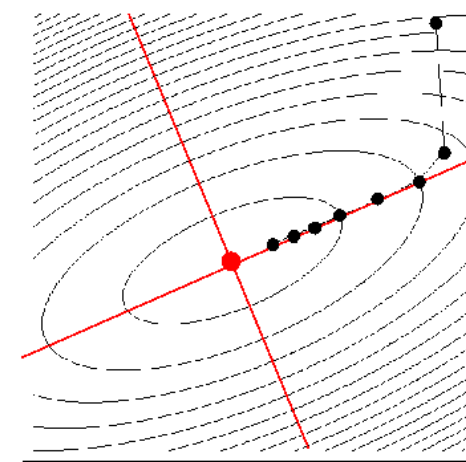
       – *upward*    for the calculation of the gradient

negative gradient gives the **direction of steepest descent** in $E$

simple gradient based minimization of $E$:

sequence $\quad \mathbf{w}_0 \rightarrow \mathbf{w}_1 \rightarrow \ldots \rightarrow \mathbf{w}_t \rightarrow \mathbf{w}_{t+1} \rightarrow \ldots$

with $\quad\quad\quad \mathbf{w}_{t+1} = \mathbf{w}_t - \eta \left.\nabla E\right|_{\mathbf{w}_t}$

approaches <u>some</u> minimum of $E$ (?)

**learning rate rate** $\quad \eta$

– controls the step size of the algorithm

– has to be small enough to ensure convergence

– should be as large as possible to facilitate fast learning

assume $E$ has a (local) minimum in $\mathbf{w}^*$, Taylor expansion in the vicinity:

$$E(\mathbf{w}) \approx E(\mathbf{w}^*) + (\mathbf{w} - \mathbf{w}^*)^T \underbrace{\nabla E|_*}_{=0} + \frac{1}{2}(\mathbf{w} - \mathbf{w}^*)^T H^* (\mathbf{w} - \mathbf{w}^*) + \dots$$

$$E(\mathbf{w}) \approx E(\mathbf{w}^*) + \frac{1}{2}(\mathbf{w} - \mathbf{w}^*)^T H^* (\mathbf{w} - \mathbf{w}^*) \qquad \nabla E|_\mathbf{w} \approx H^* (\mathbf{w} - \mathbf{w}^*)$$

with the positive definite **Hesse matrix** of second derivatives $\qquad H_{ij}^* = \left. \dfrac{\partial^2 E}{\partial w_i \, \partial w_j} \right|_*$

$H^*$ has only pos. eigenvalues $\lambda_i > 0$, orthonormal eigenvectors $\mathbf{u}_i \qquad$ (all $\lambda_i \le \lambda_{max}$)

---

gradient descent in the vicinity of $\mathbf{w}^*$: $\qquad\qquad \mathbf{w}_t - \mathbf{w}^* \equiv \boldsymbol{\delta}_t = \boldsymbol{\delta}_{t-1} - \eta \, \nabla E|_{\mathbf{w}_{t-1}}$

$$\boldsymbol{\delta}_t \approx [I - \eta H^*]\, \boldsymbol{\delta}_{t-1} \approx [I - \eta H^*]^t \, \boldsymbol{\delta}_0 \qquad\qquad \text{expansion in } \{\mathbf{u}_i\}: \; \delta_0 = \sum_i a_i \, \mathbf{u}_i$$

$$\boldsymbol{\delta}_t \approx \sum_i a_i \, [I - \eta H^*]^t \, \mathbf{u}_i = \sum_i a_i \, [1 - \eta \lambda_i]^t \, \mathbf{u}_i$$

with $\mathbf{u}_j^T \mathbf{u}_k = \delta_{jk}$ we obtain $\qquad\qquad\qquad\qquad |\delta_t|^2 = \sum_i a_i^2 \, [1 - \eta \lambda_i]^{2t}$

in detail:

$$w_t = w_{t-1} - \eta \nabla E|_{w_{t-1}}$$

$$\boldsymbol{\delta}_t = \boldsymbol{\delta}_{t-1} - \eta \nabla E|_{w_{t-1}} \approx [I - \eta H^*] \ \boldsymbol{\delta}_{t-1} \approx [I - \eta H^*]^t \boldsymbol{\delta}_o$$

$$\boldsymbol{\delta}_o = \sum_i a_i \ \mathbf{u_i} \qquad H^* \ \mathbf{u}_i = \lambda_i \ \mathbf{u_i} \qquad (1 - \eta H^*) \ \mathbf{u}_i = (1 - \eta \lambda_i) \ \mathbf{u}_i$$
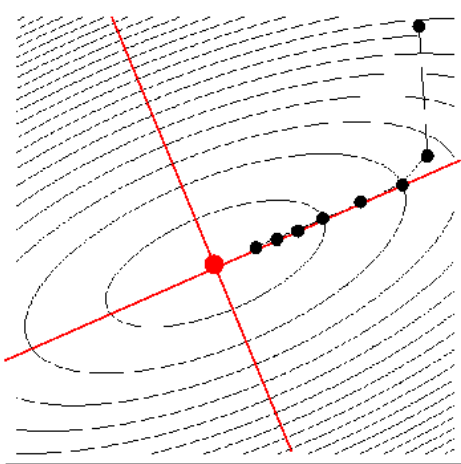
$$\boldsymbol{\delta}_t \approx [I - \eta H^*]^t \ \sum_i a_i \ \mathbf{u}_i = \sum_i a_i \ (1 - \eta \lambda_i)^t \ \mathbf{u_i}$$

$$\boldsymbol{\delta}_t^2 = \boldsymbol{\delta}_t \cdot \boldsymbol{\delta}_t = \sum_{i,j} a_i [1 - \eta \lambda_i]^t \ a_j \ [1 - \eta \lambda_j]^t \ \underbrace{\mathbf{u_i} \cdot \mathbf{u_j}}_{\substack{=1 \text{ if i=j} \\ =0 \text{ else}}}$$

$$= \sum_i a_i^2 [1 - \eta \lambda_i]^{2t}$$

78

iteration approaches the minimum, $\lim\limits_{t\to\infty} |\boldsymbol{\delta}_t| = 0,$ only if $|1 - \eta\lambda_i| < 1$ for all $i$

condition for (local) convergence: $\boldsymbol{\eta} < \boldsymbol{\eta_{max}} = \dfrac{\mathbf{2}}{\boldsymbol{\lambda_{max}}}$ (largest eigenvalue of H*)

$\eta < \dfrac{\eta_{max}}{2} = \dfrac{1}{\lambda_{max}}$ $\qquad$ $\dfrac{1}{\lambda_{max}} < \eta < \dfrac{2}{\lambda_{max}}$ $\qquad$ $\eta > \eta_{max} = \dfrac{2}{\lambda_{max}}$



$1 - \eta\lambda_{max} > 0$ $\qquad\qquad$ $1 - \eta\lambda_{max} < 0$ $\qquad\qquad$ $1 - \eta\lambda_{max} < -1$

smooth convergence $\qquad\qquad$ oscillations $\qquad\qquad$ divergence

... the above considerations

- are only valid close to the minimum

  local minima can have completely different characteristics ($\lambda_{max}$)

- do not concern global convergence properties

  e.g. the choice of the learning rate far from a minimum

**potential problems:**

- $E$ can have (many) local minima far from global optimality

- initial conditions determine which minimum will be approached

- anistropic curvatures can cause strong oscillations

- $E$ can have *saddle points* with $\nabla E = 0$ and/or *flat regions* with $\nabla E \approx 0$

  gradient learning can slow down drastically by, e.g., *plateau states*, see below

**some modifications:**

- improved gradient descent: e.g. time dependent $\eta(t)$

  **momentum:** $\quad \Delta\mathbf{w}_{t+1} = -\eta\,\nabla E + a\,\Delta\mathbf{w}_t \quad$ "keep going"

- sophisticated **optimization methods**:
  line search procedures, conjugate gradient, second order methods,

  e.g. Newton's method ("matrix update" employs $H$), ...

- **different learning rates** for different weights, examples:
  – heuristics: $\quad \eta \propto 1/N$ for input-to-hidden, $\eta \propto 1/K$ for hidden-to-output weights
  – simplified version of "matrix update" (assume $H$ is approximately diagonal):

  update each weight $\quad w_j \quad$ with a learning rate $\quad \eta_j \propto 1 \left/ \dfrac{\partial^2 E}{\partial w_j^2}\right.$

  – learning algorithms realize *descent* in $\ E\ $ as long as $\Delta\mathbf{w}\cdot\nabla E < 0$

- construction of alternative **well-behaved cost functions,** one example:

$$E = \sum_\mu \left\{ \begin{array}{ll} \gamma\,(\sigma - \tau)^2 & \text{if } \ \operatorname{sign}(\sigma) = \operatorname{sign}(\tau) \\ (\sigma - \tau)^2 & \text{if } \ \operatorname{sign}(\sigma) \neq \operatorname{sign}(\tau) \end{array}\right. \quad \text{with } \gamma \text{ increasing from } 0 \text{ to } 1.$$

small $\gamma$: emphasis on correct <u>sign</u> of the output    large $\gamma$: fine tuning of $\sigma$

**stochastic approximation**    (on-line gradient descent)

cost function   $E = \frac{1}{P} \sum_{\mu=1}^{P} e^{\mu} \equiv \overline{e^{\mu}}$   is an **empirical average** over examples

$\rightarrow$   simple approximation of $\nabla E$ by $\nabla e^{\mu}$ for one example only

- select one $\mu \in \{1, 2, \ldots, P\}$ with equal probabilty  $1/P$
- single step:  $\mathbf{w}_{t+1} = \mathbf{w}_t + \Delta \mathbf{w}_t = \mathbf{w}_t - \eta \, \nabla \, e^{\mu}|_{\mathbf{w}_t}$

– computationally cheap compared to *off-line (batch)* gradient descent

– *intrinsic noise*, fewer problems with local minima, flat regions etc.

(when) does the procedure converge?

behavior close to a (local) minimum  $\mathbf{w}^*$ of  $E$?

averaged learning step: 
$$\overline{\Delta \mathbf{w}} = -\eta \, \overline{\boldsymbol{\nabla} \, e^\mu|_{\mathbf{w}}} = -\frac{\eta}{P} \sum_{\mu=1}^{P} \boldsymbol{\nabla} \, e^\mu|_{\mathbf{w}} = -\eta \, \boldsymbol{\nabla} \, E|_{\mathbf{w}}$$

$$\overline{\Delta \mathbf{w}} = 0 \quad \text{for} \quad \mathbf{w} \to \mathbf{w}^*$$

averaged length of $\Delta \mathbf{w}$:
$$\overline{(\Delta \mathbf{w})^2} = \eta^2 \, \overline{\left(\boldsymbol{\nabla} e^\mu|_*\right)^2} > 0$$

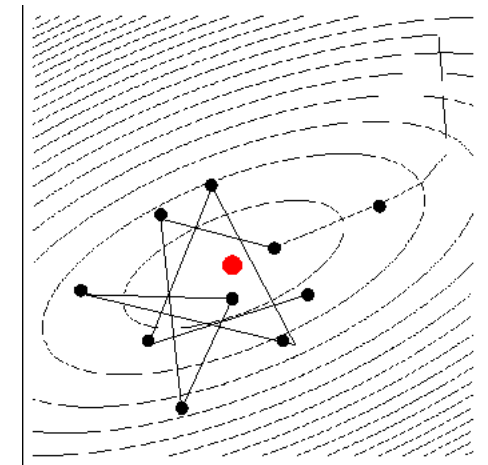($\cancel{= 0}$ ~~and~~ possible <u>if</u> all $e^{\,\mu} = 0$)

<span style="color:blue">generic behavior</span>
for constant rate $\eta > 0$:
$$\lim_{t \to \infty} \left(\Delta \mathbf{w}_t\right)^2 > 0$$

(fluctuations remain non-zero)



**convergence** in the sense of $\left(\Delta \mathbf{w}\right)^2 \to 0$ only if $\eta(t) \to 0$ for $t \to \infty$

one can show: $\lim_{t \to \infty} \sum_t \eta(t) \to \infty$ but $\lim_{t \to \infty} \sum_t \eta(t)^2 < \infty$ is required
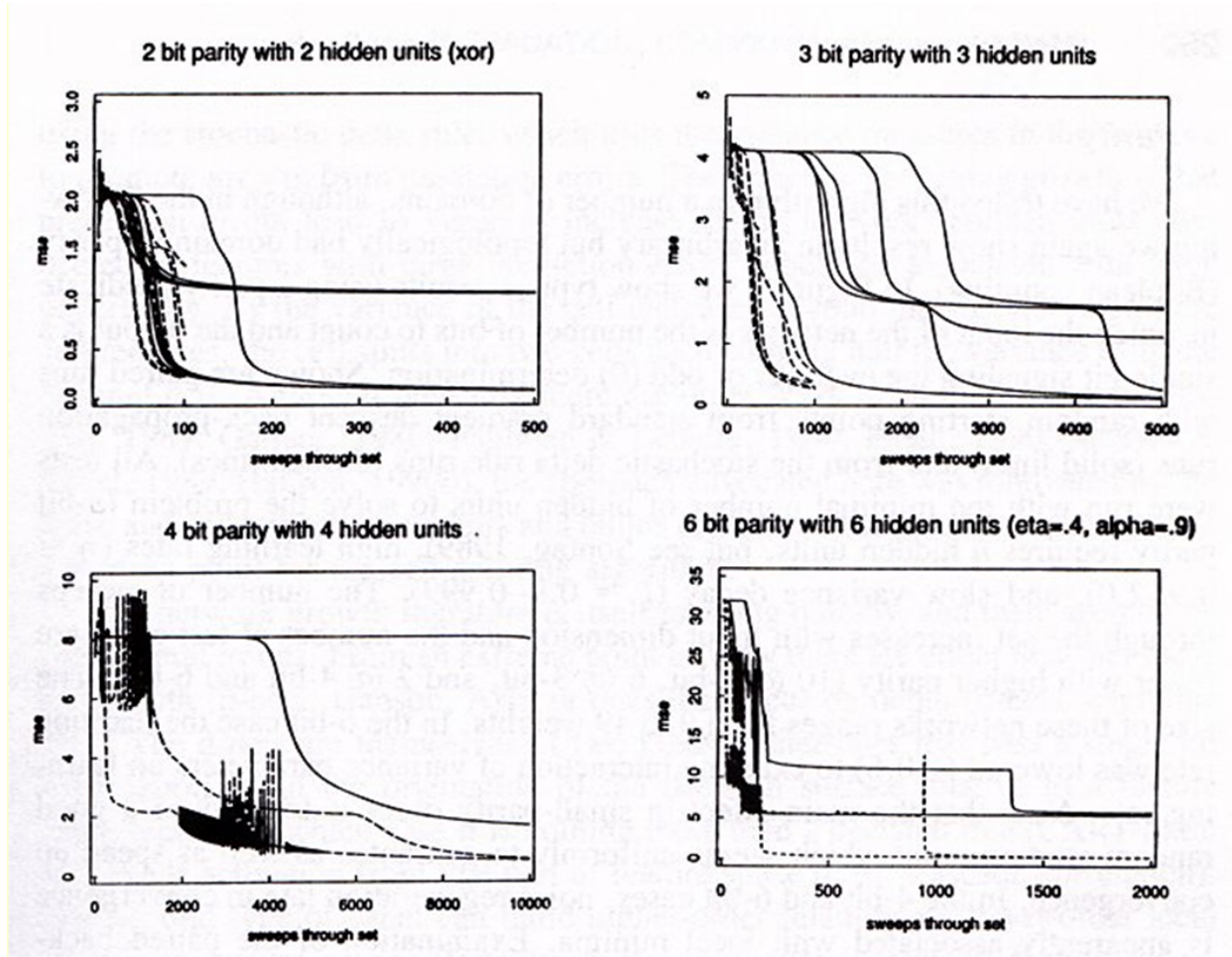
satisfied by, e.g. $\eta(t) \propto \dfrac{1}{t}$ for large $t$ **learning rate schedules,** e.g. $\eta(t) = \dfrac{a}{b+t}$

<span style="color:blue">alternative: retain finite learning rate and fluctuations, average weights over time</span>  83
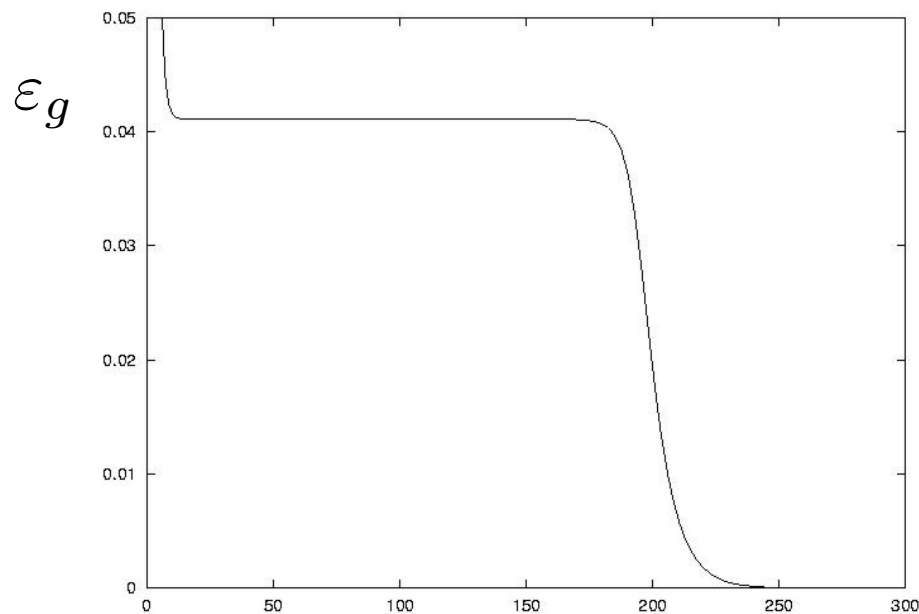
## Plateau states

frequent observation:
training of multilayer networks is delayed by *quasi-stationary plateaus*



(S.J. Hanson, in: Y. Chauvin and D. Rummelhart, *The Handbook of Backpropagation*, 1995)

example: a two-layer network trained from reliable, perfectly realizable data
by on-line gradient descent



$\varepsilon_g$

number of examples   $P/(KN)$

- fast initial decreas of $\varepsilon_g$

- fast asymptotic decrease of $\varepsilon_g \to 0$
  (here: matching complexity)

- plateau state:
  **unspecialized** h.u. with   $\mathbf{w}_k \sim \mathbf{w}_o + noise$
  have all obtained some (the same)
  information about the unknown rule

occurence of plateaus relates to symmetries:

the network output is invariant under **permutations of hidden units**

perfectly symmetric state corresponds to a flat region (saddle) in $E$

successful learning requires **specialization** and can be delayed significantly

math. analysis:  D. Saad and S. Solla (1995), M. Biehl, P.Riegler, C. Wöhler (1996)

**Remarks**

- the extension to several output units $\{\sigma^l\}_{l=1}^{L}$ is non-trivial

  the choice of cost function $E = \dfrac{1}{2} \displaystyle\sum_{l=1}^{L} \left(\sigma^l - \tau^l\right)^2$ seems plausible

  but a generalized metric $\quad E = \dfrac{1}{2} \displaystyle\sum_{k,l} \left(\sigma^k - \tau^k\right) A_{kl} \left(\sigma^l - \tau^l\right)$

  with a suitable $(L \times L)$-matrix $A$ might be more appropriate

- appoximation of a classifier by regression training:

  replace $S^\mu = \pm 1$ by $\sigma^\mu = \pm R \quad |R| < 1,$
  take $S = \text{sign}(\sigma^\mu)$ as classification output after training

  Note:

  – the decrease of $\varepsilon_g$ can be slower than *necesary*

  – for very steep activation functions or $|R| \approx 1$, the gradient becomes
    non-informative as $\nabla E \approx 0$