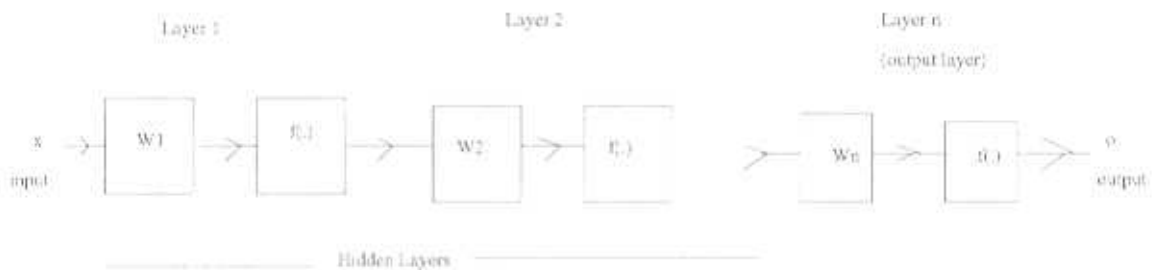


- **Back Propagation Algorithm**

- If we cascade single layer structures we obtain multilayer structures...



- $o = \Gamma(W_n(\Gamma(W_{n-1}(\dots \Gamma(W_2(\Gamma(W_1x))\dots)))$

- Usually we have a training set $\mathcal{S}_{Tr} = \{x^1, x^2, \dots, x^N\}$ and we have a set of desired outputs d^1, \dots, d^N . Note that these are **vectors**.

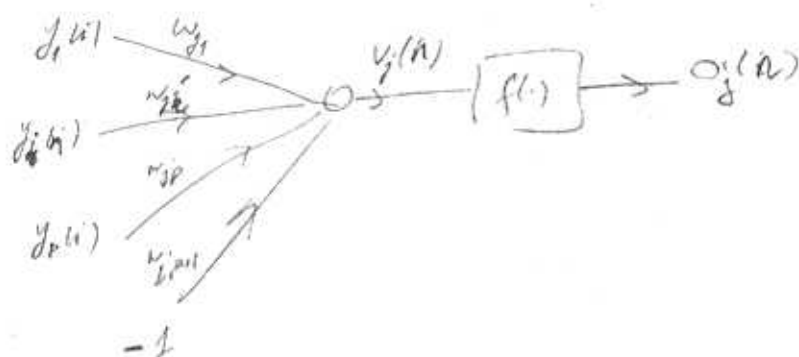
- Assume that at the output layer we have R neurons. Hence, we have $d, o \in \mathbf{R}^R$, i.e. they have R components.

- Suppose that we applied n th training set pattern, we obtain an output $o(n)$, and the corresponding desired output is $d(n) = d^n$. Componentwise we have :

- $$d(n) = \begin{pmatrix} d_1(n) \\ d_2(n) \\ \vdots \\ d_R(n) \end{pmatrix} \quad o(n) = \begin{pmatrix} o_1(n) \\ o_2(n) \\ \vdots \\ o_R(n) \end{pmatrix} \quad n = 1, 2, \dots, N$$

$$\tilde{d}_n = \begin{pmatrix} d_1(n) \\ d_2(n) \\ \vdots \\ d_R(n) \end{pmatrix}, \quad \tilde{o}_n = \begin{pmatrix} o_1(n) \\ o_2(n) \\ \vdots \\ o_R(n) \end{pmatrix}, \quad \tilde{y}_n = \begin{pmatrix} y_1(n) \\ y_2(n) \\ \vdots \\ y_P(n) \\ 1 \end{pmatrix}$$

n : pattern number
 j : output neuron



$$v_j(n) = w_{j1} y_1(n) + \dots + w_{j2} y_2(n) + \dots + w_{jp} y_p(n) + w_{j,p+1}$$

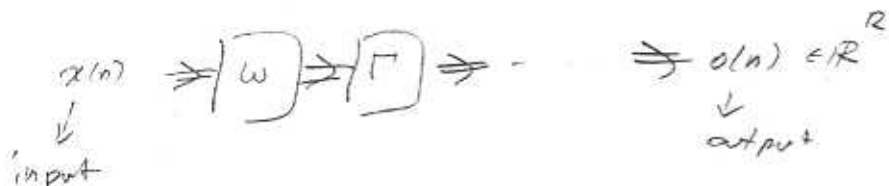
$$o_j(n) = f(v_j(n))$$

$$E(n) = \frac{1}{2} \|\tilde{d}_n - \tilde{o}_n\|^2 = \frac{1}{2} \sum_{i=1}^R (d_i(n) - o_i(n))^2$$

$e_j(n) = d_j(n) - o_j(n)$ → error node at j th output neuron for n th input pattern.

$$\rightarrow E(n) = \frac{1}{2} \sum_{j=1}^R e_j^2(n)$$

$$E_{ave} = \frac{1}{N} \sum_{n=1}^N E(n) = \frac{E(1) + E(2) + \dots + E(N)}{N}$$



update law for any weight:

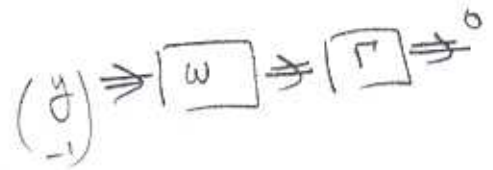
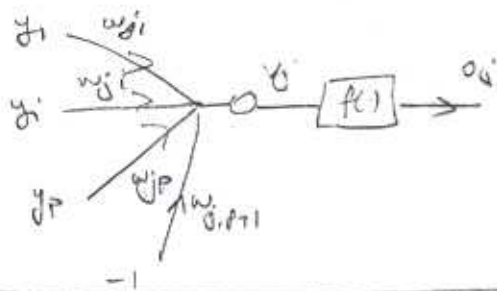
$$w \leftarrow w + \eta \frac{\partial E(n)}{\partial w}$$

Δw

FOR FREQUENTLY ASKED QUESTIONS ON NEURAL NETWORKS:

<ftp://ftp.sas.com/pub/neural/FAQ.html>

→ Neuron j is an output neuron:



$$\frac{\partial E(n)}{\partial w_{ji}} = -f'(v_j(n)) \cdot e_j(n) \cdot y_i$$

$j = 1, \dots, R$
 $i = 1, \dots, P+1$

Local gradient:

$$\delta_j(n) = f'(v_j(n)) \cdot e_j(n) \quad j = 1, 2, \dots, R$$

$$\frac{\partial E(n)}{\partial w_{ji}} = -\delta_j(n) \cdot y_i$$

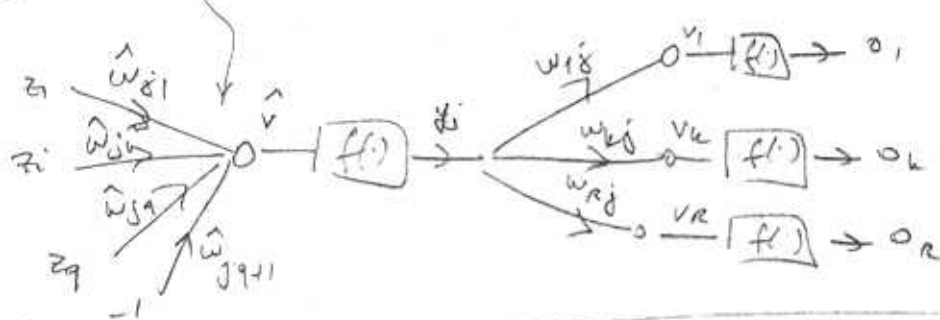
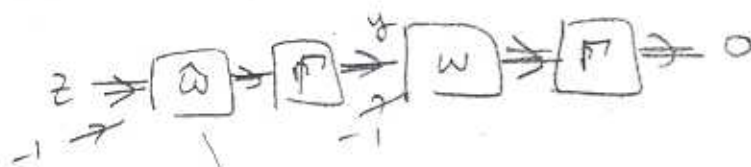
MATRIX FORMULATION:

$$\delta(n) = \begin{pmatrix} \delta_1(n) \\ \vdots \\ \delta_R(n) \end{pmatrix} = \Gamma'(v_j(n)) \cdot e(n), \quad \Gamma'(v) = \begin{bmatrix} f'(v_1) & & 0 \\ & \ddots & \\ 0 & & f'(v_R) \end{bmatrix}$$

$$W = \begin{pmatrix} w_{j1} & \dots & w_{jP} \\ \vdots & & \vdots \\ w_{R1} & \dots & w_{RP} \end{pmatrix} \quad W^e = (w, \text{threshold}) \quad y_e = \begin{pmatrix} y \\ -1 \end{pmatrix}$$

$$W^e \leftarrow W^e + \eta \delta(n) \cdot y_e^T$$

→ Neuron j is hidden layer neuron:



$\frac{\partial E(n)}{\partial z_k} = \frac{\partial E(n)}{\partial v_k} \cdot \frac{\partial v_k}{\partial z_k}$

$$\frac{\partial E(n)}{\partial \hat{w}_{ji}} = - f'(\hat{v}_j(n)) \left(\sum_{k=1}^R f'(v_k(n)) \cdot e_k(n) w_{kj} \right) z_i$$

$j=1, \dots, p$
 $i=1, \dots, q+1$

Local gradient:

$$\hat{\delta}_j = f'(\hat{v}_j(n)) \cdot \sum_{k=1}^R f'(v_k(n)) \cdot e_k(n) w_{kj}$$

$$= f'(\hat{v}_j(n)) \cdot \sum_{k=1}^R \delta_k(n) w_{kj}$$

$j=1, \dots, p$

$$\frac{\partial E(n)}{\partial \hat{w}_{kj}} = - \hat{\delta}_j(n) \cdot z_i$$

MATRIX FORMULATION

$$\hat{\delta}(n) = \begin{pmatrix} \hat{\delta}_1(n) \\ \vdots \\ \hat{\delta}_p(n) \end{pmatrix} = \Gamma'(\hat{v}_j(n)) \cdot W^T \cdot \delta(n)$$

local gradient vector of following layer.

weight matrix of following layer

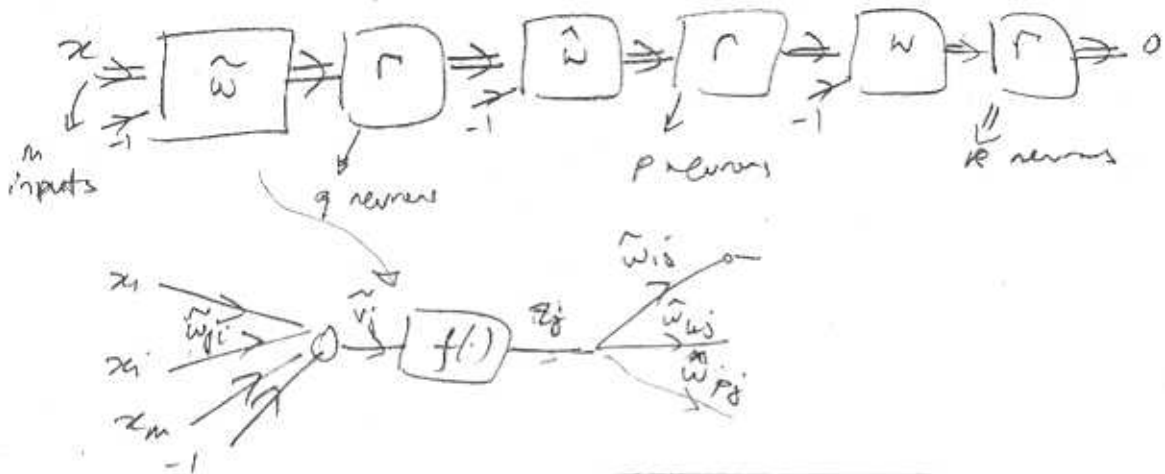
$$\hat{W} = \begin{pmatrix} \hat{w}_{11} & \dots & \hat{w}_{1q} \\ \vdots & & \vdots \\ \hat{w}_{p1} & \dots & \hat{w}_{pq} \end{pmatrix}$$

$$\hat{W}^e = (\hat{W} \text{ (threshold)})$$

$$Z_e = \begin{pmatrix} z \\ -1 \end{pmatrix}$$

$$\hat{W}^e \leftarrow \hat{W}^e + \eta \hat{\delta}(n) Z_e^T$$

How do extend to one more layer?



$$\frac{\partial E(n)}{\partial \tilde{w}_{ji}} = -f'(\tilde{v}_j(n)) \left(\sum_{k=1}^p \hat{\delta}_k(n) \hat{w}_{kj} \right) x_i$$

$j=1, \dots, q$
 $i=1, \dots, m+1$

local gradient of following layer

LOCAL GRADIENT

$$\tilde{\delta}_j = f'(\tilde{v}_j(n)) \sum_{k=1}^p \hat{\delta}_k(n) \hat{w}_{kj} \quad j=1, \dots, q$$

$$\frac{\partial E(n)}{\partial \tilde{w}_{ji}} = -\tilde{\delta}_j(n) z_i$$

MATRIX FORMULATION

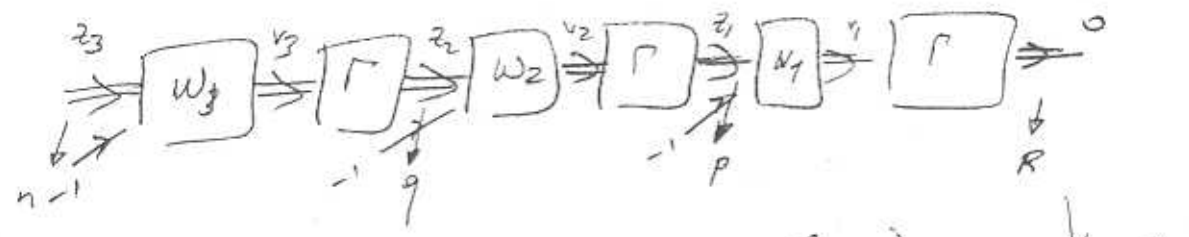
$$\tilde{\delta}(n) = \begin{pmatrix} \tilde{\delta}_1 \\ \vdots \\ \tilde{\delta}_q \end{pmatrix} = \Gamma'(\tilde{V}) \hat{W}^T \hat{\delta}(n)$$

local gradient of following layer

weight matrix of following layer

$$\tilde{W} = \begin{pmatrix} \tilde{w}_{11} & \dots & \tilde{w}_{1m} \\ \vdots & & \vdots \\ \tilde{w}_{q1} & \dots & \tilde{w}_{qm} \end{pmatrix} \quad \tilde{w}_e = \begin{pmatrix} \tilde{w} \\ \text{threshold} \end{pmatrix} \quad x_e = \begin{pmatrix} x \\ -1 \end{pmatrix}$$

$$\tilde{w}_e \leftarrow \tilde{w}_e + \eta \tilde{\delta}(n) \cdot x_e^T$$



$z_3 \in \mathbb{R}^n$
 $z_3^c = \begin{pmatrix} z_3 \\ -1 \end{pmatrix} \in \mathbb{R}^{n+1}$

$z_2 \in \mathbb{R}^q$
 $z_2^c = \begin{pmatrix} z_2 \\ -1 \end{pmatrix} \in \mathbb{R}^{q+1}$

$z_1 \in \mathbb{R}^p$
 $z_1^c = \begin{pmatrix} z_1 \\ -1 \end{pmatrix} \in \mathbb{R}^{p+1}$

$o \in \mathbb{R}^r$
 $d \in \mathbb{R}^r$

- n : # of input components (excluding threshold)
- q : # of neurons in 1st layer.
- p : # of neurons in 2nd layer.
- r : # of neurons in output layer.

$W_1 \in \mathbb{R}^{r \times p}$

$W_1^c = [W_1 \text{ threshold}] \in \mathbb{R}^{r \times (p+1)}$

$W_2 \in \mathbb{R}^{p \times q}$

$W_2^c = [W_2 \text{ threshold}] \in \mathbb{R}^{p \times (q+1)}$

$W_3 \in \mathbb{R}^{q \times n}$

$W_3^c = [W_3 \text{ threshold}] \in \mathbb{R}^{q \times (n+1)}$

present $x(n) = z_3$ at n th step, use $W_3(n), W_2(n), W_1(n)$
 Calculate everything till 0: by using current weights.

- $e(n) = d(n) - o(n)$
- $\delta_1(n) = \Gamma'(v_1(n)) e(n) \in \mathbb{R}^r = \begin{bmatrix} f'(v_{11}(n)) & & \\ & f'(v_{12}(n)) & \\ & & f'(v_{1r}(n)) \end{bmatrix} \begin{bmatrix} -e_1(n) \\ \vdots \\ e_r(n) \end{bmatrix}$
- $W_1^c(n+1) \leftarrow W_1^c(n) + \gamma \delta_1(n) (z_1^c)^T$
- $\delta_2(n) = \Gamma'(v_2(n)) W_2^T(n) \delta_1(n)$
 $p \times p \quad p \times r \quad r \times 1 \quad = p \times 1$
- $W_2^c(n+1) \leftarrow W_2^c(n) + \gamma \delta_2(n) (z_2^c)^T$
- $\delta_3(n) = \Gamma'(v_3(n)) W_3^T(n) \delta_2(n)$
 $q \times q \quad q \times p \quad p \times 1 \quad = q \times 1$
- $W_3^c(n+1) \leftarrow W_3^c(n) + \gamma \delta_3(n) (z_3^c)^T$

- pick up another pattern as x . continue these steps.

TRAINING ERRORS:

Training set: $\{z_1, \dots, z_N\}$, $o = \begin{pmatrix} o_1 \\ \vdots \\ o_R \end{pmatrix}$ $d = \begin{pmatrix} d_1 \\ \vdots \\ d_R \end{pmatrix}$

error at n^{th} pattern $\Rightarrow E(n) = \frac{1}{2} \|o(n) - d(n)\|^2 = \frac{1}{2} \sum_{i=1}^R (o_i(n) - d_i(n))^2$

cumulative error $\Rightarrow E_{\text{cum}} = \sum_{n=1}^N E(n)$

average error $\Rightarrow E_{\text{ave}} = \frac{1}{N} E_{\text{cum}}$

if # of output neurons (R) is large \Rightarrow Root-Mean-Square error:

$$E_{\text{rms}} = \frac{1}{NR} \sqrt{2 E_{\text{cum}}}$$

\Rightarrow maybe useful for comparing two NN architectures with different N and R .

In Binary case ($o_i = \pm 1$) ($d_i = \pm 1$)

$$\Rightarrow e_i = \begin{cases} 1 & \text{if output } i \text{ is incorrect (bit error)} \\ 0 & \text{if output } i \text{ is correct} \end{cases}$$

$$E(n) = \sum_{i=1}^R e_i \quad (\text{bit error at output})$$

$$E_{\text{cum}} = \sum_{n=1}^N E(n) \quad (\text{cumulative bit error})$$

$$E_{\text{ave}} = \frac{1}{N} E_{\text{cum}} \quad (\text{average error})$$

$$E_{\text{rms}} = \frac{1}{NR} \sqrt{2 E_{\text{cum}}} \quad (\text{rms error})$$

$$E_{\text{dec}} = \frac{E_{\text{cum}}}{NR} \quad (\text{decision error})$$

$$E(w+\Delta w) = E(w) + \frac{\partial E}{\partial w} \cdot \Delta w + \text{h.o.t}$$

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

- if η is low method works, but error descent may be slow.
- if η is high error descent may be high but method may not converge.

Momentum method: A method to (possibly) accelerate the descent without instability

$$\Delta w(n) = \alpha \Delta w(n-1) - \eta \frac{\partial E(n)}{\partial w}$$

α is momentum constant (forgetting factor)

- Acts like a filter:

\Rightarrow z transform:

$$\Delta w(z) = \alpha z^{-1} \Delta w(z) - \eta Z(\nabla E)$$

\leftarrow z transform of error gradient:

$$\Delta w(z) = \frac{1}{1-\alpha z^{-1}} (-\eta Z(\nabla E))$$

\hookrightarrow acts like a low pass filter.

$$-\eta \nabla E \Rightarrow \left[\frac{1}{1-\alpha z^{-1}} \right] \Rightarrow \Delta w$$

For stability: $|\alpha| < 1$!

$$\Delta w(1) = -\eta \frac{\partial E(1)}{\partial w}$$

$$\Delta w(2) = \alpha \Delta w(1) - \eta \frac{\partial E(2)}{\partial w} = -\alpha \eta \frac{\partial E(1)}{\partial w} - \eta \frac{\partial E(2)}{\partial w}$$


$$\Delta w(3) = \alpha \Delta w(2) - \eta \frac{\partial E(3)}{\partial w} = -\alpha^2 \eta \frac{\partial E(1)}{\partial w} - \alpha \eta \frac{\partial E(2)}{\partial w} - \eta \frac{\partial E(3)}{\partial w}$$

$$\Delta w(n) = -\eta \sum_{k=1}^n \alpha^{n-k} \frac{\partial E(k)}{\partial w}$$

⇒ weighted average of error gradients.

- $0 < \alpha < 1$ for stability and accelerating the descent.

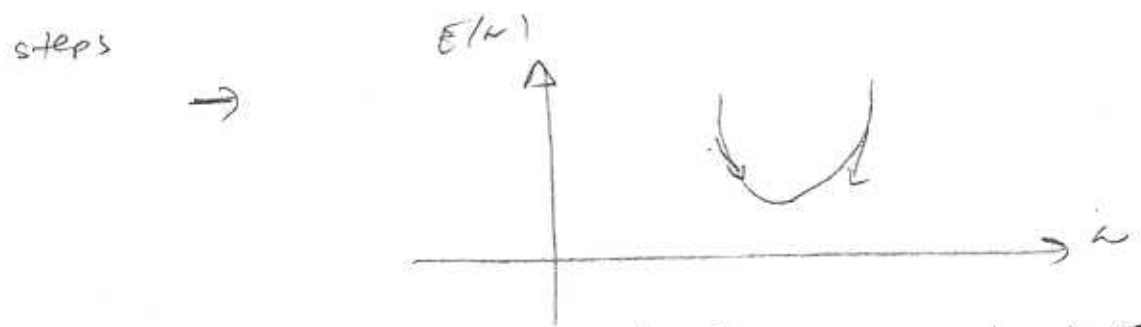
- when $\frac{\partial E(k)}{\partial w}$ has the same algebraic sign over consecutive steps ⇒



⇒ weight increment is increased

⇒ accelerating the descent.

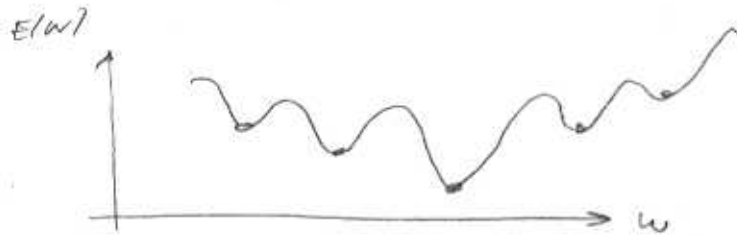
- when $\frac{\partial E(k)}{\partial w}$ has alternating signs over consecutive steps ⇒



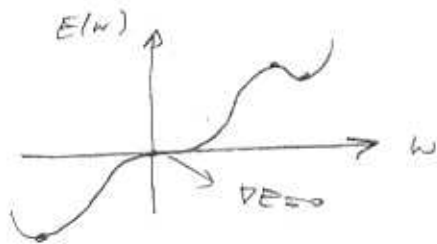
→ danger of instability (oscillations) ⇒ weight change is decreased.

- PROBLEM OF LOCAL MINIMUM.

- Convergence is never guaranteed.
- If η is sufficiently small and if $w(0)$ is sufficiently close to a local minimum, then algorithm converges to ~~the~~ local minimum.



- No guarantee to converge to a global minimum
- "sufficiently small" "sufficiently close"
- Saddle point: extremum ($\nabla E=0$) but neither minimum nor maximum.



- If in local minimum, impossible to get out by using deterministic approach.
- Start with different initial conditions.
- Introduce some randomness to the algorithm.
- Shuffle the training patterns.
- Add zero mean white noise to the training patterns.

- Weight Initialization

- Prevent the neurons from saturation.



$f'(v) = (1 - o^2) \Rightarrow$ if $o \approx \pm 1$, $f'(v)$ is small \Rightarrow small update.

- Choose the weights randomly in a small range around the origin (try to keep the output in linear range)

- $(-\frac{2.4}{F_j}, \frac{2.4}{F_j})$ a typical range for the initial weights of g^{th} neuron.

- Wrong choice of initial conditions may cause "premature saturation"

- Correct saturation, incorrect saturation.

- If a neuron is in incorrect saturation, very difficult to get out of it. \Rightarrow saddle behaviour.

- STOPPING CRITERIA.

typical:

if Error $< \epsilon$ stop

else start a new epoch.

- Error: $\rightarrow E_{cum}, E_{ave}, E_{rms}, E_{dc} \dots$

- check $\forall E$. stop when $\| \nabla E \| < \epsilon$

- Near minimum, error is stationary.

Let q be epoch index.

if $\frac{|E_{cum}(q) - E_{cum}(q-1)|}{E_{cum}(q-1)} < \epsilon$ stop