

Neural networks

Virgil Pavlu October 1, 2008

1 The perceptron

Lets suppose we are (as with regression regression) with $(\mathbf{x}_i, y_i); i = 1, \dots, m$ the data points and labels. This is a classification problem with two classes $y \in \{-1, 1\}$

Like with regression we are looking for a linear predictor (classifier)

$$h_{\mathbf{w}}(\mathbf{x}) = \mathbf{x}\mathbf{w} = \sum_{d=0}^D x^d w^d$$

(we added the $x^0 = 1$ component so we can get the free term w^0) such that $h_{\mathbf{w}}(\mathbf{x}) \geq 0$ when $y = 1$ and $h_{\mathbf{w}}(\mathbf{x}) \leq 0$ when $y = -1$.

On $y = -1$ data points: given that all \mathbf{x} and y are numerical, we will make the following transformation: when $y = -1$, we will reverse the sign of the input; that is replace \mathbf{x} with $-\mathbf{x}$ and $y = -y$. Then the condition $h_{\mathbf{w}}(\mathbf{x}) \leq 0$ becomes $h_{\mathbf{w}}(\mathbf{x}) \geq 0$ for all data points.

The perceptron objective function is a combination of the number of miss-classification points and how bad the miss-classification is

$$J(\mathbf{w}) = \sum_{\mathbf{x} \in M} -h_{\mathbf{w}}(\mathbf{x}) = \sum_{\mathbf{x} \in M} -\mathbf{x}\mathbf{w}$$

where M is the set of miss-classified data points. Note that each term of the sum is positive, since miss-classified implies $\mathbf{w}\mathbf{x} < 0$. Using gradient descent, we first differentiate J

$\mathbf{x}\mathbf{w}$

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \sum_{\mathbf{x} \in M} -\mathbf{x}^T$$

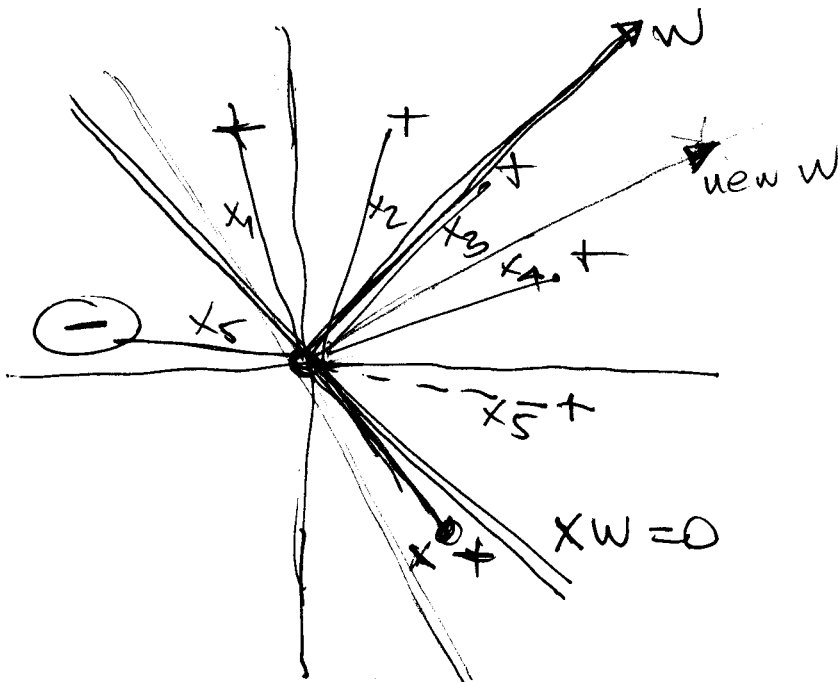
$$\frac{\partial J}{\partial w_i} = \sum_{\mathbf{x} \in M} -x_i$$

then we write down the gradient descent update rule

$$\mathbf{w} := \mathbf{w} + \lambda \sum_{\mathbf{x} \in M} \mathbf{x}^T$$

(λ is the learning rate). The batch version looks like

1. init \mathbf{w}
2. LOOP
3. get $M =$ set of missclassified data points
4. $\mathbf{w} = \mathbf{w} + \lambda \sum_{\mathbf{x} \in M} \mathbf{x}^T$
5. UNTIL $|\lambda \sum_{\mathbf{x} \in M} \mathbf{x}| < \epsilon$



x on the right side $xw > 0$
 x on the wrong side $xw < 0$

$wx < 0 \rightarrow$ move x toward x

$$w \Rightarrow w + x$$

Assume the instances are linearly separable. Then we can modify the algorithm

1. init w
2. LOOP
3. get $M =$ set of missclassified data points
4. for each $x \in M$ do $w = w + \lambda x^T$
5. UNTIL M is empty

Proof of perceptron convergence Assuming data is linearly separable, or there is a solution \bar{w} such that $x\bar{w} > 0$ for all x .

Lets call w_k the w obtained at the k -th iteration (update). Fix an $\alpha > 0$. Then

$$w_{k+1} - \alpha \bar{w} = (w_k - \alpha \bar{w}) + x_k^T$$

where x_k is the datapoint that updated w at iteration k . Then

$$\|w_{k+1} - \alpha \bar{w}\|^2 = \|w_k - \alpha \bar{w}\|^2 + 2x_k(w_k - \alpha \bar{w}) + \|x_k\|^2 \leq \|w_k - \alpha \bar{w}\|^2 - 2x_k \alpha \bar{w} + \|x_k\|^2$$

Since $x_k \bar{w} > 0$ all we need is an α sufficiently large to show that this update process cannot go on forever. When it stops, all datapoints must be classified correctly.

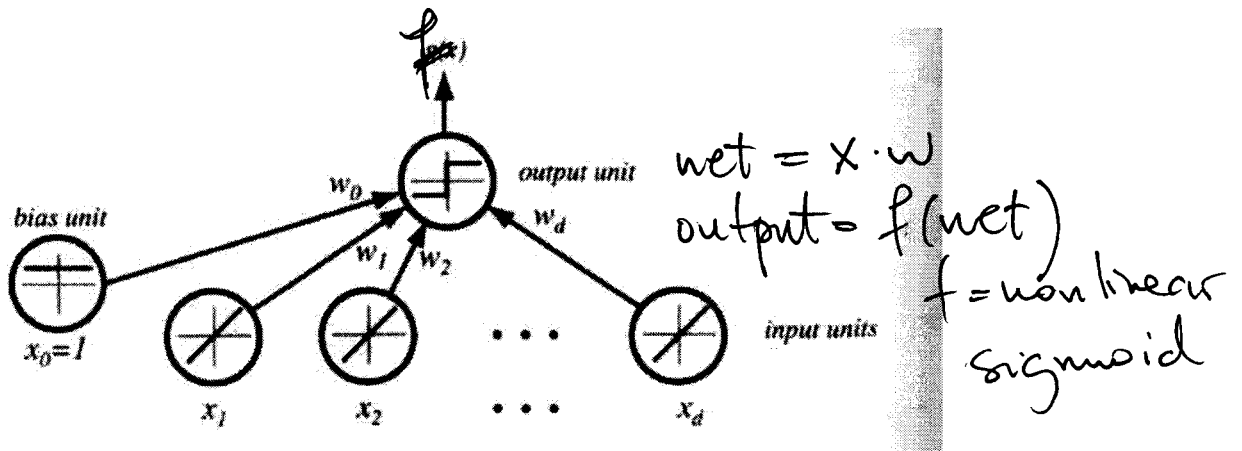


Figure 1: bias unit

2 Multilayer perceptrons

Also called *feed-forward networks*.

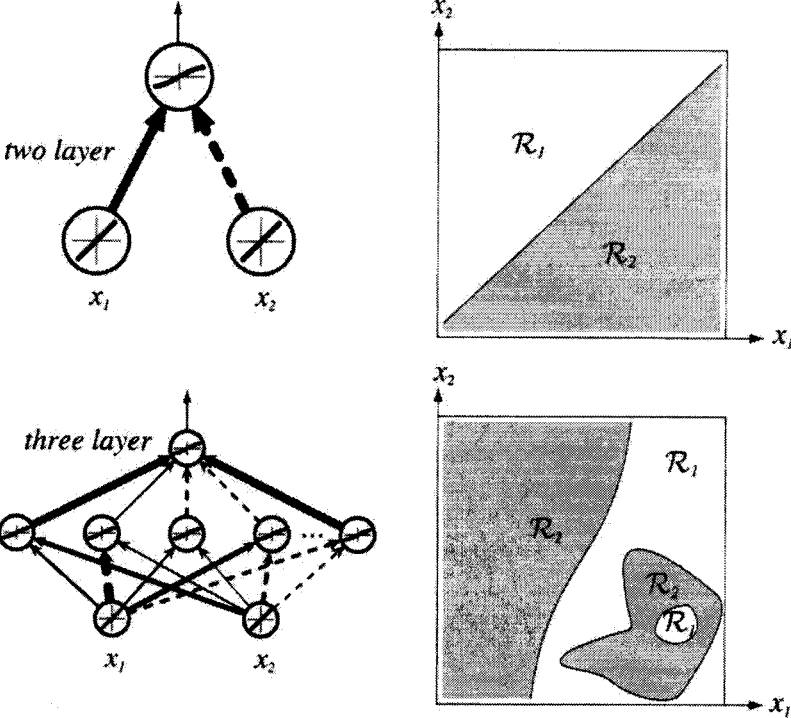


Figure 2: multilayer perceptron

2.1 More than linear functions, example: XOR

- activation functions, XOR example

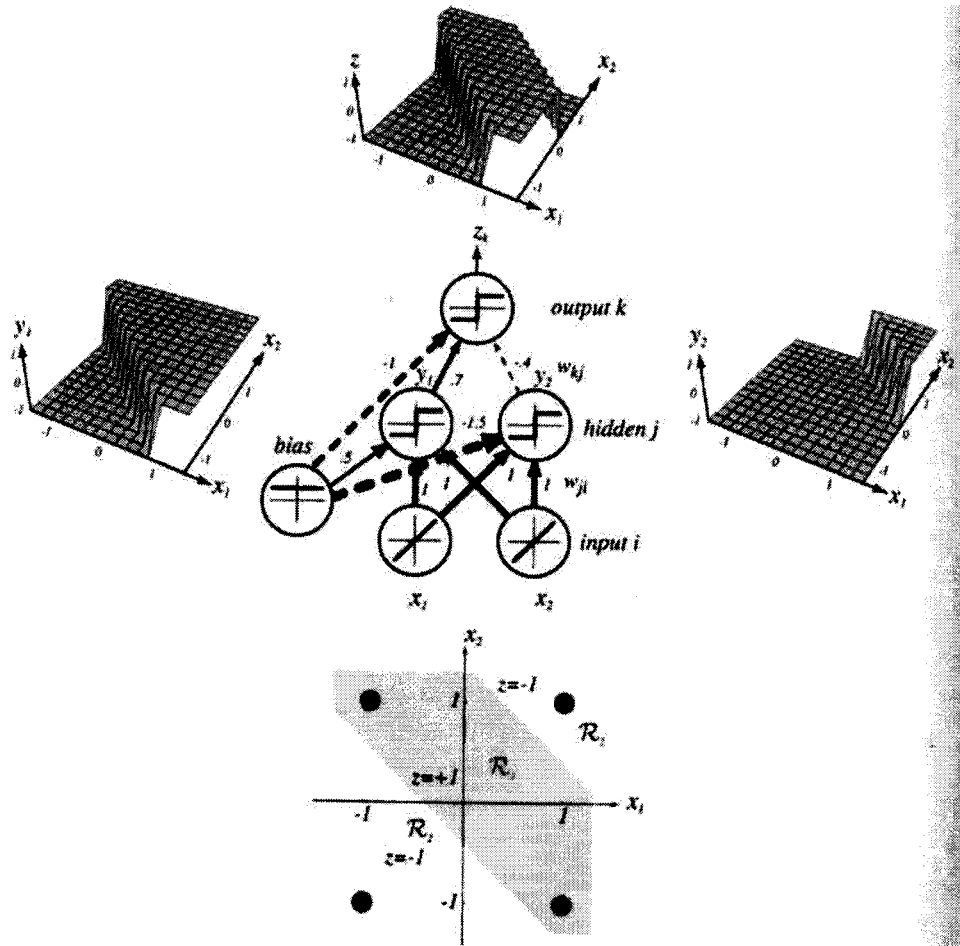


Figure 3: XOR NNNet

$$Y_1 = x_1 \text{ OR } x_2 \quad Y_2 = x_1 \text{ AND } x_2$$

$$Z = Y_1 \text{ AND NOT } Y_2$$

2.2 Construction and structure of NNets

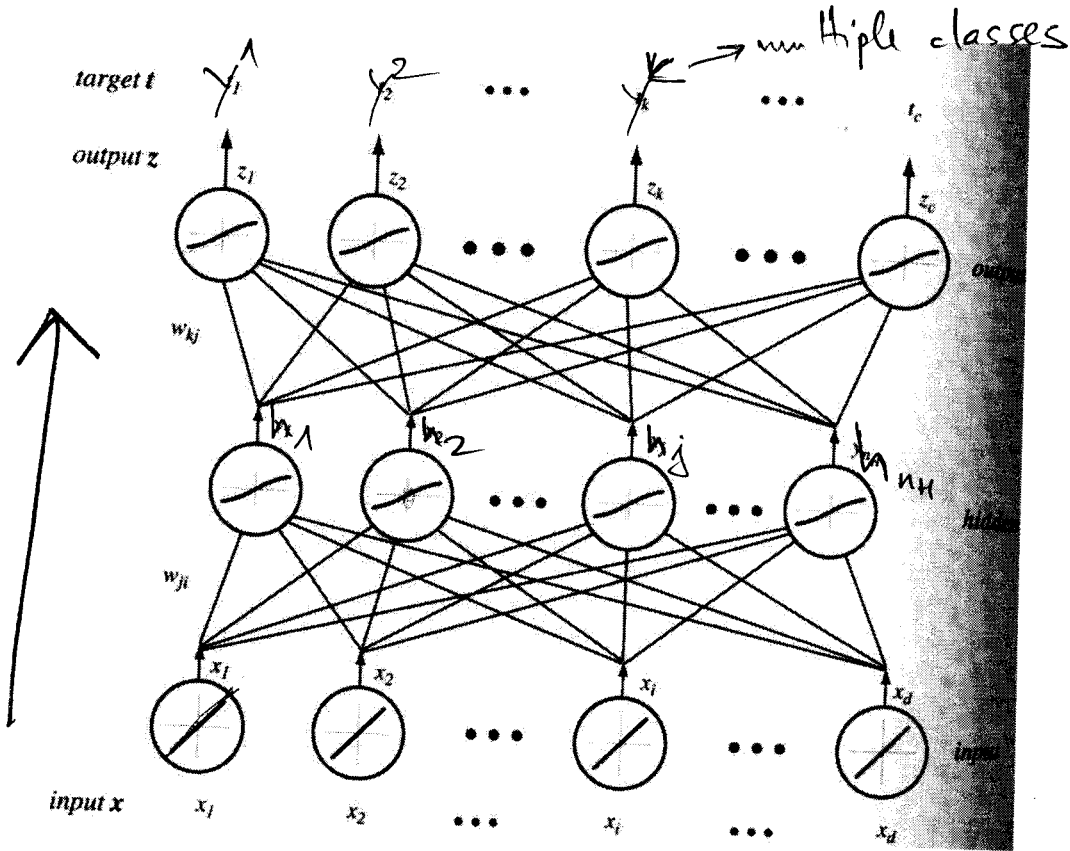


Figure 4: NNet fully connected

The output units can be written as

$$g_k(\mathbf{x}) = z_k = f \left(\sum_j w_{kj} f \left(\sum_i w_{ij} x^i + w_{j0} \right) + w_{k0} \right) = F(F(\mathbf{x}\mathbf{w}_j)\mathbf{w}_k)$$

2.3 Kolmogorov theorem, expressive power of NNet

Any function g can be written

$$g(\mathbf{x}) = \sum_j \Xi_j \left(\sum_d \Psi_{dj}(x^d) \right)$$

but there is no practical way to use this theorem in practice. Usually Ξ and Ψ are very complex and not smooth.

def net = linear combination = $\sum_i w_{ij} x^i + w_{j0}$ (first layer)

= $\sum_j w_{kj} f(\dots) + w_{k0}$ (second layer)

def output = sigmoid(net) = $f(\text{net})$

3 Training, Error backpropagation

- error

for one data point x

$$J(w) = \frac{1}{2} \sum_k (y^k - z^k)^2 = \frac{1}{2} \|y - z\|^2$$

$$\nabla_w J = -\lambda \frac{\partial J}{\partial w_{pq}}$$

- propagation to last set of weights (close to output)

$$\frac{\partial J}{\partial w_{kj}} = \frac{\partial J}{\partial \text{net}_k} \frac{\partial \text{net}_k}{\partial w_{kj}} = \delta_k \frac{\partial \text{net}_k}{\partial w_{kj}}$$

sensitivity of unit k

$$\delta_k = -\frac{\partial J}{\partial \text{net}_k} = -\frac{\partial J}{\partial z_k} \frac{\partial z_k}{\partial \text{net}_k} = (y^k - z^k) f'(\text{net}_k)$$

$$\frac{\partial \text{net}_k}{\partial w_{kj}} = h_j$$

$$w_{kj} = w_{kj} + \lambda (y^k - z^k) f'(\text{net}_k) h_j$$

$$\lambda \delta_k h_j$$

- propagation to first set of weights (close to input)

$$\frac{\partial J}{\partial w_{ji}} = \frac{\partial J}{\partial h_j} \cdot \frac{\partial h_j}{\partial \text{net}_j} \cdot \frac{\partial \text{net}_j}{\partial w_{ji}}$$

$$\begin{aligned} \frac{\partial J}{\partial h_j} &= \frac{\partial \left[\frac{1}{2} \sum_k (y^k - z^k)^2 \right]}{\partial h^j} = - \sum_k (y^k - z^k) \frac{\partial z^k}{\partial \text{net}_k} \cdot \frac{\partial \text{net}_k}{\partial h_j} = \\ &= \sum_k (y^k - z^k) f'(\text{net}_k) w_{kj} \end{aligned}$$

$$\frac{\partial h_j}{\partial \text{net}_j} = f'(\text{net}_j) \quad \left| \quad \frac{\partial \text{net}_j}{\partial w_{ji}} = x_i$$

- ~~step 1~~ batch

$$w_{ji} = w_{ji} - \lambda \left[\sum_k (y^k - z^k) f'(\text{net}_k) w_{kj} \right] \cdot f'(\text{net}_j) \cdot x_i$$

δ_j = sensitivity at unit j

$$\delta_j = f'(\text{net}_j) \sum_k w_{kj} \delta_k$$

\downarrow
 sensitivity at k

STOCHASTIC BACKP

Select x_t (randomly chosen)

$$w_{ji} = w_{ji} + \lambda \delta_j x_i$$

$$w_{kj} = w_{kj} + \lambda \delta_k h_j$$

$$\text{until } |\nabla_w J| < \epsilon$$

BATCH $J = \sum_{t=1}^m J(w, x_t)$

for each iteration

for each x_t

$$\Delta w_{ji} = \Delta w_{ji} + \lambda \delta_j x_i$$

$$\Delta w_{kj} = \Delta w_{kj} + \lambda \delta_k h_j$$

$$w_{ji} = w_{ji} + \Delta w_{ji}$$

$$w_{kj} = w_{kj} + \Delta w_{kj}$$

$$\text{until } \|\nabla J(w)\| < \epsilon$$

4 How to improve backpropagation

4.1 Activation function

- continuous, differentiable (smoothness)
 - nonlinear
 - saturation
 - monotonicity

4.2 Scale input mass (grams) vs length (meters) STANDARDIZE (0,1)

4.3 Target values for classification use 1, -1 etc. Sync with f

4.4 Noise

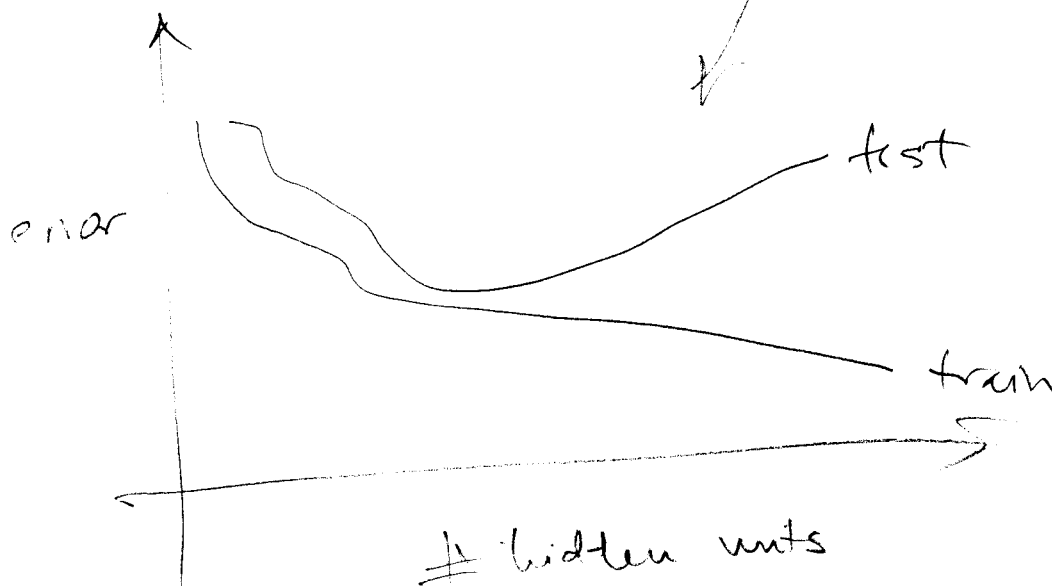
4.5 Number of hidden units = complexity, overfit, test error, accuracy, REGULARIZATION

4.6 Weights initialization → choose weights randomly from a single dist (avoid learn one class first)

4.7 Learning rates

5 Network size and structure. Regularization

6 Jacobian and Hessian



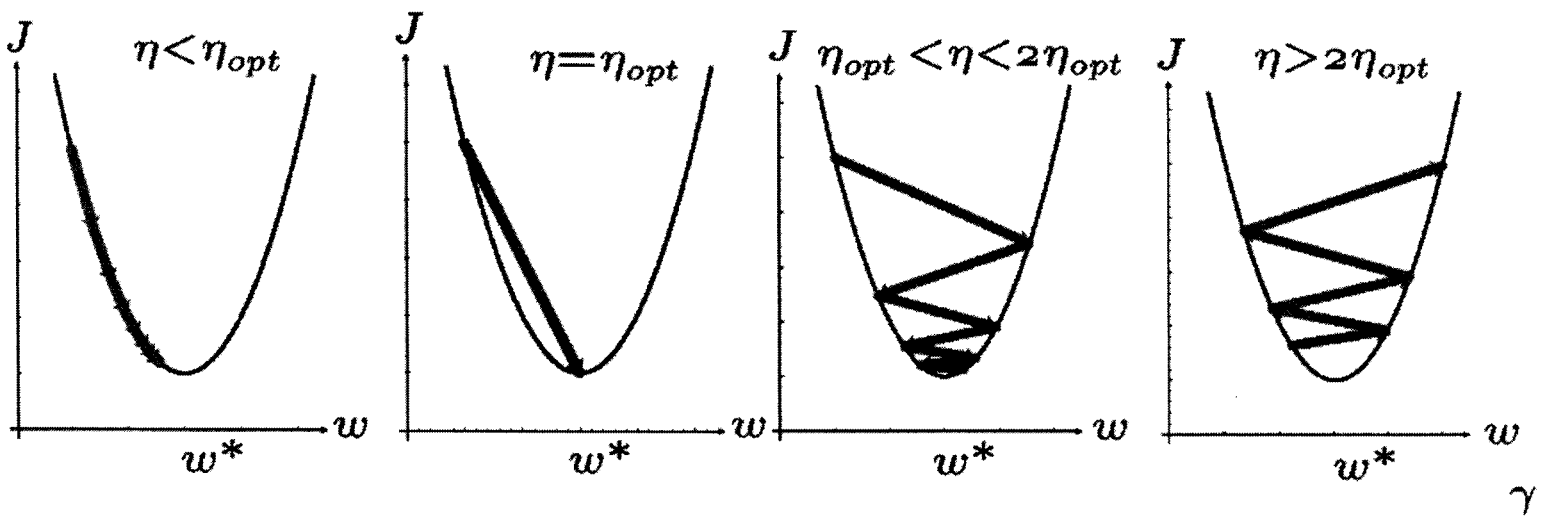


Figure 6.18: Gradient descent in a one-dimensional quadratic criterion with different learning rates. If $\eta < \eta_{opt}$, convergence is assured, but training can be needlessly slow. If $\eta = \eta_{opt}$, a single learning step suffices to find the error minimum. If $\eta_{opt} < \eta < 2\eta_{opt}$, the system will oscillate but nevertheless converge, but training is needlessly slow. If $\eta > 2\eta_{opt}$, the system diverges.