

Purpose:

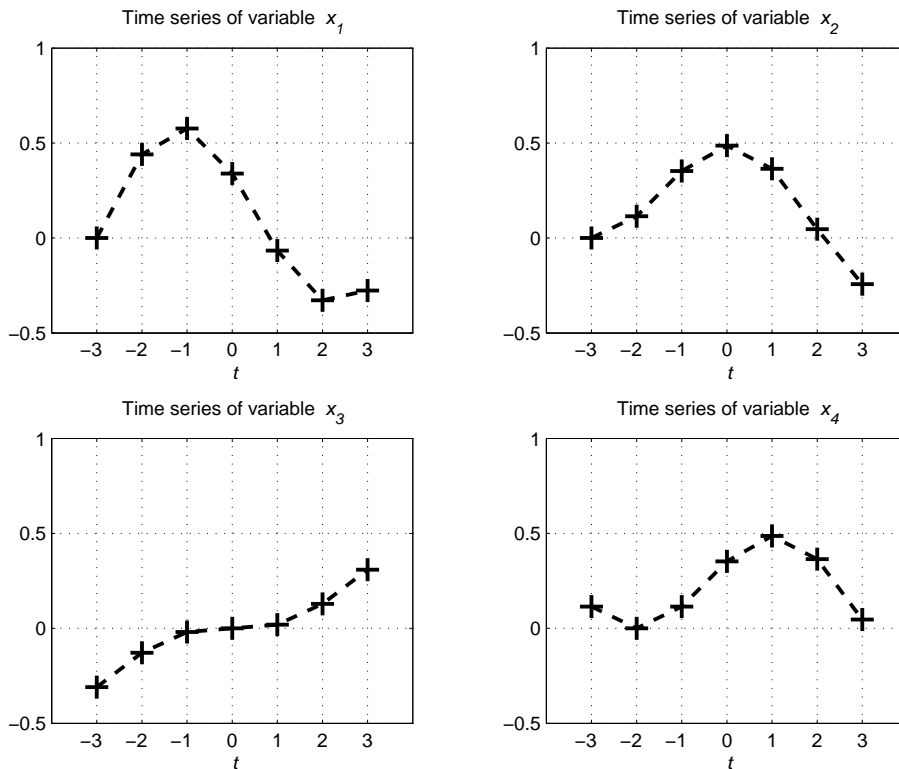
How to train an MLP neural network in MATLAB environment!

that is

For good computations,
we need good formulae
for good algorithms;
and good visualization
for good illustration
and proper testing
of good methods
and succesfull applications!

Basic Applications for MLP (I):

- time-series prediction



Derivation of learning data:

$$\begin{cases} \mathbf{x}_1 = [x_1(-1) & x_2(-1) & x_3(-1) & x_4(-1) & x_4(-2)] \\ \mathbf{x}_2 = [x_1(-2) & x_2(-2) & x_3(-2) & x_4(-2) & x_4(-3)] \\ \vdots \end{cases}, \quad \begin{cases} \mathbf{y}_1 = [x_4(0)] \\ \mathbf{y}_2 = [x_4(-1)] \\ \vdots \end{cases}$$

i.e., in general:

$$\begin{cases} \mathbf{x}_t = [x_1(-t) & x_2(-t) & x_3(-t) & x_4(-t) & x_4(-(t+1))] \\ \mathbf{y}_t = [x_4(-t+1)] \end{cases}$$

In Theory: According to *Takens' theorem* (F. TAKENS, *Detecting strange attractors in turbulence*, *Dynamical Systems and Turbulence*, 898: 336–381, 1981) there exists, for deterministic (= functional) systems, a diffeomorphism between a sufficiently large time window and underlying state of the dynamical system:

$$x(t) = g(x(t-1), x(t-2), \dots, x(t-T)).$$

Basic Applications for MLP (II):

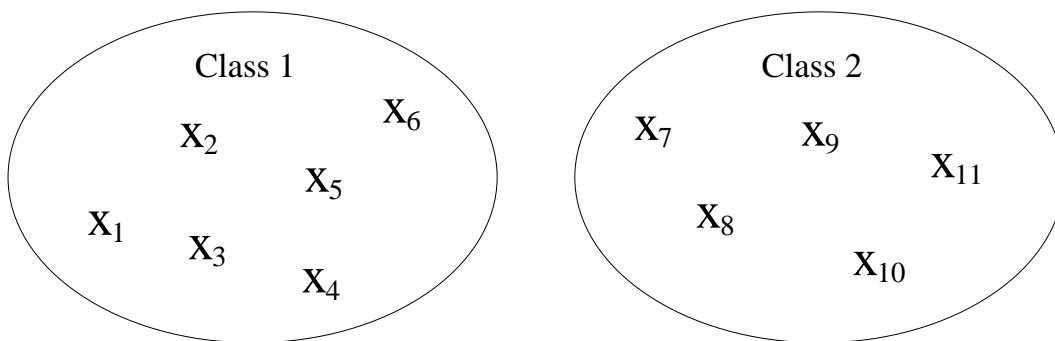
- time-series prediction (cont.)

introduction of higher-order computing (*sigma-pi*) units:

$$\mathbf{x}_t = [x1(-t) \quad x1(-t)^2 \quad x2(-t) \quad x2(-t)^2 \quad x1(-t) * x2(-t) \quad \dots] \quad t = 1, 2, \dots$$

– basic problem: size of learning problem increased significantly!

- classification



Output coding: for k th class

$$\mathbf{y}_i^T = \mathbf{e}_k^T = [0 \dots 1 \dots 0]^T$$

i.e., for the above example

$$\mathbf{y}_i = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad i = 1, \dots, 6; \quad \mathbf{y}_i = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad i = 7, \dots, 11.$$

Actual classifier after training the MLP:

$$\mathbf{x} \in c_k, \quad \text{where } k = \arg \max_j o_j,$$

i.e., for the above example $\mathbf{x} \in c_1$ if $y_1 > y_2$.

More cautious approach: for suitable $0 < \theta < 1$

- 'Class 1', if $o_1 > o_2 + \theta$,
- 'Class 2', if $o_2 > o_1 + \theta$,
- 'Unclassified', otherwise.

Basic Applications for MLP (III):

- data compression (i.e., nonlinear PCA)
 - We try to construct *coder* \mathcal{K} and *decoder* \mathcal{D} according to basic principle

$$\mathbf{x}_i \sim \mathcal{D}(\mathcal{K}(\mathbf{x}_i))$$

for all data vectors \mathbf{x}_i , $i = 1, \dots, N$.

- Multilayered Perceptron approach:

$$\mathbf{x}_i \sim \mathbf{W}^4 \hat{\mathbf{F}}^2(\mathbf{W}^3 \widehat{\mathbf{W}}^2 \hat{\mathbf{F}}^1(\mathbf{W}^1 \hat{\mathbf{x}})),$$

where $\text{size}(\mathbf{W}^2, 1) \ll \text{size}(\mathbf{X}, 2)$.

- yields very large learning problem
 - quality of coding&decoding can be evaluated from the residual of cost function
→ a way to control the choice of $\text{size}(\mathbf{W}^2, 1)$
- other applications (! = exists, ? = worth testing)
 - autopilots for moving vehicles(!)
 - functional algorithms(?): *sort, schedule* etc.
 - automatic differentiation of an optimization problem(!)
 - data inputation(?!)
 - ... (?)

Improving the performance of MLP:

Preprocessing: rescaling, PCA, normalization,...

Restricting the generality:

basic principle: try to keep the unknown weights in the neighborhood of zero, i.e. around the linear region of the activation functions

realization: penalization of weight values

$$\mathcal{J}(\{\mathbf{W}^l\}) = \frac{1}{2N} \sum_{i=1}^N \|\mathcal{N}(\{\mathbf{W}^l\})(\mathbf{x}_i) - \mathbf{y}_i\|^2 + \frac{\beta}{2} \sum_{l=1}^L \sum_{(i,j) \in I_l} |\mathbf{W}_{i,j}^l|^2,$$

where I_l is defined as

$$I_l = \begin{cases} \{(i, j) : 1 \leq i \leq n_l, 0 \leq j \leq n_{l-1}\}, & l < L, \\ \{(i, j) : 1 \leq i \leq n_l, 1 \leq j \leq n_{l-1}\}, & l = L. \end{cases}$$

- easily added into sensitivity analysis of learning problem
- choice of new free (hyper)parameter $\beta \geq 0$ creates new problem, but already a decent value can be helpful

Combination of MLP and RBFN:

- RBFN (*Radial-Basis Function Network*) approximates (uniformly) input-output mapping using local basis functions, most commonly *Gaussian kernels*:

$$\mathbf{o} = \mathbf{W} \left[\exp\left(-\frac{\|\mathbf{x} - \mu_1\|^2}{\delta_1^2}\right) \dots \exp\left(-\frac{\|\mathbf{x} - \mu_i\|^2}{\delta_i^2}\right) \dots \exp\left(-\frac{\|\mathbf{x} - \mu_m\|^2}{\delta_m^2}\right) \right]^T$$

assuming that centroids $\{\mu_i\}_{i=1}^m$ and standard deviations $\{\delta_i\}_{i=1}^m$ are given.

- SQUARE-MLP (*square unit augmented, radially extended, multilayer perceptron*, FLAKE) tries to combine good properties of MLP (global approximation) and RBFN (local representation) simply by (cf. *sigma-pi* units):

$$\begin{aligned} \mathbf{x}_i &= [(\mathbf{x}_i)_1 \quad (\mathbf{x}_i)_1^2 \quad \dots \quad (\mathbf{x}_i)_k \quad (\mathbf{x}_i)_k^2 \quad \dots \quad (\mathbf{x}_i)_{n_0} \quad (\mathbf{x}_i)_{n_0}^2] \quad i = 1, \dots, N, \\ &= [(\mathbf{x}_i)_1 \quad \dots \quad (\mathbf{x}_i)_k \quad \dots \quad (\mathbf{x}_i)_{n_0} \quad (\mathbf{x}_i)_1^2 \quad \dots \quad (\mathbf{x}_i)_k^2 \quad \dots \quad (\mathbf{x}_i)_{n_0}^2]. \end{aligned}$$