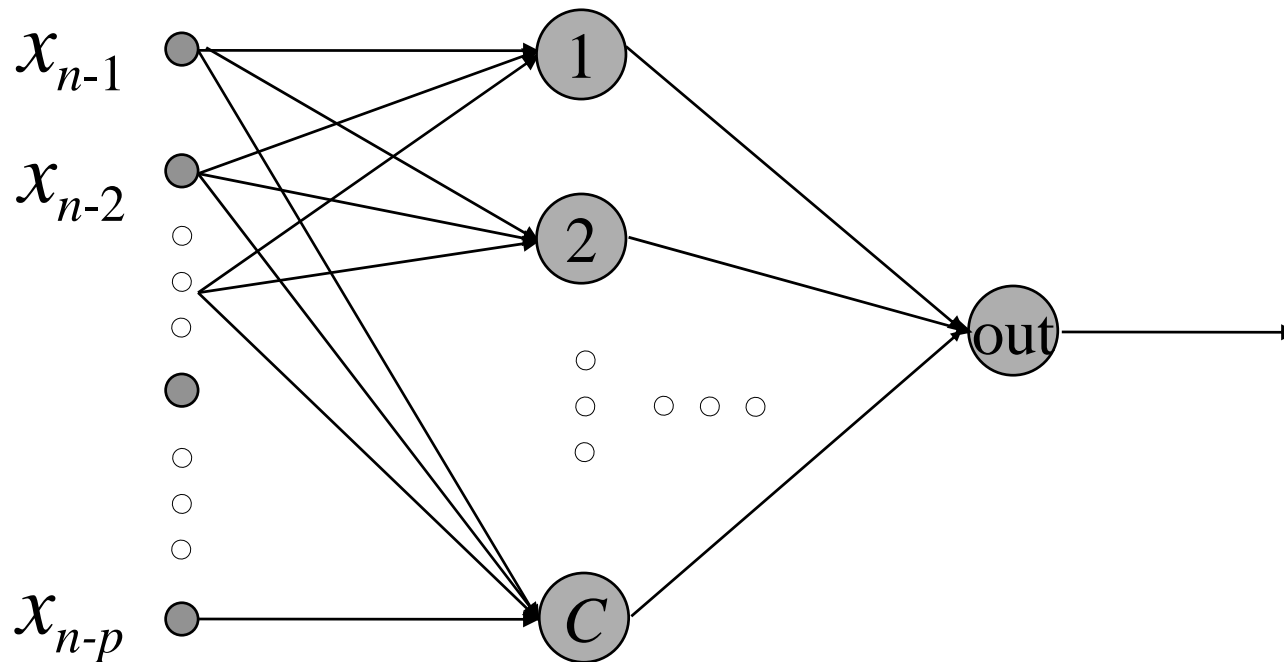


- Multi-layer perceptrons (MLP) constitute the most famous and most employed type of neural networks [1].
- At first, they were introduced in the context of classification, and it was soon recognized that they had the property of universal approximation capability.
- The introduction of the back-propagation algorithm made them flexible enough for many applications.

- Their basic structure for prediction is:



- The output of the i th neuron in the intermediate (hidden) layer is:

$$y_i = f \left(\theta_j + \sum_{j=1}^p w_{ij} x_{n-j} \right)$$

with $f(\cdot)$ a sigmoidal function such as $f(u) = \tanh(au)$.

- The final output is a weighted sum of the outputs in the hidden layer.

- As for the polynomial predictor, one considers a NAR formulation:

$$x_n = g(x_{n-1}, \dots, x_{n-p}) + \varepsilon_n$$

and one estimates $g(\cdot)$ in the least-square sense using the data at hand. That is, one uses the universal approximation capability of MLPs.

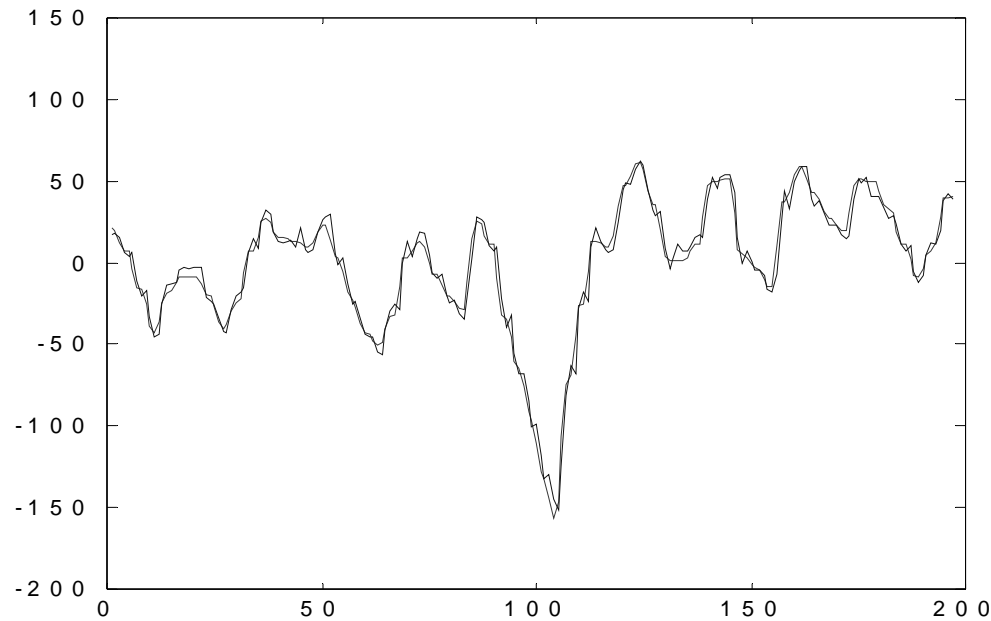
- Of course, the input vector to predict x_n is

$$\mathbf{x}_n = [x_{n-1}, x_{n-2}, \dots, x_{n-p}].$$

- It is possible to show that, when p increases, the size of the MLP increases more slowly than that of a polynomial predictor. But:
 - It is almost impossible to relate the values of the MLP parameters to the characteristic of the predicted signal (black box effect).
 - Convergence of the back-propagation algorithm is sometimes problematic.
 - Adaptive prediction is not very efficient.

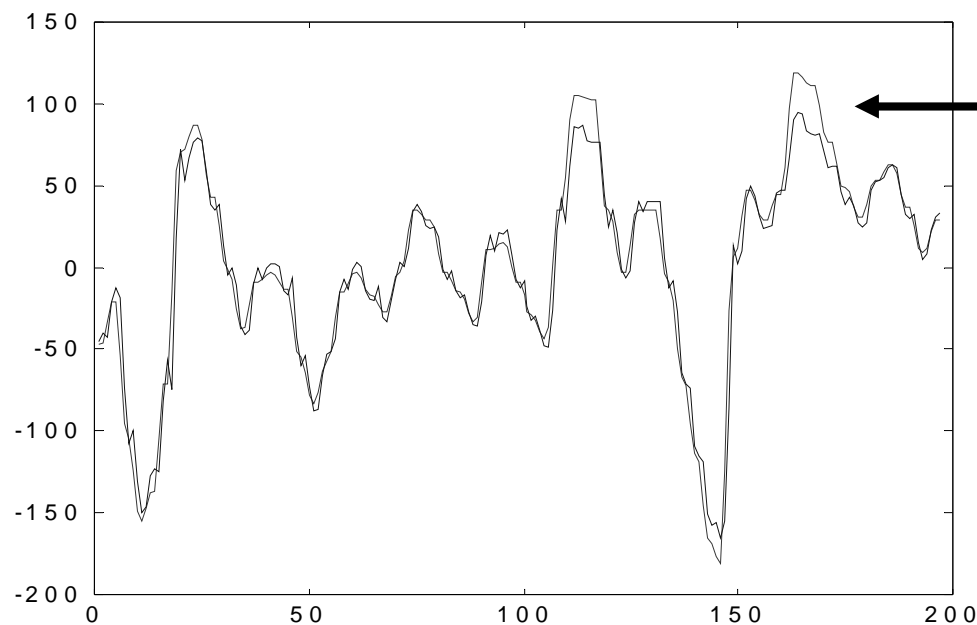
- Example: prediction of RR intervals ($p = 3, 5$ neurons in the hidden layer). ESR: - 16 dB

Signal
Prediction



- Generalization: prediction on another part of the signal. ESR = -13 dB

Signal
Prediction



Failure to
generalize

- Radial basis functions (RBF) constitute another type of neural networks, issued from interpolation theory.
- By some aspects, RBF combine the advantages of polynomial models and MLP:
 - Universal approximation property,
 - Estimation of parameters is simple,
 - There is a physical interpretation of these parameters.

- Let us suppose we have N pairs $\{\mathbf{x}_i, y_i = g(\mathbf{x}_i)\}$, $i = 1, \dots, N$, with \mathbf{x}_i p -dimensional vectors, y_i scalars, and $g(\cdot)$ a continuous function.
- the *interpolation* of $g(\cdot)$ from these N pairs is an *ill-posed* one, because there is usually no information on $g(\cdot)$ between these pairs.
- A sensible approach, called *regularization*, has been developed to tackle this type of problem.

- In this approach, one hypothesizes that $g(\cdot)$ should not oscillate erratically, and should present some degree of smoothness. This translates into defining an interpolating function $G(\cdot)$ minimizing a composite criterion:

$$C(G) = \sum_{i=1}^N [y_i - G(\mathbf{x}_i)]^2 + \rho \|PG\|^2$$

with P a differential operator.

- Parameter ρ balances the relative influence of the two terms. The first one deals with the fidelity with respect to the data, the second with interpolation smoothness.
- It is possible to show that, if P is an infinite sum of differential operators of increasing degree, and rotation and translation invariant, then $G(\cdot)$ is:

$$G(\mathbf{x}) = \sum_{i=1}^N w_i \phi(\|\mathbf{x} - \mathbf{x}_i\|) \quad \text{with} \quad \phi(s) = \exp\left(-\frac{s^2}{2\beta^2}\right)$$

- Parameter β must be specified with respect to the regularization parameter ρ . If there is a true interpolation then one must have:

$$G(\mathbf{x}_i) = y_i, \quad i = 1, \dots, N$$

and one must solve the following linear system:

$$\begin{bmatrix} \phi(\|\mathbf{x}_1 - \mathbf{x}_1\|) & \cdots & \phi(\|\mathbf{x}_1 - \mathbf{x}_N\|) \\ \vdots & \ddots & \vdots \\ \phi(\|\mathbf{x}_N - \mathbf{x}_1\|) & \cdots & \phi(\|\mathbf{x}_N - \mathbf{x}_N\|) \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_N \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

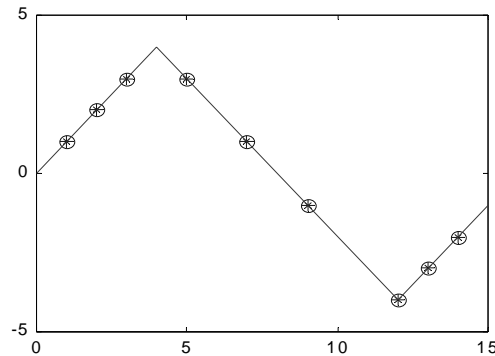
Thus a matrix equation:

$$\mathbf{A}\mathbf{w} = \mathbf{y}$$

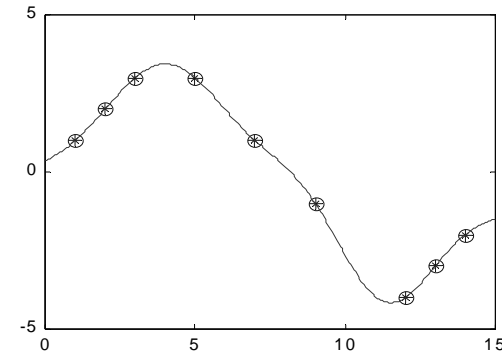
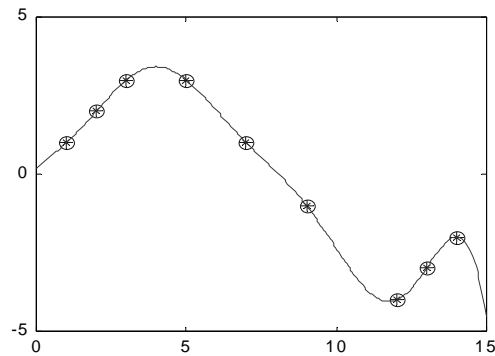
- It is possible to show that if all vectors $\{\mathbf{x}_i\}$ are different, matrix \mathbf{A} is always invertible, and it will always be possible to determine the coefficients $\{w_i\}$ of \mathbf{w} .
- The value of β defines the behavior of $G(\cdot)$ between the $\{\mathbf{x}_i\}$. The larger it is, the smoother $G(\cdot)$ is.

- Example

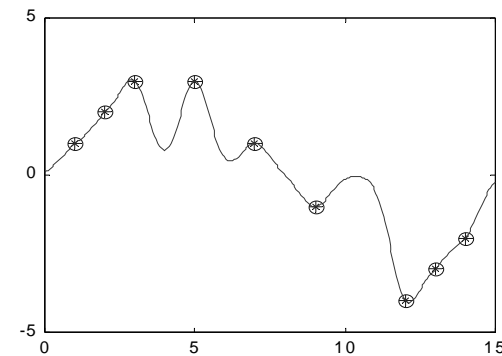
True function
(not known).



$\beta = 5$



$\beta = 2$



$\beta = 0.5$

- This approach of interpolation presents some limitations:
 - The fact that the interpolating function must include all pairs $\{\mathbf{x}_i, y_i = g(\mathbf{x}_i)\}$ makes the approach highly sensitive to noise (erroneous values) and thus prone to overfitting.
 - When the number of pairs increases $G(\cdot)$ becomes soon complex.

- The idea that has been proposed to create a new type of neural networks is quite simple. One considers a slightly different formulation:

$$G(\mathbf{x}) = \sum_{j=1}^M w_j \phi(\|\mathbf{x} - \mathbf{c}_j\|) \quad \text{avec } \phi(s) = \exp\left(-\frac{s^2}{2\beta^2}\right)$$

with $M < N$, and vectors $\{\mathbf{c}_j\}$ called *centers*, which are not constrained to be a subset of the $\{\mathbf{x}_i\}$. Function $G(\cdot)$ will now *approximate* $g(\cdot)$, typically in the least-square sense.

- To sum up, the RBF network approximates locally $g(\cdot)$ in the neighborhood of each center, and merges these local approximations to create a global one.
- The nice feature of RBF networks is that, once the centers and β have been defined, *determination of coefficients $\{w_j\}$ is simple, since $G(\cdot)$ depends linearly on them.*

- If least-square estimation is used, one must solve:

$$\begin{bmatrix} \phi(\|\mathbf{x}_1 - \mathbf{c}_1\|) & \cdots & \phi(\|\mathbf{x}_1 - \mathbf{c}_M\|) \\ \phi(\|\mathbf{x}_2 - \mathbf{c}_1\|) & \cdots & \phi(\|\mathbf{x}_2 - \mathbf{c}_M\|) \\ \vdots & & \vdots \\ \phi(\|\mathbf{x}_N - \mathbf{c}_1\|) & \cdots & \phi(\|\mathbf{x}_N - \mathbf{c}_M\|) \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_M \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} + \begin{bmatrix} e_1 \\ \vdots \\ e_N \end{bmatrix}$$

that is $\mathbf{A}\mathbf{w} = \mathbf{y} + \mathbf{e}$. \mathbf{A} is now an $N \times M$ matrix.

- One may add a column of "1" to \mathbf{A} if \mathbf{y} is significantly non zero-mean.

- When minimizing $\|\mathbf{A}\mathbf{w} - \mathbf{y}\|$, one performs of course an orthogonal projection of vector \mathbf{y} on the subspace generated by the columns of \mathbf{A} .
- Since matrix \mathbf{A} may be ill-conditioned it is better to use a robust estimation scheme, typically SVD.
- But of course centers should not be too close to each other.

- Note first there is no really optimal way to select the centers. Also, we will come back later to the problem of determining the *number* of centers.
- One of the firsts methods proposed consisted in selecting centers randomly (uniformly), in the hypercube defined by the $\{\mathbf{x}_i\}$. This is not very efficient, because the spatial distribution of the $\{\mathbf{x}_i\}$ is not taken into account.

- A second approach consists in selecting randomly the centers among the $\{\mathbf{x}_i\}$. This is already better, since statistically those centers are representative of the spatial distribution of the $\{\mathbf{x}_i\}$. That is, there are more centers in the regions of space where the $\{\mathbf{x}_i\}$ are clustered. However, the constraint remains that the centers are a subset of the $\{\mathbf{x}_i\}$.

- A solution yielding better performance is to select centers using *learning vector quantization* (LVQ), which is used in other contexts too to find representatives of vector sets. LVQ works as follows:
 1. The initial centers are chosen $\{c_j\}$ at random in the hypercube defined by the $\{x_i\}$.

2. Vectors $\{\mathbf{x}_i\}$ are presented, generally for several epochs. At iteration k , the following correction is applied to center \mathbf{c}_j :

$$\mathbf{c}_j(k+1) = \mathbf{c}_j(k) + \alpha(k) [\mathbf{x}(k) - \mathbf{c}_j(k)]$$

$$\text{if } \|\mathbf{x}(k) - \mathbf{c}_j(k)\| < \|\mathbf{x}(k) - \mathbf{c}_m(k)\|, \quad m \neq j$$

with $\mathbf{x}(k)$ the vector presented at iteration k . The net effect is to draw \mathbf{c}_j closer to $\mathbf{x}(k)$, and thus possibly to a cluster.

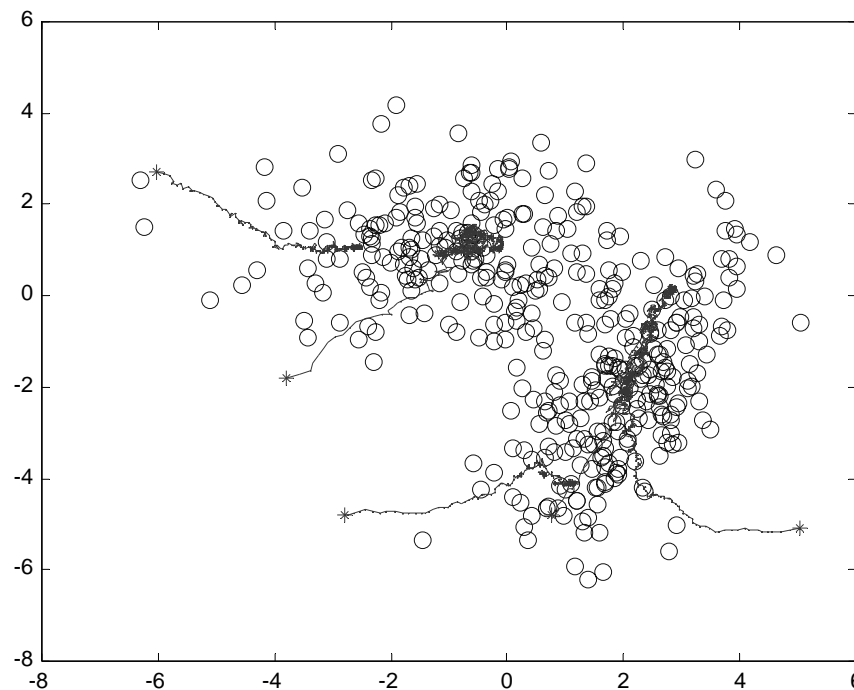
Parameter $\alpha(k)$ must decrease with k , generally using:

$$\alpha(k) = \alpha_0 \left(1 - \frac{k}{K} \right)$$

K total number of iterations, and α_0 a constant (typically $\alpha_0 = 0.05$).

- It is to be noted that LVQ is sensitive to the initial center selection.

- Example: LVQ for $M = 5$ centers, vectors $\{\mathbf{x}_i\}$ drawn from two bi-dimensional Gaussian pdf.



- Another scheme, less sensitive to initial center choice, is *Lloyd-Max vector quantization algorithm*.
- Centers $\{\mathbf{c}_j\}$ (sometimes called in this context *centroids*) are chosen such that they minimize a distortion criterion:

$$D = \sum_{j=1}^M \sum_{\mathbf{x}_i \in V_j} \|\mathbf{x}_i - \mathbf{c}_j\|^2$$

- Subsets V_j contain vectors \mathbf{x}_j such that \mathbf{c}_j is the center closest to them. So to speak, \mathbf{c}_j must be a good representative of V_j . Lloyd-Max algorithm works as follows

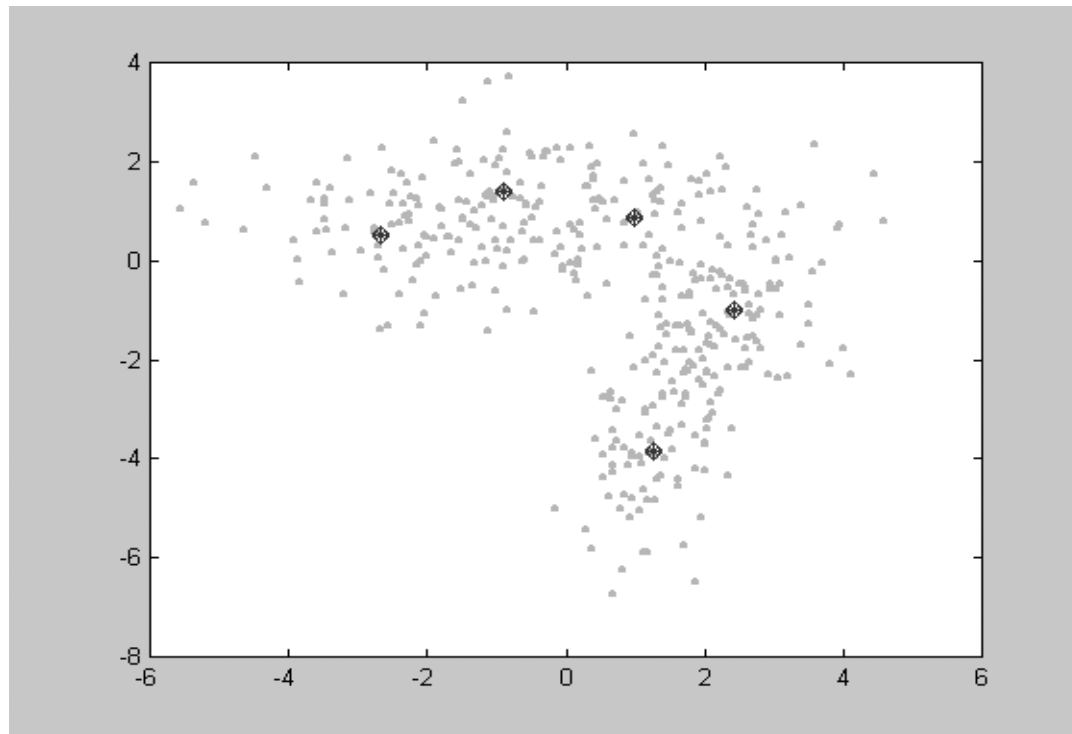
1. Initial random selection of $\{\mathbf{c}_j\}$ among the $\{\mathbf{x}_i\}$
2. Update:

$$\mathbf{c}_j \rightarrow \mathbf{c}_j = \text{mean}(\mathbf{x}_i \in V_j)$$

3. If distortion D stabilizes, stop, otherwise back to 2.

- Example: Lloyd-Max for $M = 5$ centers, vectors $\{\mathbf{x}_i\}$ drawn from two bi-dimensional Gaussian pdf.

- data
- centroids



- Once the centers $\{\mathbf{c}_j\}$ have been chosen, parameter β can be determined.
- This is the same as for interpolation: too large a β means too strong an interaction (overlap) between the local approximations. Too small a β means a “bumpy” approximation.
- A good empirical rule is:

$$\beta = \frac{d}{\sqrt{M}} \quad \text{with } d = \max_{1 \leq i, j \leq M} \|\mathbf{c}_i - \mathbf{c}_j\|$$

- Nothing imposes to use the same value for β for each RBF.
- It is even intuitively appealing that β be smaller in regions with many centers, since the spatial influence of the RBFs corresponding to these centers should be smaller
- A good empirical rule for β_j corresponding to center \mathbf{c}_j is:

$$\beta_j = \frac{1}{N\sqrt{8}} \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{c}_j\|$$

- As for MLP, one considers a NAR formulation:

$$x_n = g(x_{n-1}, \dots, x_{n-p}) + \varepsilon_n$$

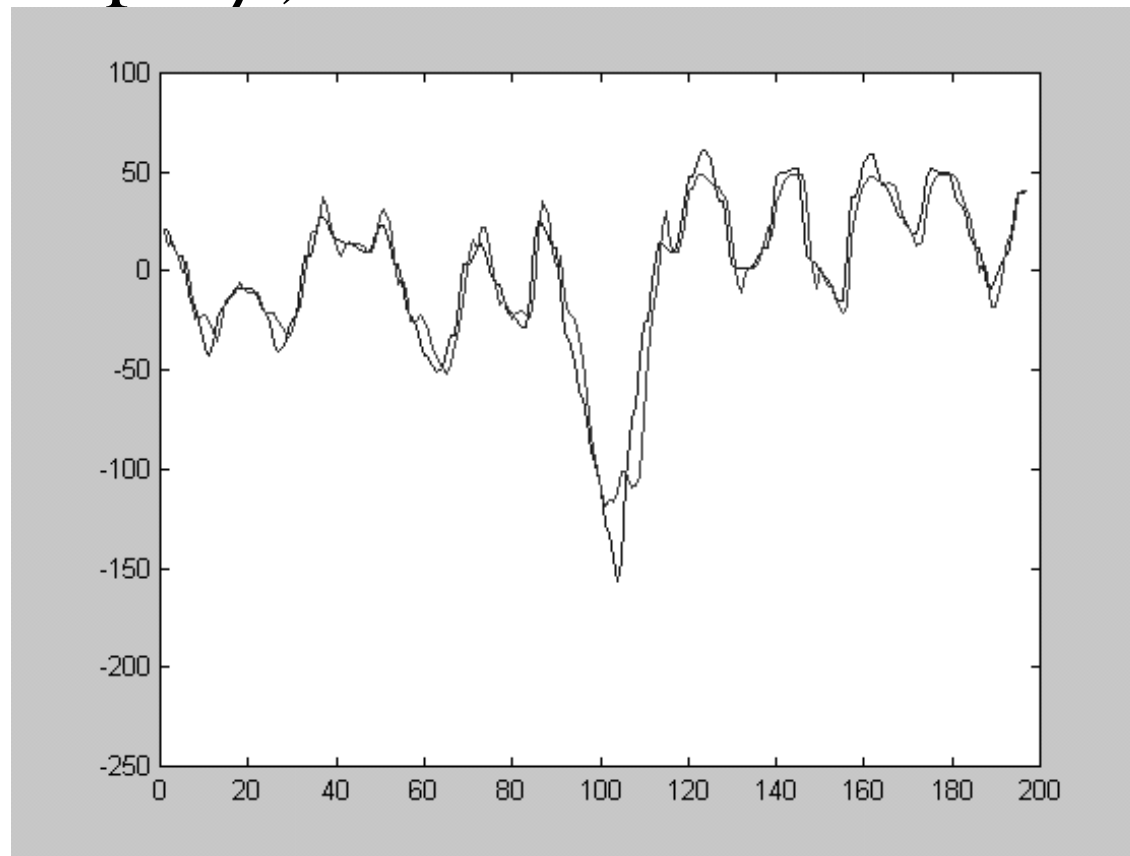
and one estimates $g(\cdot)$ in the least-square sense using the data at hand. That is, one uses the universal approximation capability of RBF networks.

- Of course, the input vector to predict x_n is

$$\mathbf{x}_n = [x_{n-1}, x_{n-2}, \dots, x_{n-p}].$$

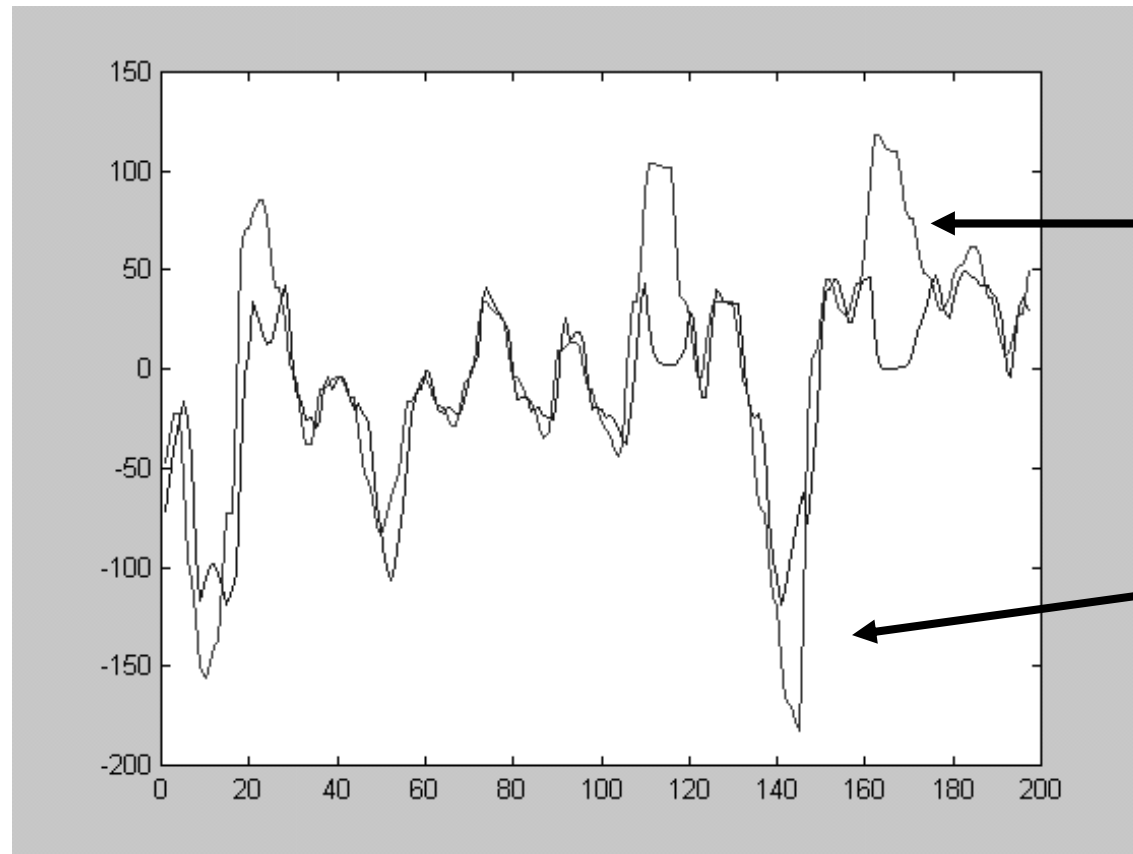
- Example: prediction of RR intervals ($p = 3, 20$ RBF, multiple β). ESR: - 10 dB

Signal
Prediction



- Generalization: prediction on another part of the signal. ESR = -4 dB

Signal
Prediction



failure to
generalize

- Center selection, which conditions the choice of parameter β also, has an important role in RBF network performance.
- LVQ and Lloyd-Max algorithms are intuitively appealing. There are however other interesting schemes [3].
- Again, we note that solving $\mathbf{A}\mathbf{w} = \mathbf{y} + \mathbf{e}$ in the least squares sense means performing an *orthogonal projection of \mathbf{y}* on the linear subspace generated by the *columns of \mathbf{A}* .

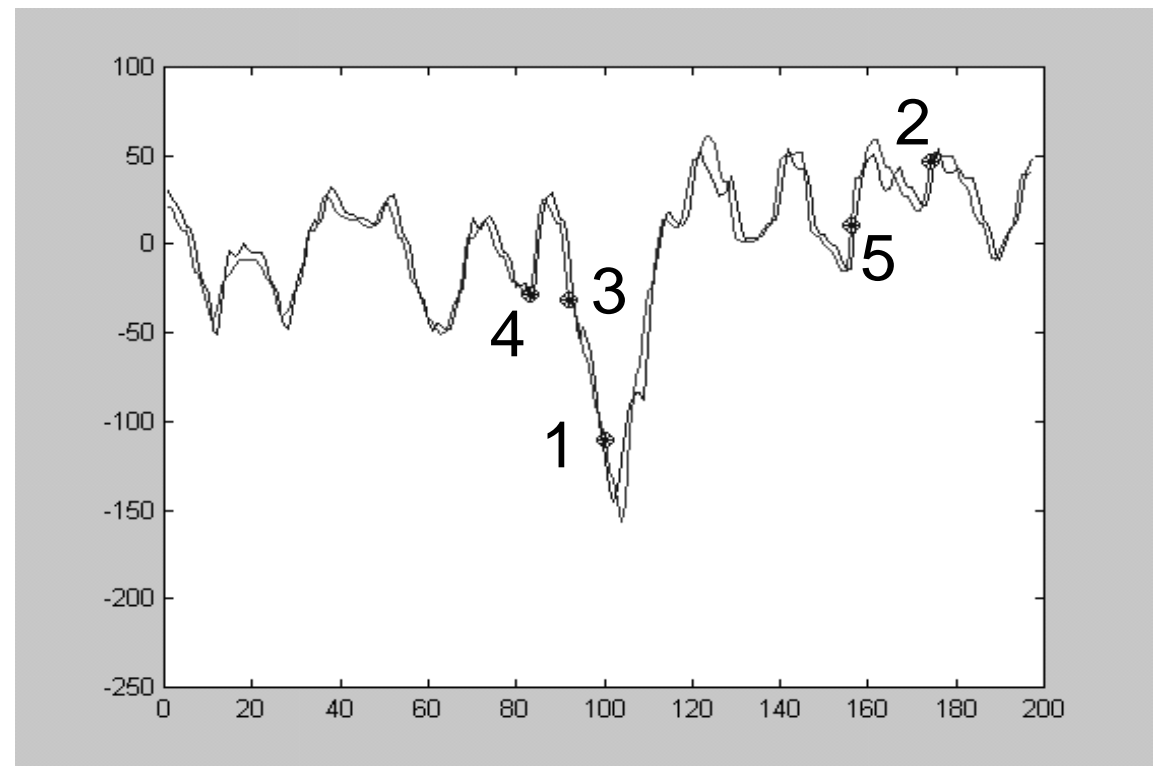
- One can apply the same procedure as for term selection in polynomial prediction, that is a Gram-Schmidt orthogonalization on selected columns (*orthogonal least squares*, OLS).
- At each iteration, one adds in \mathbf{A} the column yielding maximum error reduction. This column is defined by a pair (\mathbf{c}_j, β_j) , and center \mathbf{c}_j is selected among the vectors $\{\mathbf{x}_i\}$.

- One proceeds as follows:
 - 1) The first center is selected by examining all the $\{\mathbf{x}_i\}$, and keeping the smallest least squares error.
 - 2) Successive centers are chosen among the remaining $\{\mathbf{x}_i\}$, with an orthogonalization of $[\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_N)]^T$ with respect to the terms already selected. In this way the center giving maximum error reduction is spotted.

- Example: prediction of RR intervals ($p = 3, 5$ RBF, multiple β). ESR: - 11 dB

Signal Prediction

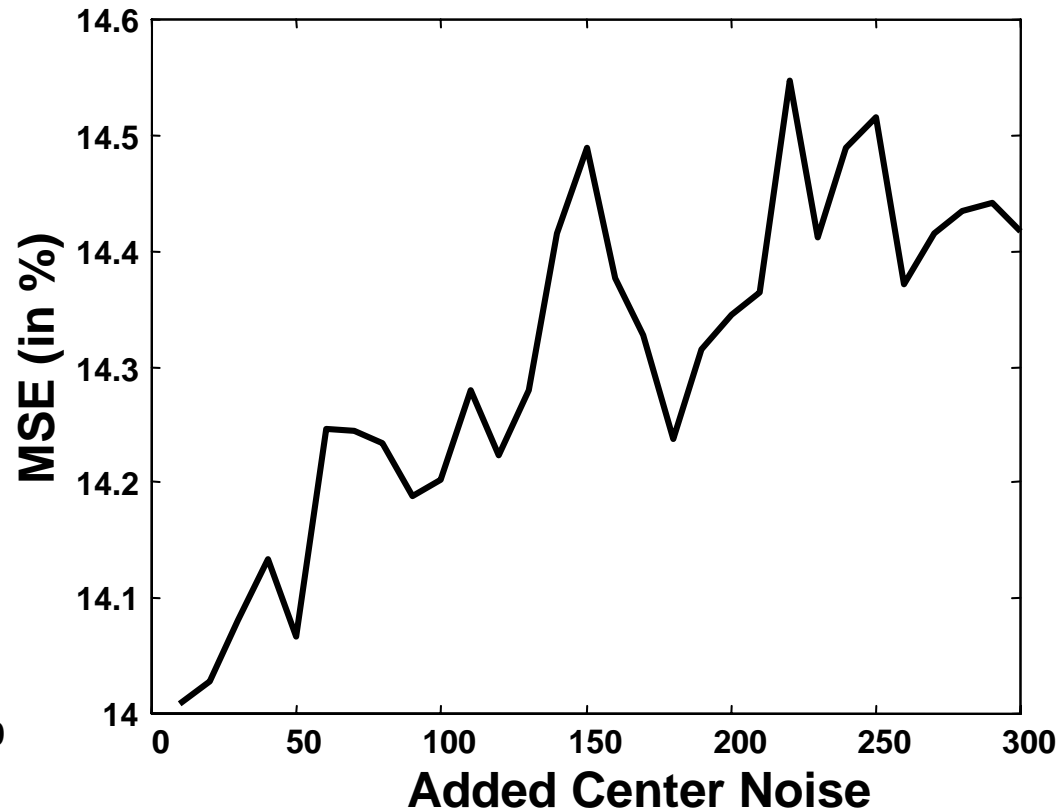
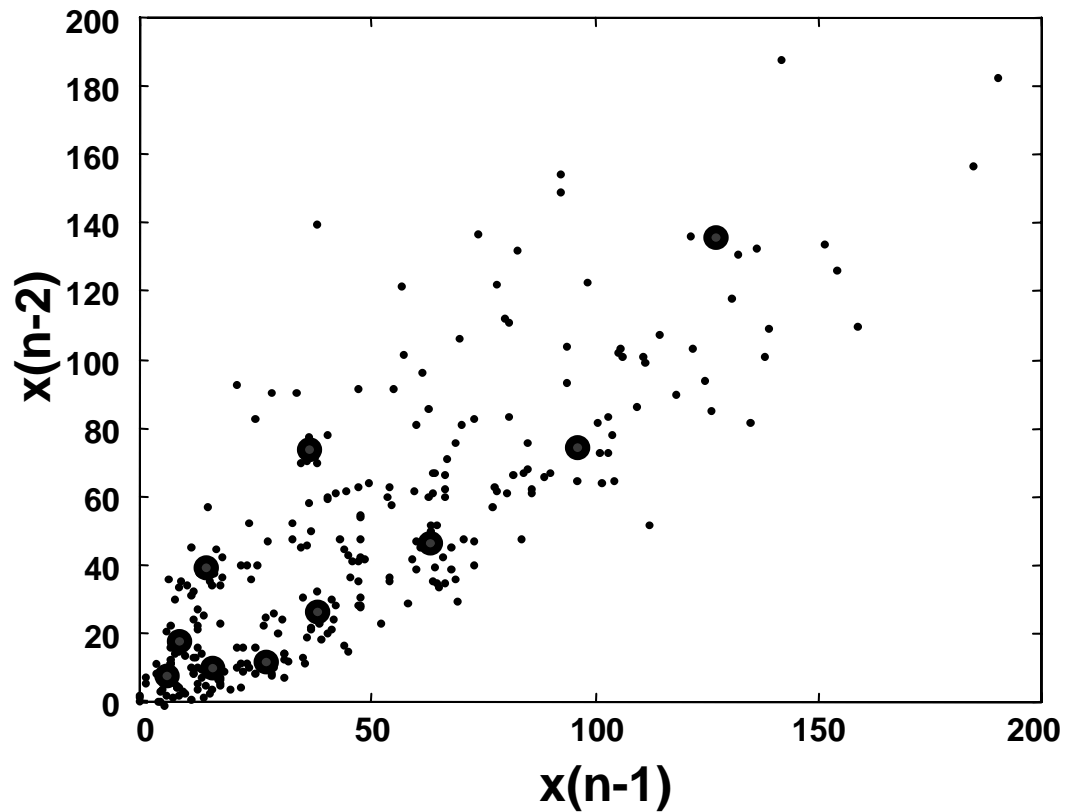
- samples corresponding to selected \mathbf{x}_i

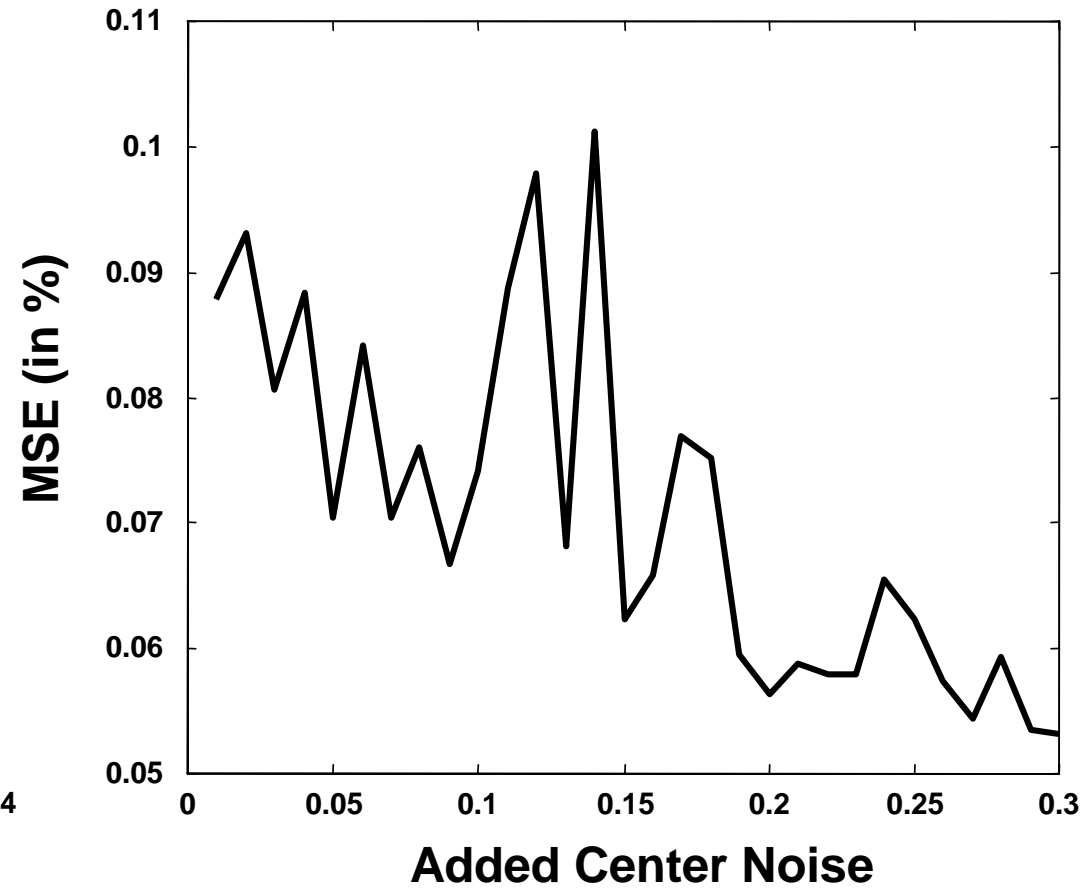
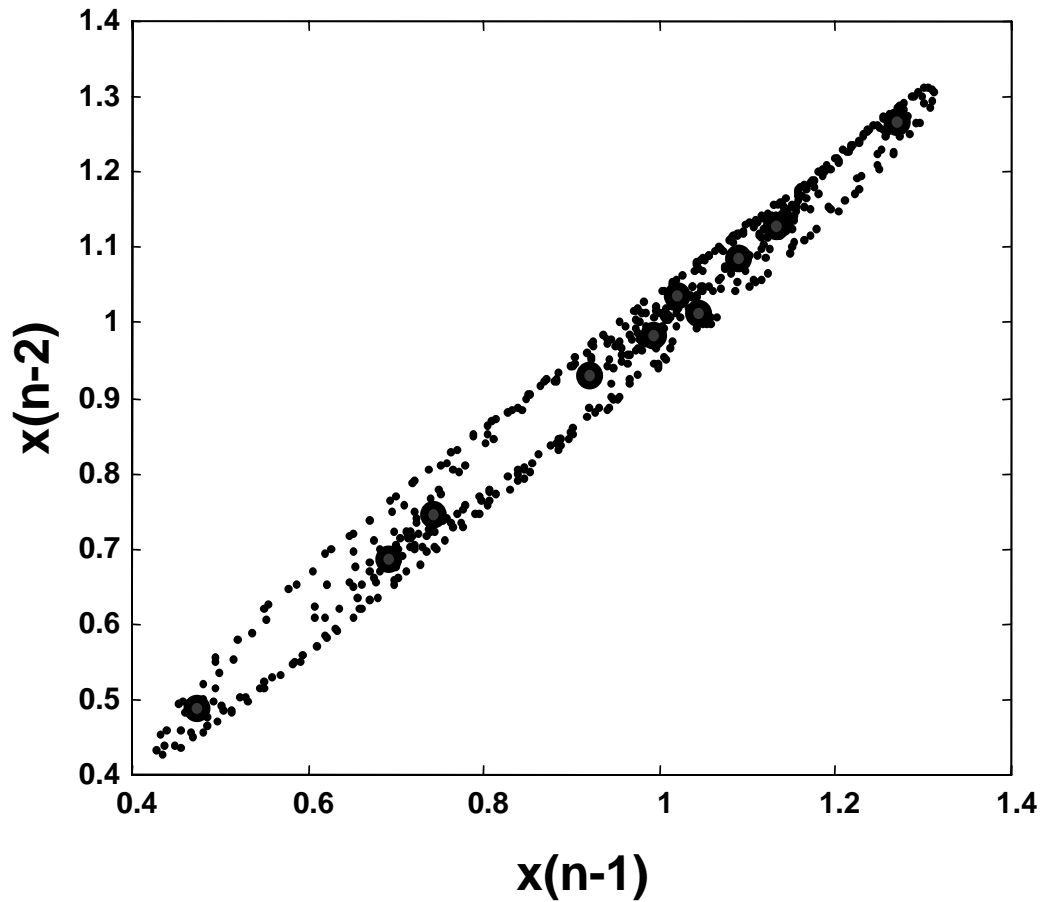


$$MDL = \underbrace{N \cdot \ln \sigma_e^2}_{\text{Data}} + \underbrace{M \cdot \ln N}_{\text{Model}} + \underbrace{?}_{\text{Model}}$$

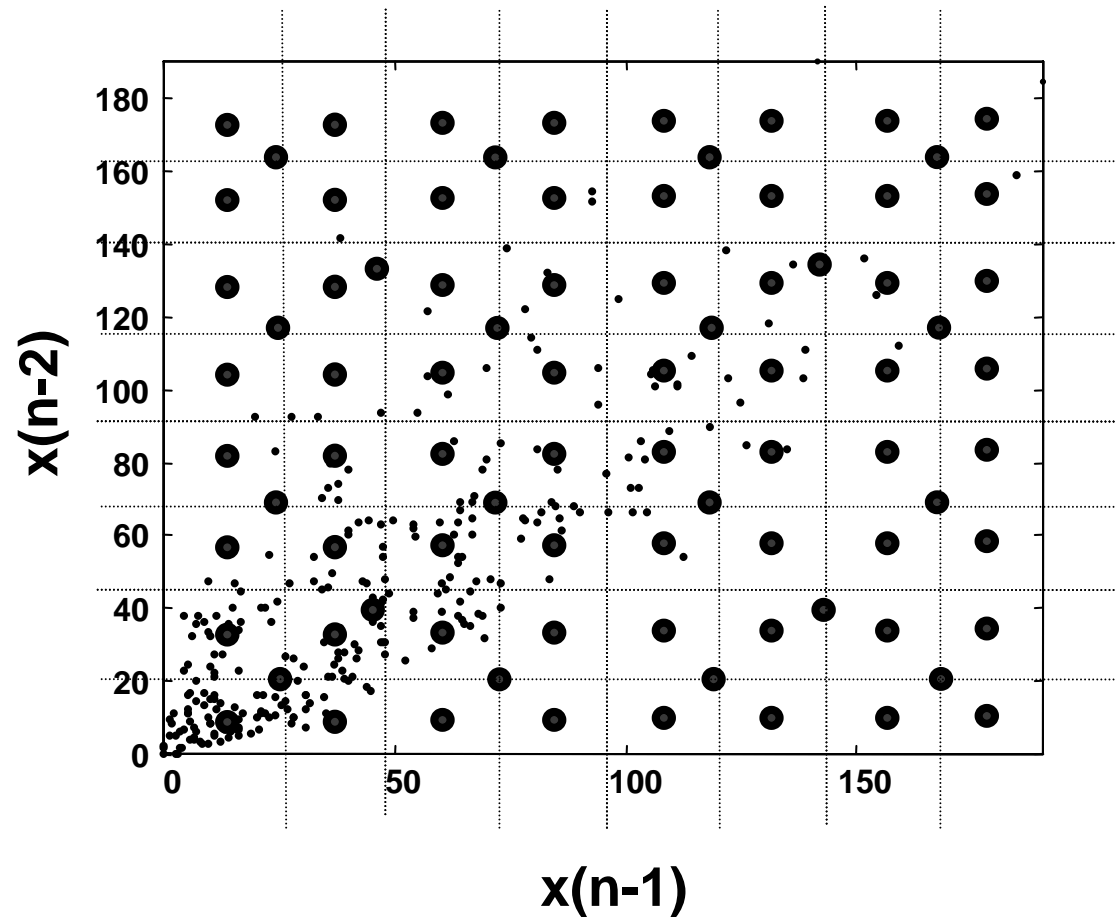
linear parameters only !

nonlinear parameters

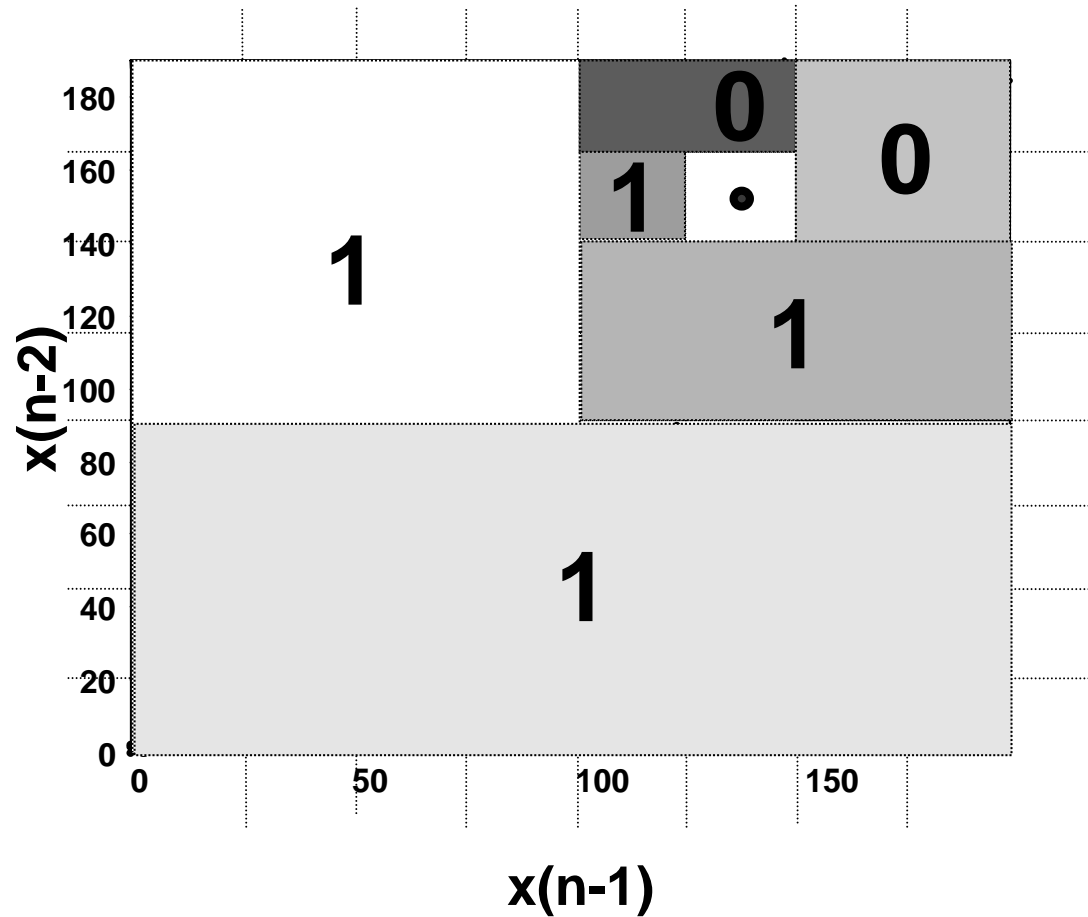




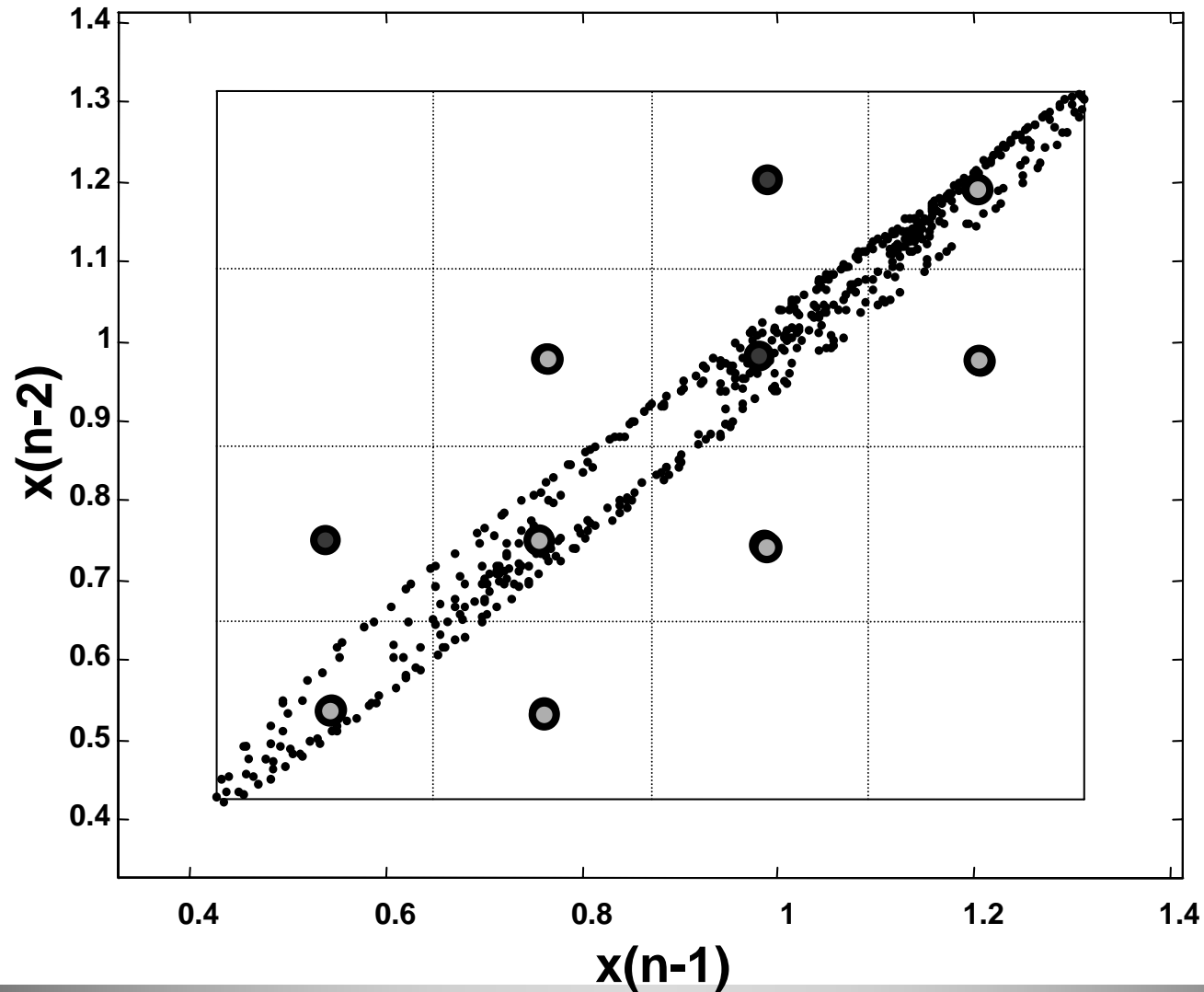
2 bits
4 bits
6 bits



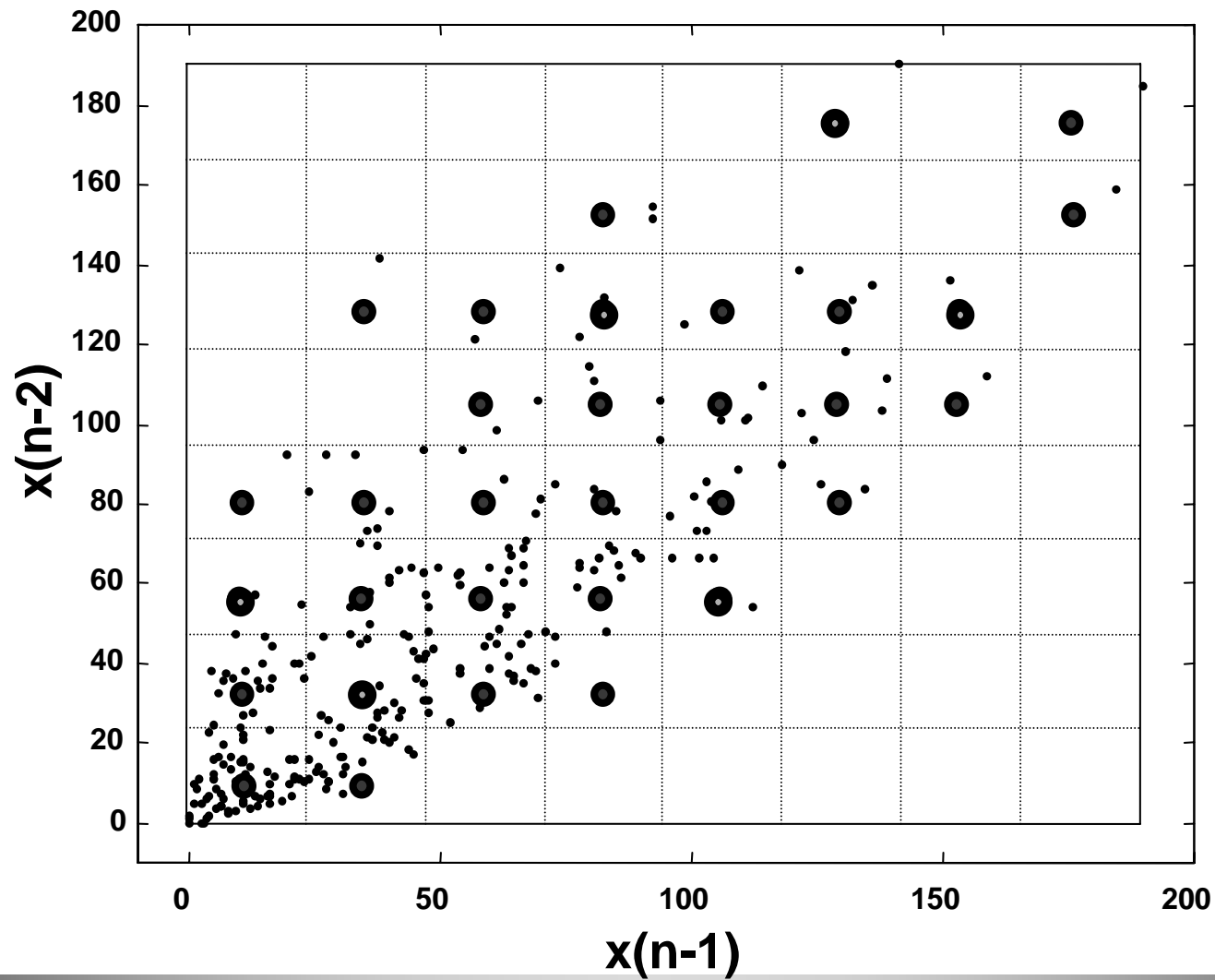
1 1 0 1 1 0



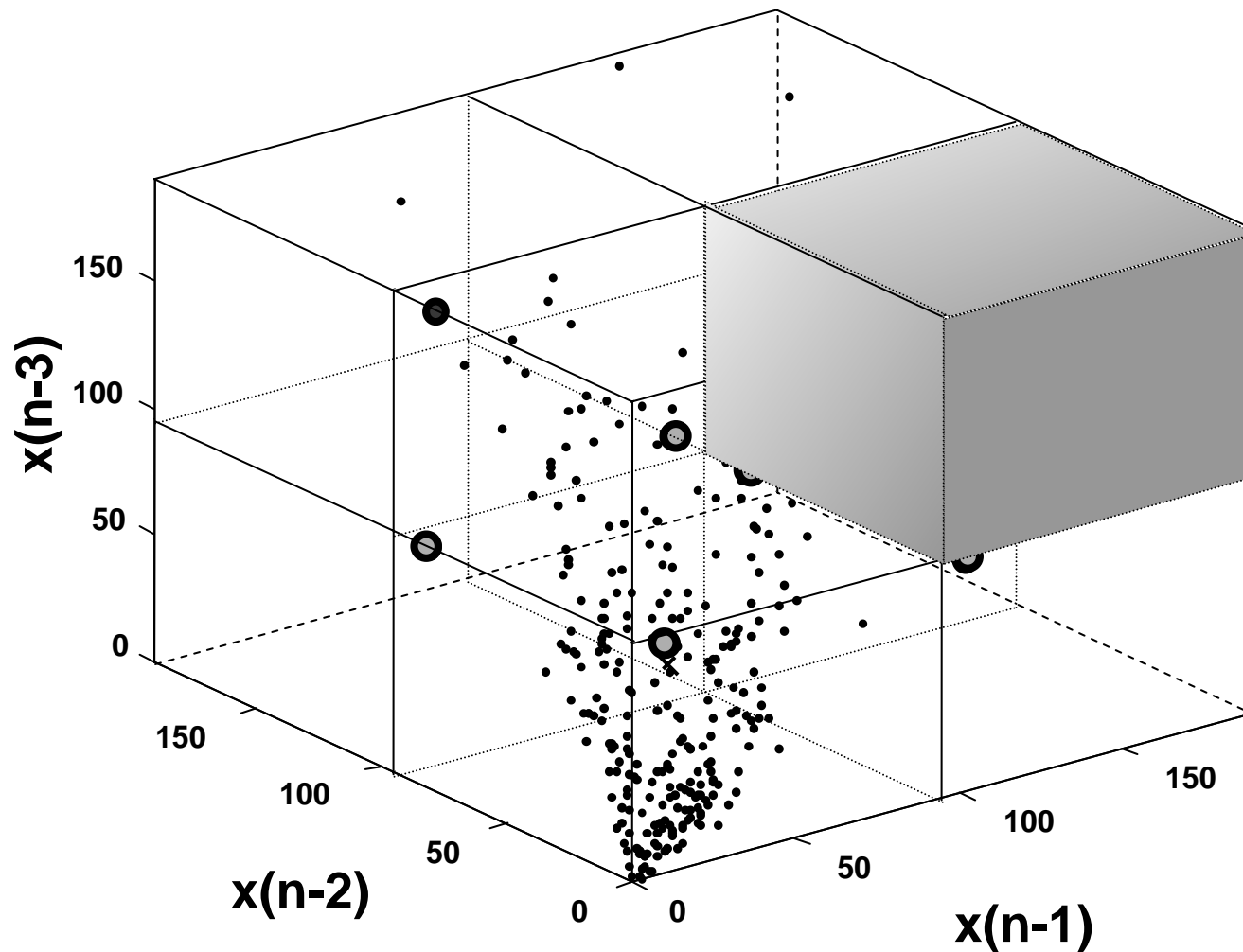
The Best Subset of Centers (4 bits)



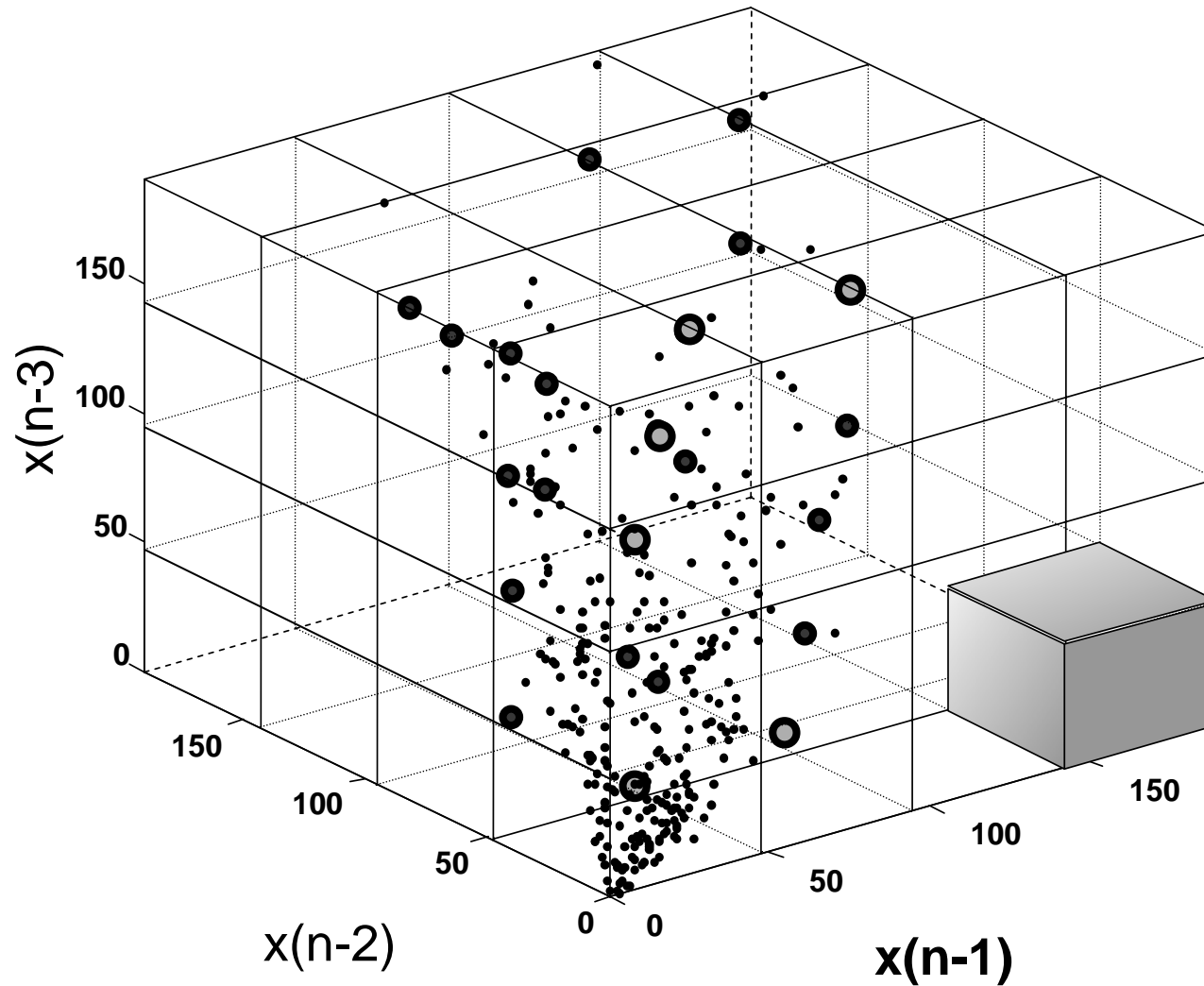
The Best Subset of Centers (6 bits)



3-D State Space Quantization (3 bits)

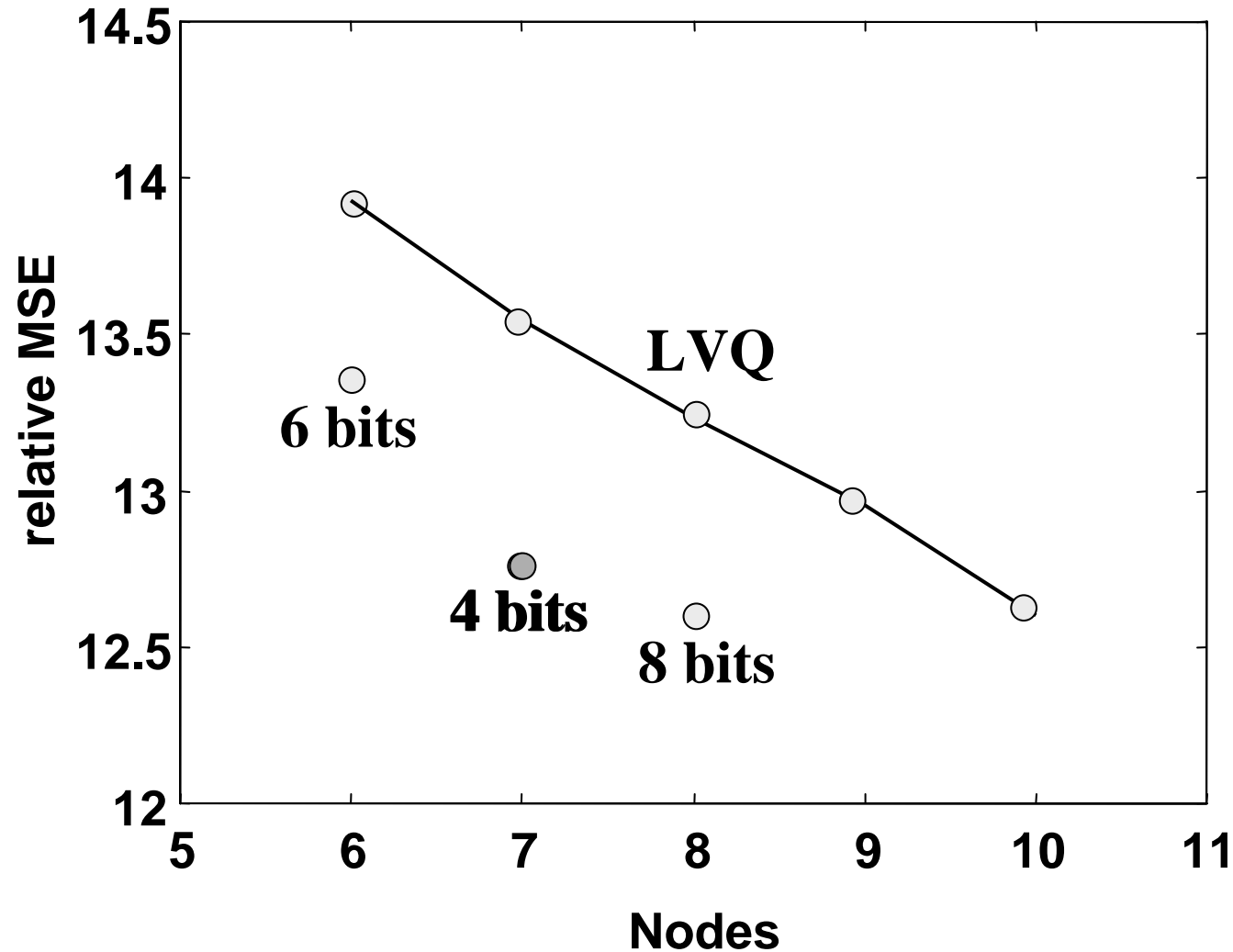


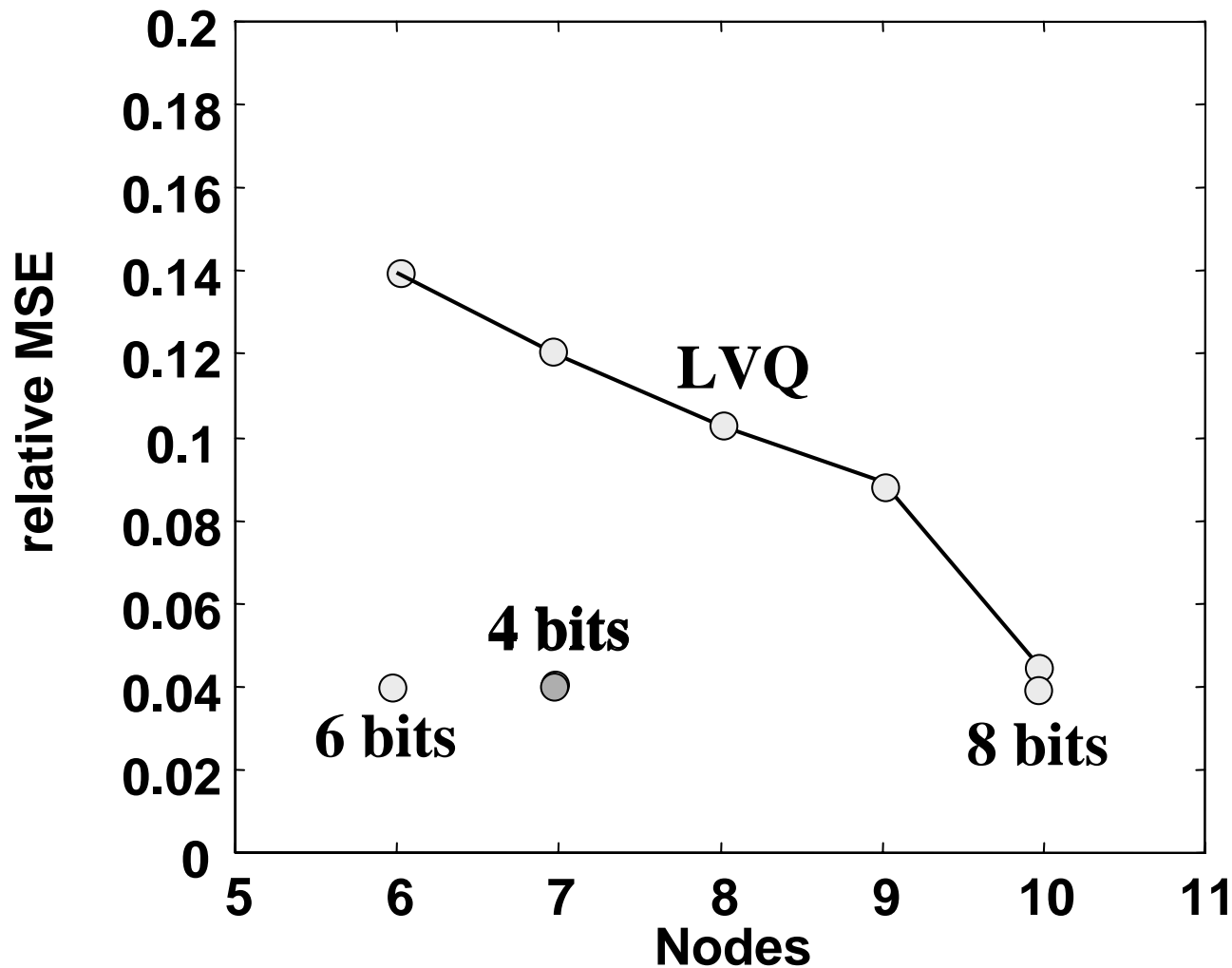
3-D State Space Quantization (6 bits)



$$MDL = N \cdot \ln \sigma_e^2 + M \cdot (\ln N + B \cdot \ln 2)$$

B: Number of bits to encode one center





- It is of course possible to derive an adaptive version of RBF network prediction.
- In what follows, reference to time index k will not be explicit.

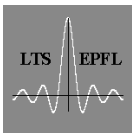
$w(i)$ i th component of coefficient vector \mathbf{w}

$\phi(i)$ i th RBF

ϕ RBF vector

- instantaneous square error is thus

$$\varepsilon^2 = [y - \mathbf{G}(\mathbf{x})]^2 = [y - \mathbf{w}^T \phi]^2$$



- In LMS, the updates are:

$$\begin{aligned} \mathbf{w} &\rightarrow \mathbf{w} - \mu_1 \frac{\partial \varepsilon^2}{\partial \mathbf{w}} = \mathbf{w} - 2\mu_1 \varepsilon \frac{\partial \varepsilon}{\partial \mathbf{w}} = \mathbf{w} + 2\mu_1 \varepsilon \frac{\partial G(\mathbf{x})}{\partial \mathbf{w}} \\ &= \mathbf{w} + 2\mu_1 \varepsilon \phi \end{aligned}$$

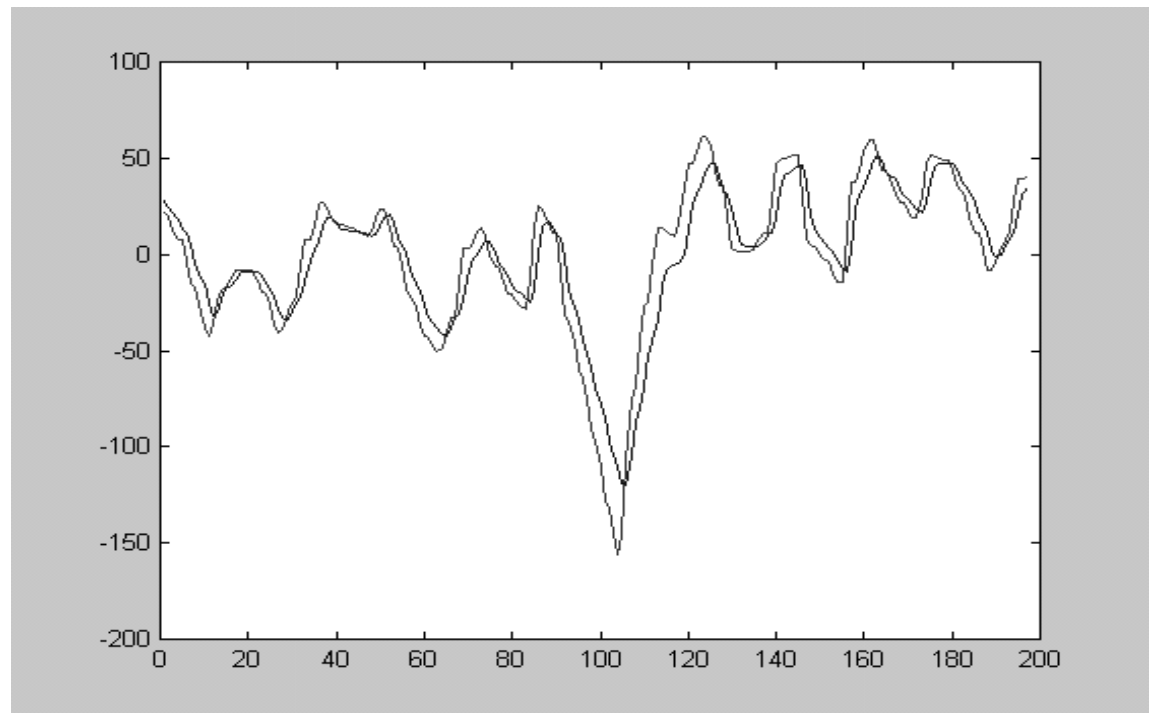
$$\mathbf{c}_j \rightarrow \mathbf{c}_j + 2\mu_2 \varepsilon \frac{\partial G(\mathbf{x})}{\partial \mathbf{c}_j} = \mathbf{c}_j + 2\mu_2 \varepsilon \mathbf{w}(j) \phi(j) (\mathbf{x} - \mathbf{c}_j)$$

$$\beta \rightarrow \beta + 2\mu_3 \varepsilon \frac{\partial G(\mathbf{x})}{\partial \beta} = \beta + \frac{2\mu_3 \varepsilon}{\beta^3} \sum_{i=1}^M \|\mathbf{x} - \mathbf{c}_i\|^2 \phi(i)$$

- Example: prediction of RR intervals ($p = 3, 5$ RBF, unique β). ESR: - 8 dB

Signal
Prediction

*Multiple epochs
before convergence*



- Remarks:
 - Better than batch algorithm, but not impressive
 - Tuning of adaptation parameters is uneasy (except for w for which it is possible to normalize by $\phi^T \phi$)
- *But maybe the gradient approach is not suited to center and β parameter update. For instance, nothing guarantees that the mean square error has a unique minimum.*

- Why not *decouple* the update of \mathbf{w} , for which LMS should work well, from that of the centers and β ?
- One might use an LVQ-style update for the centers, then modify β by using the empirical rule based on the distance between centers.
- But it is necessary to define a link between the two update schemes.

1. One first computes the modification on \mathbf{c}_i , the center closest to \mathbf{x} with LVQ (but α must remain constant to preserve adaptability):

$$\mathbf{c}_i \rightarrow \mathbf{c}_i + \Delta\mathbf{c}_i = \mathbf{c}_i + \alpha(\mathbf{x} - \mathbf{c}_i)$$

2. If the maximum distance between centers changes, β should become $\beta + \Delta\beta$, by using the empirical rule for unique β .

The output changes from $G(\mathbf{x})$ to $G(\mathbf{x}) + \Delta G(\mathbf{x})$.

- But:
$$\Delta G(\mathbf{x}) \approx (\Delta \mathbf{c}_i)^\top \frac{\partial G(\mathbf{x})}{\partial \mathbf{c}_i} + \Delta \beta \frac{\partial G(\mathbf{x})}{\partial \beta}$$
- Thus:

$$\frac{\partial (G(\mathbf{x}) + \Delta G(\mathbf{x}))}{\partial \mathbf{w}} \approx \frac{\partial G(\mathbf{x})}{\partial \mathbf{w}} + (\Delta \mathbf{c}_i)^\top \frac{\partial G(\mathbf{x})}{\partial \mathbf{w} \partial \mathbf{c}_i} + \Delta \beta \frac{\partial G(\mathbf{x})}{\partial \mathbf{w} \partial \beta}$$

- And the instantaneous square error becomes:

$$\varepsilon^2 = [y - G(\mathbf{x}) - \Delta G(\mathbf{x})]^2$$

The coefficient vector \mathbf{w} is updated with:

$$\mathbf{w} \rightarrow \mathbf{w} - \mu_1 \frac{\partial \varepsilon^2}{\partial \mathbf{w}} = \mathbf{w} + 2\mu\varepsilon \frac{\partial (G(\mathbf{x}) + \Delta G(\mathbf{x}))}{\partial \mathbf{w}}$$

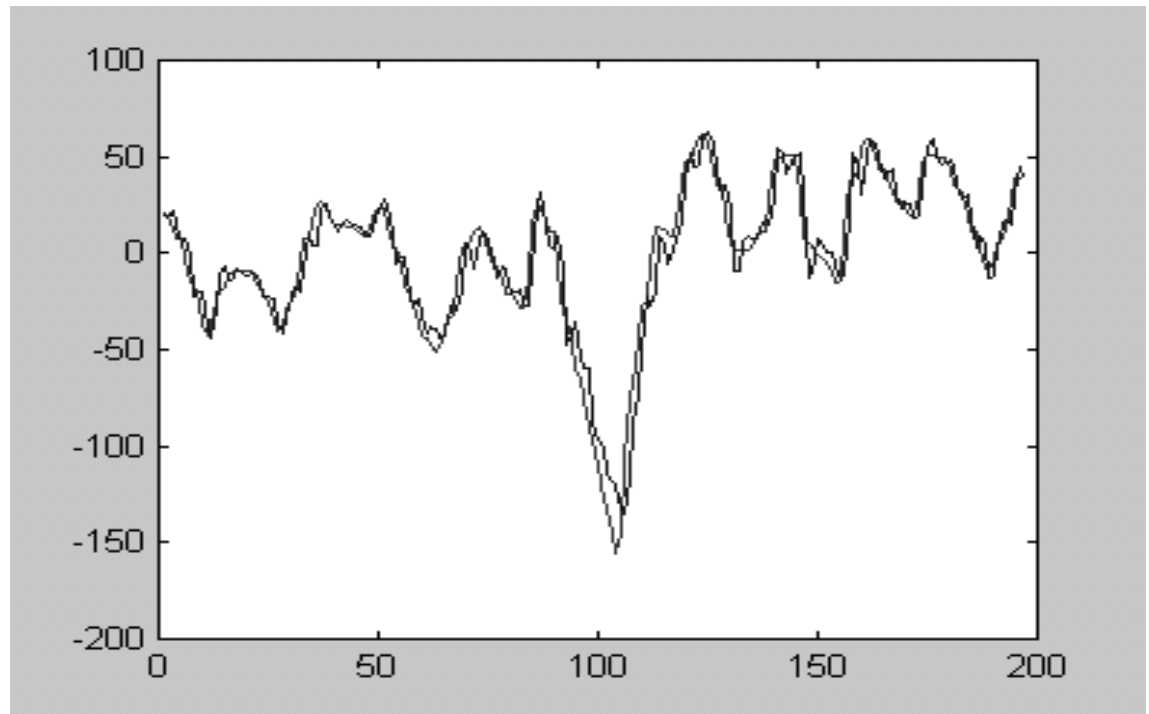
- With some computations:

$$(\Delta \mathbf{c}_i)^\top \frac{\partial G(\mathbf{x})}{\partial \mathbf{w} \partial \mathbf{c}_i} = \frac{\alpha}{\beta^2} \|\mathbf{x} - \mathbf{c}_i\|^2 \phi(i) \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \quad \leftarrow \begin{array}{l} \text{ith} \\ \text{position} \end{array}$$

$$\frac{\partial G(\mathbf{x})}{\partial \mathbf{w} \partial \beta} = \frac{1}{\beta^3} \begin{bmatrix} \|\mathbf{x} - \mathbf{c}_1\|^2 \phi(1) \\ \vdots \\ \|\mathbf{x} - \mathbf{c}_M\|^2 \phi(M) \end{bmatrix}$$

- Example: prediction of RR intervals ($p = 3, 5$ RBF, unique β). ESR: - 10.7 dB

Signal
Prediction



*Multiple epochs
before convergence*

- Remarques:
 - Prediction performance improves in the non-stationary intervals.
 - It is a lot easier to tune μ for w and α for the centers.
 - This approach is easily modified for the multiple β case.

1. S. Haykin, *Neural Networks, a Comprehensive Foundation*, McMillan, NY, 1994.
2. J. C. Principe, N. R. Euliano, and W. Curt Lefebvre, *Neural and Adaptive Systems*, J. Wiley, NY, 2000.
3. S. Chen, C.F.N. Cowan, and P.M. Grant, “Orthogonal least squares learning algorithm for radial basis function networks,” *IEEE Trans. Neural Networks*, vol. 2, no. 2, 1991, pp. 302-309.
4. B. Mulgrew, “Applying radial basis functions,” *IEEE Signal Process. Mag.*, vol. 13, no. 2, 1996, pp. 50-65.