

# Fast Bootstrap for model selection

M. Verleysen

Université catholique de Louvain (Belgium)

Machine Learning Group

<http://www.dice.ucl.ac.be/mlg/>

# Outline

- Model selection
  - estimation of generalization error
  - resampling methods
- Bootstrap
  - plug-in principle
  - computational load
- Fast Bootstrap
  - idea and hypothesis
  - reduced number of experiments
- Experiments
  - artificial regression example
  - Santa Fe A time series prediction

# Outline

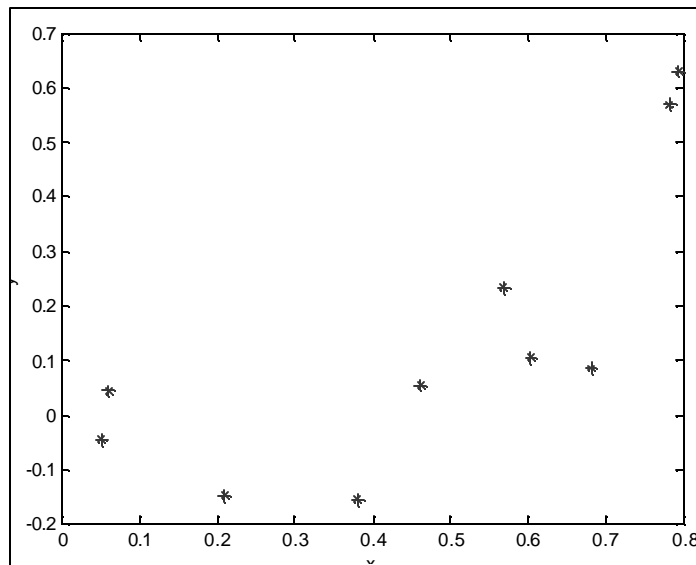
- Model selection
  - estimation of generalization error
  - resampling methods
- Bootstrap
  - plug-in principle
  - computational load
- Fast Bootstrap
  - idea and hypothesis
  - reduced number of experiments
- Experiments
  - artificial regression example
  - Santa Fe A time series prediction

# Model selection is necessary

- Database to model

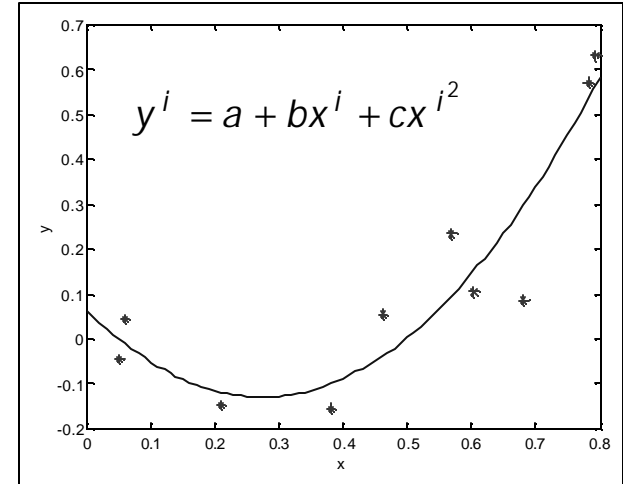
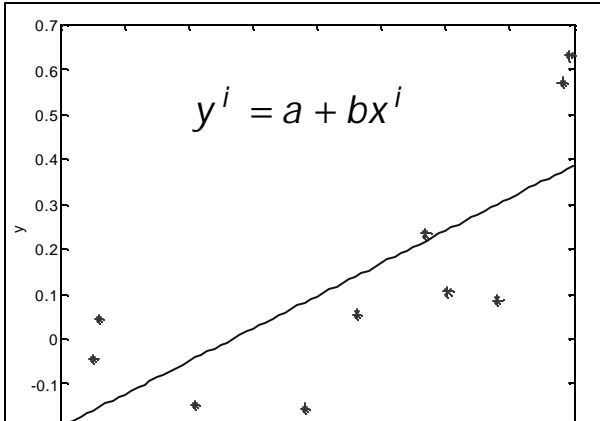
$$\{\mathbf{x}^i, y^i\}, \text{ with } \mathbf{x}^i \in \mathfrak{R}^D, y^i \in \mathfrak{R}, 1 \leq i \leq N$$

- Example ( $d=1$ ):

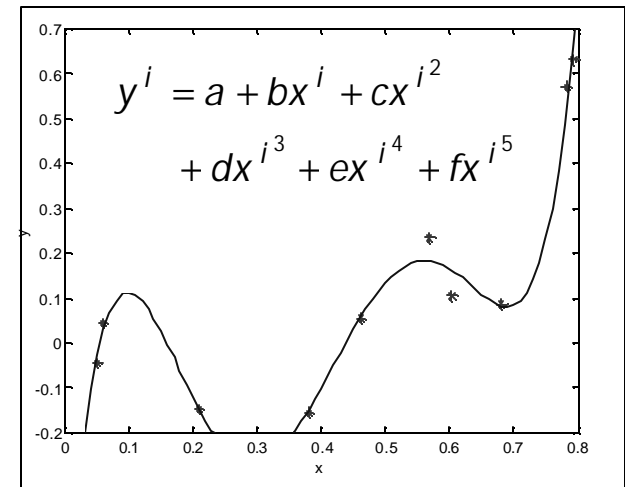
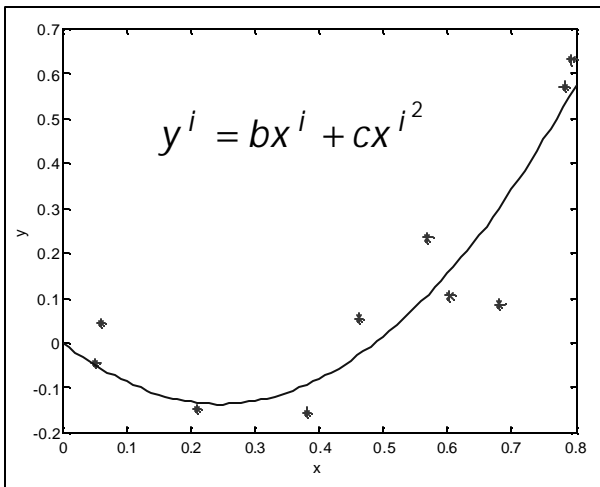


Which (polynomial) model structure to use?

# Some possible model structures...



The model structure used to generate the data (with noise)



# Model structure selection is performed according to an error criterion

## ■ Notations

$$\mathbf{x}_t \in R^d, y_t \in R$$

$$\hat{y}_t = g(x_t, \mathbf{q})$$

## ■ Generalization error

$$E_{gen}(\mathbf{q}) = \lim_{N \rightarrow \infty} \sum_{t=1}^N \frac{(g(\mathbf{x}_t, \mathbf{q}) - y_t)^2}{N}$$



$$\hat{E}_{gen}(\mathbf{q}) = \sum_{t=1}^N \frac{(\hat{y}_t - y_t)^2}{N} \quad (= \hat{E}_{gen})$$

# How to estimate the generalization error ?

$$\hat{E}_{gen}(\mathbf{q}) = \sum_{t=1}^N \frac{(\hat{y}_t - y_t)^2}{N} \quad (= \hat{E}_{gen})$$

- Problems:
  - The test sample cannot be used for learning
  - Need for 3 sets:
    - Learning
    - Validation (model selection)
    - Test (estimate of model performances)
  - The available sample is always too small in practise...
  
- Validation is learning (of hyperparameters), not test!

# Resampling methods

- Small sample: asymptotic results are difficult to use
- **Empirical** model structure selection !
- Possible methods
  - Validation
  - *K*-fold cross-validation
  - Leave-One-Out
  - Bootstrap



# Outline

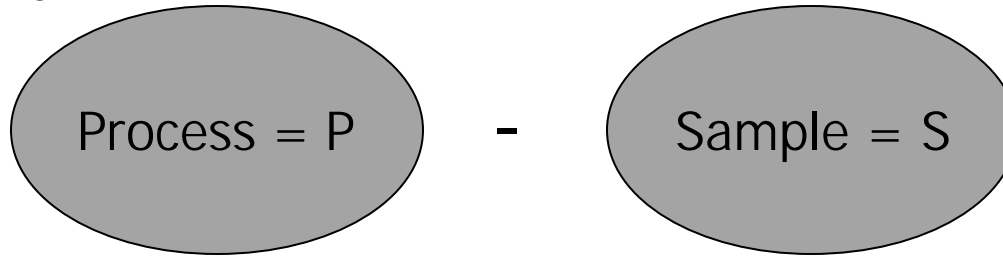
- Model selection
  - estimation of generalization error
  - resampling methods
- Bootstrap
  - plug-in principle
  - computational load
- Fast Bootstrap
  - idea and hypothesis
  - reduced number of experiments
- Experiments
  - artificial regression example
  - Santa Fe A time series prediction

# Why the bootstrap ?

- Experimentally: its variance is lower, for a fixed number of experiments
- It is sound on a statistical point of view (both other methods are too!)
- It makes the following approximation possible...

# The bootstrap is rather intuitive

We have

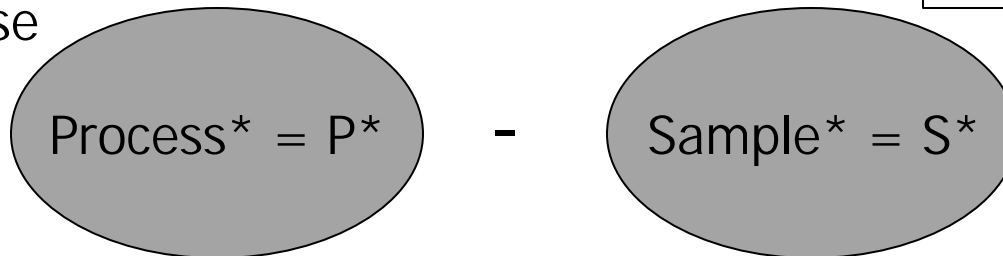


= optimism

**Bootstrap hypothesis**

=

We use



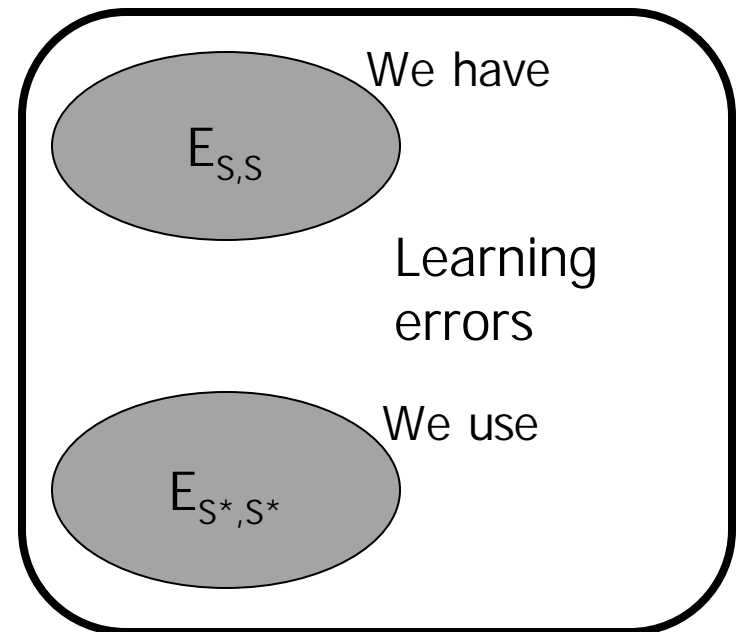
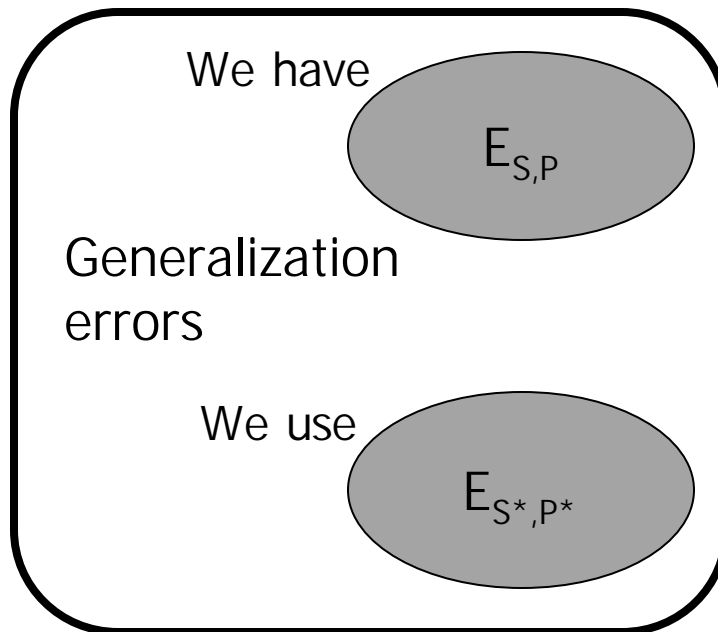
= difference

Key ideas: plug-in principle and drawing with replacement

Bootstrap – plug-in

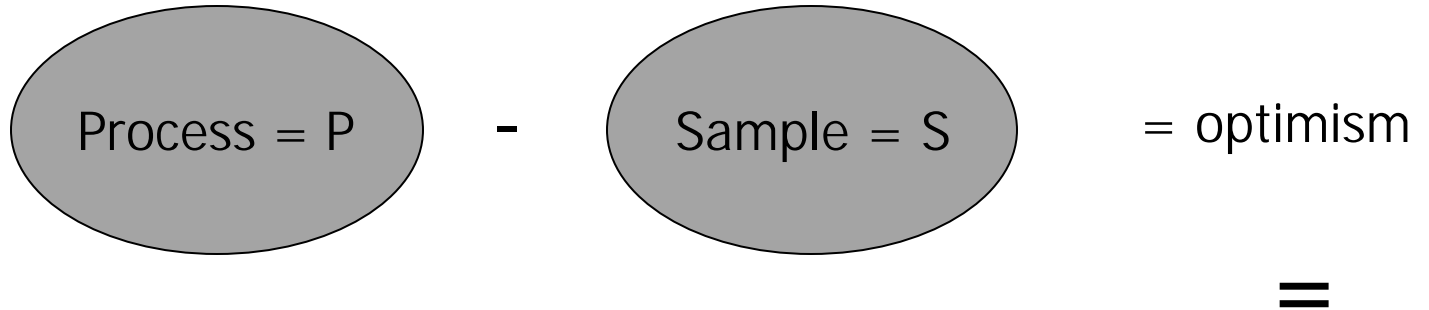
# The plug-in principle is the main idea of the bootstrap

- The plug-in principle is often used in statistics
- Here, this principle gives:  $\text{Process}^* = \text{Sample}$
- All errors are denoted here as  $E_{\text{learning set, test set}}$ 
  - $E_{S,S}$  is a **learning error**
  - $E_{S,P}$  is a **generalization error** ( $S \cap P = \emptyset$ )

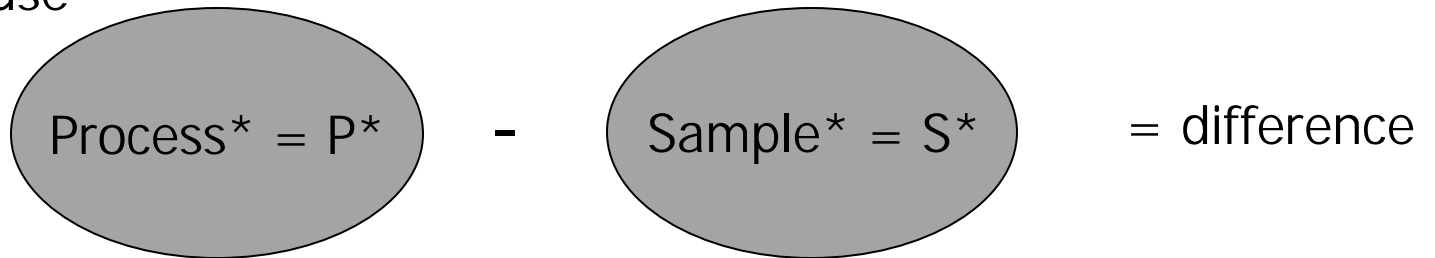


# What are the samples in the bootstrap ?

instead of



If we use

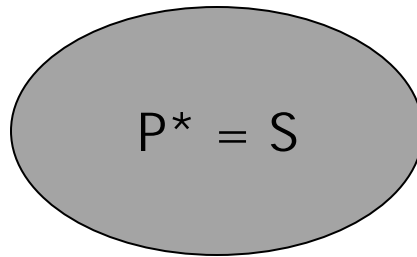


then we must have  $S^*$  and  $P^*$  !

- Bootstrap:  $P^* = S$  (this is all what we know...)
- What about  $S^*$  ?

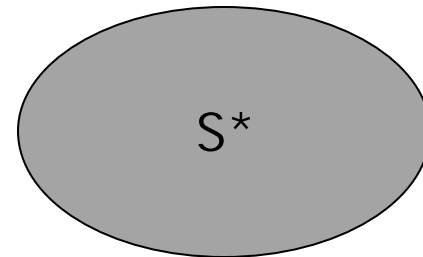
# What about the new learning sample $S^*$ ?

- The new learning sample  $S^*$  must be randomly drawn from what we know (i.e.  $P^* = S$ )
- To keep the same size: draw with replacement



1  
2  
3  
4  
5  
6  
7  
8  
9  
10

Same size !



4  
6  
3  
7  
3  
10  
5  
3  
9  
10

# The intuitive idea can be easily rewritten

We have

$$E_{S,P} - E_{S,S} = \text{optimism}$$

=

We use

$$E_{S^*,P^*} - E_{S^*,S^*} = \text{difference}$$

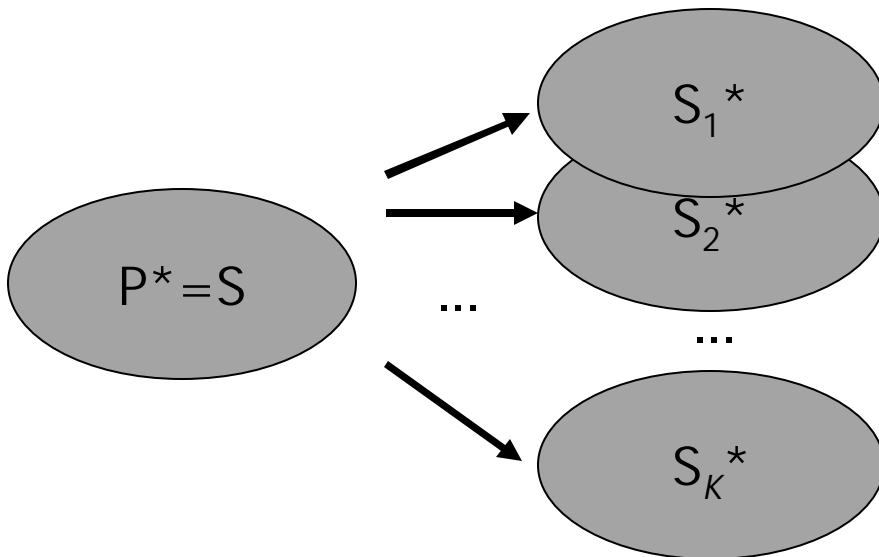
$$\begin{aligned} \hat{E}_{gen} &= E_{S,P} + E_{S,S} - E_{S,S} \\ &\stackrel{\text{Plug-in}}{=} E_{S,S} + (E_{S,P} - E_{S,S}) \\ &= E_{S,S} + (E_{S^*,P^*} - E_{S^*,S^*}) \\ &= E_{S,S} + \text{optimism} \end{aligned}$$

- Last step: to estimate

# Estimation of the optimism

$$\text{optimism} = E_{S^*, P^*} - E_{S^*, S^*}$$

- One estimation is not enough; fluctuations due to
  - sample  $S^*$
  - learning (initial conditions, local minima,...)
  - other numerical problems
- Then: repeat !

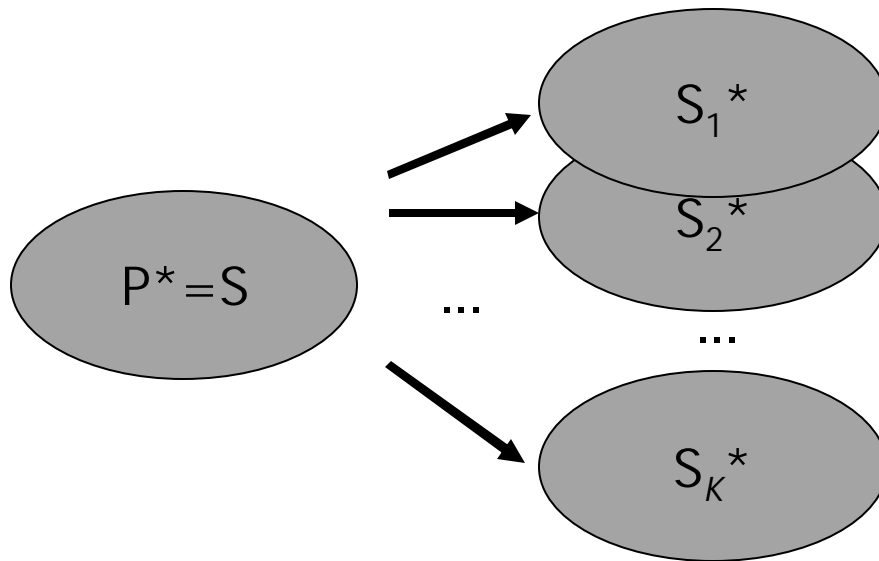


$$\text{optimism} = \frac{1}{K} \sum_{k=1}^K (E_{S^{*k}, P^*} - E_{S^{*k}, S^{*k}})$$



# In summary...

- This is what we know:



- We draw
- Then we compute

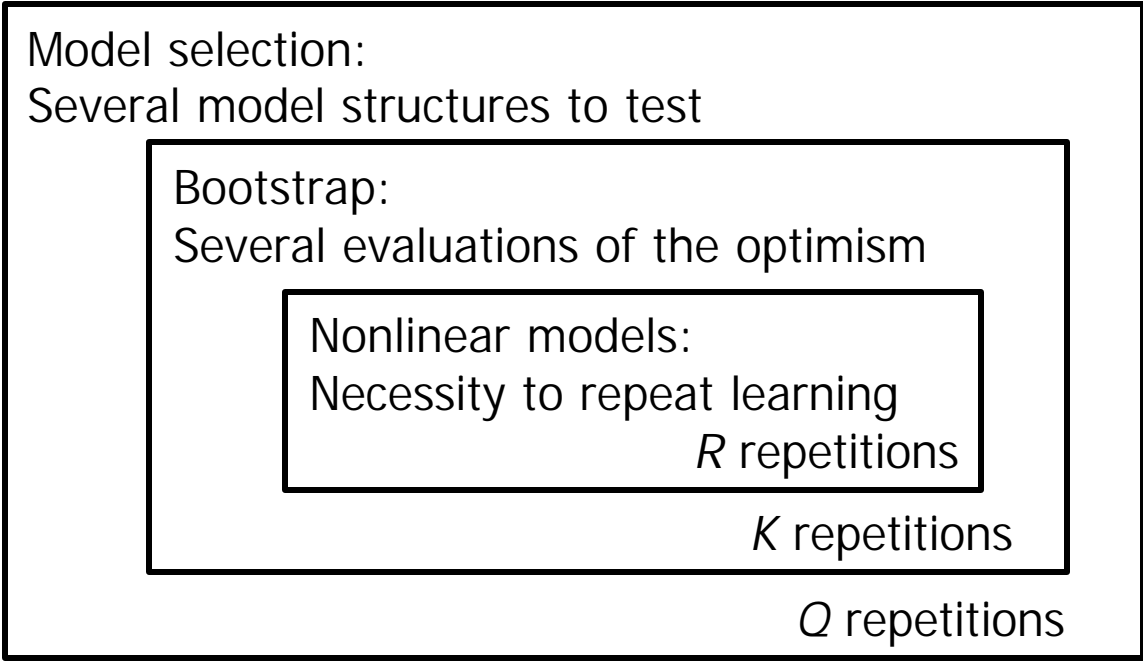
$$\hat{E}_{gen} = E_{S,S} + \frac{1}{K} \sum_{k=1}^K (E_{S_k^*, P^*} - E_{S_k^*, S_k^*})$$

K+1 learnings

2K+1 evaluations

# It seems very nice, but...

- The number of experiments grows dramatically!
- Context: model selection ( $Q$  models to test,  $1 = q = Q$ )



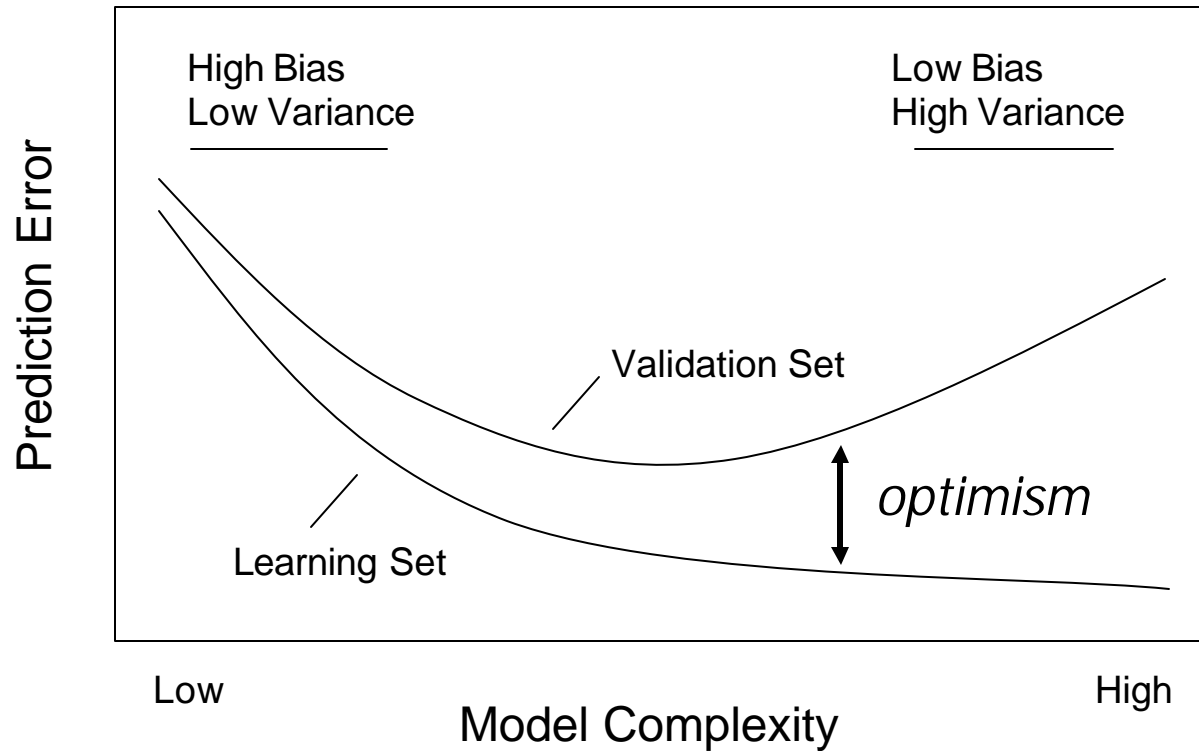
- $R \times K \times Q$  : too much!

# Outline

- Model selection
  - estimation of generalization error
  - resampling methods
- Bootstrap
  - plug-in principle
  - computational load
- Fast Bootstrap
  - idea and hypothesis
  - reduced number of experiments
- Experiments
  - artificial regression example
  - Santa Fe A time series prediction

# Fast bootstrap

- Attempt to decrease  $K \times Q$



- Idea: *optimism* is a very smooth function
  - A polynomial of order  $o$

# Is it justified ?

- Yes, experimentally
- Yes, see AIC and BIC:

– AIC: 
$$\hat{E}_{gen}(\mathbf{q}) = \sum_{t=1}^N \frac{(g(x_t, \mathbf{q}) - y_t)^2}{N} + \frac{2}{N} \dim(\mathbf{q})$$

– BIC: 
$$\hat{E}_{gen}(\mathbf{q}) = \sum_{t=1}^N \frac{(g(x_t, \mathbf{q}) - y_t)^2}{N} + \frac{\ln(N)}{N} \dim(\mathbf{q})$$

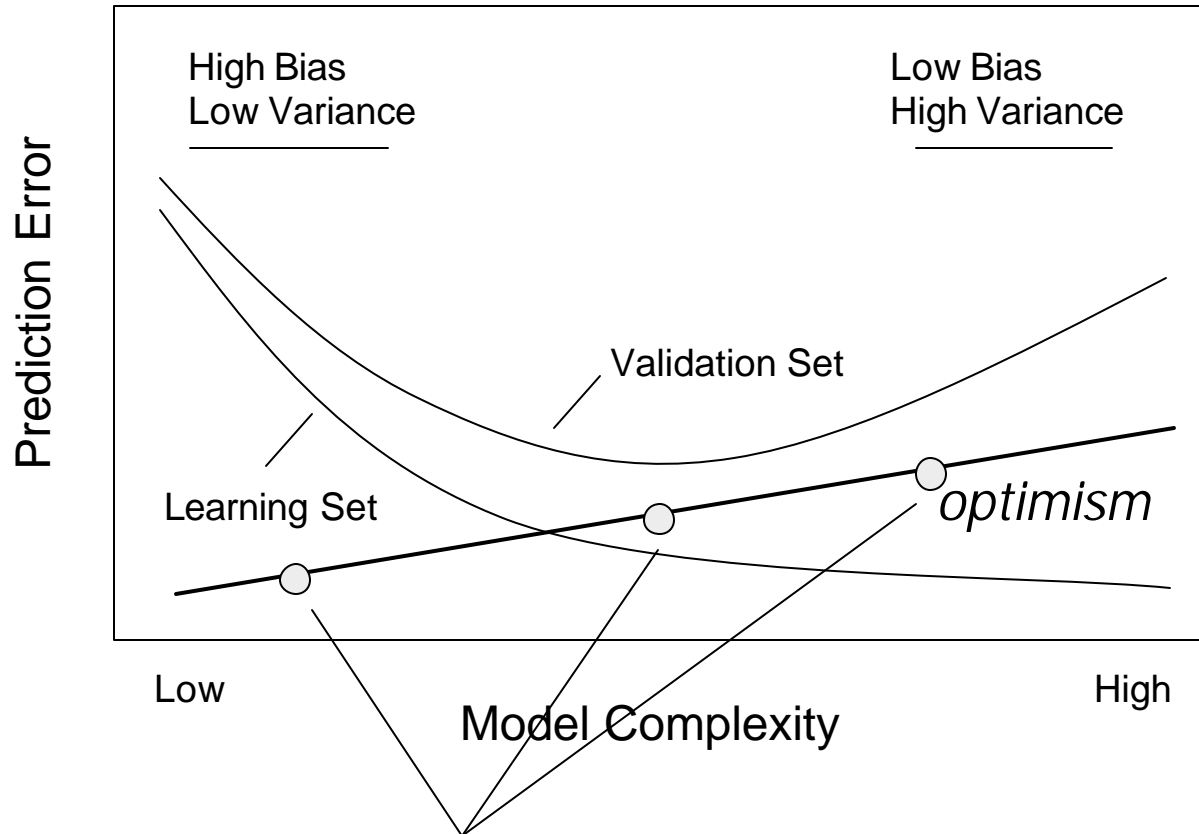
– Extension to nonlinear models

$\dim(\mathbf{q})$  is replaced by  $\dim(\mathbf{q})\mathbf{s} = \dim(\mathbf{q}) \frac{\sum_{t=1}^N (\hat{y}_t - y_t)^2}{N - \dim(\mathbf{q})}$

- In all cases: *optimism* proportional to  $\dim(\theta)$ 
  - Polynom of order  $o=1$  !

# Where is the lower number of experiments ?

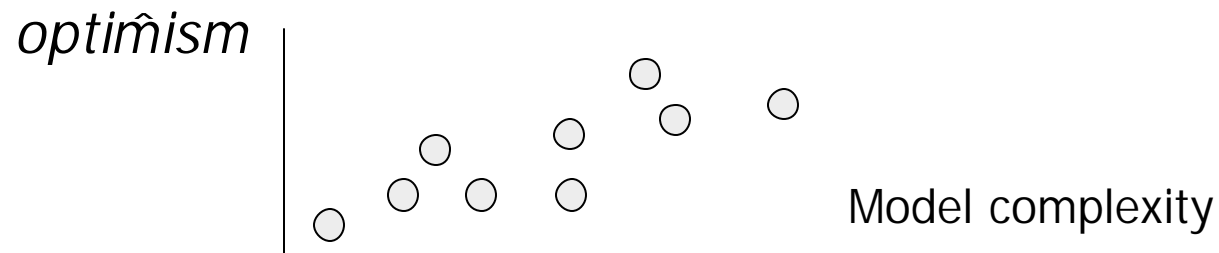
Fast Bootstrap – reduced number of experiments



- Lower number  $Q$  of models to test
- Each model to test: possibility to decrease  $K$  (evaluation noise averaged in fitting)

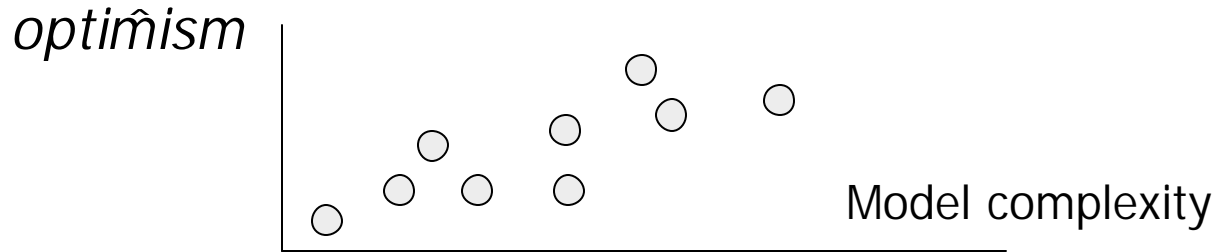
# Yes, but... how many experiments ?

- How to choose  $Q \times K$  in practise?
- Try... and test
- Example:
  - $Q$  models are tested (each one needs  $K$  bootstrap samples)



- Is it linear ? Use Fisher's statistics !

# Yes, but... how many experiments ? 2/



- Test  $o_1=1$  (linear) against  $o_2=2$

$$F_{o_2-o_1, Q-o_2-1} = \frac{\frac{SR_1 - SR_2}{o_2 - o_1}}{\frac{SR_2}{Q - o_2 - 1}} \quad \left( SR = \sum_{t=1}^N (\hat{y}_t - y_t)^2 \right)$$

- If test passes:
  - It is linear
  - We have enough experiments !



# Outline

- Model selection
  - estimation of generalization error
  - resampling methods
- Bootstrap
  - plug-in principle
  - computational load
- Fast Bootstrap
  - idea and hypothesis
  - reduced number of experiments
- Experiments
  - artificial regression example
  - Santa Fe A time series prediction

# And now... the experiments

- Two datasets:
  - Artificial function approximation problem
  - Santa Fe A time series forecasting
- Three models
  - Multi-Layer Perceptron
  - Radial-Basis Function Network
  - Least-Square Support Vector Machine

# The model structure parameter

## ■ Multi-Layer Perceptron

$$\hat{y}_t = \sum_{i=1}^M w_i \tanh \left( \sum_{j=1}^D w_{ij} x_{tj} \right)$$

## ■ Radial-Basis Function Networks

$$\hat{y}_t = \sum_{i=1}^M w_i \exp \left( - \frac{\|\mathbf{x}_t - \mathbf{c}_i\|}{2s_i^2} \right)$$

Model structure parameters

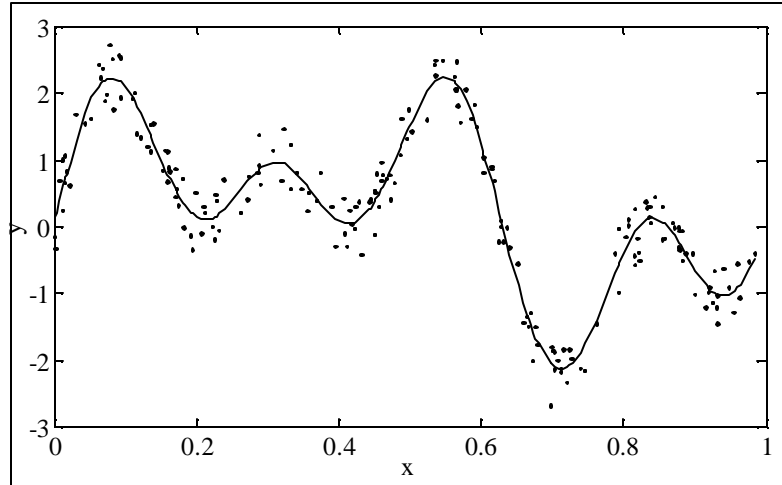
## ■ Least-Square Support Vector Machine

$$\hat{y}_t = \mathbf{?}^T \mathbf{j}(\mathbf{x}_t)$$

- But optimization is regularized

$$\min_{\mathbf{?}, \mathbf{e}} J(\mathbf{?}, \mathbf{e}) = \mathbf{?}^T \mathbf{?} + \mathbf{g} \sum_{i=1}^N (y_t - \hat{y}_t)$$

# Artificial regression example

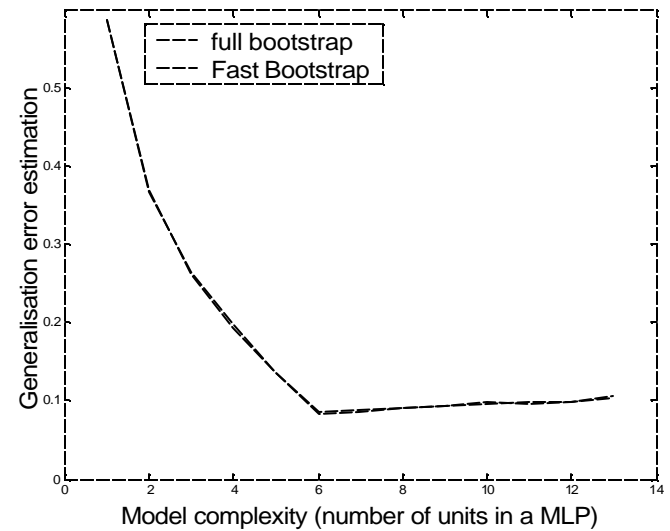
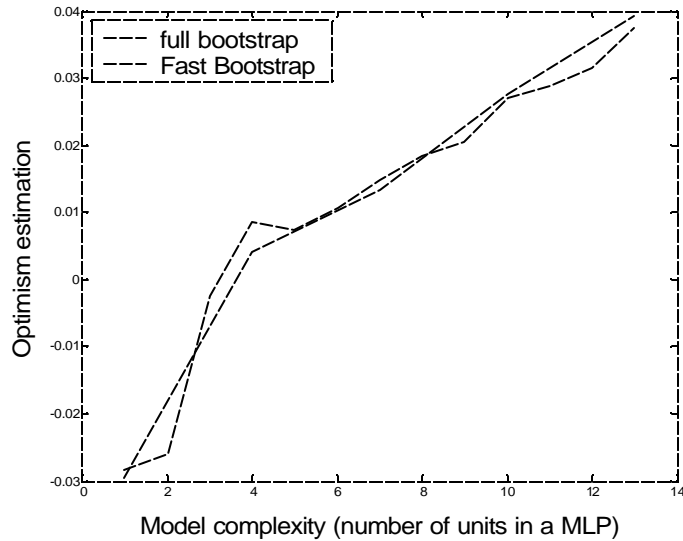


$$y_t = \sin(5x_t) + \sin(15x_t) + \sin(25x_t) + \mathbf{e}_t$$

- 200 sample
- distribution of  $\varepsilon_t$  : i.i.d. uniformly in  $[-0.5, 0.5]$

# MLP on artificial regression example

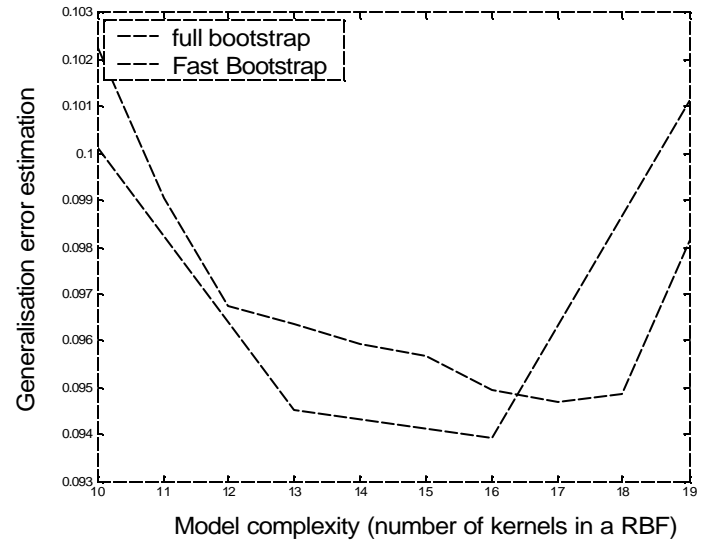
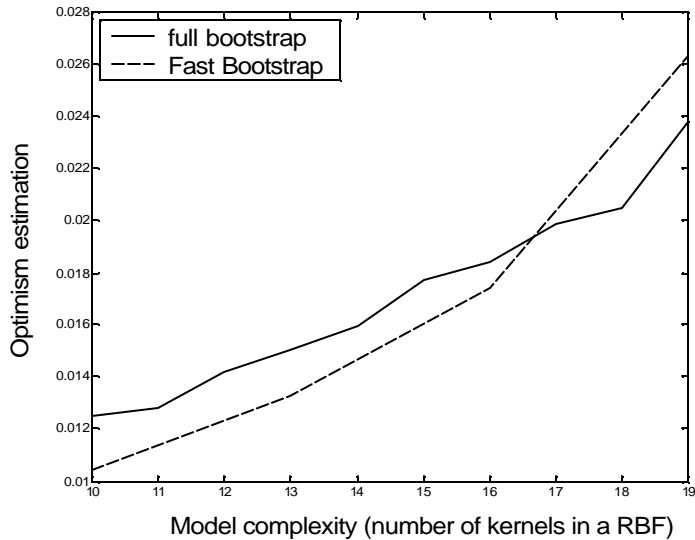
Experiments – artificial regression problem



<b>MLP</b>	<i>Number of hidden neurons</i>	<i>Bootstrap replications</i>	<i>Number of experiments</i>	<i>Gain</i>	$F_{1,2}$
<i>Bootstrap</i>	1-13 by steps of 1	100	1300		
<i>Fast Bootstrap</i>	1-13 by steps of 3	10	50	96.2%	3.25

# RBFN on artificial regression example

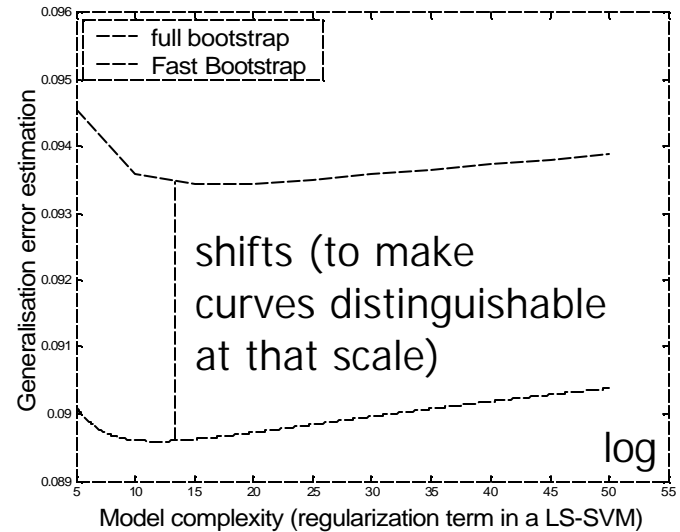
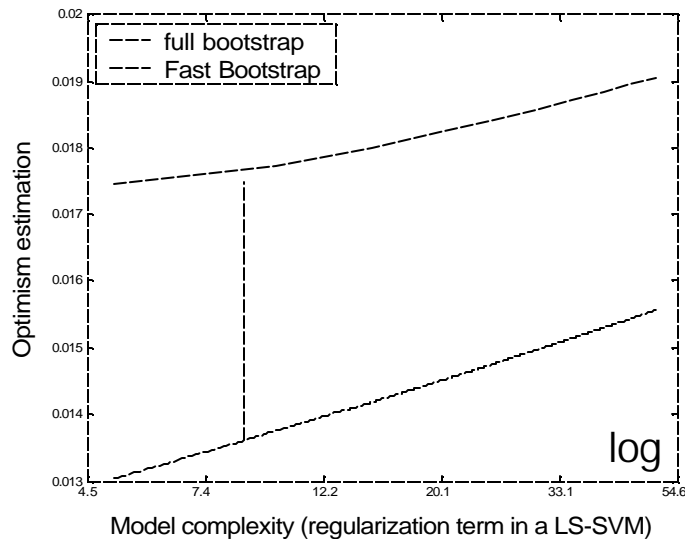
Experiments – artificial regression problem



<b>RBFN</b>	<i>Number of kernels</i>	<i>Bootstrap replications</i>	<i>Number of experiments</i>	<i>Gain</i>	$F_{1,1}$
<i>Bootstrap</i>	10-19 by steps of 1	100	1000		
<i>Fast Bootstrap</i>	10-19 by steps of 3	10	40	96%	16.31

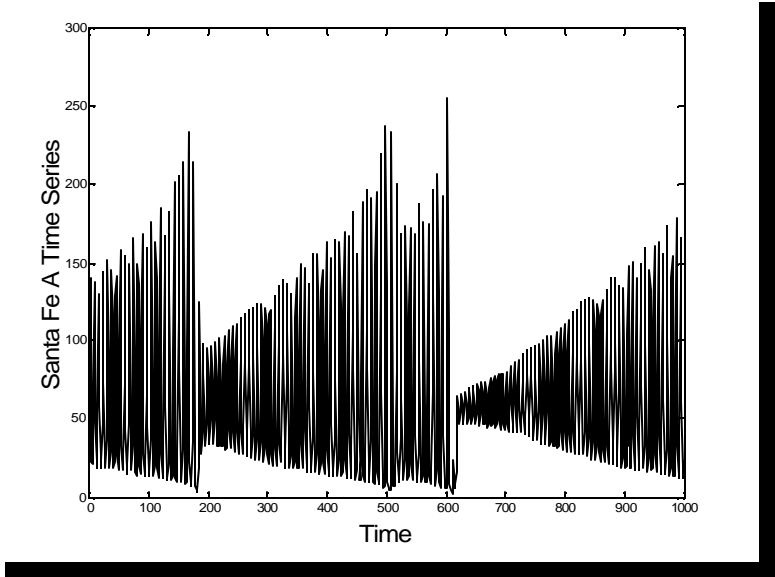
# LS-SVM on artificial regression example

Experiments – artificial regression problem



<b>LS-SVM</b>	<i>Regularization <math>g</math></i>	<i>Bootstrap replications</i>	<i>Number of experiments</i>	<i>Gain</i>	$F_{1,7}$
<i>Bootstrap</i>	5-50 by steps of 0.1	100	45100		
<i>Fast Bootstrap</i>	5-50 by steps of 5	10	100	99.8%	0

# Santa Fe A time series forecasting



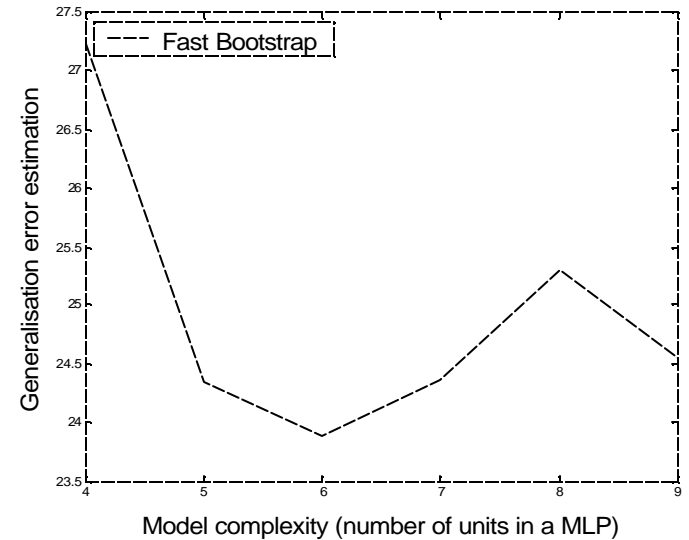
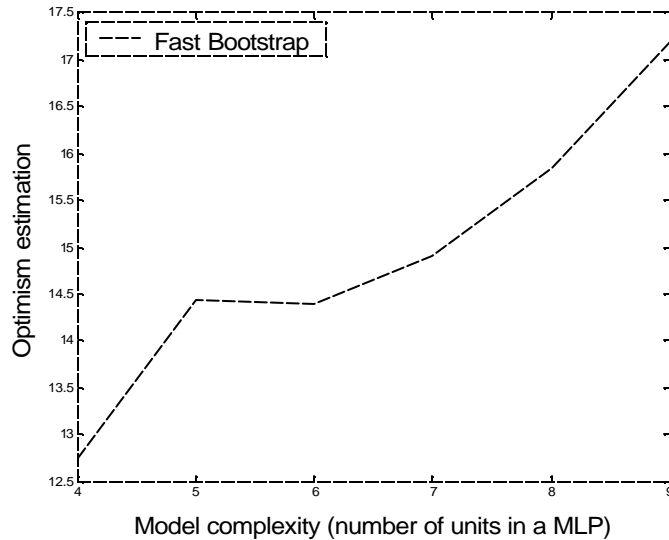
- Far-infrared laser in a chaotic state

$$\hat{y}_t = g \left( \underbrace{y_t, y_{t-1}, y_{t-2}, y_{t-3}, y_{t-4}, y_{t-5}}_{\mathbf{x}_t}, \mathbf{q}(q) \right)$$

MLP, RBFN, LS-SVM                      parameters                      model structure parameter



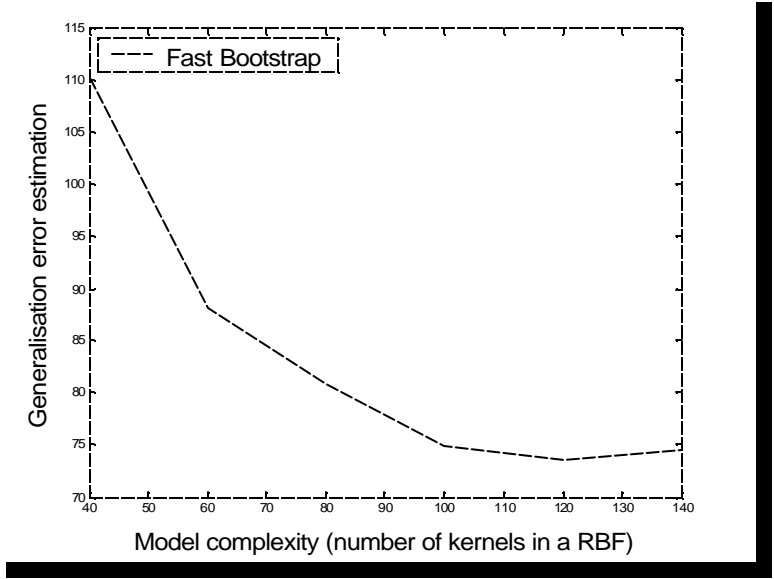
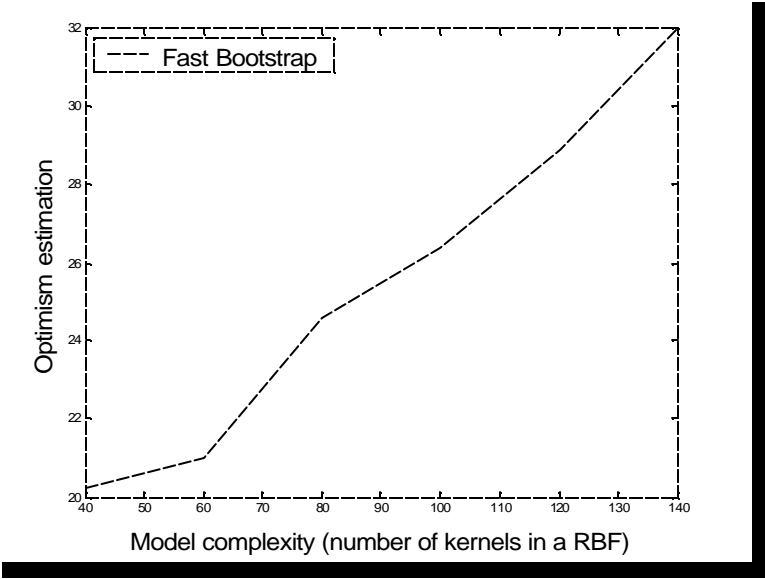
# MLP on Santa Fe A time series example



<b>MLP</b>	<i>Number of hidden neurons</i>	<i>Bootstrap replications</i>	<i>Number of experiments</i>	$F_{1,3}$
<i>Fast Bootstrap</i>	4-9 by steps of 1	10	60	0.2002

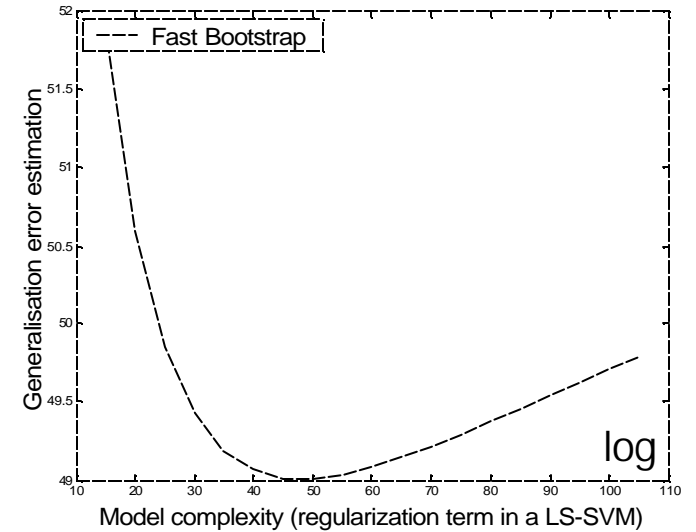
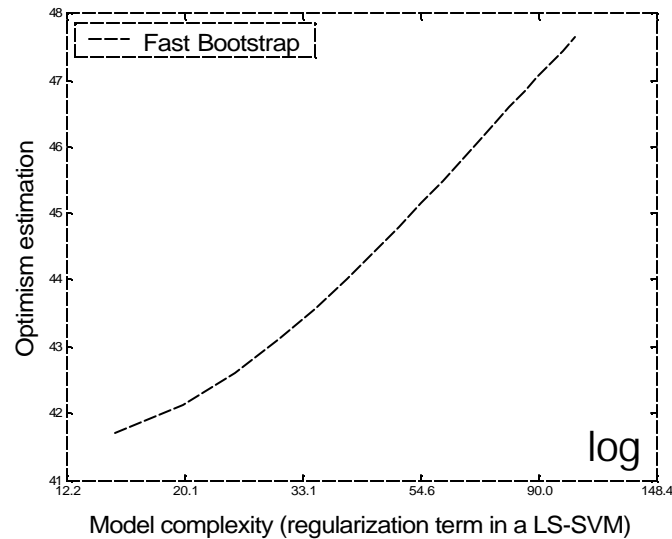
# RBFN on Santa Fe A time series example

Experiments – Santa Fe A forecasting



<b>RBFN</b>	<i>Number of kernels</i>	<i>Bootstrap replications</i>	<i>Number of experiments</i>	$F_{1,2}$
<i>Fast Bootstrap</i>	60-140 by steps of 20	20	100	1.6472

# LS-SVM on Santa Fe A time series example



<b>LS-SVM</b>	<i>Regularization <math>g</math></i>	<i>Bootstrap replications</i>	<i>Number of experiments</i>	$F_{1,16}$
<i>Fast Bootstrap</i>	15- 105 by steps of 5	10	190	0

# To conclude

- Bootstrap: efficient model selection method
- But: computationally intensive!
- Solution: Fast Bootstrap
  - uses smoothness hypothesis on *optimism*
  - hypothesis seems reasonable
  - makes the number of experiments decrease by 1-2 orders of magnitude!
- Further work:
  - theoretical insights about hypothesis
  - Moody's Generalized Prediction Error could explain the log in LS-SVM ?

# Thanks to...

- Amaury Lendasse and Geoffroy Simon

for

- the initial idea
- the experiments !