

# Model generation and selection using coherent Bayesian inference

Vadim Strijov

Visiting Professor at Laboratoire d'Informatique de Grenoble,  
Apprentissage : modeles et algorithmes

April 22<sup>nd</sup>, 2015

## Problem significance

To get an accurate and stable forecast we develop the methods of model selection from the set of admissible basic models.

## Our approach

Optimization of parameters for an arbitrary model is a non-trivial optimization problem. Our approach is to simplify the problem by considering sets of the successively generated stable models of given complexity.

## We solve a regression problem:

estimate the conditional expectation  $E(Y|\mathbf{x}) = f(\mathbf{w}_0, \mathbf{x})$ .

The sample:  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}$ ,  $i \in \mathcal{I} = \{1, \dots, m\}$ . The set  $\mathcal{G}$  is a set of parametric basic functions  $g(\mathbf{b}, \mathbf{x}')$ .

## Regression model

$$f = f(\mathbf{w}, \mathbf{x}) = g_1(\mathbf{b}_1, \mathbf{x}'_1) \circ \dots \circ g_r(\mathbf{b}_r, \mathbf{x}'_r)(\mathbf{x}),$$

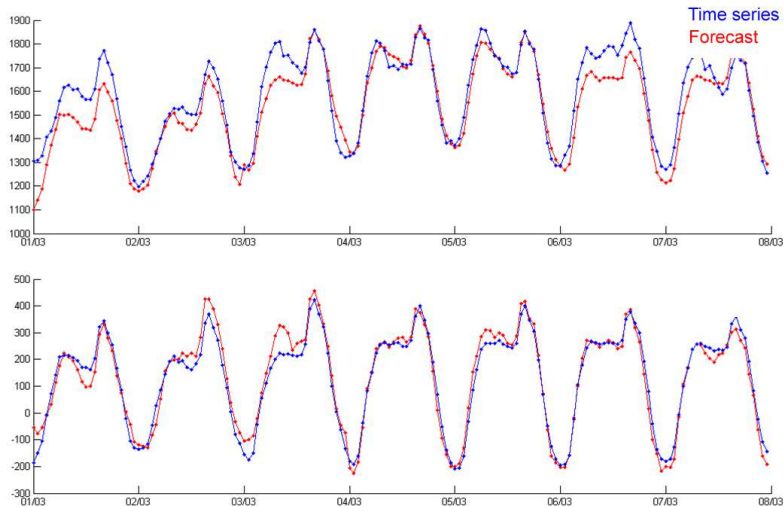
$$f : \mathbb{W} \times \mathbb{X} \rightarrow \mathbb{Y}, \quad \text{or elementwise: } f : (\mathbf{w}, \mathbf{x}) \mapsto y,$$

is chosen from the successively generated set  $\mathfrak{F}$ .

We find the regression function, the restriction of the model over the set of parameters

$$\hat{f}|_{\mathbb{W} \ni \mathbf{w} = \mathbf{w}_0} : \mathbb{X} \rightarrow \mathbb{Y}.$$

# Energy consumption one-week forecast, an example



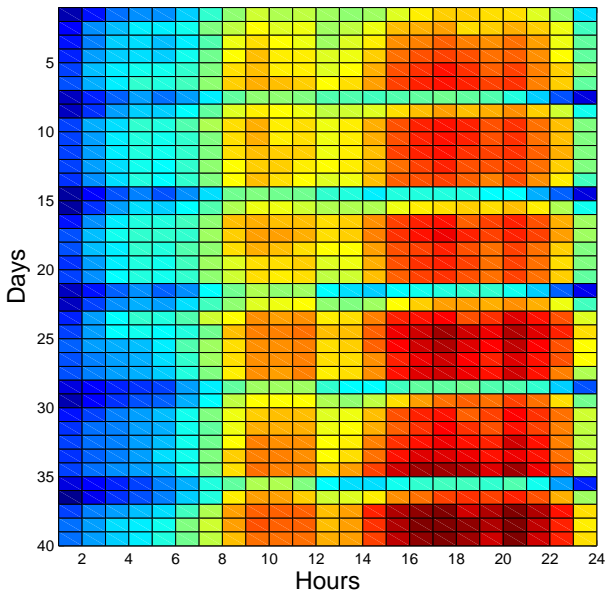
## The time series:

- energy price,
- consumption,
- daytime,
- temperature,
- humidity,
- wind force,
- holiday schedule.

## Periods:

- one year seasons  
(temperature, daytime),
- one week,
- one day (working day,  
week-end),
- a holiday,
- aperiodic events.

# The autoregressive matrix, five week-ends



$$\mathbf{X}^*_{(m+1) \times (n+1)} = \left( \begin{array}{c|ccc} S_T & S_{T-1} & \dots & S_{T-\kappa+1} \\ \hline S_{(m-1)\kappa} & S_{(m-1)\kappa-1} & \dots & S_{(m-2)\kappa+1} \\ \dots & \dots & \dots & \dots \\ S_{n\kappa} & S_{n\kappa-1} & \dots & S_{n(\kappa-1)+1} \\ \dots & \dots & \dots & \dots \\ S_\kappa & S_{\kappa-1} & \dots & S_1 \end{array} \right) .$$

In a nutshell,

$$\mathbf{X}^* = \left[ \begin{array}{c|c} S_T & \mathbf{x}_{m+1} \\ \hline \mathbf{y} & \mathbf{X} \end{array} \right] .$$

$\begin{matrix} 1 \times 1 & 1 \times n \\ m \times 1 & m \times n \end{matrix}$

In terms of linear regression:

$$\mathbf{y} = \mathbf{X}\mathbf{w},$$

$$y_{m+1} = S_T = \mathbf{w}^\top \mathbf{x}_{m+1}^\top .$$

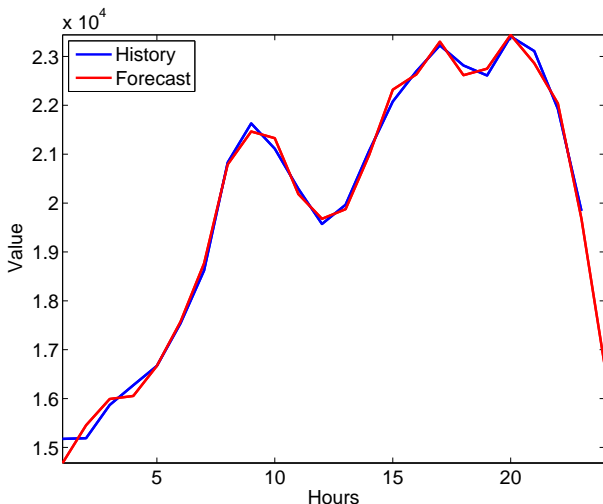
Introduce a set of the primitive functions  $G = \{g_1, \dots, g_r\}$ ,  
 for example  $g_1 = 1$ ,  $g_2 = \sqrt{x}$ ,  $g_3 = x$ ,  $g_4 = x\sqrt{x}$ , etc.

The generated set of features  $\mathbf{X} =$

$$\left( \begin{array}{ccc|ccc} g_1 \circ s_{T-1} & \dots & g_r \circ s_{T-1} & \dots & g_1 \circ s_{T-\kappa+1} & \dots & g_r \circ s_{T-\kappa+1} \\ g_1 \circ s_{(m-1)\kappa-1} & \dots & g_r \circ s_{(m-1)\kappa-1} & \dots & g_1 \circ s_{(m-2)\kappa+1} & \dots & g_r \circ s_{(m-2)\kappa+1} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ g_1 \circ s_{n\kappa-1} & \dots & g_r \circ s_{n\kappa-1} & \dots & g_1 \circ s_{n(\kappa-1)+1} & \dots & g_r \circ s_{n(\kappa-1)+1} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ g_1 \circ s_{\kappa-1} & \dots & g_r \circ s_{\kappa-1} & \dots & g_1 \circ s_1 & \dots & g_r \circ s_1 \end{array} \right).$$



## The one-day forecast (an example)



The function  $y = f(\mathbf{x}, \mathbf{w})$  could be a linear model, neural network, deep NN, SVN, ...

Assume we have hourly data on price/consumption for three years.

Then the matrix  $\mathbf{X}^*$  is  
 $(m+1) \times (n+1)$

$156 \times 168$ , in details:  $52w \cdot 3y \times 24h \cdot 7d$ ;

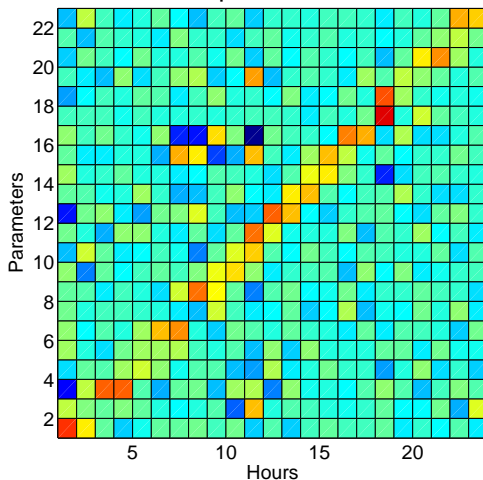
- for 6 time series the matrix  $\mathbf{X}$  is  $156 \times 1008$ ,
- for 4 primitive functions it is  $156 \times 4032$ ,

$$m \ll n.$$

The autoregressive matrix could be considered as *ill-conditioned* and *multi-correlated*. The model selection procedure is required.

# How many parameters must be used to forecast?

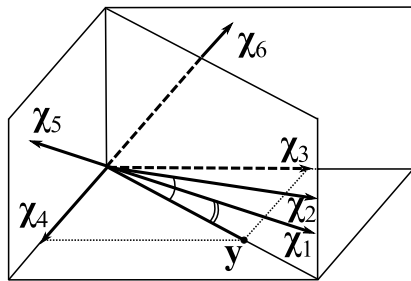
The color shows the value of a parameter for each hour.



Estimate parameters  $\mathbf{w}(\tau) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ , then calculate the sample  $s(\tau) = \mathbf{w}^T(\tau) \mathbf{x}_{m+1}$  for each  $\tau$  of the next ( $m+1$ -th) period.

## Selection of a stable set of features of restricted size

The sample contains multicollinear  $\chi_1, \chi_2$  and noisy  $\chi_5, \chi_6$  features, columns of the design matrix  $\mathbf{X}$ . We want to select two features from six.

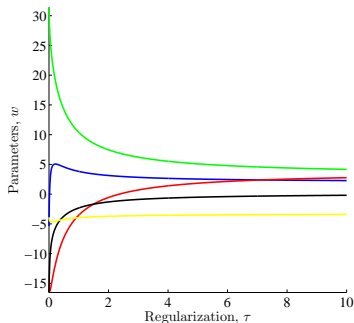


### Stability and accuracy for a fixed complexity

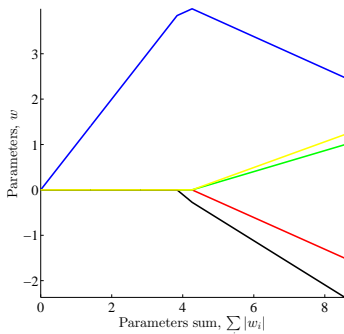
The solution:  $\chi_3, \chi_4$  is an orthogonal set of features minimizing the error function.

Algorithms: GMDH, Stepwise, Ridge, Lasso, Stagewise, FOS, LARS, Genetics, ...

Vector-function  $\mathbf{f} = \mathbf{f}(\mathbf{w}, \mathbf{X}) = [f(\mathbf{w}, \mathbf{x}_1), \dots, f(\mathbf{w}, \mathbf{x}_m)]^T \in \mathbb{Y}^m$ .

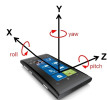


$$S(\mathbf{w}) = \|\mathbf{f}(\mathbf{w}, \mathbf{X}) - \mathbf{y}\|^2 + \gamma^2 \|\mathbf{w}\|^2$$

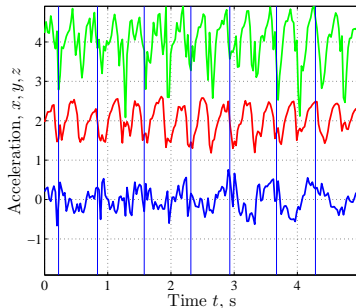
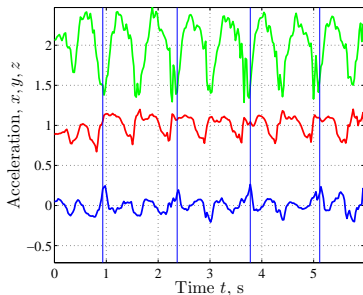


$$S(\mathbf{w}) = \|\mathbf{f}(\mathbf{w}, \mathbf{X}) - \mathbf{y}\|^2, \\ T(\mathbf{w}) \leq \tau$$

# Classification of accelerometric time series



Examples of accelerometric time series for slow walking and jogging:



3-dimensional time series of acceleration projections to spatial axis

$$\mathbf{x} = \{acc_x(t); acc_y(t); acc_z(t)\}_{t=1}^n \mapsto \mathbf{y} \in \mathbb{R}^S.$$

Class labels  $y_i$  correspond to one of  $S = 6$  types of activity: Jogging, Walking, Upstairs, Downstairs, Sitting, Standing.

Construct a classifier

$$\mathbf{f} = \mathbf{a}(\mathbf{h}_N(\dots \mathbf{h}_1(\mathbf{x}))),$$

where  $\mathbf{h}_k$  are autoencoding blocks of the form

$$\mathbf{h}_k(\mathbf{x}) = \sigma(\mathbf{W}_k \mathbf{x} + \mathbf{b}_k),$$

and  $\mathbf{a}$  is multinomial logistic regression classifier

$$\mathbf{a}(\mathbf{x}) = \mathbf{W}_2^T \tanh(\mathbf{W}_1^T \mathbf{x}).$$

Vectorize matrices  $\mathbf{W}_1 \in \mathbb{R}^{n \times N_h}$ ,  $\mathbf{W}_2 \in \mathbb{R}^{N_h \times S}$  of parameters of each layer to obtain vector of model parameters

$$\mathbf{w} = \text{vec}(\mathbf{W}_1^T | \mathbf{W}_2^T) \in \mathbb{R}^k.$$

Here number  $N_h$  of neurons in the hidden layer — the *structure parameter* of the model — is fixed.

## Model structure

Parameter  $w_j$  of model  $\mathbf{f}$  is called *active*, if  $w_j \neq 0$ .

The set of active indices  $\mathcal{A} = \{j : w_j \neq 0\} \subseteq \mathcal{J}$  is called *structure*  $\mathcal{A}$  of model  $\mathbf{f}$ .

Each structure  $\mathcal{A} \subseteq \mathcal{J}$  defines a model  $\mathbf{f}_{\mathcal{A}}$

$$\mathbf{f}_{\mathcal{A}} : \hat{\mathbf{w}}_{\mathcal{A}} \in \mathbb{R}^k,$$

where  $\hat{\mathbf{w}}_{\mathcal{A}} \in \mathbb{R}^k$  is an optimal parameter vector of  $\mathbf{f}_{\mathcal{A}}$  which minimizes error function

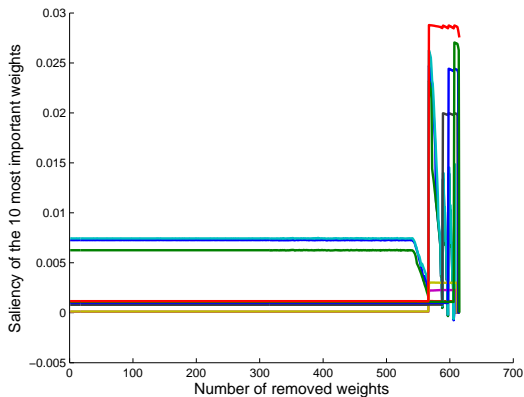
$$S(\mathbf{w}|\mathcal{L}) = - \sum_{i \in \mathcal{K}} \sum_{\xi=1}^S t_{i\xi} \ln(p_{\xi}(\mathbf{x}_i, \mathbf{w})), \quad \mathbf{p}(\mathbf{x}) = \frac{\exp(\mathbf{a}(\mathbf{x}))}{\sum_j \exp(a_j(\mathbf{x}))},$$

computed at learning subset of  $\mathcal{D}$ , defined by set of indices  $\mathcal{L}$ .

We chose optimal model  $\hat{\mathbf{f}}_{\mathcal{A}}$  from a set  $\mathfrak{F}$  of *admissible models*:

$$\mathfrak{F} = \bigcup_{\mathcal{A} \subseteq \mathcal{J}} \{\mathbf{f}_{\mathcal{A}}\}.$$





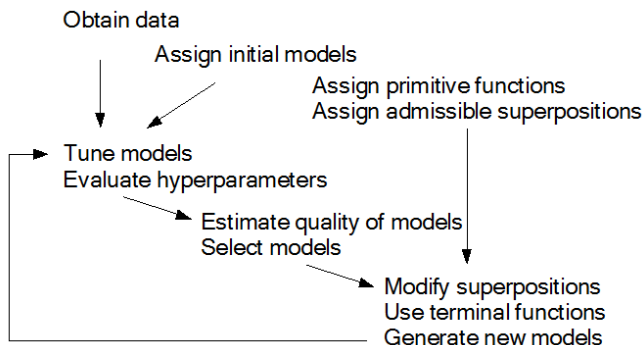
Dependency of a saliency  $L_j = \frac{w_j^2}{2\mathbf{H}_{jj}^{-1}}$  from a number of removed parameters.

## The basic goal of research

To develop a methodology for selection of successively generated models for regression and classification problems.

## The approach

- a) we successively generate a set of regression models,
- b) we investigate space of model parameters,
- c) we compare model elements by estimating a covariance matrix and its parameters,
- d) we choose the model according to the MDL principle.



- |  |                            |
|--|----------------------------|
| ① Stepwise method of model selection     | M. A. Efroimson, 1960.     |
| ② Regularization for the inverse problem | A. N. Tikhonov, 1963.      |
| ③ Group method of data handling          | A. G. Ivakhnenko, 1971.    |
| ④ Optimal brain damage                   | Y. LeCun, 1999.            |
| ⑤ Model hyperparameters estimation       | Y. Nabney, 2004.           |
| ⑥ Symbol regression                      | I. Zelinka, D. Koza, 2004. |
| ⑦ Least angle regression                 | B. Efron, T. Hastie, 2002. |
| ⑧ Entropy methods for MDL                | P. Gruenwald, 2006.        |
| ⑨ MDL principle in regression            | J. Rissanen, 2009.         |
| ⑩ Learning of Bayesian network structure | T. Jaakkola, 2012.         |

## Data and parameters generation assumption

Distribution of the dependent random variable  $\mathbf{y} = \boldsymbol{\mu}^{-1}(\mathbf{X}, \mathbf{w})$  belongs to the *exponential family*

$$p(\mathbf{y}|\boldsymbol{\eta}) = h(\mathbf{y})g(\boldsymbol{\eta}) \exp(\boldsymbol{\eta}^T \mathbf{u}(\mathbf{y})) \quad (\text{ED})$$

with a vector  $\boldsymbol{\eta}$  of parameters. The special cases: normal (ND) and binomial (BD) distributions:

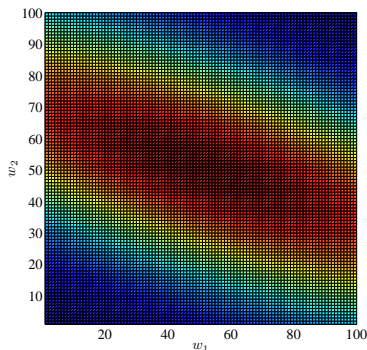
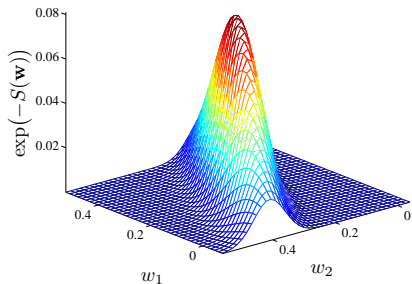
$$p(\mathcal{D}|\mathbf{B}, \mathbf{w}, \mathbf{f}) = (2\pi)^{-\frac{m}{2}} |\mathbf{B}^{-1}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{y} - \mathbf{f})^T \mathbf{B}(\mathbf{y} - \mathbf{f})\right), \quad (\text{ND})$$

$$p'(\mathcal{D}|\mathbf{w}, \mathbf{f}) = \prod_{i \in \mathcal{I}} f_i^{y_i} (1 - f_i)^{1 - y_i}. \quad (\text{BD})$$

### Distributions $p(\mathcal{D}|\mathbf{B}, \mathbf{w}, \mathbf{f})$ and $p(\mathbf{w}|\mathbf{A}, \mathbf{f})$ : different cases

Dependent variable $\mathbf{y}$	Model parameters $\mathbf{w}$
$\mathbf{y} \sim \mathcal{N}(\mathbf{f}, \sigma_{\mathbf{y}}^2 \mathbf{I}) \stackrel{\text{def}}{=} \mathcal{N}(\mathbf{f}, \beta^{-1} \mathbf{I})$	$\mathbf{w} \sim \mathcal{N}(\mathbf{w}_0, \sigma_{\mathbf{w}}^2 \mathbf{I}) \stackrel{\text{def}}{=} \mathcal{N}(\mathbf{0}, \alpha^{-1} \mathbf{I})$
$\mathbf{y} \sim \mathcal{N}(\mathbf{f}, \text{diag}^{-1}(\beta_1, \dots, \beta_m) \mathbf{I})$	$\mathbf{w} \sim \mathcal{N}(\mathbf{w}_0, \text{diag}^{-1}(\alpha_1, \dots, \alpha_n) \mathbf{I})$
$\mathbf{y} \sim \mathcal{N}(\mathbf{f}, \mathbf{B}^{-1})$	$\mathbf{w} \sim \mathcal{N}(\mathbf{w}_0, \mathbf{A}^{-1})$

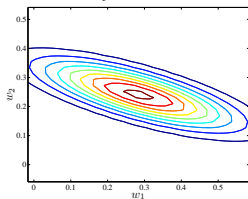
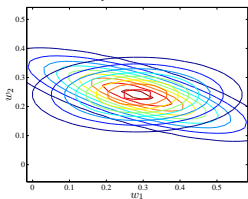
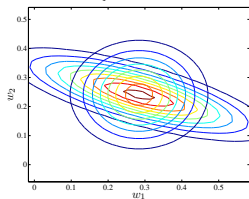
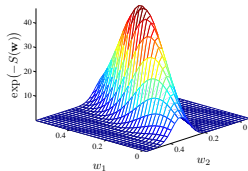
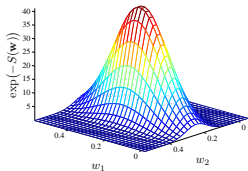
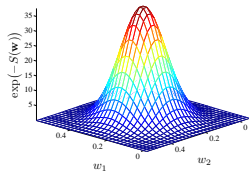
There given a sample  $\{\mathbf{w}_1, \dots, \mathbf{w}_K\}$  of realizations of the m.r.v.  $\mathbf{w}$  and an error function  $S(\mathbf{w}|\mathcal{D}, \mathbf{f})$ . Consider the set of points  $\{s_k = \exp(-S(\mathbf{w}_k|\mathcal{D}, \mathbf{f})) | k = 1, \dots, K\}$ .



x- and y-axis: parameters  $\mathbf{w}$ , z-axis:  $\exp(-S(\mathbf{w}))$ .

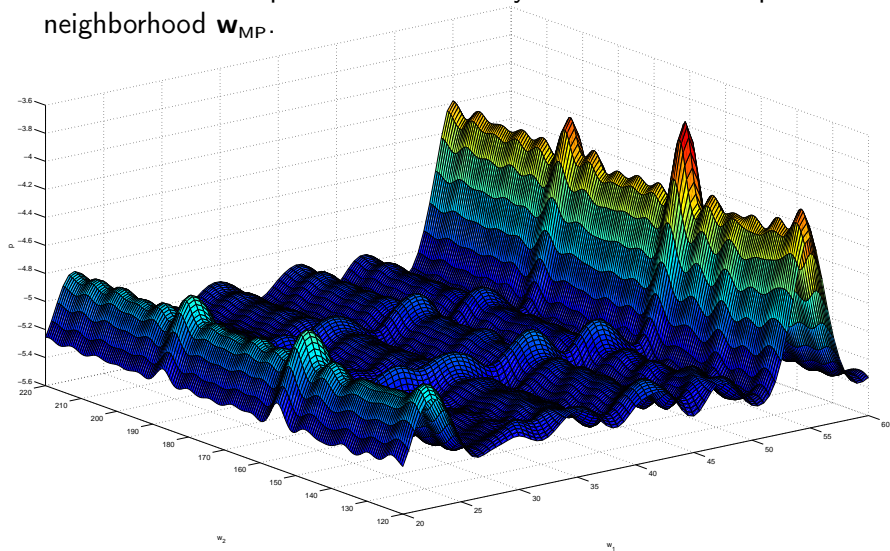
Approximate the set of points  $\{s_k\}$  by a function  $p(\mathbf{w}|\mathbf{A})$  (ND), considering assumptions about the covariance matrix  $\mathbf{A}^{-1}$  type:

$$\mathbf{A} = \alpha \mathbf{I}, \quad \alpha \geq 0; \quad \mathbf{A} = \text{diag}(\alpha_1, \dots, \alpha_n); \quad \mathbf{A}, \quad \mathbf{w}^T \mathbf{A} \mathbf{w} \geq 0.$$



x- and y-axis: parameters  $\mathbf{w}$ , z-axis:  $\exp(-S(\mathbf{w}))$ .

Distribution of parameters  $\mathbf{w}$  beyond the most probable neighborhood  $\mathbf{w}_{MP}$ .





## Posterior parameter distribution

for the given sample  $\mathcal{D}$ , model  $\mathbf{f} = \mathbf{f}(\mathbf{w}, \mathbf{X})$  and matrices  $\mathbf{A}, \mathbf{B}$ :

$$p(\mathbf{w}|\mathcal{D}, \mathbf{A}, \mathbf{B}, \mathbf{f}) = \frac{p(\mathcal{D}|\mathbf{w}, \mathbf{B}, \mathbf{f})p(\mathbf{w}|\mathbf{A}, \mathbf{f})}{p(\mathcal{D}|\mathbf{A}, \mathbf{B}, \mathbf{f})}.$$

The elements of this expression and the corresponding parameters:

$p(\mathbf{w}|\mathcal{D}, \mathbf{A}, \mathbf{B}, \mathbf{f})$  — posterior parameter distribution,

$\mathbf{w}_{\text{MP}} = \arg \max p(\mathbf{w}|\mathcal{D}, \mathbf{A}, \mathbf{B}, \mathbf{f})$  — most probable parameters,

$p(\mathcal{D}|\mathbf{w}, \mathbf{B}, \mathbf{f})$  — data likelihood,

$\mathbf{w}_{\text{ML}} = \arg \max p(\mathcal{D}|\mathbf{w}, \mathbf{B}, \mathbf{f})$  — most plausible parameters,

$p(\mathbf{w}|\mathbf{A}, \mathbf{f})$  — prior distribution,

$p(\mathcal{D}|\mathbf{A}, \mathbf{B}, \mathbf{f})$  — model likelihood.

For a set of models  $\mathfrak{F} = \{f_1, \dots, f_K\}$  to approximate  $\mathcal{D}$

$$p(f_k|D) = \frac{p(D|f_k)p(f_k)}{\sum_{q=1}^K p(D|f_q)p(f_q)}$$

$p(f_k)$  — prior probability,

$p(D|f_k)$  — model evidence,

$p(f_k|D)$  — posterior probability.

Select the most evident model by comparison

$$\frac{p(f_k|D)}{p(f_q|D)} = \frac{p(D|f_k)p(f_k)}{p(D|f_q)p(f_q)}$$

since the denominator does not depend on the model.

Assuming equal prior probability of the models from the set  $\mathfrak{F}$ ,

$$p(f_k) = p(f_q)$$

maximize the model evidence.

## Error function of the general form

Writing the error function  $S(\mathbf{w})$  in the following form,

$$S(\mathbf{w}) = -\ln p(\mathcal{D}|\mathbf{w}, \mathbf{B}, \mathbf{f})p(\mathbf{w}|\mathbf{A}, \mathbf{f}) = E_{\mathbf{w}} + E_{\mathcal{D}},$$

we obtain the following posterior distribution:

$$p(\mathbf{w}|\mathcal{D}, \mathbf{A}, \mathbf{B}, \mathbf{f}) \propto \frac{\exp(-S(\mathbf{w}))}{Z_S}.$$

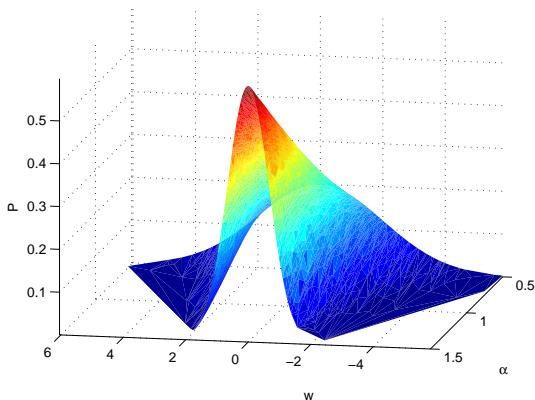
**The case of normal distribution for the dependent variable (ND)**

$$S(\mathbf{w}) = \frac{1}{2}(\mathbf{w} - \mathbf{w}_0)^T \mathbf{A}(\mathbf{w} - \mathbf{w}_0) + \frac{1}{2}(\mathbf{y} - \mathbf{f})^T \mathbf{B}(\mathbf{y} - \mathbf{f}).$$

**The case of binomial distribution for the dependent variable (BD)**

$$S(\mathbf{w}) = E_{\mathbf{w}} + \sum_{i \in \mathcal{I}} (y_i \ln f_i + (1 - y_i) \ln(1 - f_i)).$$

## Posterior parameter distribution with $\mathbf{A} = \alpha \mathbf{I}$



x-axis:  $w$  is a model parameter.

y-axis:  $\alpha$  is an inverted covariance,

z-axis:  $p(\mathbf{w} | \mathcal{D}, \mathbf{A}, \mathbf{B}, \mathbf{f})$  is a distribution of parameters.

There is given a sample  $\mathcal{D}$ , a set of models  $\mathfrak{F} = \{f_k\}$ ,  $k \in \mathcal{K}$  and prior probabilities  $p(f_k)$ .

The problem is to find the most plausible model  $f_k$ :

$$\hat{k} = \arg \max_{k \in \mathcal{K}} p(f_k | \mathcal{D}) =$$
$$\arg \max_{k \in \mathcal{K}} \int_{\mathbf{w} \in \mathbb{W}_k} p(\mathcal{D} | \mathbf{w}, \mathbf{B}_k, \mathbf{f}_k) p(\mathbf{w} | \mathbf{A}_k, \mathbf{f}_k) d\mathbf{w}.$$

Posterior model probability

$$p(f_k | \mathcal{D}) = \frac{1}{p(\mathcal{D})} p(\mathcal{D} | f_k) p(f_k),$$

where the function  $p(\mathcal{D} | f_k)$  of the sample  $\mathcal{D}$ , with a fixed model  $f_k$  is a model likelihood. The normalized coefficient doesn't depend on the model.

There is given a sample  $\mathcal{D}$ , a model  $\mathbf{f} = \mathbf{f}(\mathbf{w}, \mathbf{x})$ , a data generation assumption, and an error function

$$S(\mathbf{w}|\mathcal{D}, \hat{\mathbf{A}}, \hat{\mathbf{B}}, \mathbf{f}) = -\ln(p(\mathcal{D}|\mathbf{w}, \mathbf{B}, \mathbf{f})p(\mathbf{w}|\mathbf{A}, \mathbf{f})).$$

The goal is to find parameters  $\mathbf{w}_{\text{MP}}$  of the model  $\mathbf{f}$

$$\mathbf{w}_{\text{MP}} = \arg \min_{\mathbf{w} \in \mathbb{W}} S(\mathbf{w}|\mathcal{D}, \hat{\mathbf{A}}, \hat{\mathbf{B}}, \mathbf{f}).$$

The covariance matrix estimation

$$(\hat{\mathbf{A}}, \hat{\mathbf{B}}) = \arg \max_{\mathbf{A} \in \mathbb{R}^{n^2}, \mathbf{B} \in \mathbb{R}^{m^2}} \int_{\mathbf{w} \in \mathbb{W}} p(\mathcal{D}|\mathbf{w}, \mathbf{B}, \mathbf{f})p(\mathbf{w}|\mathbf{A}, \mathbf{f})d\mathbf{w}.$$

## Theorem (2014)

The linear model likelihood for the data generation assumption (ND) has the form

$$p(\mathcal{D}|\mathbf{A}, \mathbf{B}) = \frac{|\mathbf{B}|^{\frac{1}{2}}|\mathbf{A}|^{\frac{1}{2}}}{(2\pi)^{\frac{m}{2}}|\mathbf{K}|^{\frac{1}{2}}} \exp\left(\frac{1}{2}\mathbf{y}^T(\mathbf{C}^T\mathbf{K}\mathbf{C} - \mathbf{B})\mathbf{y}\right),$$

and its logarithm has the form  $\ln p(\mathcal{D}|\mathbf{A}, \mathbf{B}) =$

$$= -\frac{1}{2}(\ln |\mathbf{K}| + m \ln 2\pi - \ln |\mathbf{B}| - \ln |\mathbf{A}| - \mathbf{y}^T(\mathbf{C}^T\mathbf{K}\mathbf{C} - \mathbf{B})\mathbf{y}).$$

Here

$$\mathbf{K} = \mathbf{X}^T\mathbf{B}\mathbf{X} + \mathbf{A}, \quad \mathbf{C} = \mathbf{K}^{-1}\mathbf{X}^T\mathbf{B}.$$

## Theorem (2013)

For the data generation assumption(ND) with the fixed covariance matrices  $\mathbf{A}^{-1}$ ,  $\mathbf{B}^{-1}$  the iterative algorithm of parameters estimation,

$$\Delta \mathbf{w}_{k+1} = (\mathbf{J}^T \mathbf{J})^{-1} \left( \mathbf{J}^T (\mathbf{y} - \mathbf{f}(\mathbf{w}, \mathbf{X})) - \frac{1}{\beta} \mathbf{A}^{-1} \mathbf{w}_k \right),$$

finds a minimum of the error function of general form  $S(\mathbf{w}|\mathcal{D}, \mathbf{A}, \mathbf{B}, \mathbf{f})$  with the convergence of vectors sequence  $\mathbf{w}_k$ .

## Remark

The iterative algorithm  $\mathbf{w}_{k+1} = \Delta \mathbf{w}_{k+1} + \mathbf{w}_k$  requires the initial value  $\mathbf{w}_0$ . The sequence  $\|\mathbf{w}_{k+1} - \mathbf{w}_k\|^2$  monotonically decreases due to increase of the step  $k$ .



## Theorem (2013)

For the data generation assumption (BD) with the fixed covariance matrices  $\mathbf{A}^{-1}$ ,  $\mathbf{B}^{-1}$  the iterative algorithm of parameters estimation for the generalized linear model,

$$\Delta \mathbf{w}_{k+1} = (\mathbf{X}^T \mathbf{B} \mathbf{X} + \mathbf{A})^{-1} \mathbf{X}^T \mathbf{B}^T \mathbf{y} - \mathbf{w}_k, \quad \text{variant:}$$

$$\Delta \mathbf{w}_{k+1} = (\mathbf{X}^T \mathbf{B} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{B} (\mathbf{X} \mathbf{w}_k - \mathbf{B}^{-1} (\mathbf{f} - \mathbf{y})) + \frac{1}{2} \mathbf{w}_k^T \mathbf{A} \mathbf{w}_k,$$

finds a local minimum of the error function of general form with the convergence of vectors sequence  $\mathbf{w}_k$ .

Let the vector of parameters  $\mathbf{w}_0 = [w_{1(0)}, \dots, w_{n(0)}]^T$  be fixed.

## Theorem (2013)

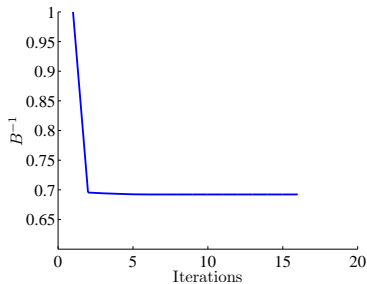
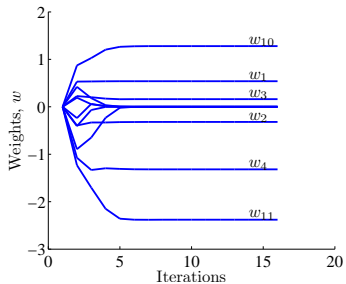
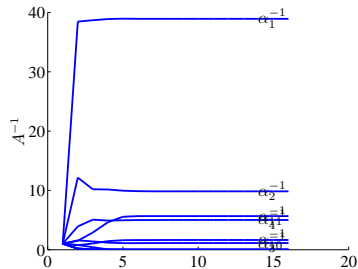
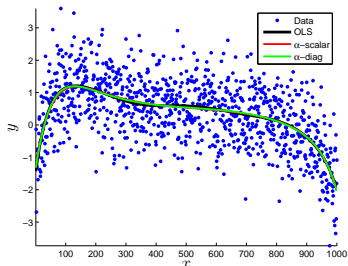
In a neighborhood of the parameters  $\mathbf{w}_0$  the covariance matrix estimations  $\mathbf{A}^{-1}$ ,  $\mathbf{B}^{-1}$  for the data generation assumption (ND) has the form

$$\alpha_i = \frac{1}{2} \lambda_i \left( \sqrt{1 + \frac{4}{(w_i - w_{i(0)})^2 \lambda_i}} - 1 \right), \text{ where } \lambda_i = \beta \mathbf{diag}(h_i),$$

$$\beta = \frac{m - \gamma}{2(\mathbf{f} - \mathbf{y})^T \mathbf{B}'(\mathbf{f} - \mathbf{y})}, \text{ где } \gamma = \sum_{j=1}^W \frac{\lambda_j}{\lambda_j + \alpha_j}.$$

The sequences  $\|\mathbf{A}_{k+1} - \mathbf{A}_k\|^2$  and  $\|\beta_{k+1} - \beta_k\|^2$  monotonically decrease due to increase of the step  $k$ .

# Estimation of parameters and covariance matrices



There is given a set  $\mathfrak{G} = \{\text{id}, g_1, \dots, g_l \mid g = g(\mathbf{b}, \mathbf{x}')\}$ , that is, there are given

- 1) the function  $g : (\mathbf{b}, \mathbf{x}') \mapsto \mathbf{x}''$ ,
- 2) its parameters  $\mathbf{b}$ ,
- 3) arity  $v(g)$  of the function  $g$  and an order of arguments,
- 4) a domain  $\text{dom}(g)$  and a codomain  $\text{cod}(g)$ .

Consider a model  $f(\mathbf{w}, \mathbf{x})$  given by a superposition

$$f(\mathbf{w}, \mathbf{x}) = (g_{i(1)} \circ \dots \circ g_{i(K)})(\mathbf{x}), \text{ где } \mathbf{w} = [\mathbf{b}_{i(1)}^T, \dots, \mathbf{b}_{i(K)}^T]^T.$$

### An admissible superposition $f$

is a superposition such that

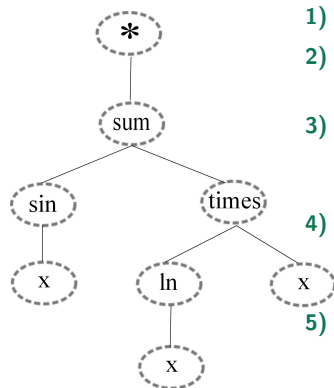
$$\text{cod}(g_{i(k+1)}) \subseteq \text{dom}(g_{i(k)}), \text{ для всех } k = 1, \dots, K - 1.$$

## To generate the models we use

- 1) the set  $\text{dom}(\mathbf{x})$ ,
- 2) the set of basic functions  $\mathfrak{G} = \{\text{id}, g\}$ ,  $g : \mathbf{x} \mapsto \mathbf{x}'$ ,
- 3) the set  $\text{Gen}$  of rules for superposition generation,
- 4) the set  $\text{Rem}$  of rules for isomorphic superpositions simplification and estimation.

We propose the following basic methods for the superpositions generation:

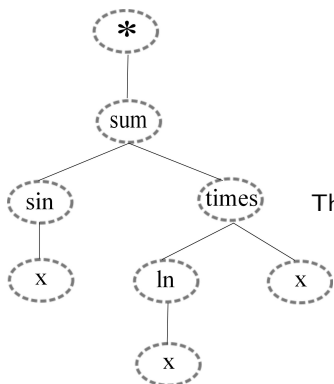
- inductive generation,
- structure learning,
- direct search.



$$f = \sin(x) + (\ln x)x$$

- 1) the root  $*$  of the tree  $\Gamma_f$  has the single vertex,
- 2) other vertices  $V_i$  correspond to the functions  $g_r \in \mathfrak{G}: V_i \mapsto g_r$ ,
- 3) the number of children  $V_j$  of the vertex  $V_i$  equals to an arity of the corresponding function  $g_r$ :  $\text{val}(V_j) = v(g_{r(i)})$ ,
- 4) the domain of the function  $g_{r(i)}$  of a child  $V_j$  contains the codomain of the function  $g_{r(j)}$  of the parent  $V_i$ :  $\text{dom}(g_{r(i)}) \supseteq \text{cod}(g_{r(j)})$ ,
- 5) an order of vertices traversal with a parent vertex  $V_i$  corresponds to the order of arguments of the corresponding function  $g_{r(i)}$ ,
- 6) the leaves  $\Gamma_f$  correspond to the independent variables, elements of the vector  $x$ .

The link matrix  $Z_f$  for the tree  $\Gamma_f$



$$f = \sin(x) + (\ln x)x$$

	sum	times	ln	sin	x
*	1	0	0	0	0
sum	0	1	1	0	0
times	0	0	0	1	1
ln	0	0	0	0	1
sin	0	0	0	0	1

The link probability matrix  $P_f$  for the tree  $\Gamma_f$

	sum	times	ln	sin	x
*	0.7	0.1	0.1	0.1	0.2
sum	0.2	0.7	0.8	0.1	0.2
times	0.1	0.3	0	0.8	0.8
ln	0.2	0.1	0.3	0.1	0.9
sin	0.1	0.2	0.1	0	0.8

$\mathfrak{J}$  is a set of matrices corresponding to the superpositions from  $\mathfrak{F}$ .

There is given a sample  $\mathcal{D} = \{(\mathbf{D}_k, f_k)\}$  where the element  $\mathbf{D}_k = \begin{pmatrix} \mathbf{X} & \mathbf{y} \\ m \times n & m \times 1 \end{pmatrix}$ , there given  $\mathfrak{S}$  and  $\mathfrak{F} = \{f_s \mid \mathbf{f}_s : (\hat{\mathbf{w}}_k, \mathbf{X}) \mapsto \mathbf{y}, s \in \mathbb{N}\}$ .

### The goal

to find an algorithm  $a : \mathbf{D}_k \mapsto f_s$  following the condition

$$\mathbf{z}_{f_s} = \arg \max_{\mathbf{z} \in \mathfrak{F}} \sum_{i,j} P_{ij} \times Z_{i,j}.$$

The index  $\hat{s}$ , что  $f_{\hat{s}}$  provides a minimum for the error function  $S$ :

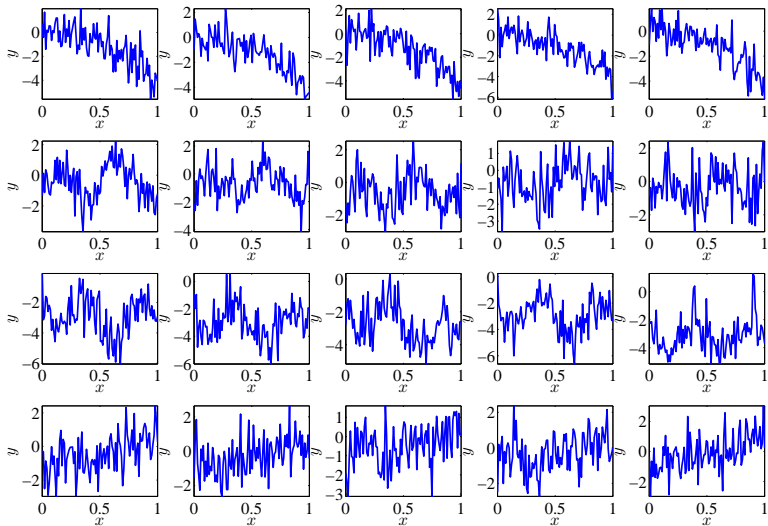
$$\hat{s} = \arg \min_{s \in \{1, \dots, |\mathfrak{F}|\}} S(f_s \mid \hat{\mathbf{w}}_k, \mathbf{D}_k),$$

where  $\hat{\mathbf{w}}_k$  is an optimal vector of parameters  $f_s$  for each  $f_s \in \mathfrak{F}$  with the fixed  $\mathbf{D}_k$ :

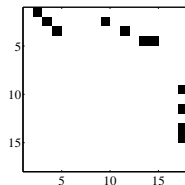
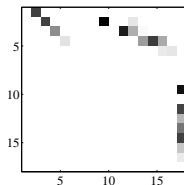
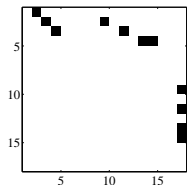
$$\hat{\mathbf{w}}_k = \arg \min_{\mathbf{w} \in \mathbb{W}_s} S(\mathbf{w} \mid f_s, \mathbf{D}_k).$$



# An example of the time series sample for physical activity monitoring



# Initial and forecasted superposition



$$f = w_1 \cos(w_2 x + w_3) + w_4 x + w_5 \ln(w_6 x + w_7) + w_8,$$

$$f = \cos(x) + x + \ln(x), \quad \mathbf{w} = [1, 1, 0, 1, 1, 1, 0, 0]^T.$$

The set  $\mathcal{A}$  uniquely defines a model  $f_{\mathcal{A}} \in \mathfrak{F}$ .

## The successive modification procedure

**Add:** to add an index  $j$  to the set  $\mathcal{A}_k = \mathcal{A}_{k-1} \cup \{j\}$ , that corresponds to the maximum value of the model likelihood

$$\hat{j} = \arg \max_{j \in \mathcal{J} \setminus \mathcal{A}_k} p(f_{\mathcal{A}_k} | \mathbf{w}_{\text{MP}}, \mathbf{A}, \mathbf{B}, \mathcal{D}).$$

**Del:** to remove an index  $j$  from the set  $\mathcal{A}_k = \mathcal{A}_{k-1} \setminus \{j\}$  to maximum increase the stability,  $\hat{j} = \arg \max_{j \in \mathcal{A}_k} Q(f_{\mathcal{A}_k} | \mathbf{w}_{\text{MP}}, \mathbf{A}, \mathbf{B}, \mathcal{D})$ :

$$\hat{j} = \arg \max_{j \in \mathcal{A}_{k-1}} \sum_{g=t-\hat{i}+1}^t q_g^j, \quad \text{где} \quad \hat{i} = \sum_{g=1}^t [\eta_g^2 > \eta_t].$$

The stages Add and Del repeated independently such that the inequality holds on each stage:  $\max_{\text{Add-Del } k \in \mathbb{N}} (\mathcal{E}(f_{\mathcal{A}'_k})) - \mathcal{E}(f_{\mathcal{A}_k}) \leq \Delta \mathcal{E}$ .

The algorithm is repeated while the expectation of the likelihood function  $E\mathcal{E}(f_{\mathcal{A}_k})$  remains constant.

We approximate error function  $S$  by

$$\Delta S = S(\mathbf{w}_0 + \Delta \mathbf{w}) - S(\mathbf{w}_0) = \frac{1}{2} \Delta \mathbf{w}^\top \mathbf{H} \Delta \mathbf{w}$$

near its local optimum  $\mathbf{w}_0$ . Here  $\Delta \mathbf{w} = \mathbf{w} - \mathbf{w}_0$  and  $\mathbf{H}$  stands for Hessian matrix of  $S$ .

Since

$$w_j = 0 \equiv \mathbf{e}_j^\top \Delta \mathbf{w} + w_j = 0,$$

we specify Lagrange function

$$L = \frac{1}{2} \Delta \mathbf{w}^\top \mathbf{H} \Delta \mathbf{w} - \lambda_j (\mathbf{e}_j^\top \Delta \mathbf{w} + w_j)$$

for conditional optimization  $\Delta S \rightarrow \min$ ,  $\mathbf{e}_j^\top \Delta \mathbf{w} + w_j = 0$ .

The optimal pruning criterion is then given by

$$\hat{j} = \operatorname{argmin}_{j \in \mathcal{A}} L_j, \text{ where } L_j = \frac{w_j^2}{2[\mathbf{H}^{-1}]_{j,j}}.$$

## Decomposition of the covariance matrix $\mathbf{A}^{-1}$

Consider the condition numbers  $\eta_j = \frac{\lambda_{\max}}{\lambda_j}$  in the singular decomposition of the covariance matrix  $\mathbf{A}^{-1}\mathbf{V} = \mathbf{V}\mathbf{\Lambda}^2$ . Find covariance of the parameters  $\mathbf{w}$

$$\mathbf{Var}(\mathbf{w}) = \frac{1}{\beta}(\mathbf{V}^T)^{-1}\mathbf{\Lambda}^{-2}\mathbf{V}^{-1} = \frac{1}{\beta}\mathbf{V}\mathbf{\Lambda}^{-2}\mathbf{V}^T,$$

where  $\beta$  is an inverse covariance of the residuals, and the covariance of the parameter  $w_j$  is a  $j$ -th diagonal element  $\mathbf{Var}(\mathbf{w})$ .

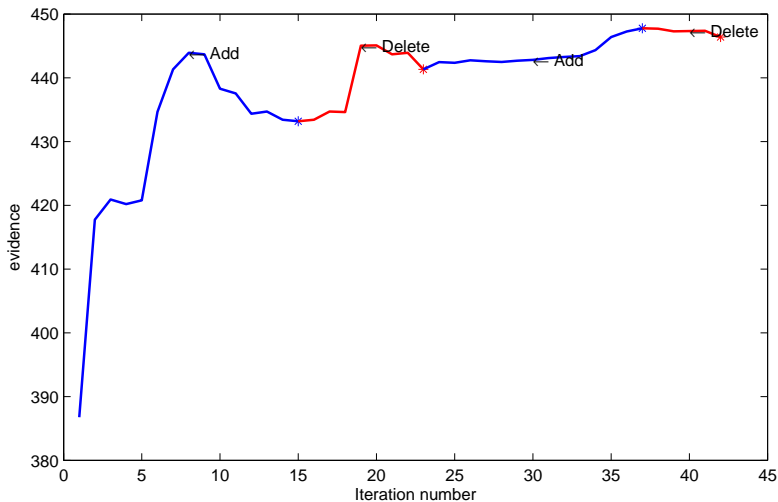
### Removal of the index $\hat{j}$ from the set $\mathcal{A}_k = \mathcal{A}_{k-1} \setminus \{\hat{j}\}$

$$\hat{j} = \arg \max_{j \in \mathcal{A}_{k-1}} \sum_{g=t-\hat{i}+1}^t q_g^j, \quad \text{where} \quad \hat{i} = \sum_{g=1}^t [\eta_g^2 > \eta_t], \text{ where}$$

$$\beta \mathbf{var}(w_i) = \sum_{j=1}^n \frac{v_{ij}^2}{\lambda_j^2} = (q_{i1} + q_{i2} + \dots + q_{in})$$

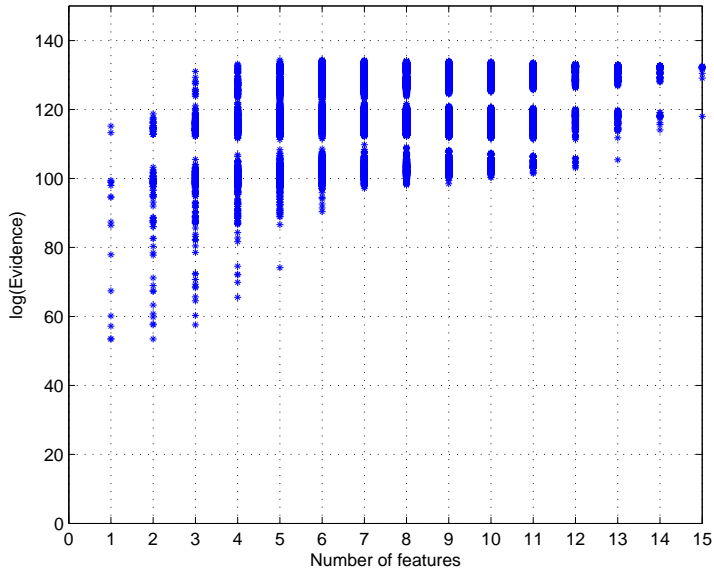
makes maximum increase the model stability  $f_{\mathcal{A}_k}$  on the pair of steps  $k, k-1$ .

# Likelihood maximization during the successive model modification

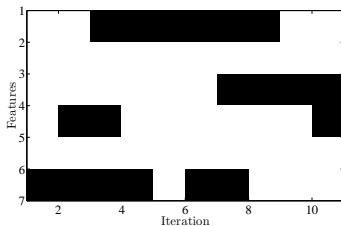
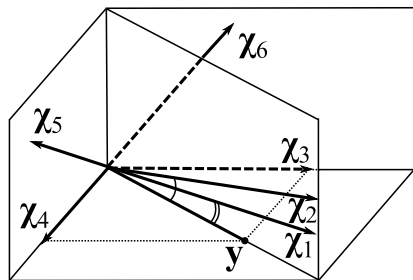


x-axis: iterations  $k$ , y-axis: likelihood  $p(f_{A_k} | \mathbf{w}_{MP}, \mathbf{A}, \mathbf{B}, \mathcal{D})$ .

# Change of likelihood at the arbitrary modification



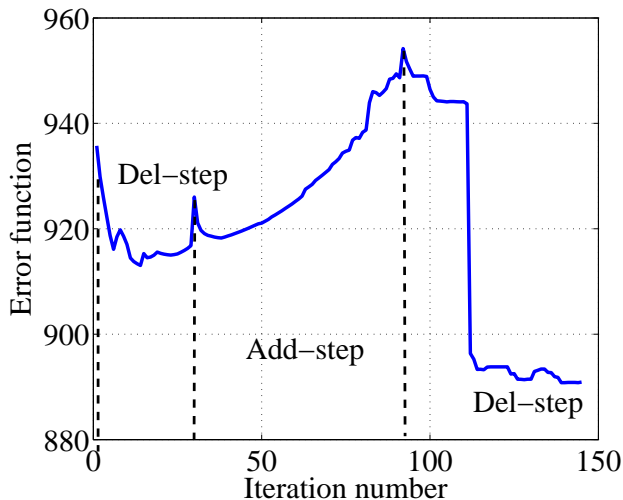
# Choice of the most plausible and stable model



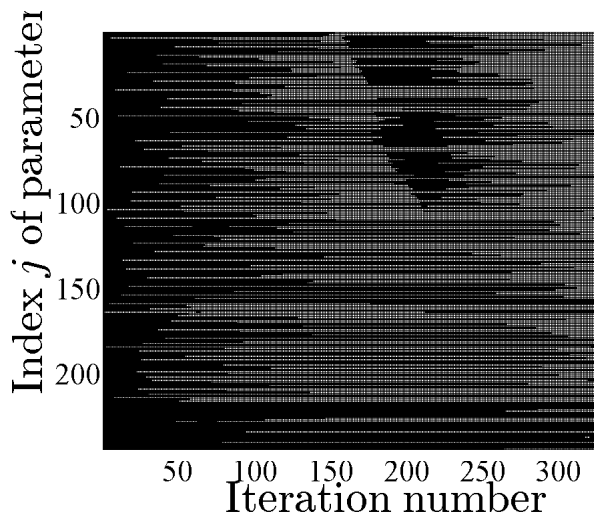
x-axis: the iterations  $k$ , y-axis: the indices of the elements  $j$ , the black rectangle: the index  $j$  added to the set  $\mathcal{A}_k$ .



## Modification procedure: adding and deleting connections



Modification procedure runs until the process stabilizes. The termination criterion



Black cells denote active parameters  $w_j$   $j \in \mathcal{A}$ .

To compare models we use three quality criteria for model  $\mathbf{f}$  with parameter vector  $\mathbf{w}$ : complexity, robustness and precision.

*Complexity*  $C$  is the size of the set  $\mathcal{A}$  of active parameters:

$$C(\mathbf{w}) = \sum_{i=1}^k [w_i \neq 0].$$

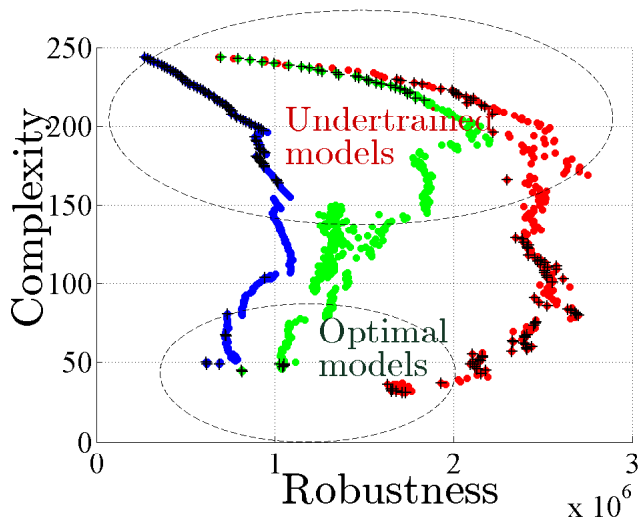
*Robustness*  $\eta = \eta(\hat{\mathbf{w}})$  is equal to the condition number of inverse covariance matrix  $\mathbf{A}$  of  $\mathbf{w}$ :

$$\eta(\hat{\mathbf{w}}) = \frac{\lambda_{\max}}{\lambda_{\min}},$$

where  $\lambda$ . stand for eigenvalues of  $\mathbf{A}$ .

Robustness increases with decrease of  $\eta$ : the best case is  $\lambda_{\min} = \lambda_{\max}$ ,  $\eta = 1$ .

*Precision*  $S$  is measured as error  $S(\mathbf{w}|\mathcal{L})$ .



Generated models in complexity-robustness coordinates.

Dataset: Energy consumption, an example

Algorithm	$S_{\mathcal{L}}$	$S_{\mathcal{C}}$	AIC	BIC	$C_p$	$\lg \kappa$	$k$
Genetics	0,073	0,107	-1152	-1072	337	13	26
GMDH	0,146	0,194	-1076	-1045	745	6	10
Stepwise	0,128	0,154	-1092	-1055	644	7	12
Ridge	0,111	0,146	-819	-330	832	33	160
Lasso	0,121	0,147	-1089	-1034	611	5	18
Stagewise	0,071	0,096	-1157	-1077	324	9	26
FOS	0,106	0,135	-1105	-1044	527	7	20
LARS	0,098	0,095	-1102	-1017	492	7	28
Consequent	0,097	0,123	-1118	-1054	469	5	21