

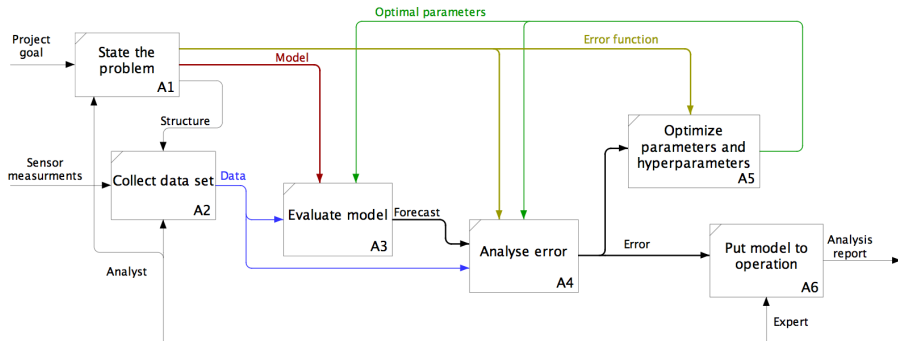
Постановка задач и выбор моделей в машинном обучении

Вадим Викторович Стрижов

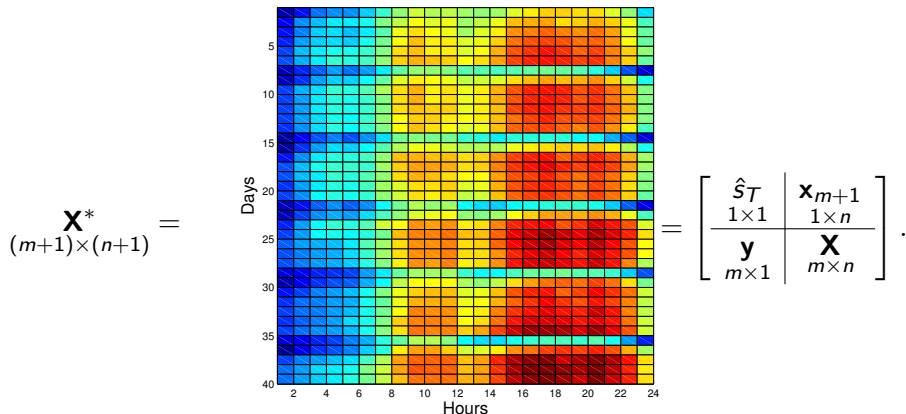
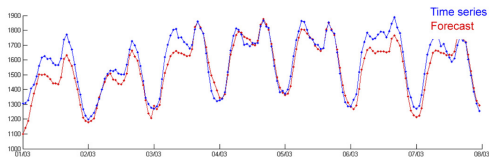
Московский физико-технический институт

Осенний семестр 2019

Analyst creates a model for expert to put it to operation



The autoregressive design matrix and the model



In terms of regression: $\hat{\mathbf{y}} = \mathbf{f}(\mathbf{X}, \mathbf{w}) = \mathbf{X}\mathbf{w}$, $\hat{y}_{m+1} = \hat{\mathbf{S}}_T = \langle \mathbf{x}_{m+1}, \hat{\mathbf{w}} \rangle$.

Model generation

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

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

$$\left(\begin{array}{ccc|ccc} g_1 \circ S_{T-1} & \dots & g_r \circ S_{T-1} & \dots & g_1 \circ S_{T-\kappa+1} & \dots & g_r \circ S_{T-\kappa+1} \\ g_1 \circ S_{(m-1)\kappa-1} & \dots & g_r \circ S_{(m-1)\kappa-1} & \dots & g_1 \circ S_{(m-2)\kappa+1} & \dots & g_r \circ S_{(m-2)\kappa+1} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ g_1 \circ S_{n\kappa-1} & \dots & g_r \circ S_{n\kappa-1} & \dots & g_1 \circ S_{n(\kappa-1)+1} & \dots & g_r \circ S_{n(\kappa-1)+1} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ g_1 \circ S_{\kappa-1} & \dots & g_r \circ S_{\kappa-1} & \dots & g_1 \circ S_1 & \dots & g_r \circ S_1 \end{array} \right)$$

Kolmogorov-Gabor polynomial as a variant for model generation

$$y = w_0 + \sum_{i=1}^{UV} w_i x_i + \sum_{i=1}^n \sum_{j=1}^n w_{ij} x_i x_j + \dots + \sum_{i=1}^n \dots \sum_{z=1}^n w_{i\dots z} x_i \dots x_z,$$

where the coefficients

$$\mathbf{w} = (w_0, w_i, w_{ij}, \dots, w_{i\dots z})_{i,j,\dots,z=1,\dots,n}$$

Examples of nonparametric transformation functions

► Univariate

| Formula | Output dimension |
|-------------|------------------|
| \sqrt{x} | 1 |
| $x\sqrt{x}$ | 1 |
| $\arctan x$ | 1 |
| $\ln x$ | 1 |
| $x \ln x$ | 1 |

► Bivariate

| | |
|----------|-------------------|
| Plus | $x_1 + x_2$ |
| Minus | $x_1 - x_2$ |
| Product | $x_1 \cdot x_2$ |
| Division | $\frac{x_1}{x_2}$ |
| | $x_1 \sqrt{x_2}$ |
| | $x_1 \ln x_2$ |

Nonparametric aggregation: sample statistics

Nonparametric transformations include basic data statistics:

- ▶ Sum or average value of each row \mathbf{x}_i , $i = 1, \dots, m$:

$$\phi_i = \sum_{j=1}^n x_{ij}, \text{ or } \phi'_i = \frac{1}{n} \sum_{j=1}^n x_{ij}.$$

- ▶ Min and max values: $\phi_i = \min_j x_{ij}$, $\phi'_i = \max_j x_{ij}$.
- ▶ Standard deviation:

$$\phi_i = \frac{1}{n-1} \sqrt{\sum_{j=1}^n (x_{ij} - \text{mean}(\mathbf{x}_i))^2}.$$

- ▶ Data quantiles: $\phi_i = [X_1, \dots, X_K]$, where

$$\sum_{j=1}^n [X_{k-1} < x_{ij} \leq X_k] = \frac{1}{K}, \text{ for } k = 1, \dots, K.$$

Nonparametric transformations: Haar's transform

Applying Haar's transform produces multiscale representations of the same data.

Assume that $n = 2^K$ and init $\phi_{i,j}^{(0)} = \phi'_{i,j}^{(0)} = x_{ij}$ for $j = 1, \dots, n$.

To obtain coarse-graining and fine-graining of the input feature vector \mathbf{x}_i , for $k = 1, \dots, K$ repeat:

- ▶ data averaging step

$$\phi_{i,j}^{(k)} = \frac{\phi_{i,2j-1}^{(k-1)} + \phi_{i,2j}^{(k-1)}}{2}, \quad j = 1, \dots, \frac{n}{2^k},$$

- ▶ and data differencing step

$$\phi'_{i,j}^{(k)} = \frac{\phi_{i,2j}^{(k-1)} - \phi_{i,2j-1}^{(k-1)}}{2}, \quad j = 1, \dots, \frac{n}{2^k}.$$

The resulting multiscale feature vectors are $\phi_i = [\phi_i^{(1)}, \dots, \phi_i^{(K)}]$ and $\phi'_i = [\phi'_i^{(1)}, \dots, \phi'_i^{(K)}]$.

Examples of parametric transformation functions

| Function name | Formula | Output dim. | Num. of args | Num. of pars |
|----------------------|--|-------------|--------------|--------------|
| Add constant | $x + w$ | 1 | 1 | 1 |
| Quadratic | $w_2x^2 + w_1x + w_0$ | 1 | 1 | 3 |
| Cubic | $w_3x^3 + w_2x^2 + w_1x + w_0$ | 1 | 1 | 4 |
| Logarithmic sigmoid | $1/(w_0 + \exp(-w_1x))$ | 1 | 1 | 2 |
| Exponent | $\exp x$ | 1 | 1 | 0 |
| Normal | $\frac{1}{w_1\sqrt{2\pi}} \exp\left(\frac{(x-w_2)^2}{2w_1^2}\right)$ | 1 | 1 | 2 |
| Multiply by constant | $x \cdot w$ | 1 | 1 | 1 |
| Monomial | $w_1x^{w_2}$ | 1 | 1 | 2 |
| Weibull-2 | $w_1w_2x^{w_2-1} \exp -w_1x^{w_2}$ | 1 | 1 | 2 |
| Weibull-3 | $w_1w_2x^{w_2-1} \exp -w_1(x - w_3)^{w_2}$ | 1 | 1 | 3 |
| ... | ... | ... | ... | ... |

Ill-conditioned matrix, or curse of dimensionality

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

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

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

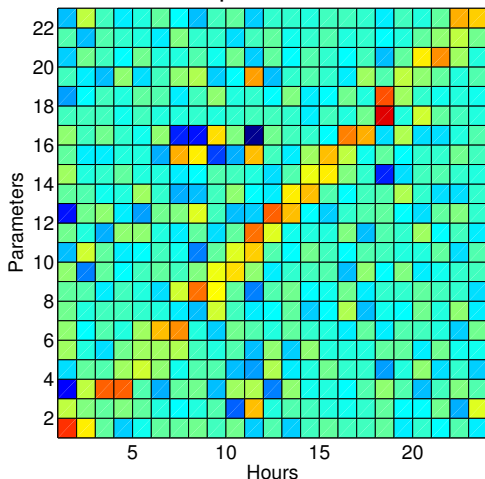
- ▶ for 6 time series the matrix \mathbf{X} is 156×1008 ,
- ▶ for 4 primitive functions it is 156×4032 ,

$$m \ll n.$$

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 τ of the next $(m+1)$ -th period.

произвольный вектор $X\mathbf{v}$, ортогональный вектору регрессионных остатков $X\mathbf{w} - \mathbf{y}$:

$$(X\mathbf{v})^T(X\mathbf{w} - \mathbf{y}) = \mathbf{v}^T(X^T X\mathbf{w} - X^T \mathbf{y}) = 0.$$

Так как это равенство должно быть справедливо для произвольного вектора \mathbf{v} , то $X^T X\mathbf{w} - X^T \mathbf{y} = 0$, см. рис. 6. Если столбцы матрицы X линейно независимы, то матрица $X^T X$ обратима и уравнение имеет единственное решение относительно параметров

$$\mathbf{w} = (X^T X)^{-1} X^T \mathbf{y}. \quad (38)$$

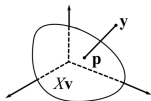
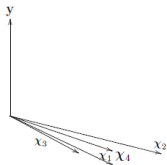


Рис. 6. Проекция вектора зависимой переменной на пространство столбцов матрицы плана.

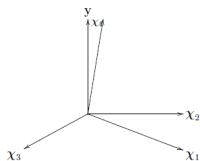
Проекция вектора \mathbf{y} на пространство столбцов матрицы X имеет вид

$$\mathbf{p} = X\mathbf{w} = X(X^T X)^{-1} X^T \mathbf{y} = P\mathbf{y}.$$

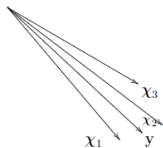
Multicollinear features to forecast: possible configurations



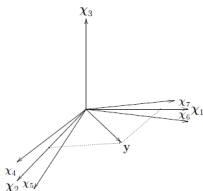
Inadequate and correlated



Adequate and random



Adequate and redundant

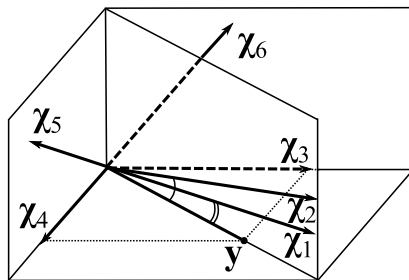


Adequate and correlated

Katrutsa A.M., Strijov V.V. Stresstest procedure for feature selection algorithms // Chemometrics and Intelligent Laboratory Systems, 2015, 142 : 172-183.

Selection of a stable set of features of restricted size

The sample contains multicollinear χ_1, χ_2 and noisy χ_5, χ_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: χ_3, χ_4 is an orthogonal set of features minimizing the error function.

Exhaustive search algorithm

The basic linear model includes all independent variables

$$y = w_0 + a_1 w_1 x_1 + a_2 w_2 x_2 + \dots + a_n w_n x_n$$

structure

The hyperparameter $a \in \{0, 1\}$ is included for the model. The exhaustive search

| a_1 | a_2 | ... | a_n |
|-------|-------|-----|-------|
| 1 | 0 | ... | 0 |
| 0 | 1 | ... | 0 |
| ... | ... | ... | ... |
| 1 | 1 | ... | 1 |

$$A \subseteq Y = \{1, \dots, n\}$$

Denote by vector a the indicator function.

$$S(w|D_A) = \|y - f_a\|_2^2 \rightarrow \min_A$$

Add (append a feature)

Step 0.

The active set $\mathcal{A}_0 = \emptyset$, and \mathcal{W} is the set of feature indices, $P = |\mathcal{W}|$.

Step $k = 1, \dots, \mathcal{A}$

Select the next best feature index

$$\hat{j} = \arg \min_{j \in P \setminus \mathcal{A}_k} \min_{\mathbf{w} \in \mathbb{W}_k} \|(X_{\mathcal{A}_k} \mathbf{x}_j) \mathbf{w} - \mathbf{y}\|_2^2,$$

then

$$\mathcal{A}_{k+1} = \mathcal{A}_k \cup \hat{j}.$$

Assume the following

The column vectors

$$\mathbf{x}^j = \{x_i^j | i \in 1, \dots, \ell\} \quad \text{and} \quad \mathbf{y} = \{y_i | i \in 1, \dots, \ell\}.$$

The model

$$\mathbf{y} = w_1 \mathbf{x}^1 + \dots + w_p \mathbf{x}^p + \varepsilon,$$

in the other words,

$$\mathbf{y} = X\mathbf{w} + \varepsilon.$$

Assume for all $j \in \mathcal{N}$

$$\|\mathbf{x}^j\|_1 = 0, \quad \|\mathbf{x}^j\|_2 = 1 \quad \text{and} \quad \|\mathbf{y}\|_1 = 0, \quad \|\mathbf{y}\|_2 = 1.$$

For all $j, k \in \mathcal{N}$, $j \neq k$ the vectors $\mathbf{x}^j, \mathbf{x}^k$ are linear independent.
Then the vector of correlation coefficients

$$\mathbf{b} = X^T \mathbf{y}.$$

Fast orthogonal search

Step 0.

The residuals $\varepsilon_0 = \mathbf{0}$, the active set $\mathcal{A}_0 = \emptyset$.

Step $k = 1, \dots, P$.

$$\mathcal{A}_k = \mathcal{A}_{k-1} \cup \hat{j},$$

where \hat{j} — feature, which has maximum correlation with ε_k :

$$\hat{j} = \arg \max_{j \in \{N \setminus \mathcal{A}_k\}} \frac{\langle \mathbf{w}, \mathbf{x}^j \rangle}{\|\mathbf{x}^j\| \|\varepsilon_k\|},$$

and

$$\varepsilon_k = X_{\mathcal{A}} \mathbf{w}_{\mathcal{A}} - \varepsilon_{k-1}.$$

Least angle regression, LARS

Denote $\boldsymbol{\mu} = X\mathbf{w}$.

Step 0.

$\boldsymbol{\mu}_0 = \mathbf{0}$, residual vector $\boldsymbol{\varepsilon}_0 = \mathbf{y} - \boldsymbol{\mu}_0$.

Step 1.

Let \mathbf{y} has greater correlation with \mathbf{x}^1 than with \mathbf{x}^2 . Then the new value of $\boldsymbol{\mu}_1 = \boldsymbol{\mu}_0 + w_1 \mathbf{x}^1$, where w_1 is chosen so, that the vector $\mathbf{y}_2 - \boldsymbol{\mu}$ is a bisector for the vectors $\mathbf{x}^1, \mathbf{x}^2$.

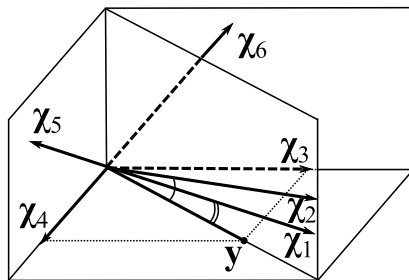
Step 2.

For the unit bisector \mathbf{u}_2 calculate w_2 :

$$\boldsymbol{\mu}_2 = \boldsymbol{\mu}_1 + w_2 \mathbf{u}_2 = \mathbf{y}_2 \quad \text{for } \mathbb{R} = 2.$$

Selection of a stable set of features of restricted size

The sample contains multicollinear χ_1, χ_2 and noisy χ_5, χ_6 features, columns of the design matrix \mathbf{X} . We want to select two features from six.

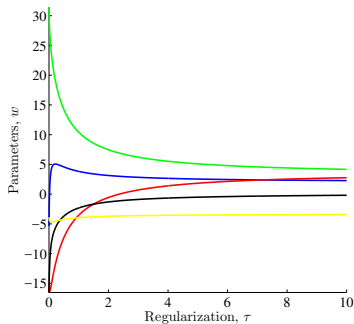


Stability and accuracy for a fixed complexity

The solution: χ_3, χ_4 is an orthogonal set of features minimizing the error function.

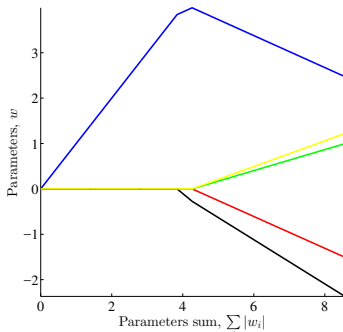
Model parameter values with regularization

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



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

$+ \delta^2 |w|$ elastic



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

Minimize number of similar and maximize number of relevant features

Introduce a feature selection method QP(Sim, Rel) to solve the optimization problem

$$\mathbf{a}^* = \arg \min_{\mathbf{a} \in \mathbb{B}^n} \mathbf{a}^T \mathbf{Q} \mathbf{a} - \mathbf{b}^T \mathbf{a},$$

where matrix $\mathbf{Q} \in \mathbb{R}^{n \times n}$ of pairwise similarities of features χ_i and χ_j is

$$\mathbf{Q} = [q_{ij}] = \text{Sim}(\chi_i, \chi_j) = \left| \frac{\text{Cov}(\chi_i, \chi_j)}{\sqrt{\text{Var}(\chi_i)\text{Var}(\chi_j)}} \right|$$

and vector $\mathbf{b} \in \mathbb{R}^n$ of feature relevances to the target is

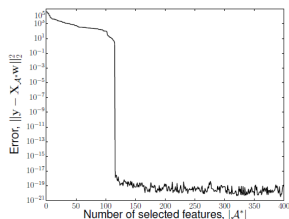
$$\mathbf{b} = [b_i] = \text{Rel}(\chi_i),$$

where elements b_i equal absolute values of the sample correlation coefficient between feature χ_i and the target vector \mathbf{y} .

Number of correlated features Sim \rightarrow min, number of correlated to the target Rel \rightarrow max.

Evaluation criteria for the NIR spectra data set

| Method | C_p | RSS | $\ln \frac{\lambda_1}{\lambda_n}$ SVD | VIF | BIC |
|-------------------------|----------------------|-----------------------|---------------------------------------|----------------------|-------------------|
| QP ($\tau = 10^{-9}$) | -110 | $1.37 \cdot 10^{-18}$ | -25.7 | $6.43 \cdot 10^6$ | 548.38 |
| Genetic | -110.88 | $7.68 \cdot 10^{-30}$ | -24 | $8.13 \cdot 10^5$ | 534.19 |
| LARS | $3.22 \cdot 10^{21}$ | $2.07 \cdot 10^{-7}$ | -28.3 | $7.94 \cdot 10^7$ | 529.47 |
| Lasso | $2.5 \cdot 10^{28}$ | 1.61 | -27.72 | $1.03 \cdot 10^{21}$ | 1712.92 |
| ElasticNet | $2.51 \cdot 10^{28}$ | 1.61 | -27.72 | $1.03 \cdot 10^{21}$ | 1712.92 |
| Stepwise | $3.66 \cdot 10^{29}$ | 23.56 | -36.78 | $1.94 \cdot 10^{22}$ | 1919.23 |
| Ridge | $1.59 \cdot 10^{28}$ | 1.02 | -36.22 | $1.07 \cdot 10^{22}$ | $1.79 \cdot 10^3$ |



Dependence of residual norm on the number of selected features QP(Sim, Rel).

Katrutsa A.M., Strijov V.V. Comprehensive study of feature selection methods to solve multicollinearity problem according to evaluation criteria // Expert Systems with Applications, 2017, 76 : 1-11.

- 1 There are set of binary vectors $\{\mathbf{a}_1, \dots, \mathbf{a}_P\}$, $\mathbf{a} \in \{0, 1\}^n$;
- 2 get two vectors $\mathbf{a}_p, \mathbf{a}_q$, $p, q \in \{1, \dots, P\}$;
- 3 chose random number $\nu \in \{1, \dots, n - 1\}$;
- 4 split both vectors and change their parts:

$$[a_{p,1}, \dots, a_{p,\nu}, a_{q,\nu+1}, \dots, a_{q,n}] \rightarrow \mathbf{a}'_p,$$

$$[a_{q,1}, \dots, a_{q,\nu}, a_{p,\nu+1}, \dots, a_{p,n}] \rightarrow \mathbf{a}'_q;$$

- 5 choose random numbers $\eta_1, \dots, \eta_Q \in \{1, \dots, n\}$;
- 6 invert positions η_1, \dots, η_Q of the vectors $\mathbf{a}'_p, \mathbf{a}'_q$;
- 7 repeat items 2-6 $P/2$ times;
- 8 evaluate the obtained models.

Repeat R times; here P, Q, R are the parameters of the algorithm and n is the number of the corresponding model features.

- 1 There are set of binary vectors $\{\mathbf{a}_1, \dots, \mathbf{a}_P\}$, $\mathbf{a} \in \{1, \dots, k\}^n$;
- 2 get two vectors $\mathbf{a}_p, \mathbf{a}_q$, $p, q \in \{1, \dots, P\}$;
- 3 chose random number $\nu \in \{1, \dots, n-1\}$;
- 4 split both vectors and change their parts:

$$[a_{p,1}, \dots, a_{p,\nu}, a_{q,\nu+1}, \dots, a_{q,n}] \rightarrow \mathbf{a}'_p,$$

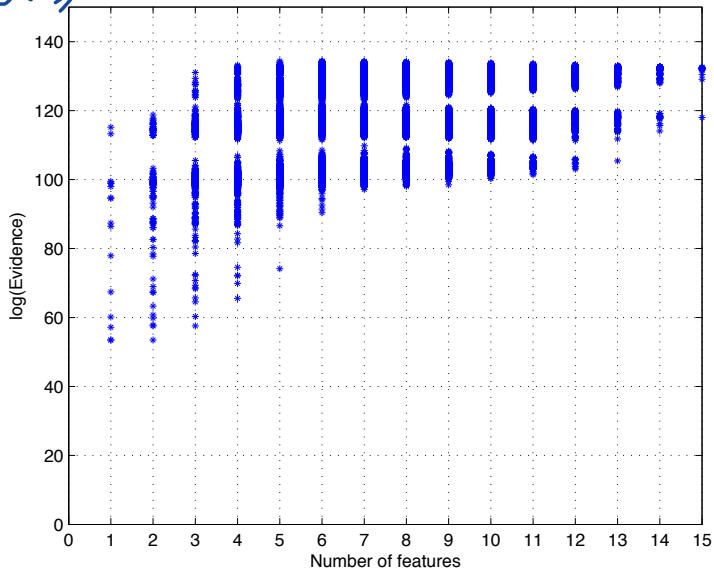
$$[a_{q,1}, \dots, a_{q,\nu}, a_{p,\nu+1}, \dots, a_{p,n}] \rightarrow \mathbf{a}'_q;$$

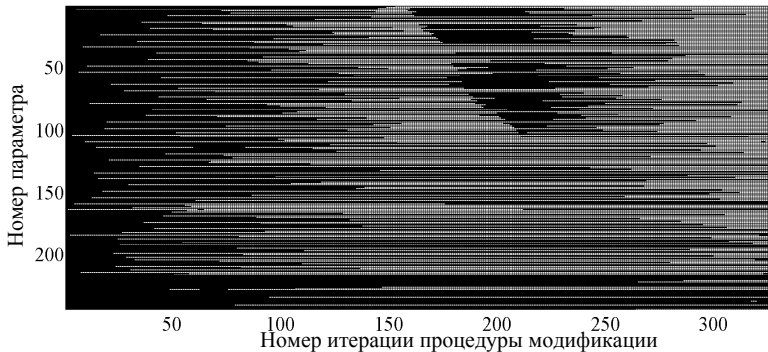
- 5 choose random numbers $\eta_1, \dots, \eta_Q \in \{1, \dots, n\}$;
- 6 replace values in positions η_1, \dots, η_Q of the vectors $\mathbf{a}'_p, \mathbf{a}'_q$ for random values from $\{1, \dots, k\}$;
- 7 repeat items 2-6 $P/2$ times;
- 8 evaluate the obtained models.

Repeat R times; here P, Q, R are the parameters of the algorithm and k is desired number of categories.

Change of likelihood at the arbitrary modification

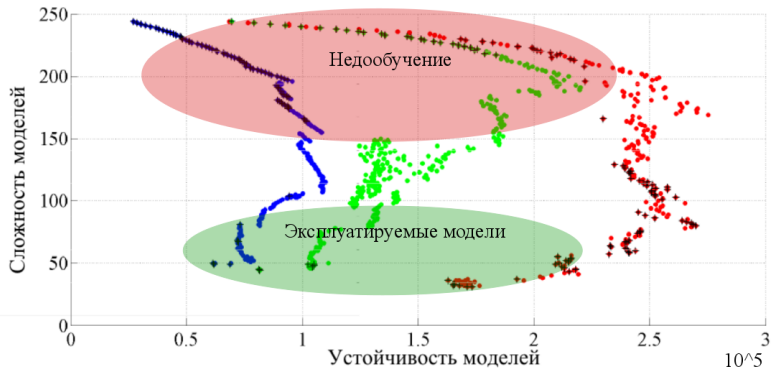
$$\exp(-S(w))$$





- каждый столбец описывает структуру модели на данной итерации;
- белая клетка – параметр неактивен;
- черная клетка – параметр активен.

Эксперимент: результаты



| Стратегия | C | S | η |
|-----------------------------|----|-----|------------------|
| Оптимальное прореж-ние | 50 | 877 | $2.2 \cdot 10^5$ |
| Последовательное прореж-ние | 50 | 870 | $1.2 \cdot 10^5$ |
| Устойчивое прореж-ние | 50 | 866 | $6 \cdot 10^4$ |