

# Выбор устойчивых прогностических моделей в задачах нелинейного регрессионного анализа

А. А. Токмакова

Научный руководитель: В. В. Стрижов  
Московский физико-технический институт  
Факультет управления и прикладной математики  
Кафедра интеллектуальных систем

10 июня 2014 г.  
Москва

## Цель исследования

Предложить стратегию выбора устойчивых моделей в условиях мультиколлинеарности выборки в задачах нелинейного регрессионного анализа.

## В работе предложена

Процедура получения модели с оптимальным числом параметров, использующая методы последовательного добавления и удаления элементов модели.

## План презентации

- 1 Вид функции ошибки
- 2 Метод оценки ковариационной матрицы
- 3 Описание процедуры Add-Del

Выборка:  $\mathcal{D} = \{\mathbf{x}_i, y_i\}, i \in \mathcal{I} = \{1, \dots, m\}$ .

Регрессионная модель:  $\mathbf{f} : (\mathbf{w}, \mathbf{X}) \mapsto \mathbf{y}$ , где  $\mathbf{X} \in \mathbb{R}^{m \times n}$ ,  
 $\mathbf{w} = [w_1, \dots, w_j, \dots, w_t], j \in \mathcal{J} = \{1, \dots, t\}$ .

Нейронная сеть:  $\mathbf{f}(\mathbf{x}, \mathbf{w}) = \mathbf{W}^{(2)}\mathbf{a}$ ,  $\mathbf{a} = \tanh(\mathbf{W}^{(1)}\mathbf{x})$ , где  $\mathbf{W}^{(1)}$ ,  
 $\mathbf{W}^{(2)}$  — матрицы весов первого и второго слоёв нейронной  
сети,  $\mathbf{w} = \text{vec}(\mathbf{W}^{(1)}|\mathbf{W}^{(2)})$ .

Требуется найти такое множество индексов  $\mathcal{A}^* \subseteq \mathcal{J}$ , что:

$$\mathcal{A}^* = \underset{\mathcal{A} \subseteq \mathcal{J}}{\text{argmin}} S(\mathbf{w}_{\mathcal{A}}^* | \mathcal{D}, \mathbf{f}), \quad \mathbf{w}_{\mathcal{A}}^* = \underset{\mathbf{w} \in \mathbb{R}^t}{\text{argmin}} S(\mathbf{w} | \mathcal{D}, \mathbf{f}),$$

где  $S(\mathbf{w})$  — функция ошибки модели.

Пусть многомерная случайная величина  $\mathbf{y} \sim \mathcal{N}(\mathbf{f}, \mathbf{B})$  имеет нормальное распределение

$$p(\mathbf{y}|\mathbf{X}, \mathbf{B}) = \frac{1}{(2\pi)^{\frac{m}{2}} \det^{\frac{1}{2}}(\mathbf{B}^{-1})} \exp\left(-\frac{1}{2}(\mathbf{y} - \mathbf{f})^T \mathbf{B}(\mathbf{y} - \mathbf{f})\right) = \frac{\exp(-E_D)}{Z_D(\mathbf{B})}.$$

Вектор параметров  $\mathbf{w}$  модели  $\mathbf{f}$  — многомерная случайная величина с математическим ожиданием  $\mathbf{w}_0$  и ковариационной матрицей  $\mathbf{A}^{-1}$

$$p(\mathbf{w}|\mathbf{A}) = \frac{1}{(2\pi)^{\frac{t}{2}} \det^{\frac{1}{2}}(\mathbf{A}^{-1})} \exp\left(-\frac{1}{2}(\mathbf{w} - \mathbf{w}_0)^T \mathbf{A}(\mathbf{w} - \mathbf{w}_0)\right) = \frac{\exp(-E_w)}{Z_w(\mathbf{A})}.$$

$$p(\mathcal{D}|\mathbf{w}, \mathbf{B}) = \frac{\exp(-E_D)}{Z_D(\mathbf{B})}, \quad p(\mathbf{w}|\mathbf{A}) = \frac{\exp(-E_w)}{Z_w(\mathbf{A})}.$$

Апостериорное распределение параметров модели для заданных ковариационных матриц  $\mathbf{A}^{-1}$  и  $\mathbf{B}^{-1}$  имеет вид:

$$p(\mathbf{w}|\mathcal{D}, \mathbf{A}, \mathbf{B}) = \frac{p(\mathcal{D}|\mathbf{w}, \mathbf{B})p(\mathbf{w}|\mathbf{A})}{p(\mathcal{D}|\mathbf{A}, \mathbf{B})} = \frac{\exp(-(E_D + E_w))}{Z_D Z_w p(\mathcal{D}|\mathbf{A}, \mathbf{B})}.$$

Определим функцию ошибки  $S(\mathbf{w})$ :

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

Ковариационная матрица случайной величины  $\mathbf{w} \in \mathbb{R}^t$ :

$$\mathbf{A}^{-1} = \mathbb{E} \left[ (\mathbf{w} - \mathbb{E}\mathbf{w})(\mathbf{w} - \mathbb{E}\mathbf{w})^\top \right] = \mathbb{E}[\mathbf{w}\mathbf{w}^\top] - \mathbb{E}[\mathbf{w}] \cdot \mathbb{E}[\mathbf{w}^\top].$$

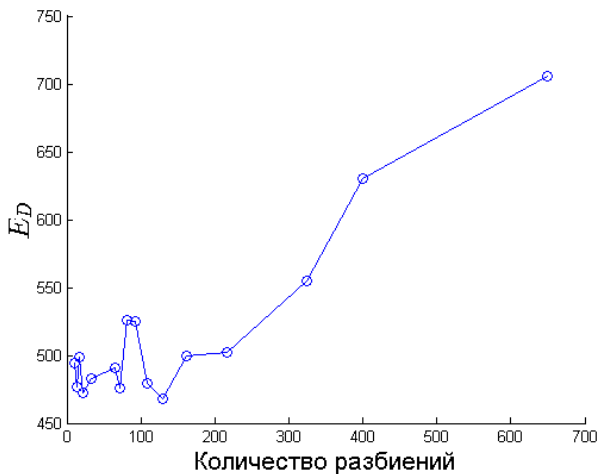
Разобьем множество индексов выборки  $\mathcal{I}$  на  $b$  подмножеств, на каждой подвыборке оценим вектор параметров  $\mathbf{w}$  модели  $\mathbf{f}$ :

$$\mathbf{W} = [\tilde{\mathbf{w}}_1, \dots, \tilde{\mathbf{w}}_b] \in \mathbb{R}^{t \times b}.$$

Центрируем каждую строку  $\mathbf{W}_j - \mathbb{E}\mathbf{W}_j \mapsto \mathbf{W}_j$  и оценим ковариационную матрицу параметров модели  $\mathbf{f}$ :

$$\mathbf{A}^{-1} = \frac{1}{b} \mathbf{W}\mathbf{W}^\top.$$

Зависимость функции  $E_D$  от  $b$  (данные winequality, UCI):



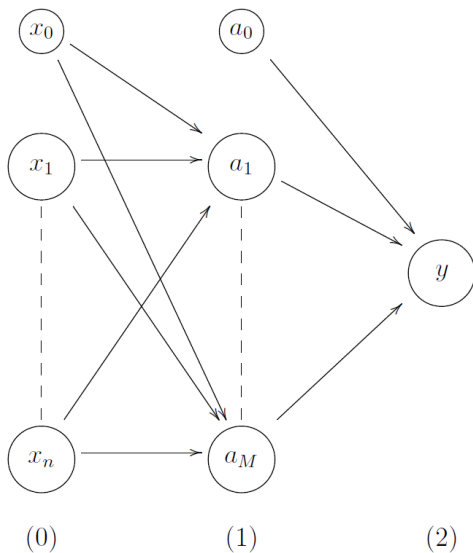
Определим структуру модели как ациклический направленный граф с тремя множествами вершин фиксированной максимальной валентности:

- (0) — множество вершин свободных переменных размера  $n + 1$ ;
- (1) — множество вершин функции активации первого слоя размера  $M + 1$ ;
- (2) — множество вершин функции активации выходного слоя размера 1.

Ребра графа могут идти только из (0)  $\rightarrow$  (1) или (1)  $\rightarrow$  (2).



Представление нейронной сети в виде графа:



Зададим вектор  $\mathbf{z} \in \{0, 1\}^t$ :

$$z_j = \begin{cases} 0, & \text{если } w_j + \Delta w_j = 0, \text{ т. е. ребра нет в графе;} \\ 1, & \text{иначе.} \end{cases}$$

$\mathbf{z}$ ,  $n$  и  $M$  однозначно задают вид графа, а значит и структуру модели.

Представим процедуру последовательной модификации структуры модели в виде пути внутри  $t$ -мерного гиперкуба  $\mathcal{Z}$ , каждая вершина которого является бинарным вектором  $\mathbf{z} = \{0, 1\}^t$ .

Соседние вершины гиперкуба  $\mathcal{Z}$  отличаются единственной компонентой вектора  $\mathbf{z}$ .

Стратегия задается следующими математическими объектами:

- набором критериев качества  $\{Q\}$ ;
- набором ограничений на структуру и параметры модели  $\mathcal{A} \in \mathcal{J}, \mathbf{w} = \mathbf{w}_{\mathcal{A}}^*$ ;
- критериями останова шагов удаления и добавления структурных единиц в модель;
- критерием останова процедуры выбора модели.

То есть, действуя согласно стратегии, будем изменять структуру модели, удаляя из неё элементы и добавляя их до тех пор, пока значение критериев качества не стабилизируется.

Пусть на шаге  $q$  известна правдоподобная нелинейная модель субоптимальной сложности, то есть известны оценки её параметров  $\mathbf{w}_{\mathcal{A}_q}^*$  и гиперпараметров  $\mathbf{w}_{0q}^*$  и  $\mathbf{A}_q^*$ :

$$\mathbf{f}_{\mathcal{A}_q} \big|_{\mathbf{w}=\mathbf{w}_{\mathcal{A}_q}^*} : \mathbf{X} \mapsto \mathbf{y}.$$

- ① Найдем такой элемент множества  $\mathcal{J} \setminus \mathcal{A}_q$ , что:

$$j^* = \operatorname{argmin}_{j \in \mathcal{J} \setminus \mathcal{A}_q} S(\mathbf{w}_{\mathcal{A}_q \cup \{j\}} \big| \mathcal{D}, \mathbf{f}_{\mathcal{A}_q \cup \{j\}}).$$

- ② Добавим новый элемент  $j^*$  к текущему набору

$$\mathcal{A}_{q+1} = \mathcal{A}_q \cup \{j^*\}$$

и будем повторять эту процедуру, пока

$$\left| S(\mathbf{w}_{\mathcal{A}_{q+1}}^* \big| \mathcal{D}, \mathbf{f}_{\mathcal{A}_{q+1}}) - S(\mathbf{w}_{\mathcal{A}_q}^* \big| \mathcal{D}, \mathbf{f}_{\mathcal{A}_q}) \right| \leq \Delta S_1.$$

Представим оценку ковариационной матрицы на шаге  $q$  как  $\mathbf{A}_q^{-1} = (\mathbf{L}\mathbf{L}^T)^{-1}$ . Запишем сингулярное разложение матрицы  $\mathbf{L}$ :

$$\mathbf{L} = \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T,$$

где  $\mathbf{U}$  и  $\mathbf{V}$  — ортогональные матрицы, а  $\mathbf{\Lambda}$  — диагональная матрица с собственными значениями  $\lambda_j$  на диагонали, такими что

$$\lambda_1 > \lambda_2 > \dots > \lambda_q > 0.$$

Индексом обусловленности с индексом  $j$  будем называть:

$$\eta_j = \frac{\lambda_{\max}}{\lambda_j}.$$

Запишем матрицу  $\mathbf{A}_q^{-1}$  в виде:

$$\mathbf{A}_q^{-1} = (\mathbf{L}\mathbf{L}^T)^{-1} = (\mathbf{V}\mathbf{\Lambda}^T \mathbf{U}^T \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T)^{-1} = \mathbf{V}\mathbf{\Lambda}^{-2}\mathbf{V}^{-1}.$$

Оценками дисперсии параметров будут диагональные элементы матрицы  $\mathbf{A}_q^{-1}$ :

$$\alpha_{ii} = \sum_{j=1}^t \frac{v_{ij}^2}{\lambda_j^2}.$$

**Дисперсионной долей**  $r_{ij}$  будем называть вклад  $j$ -го признака в дисперсию  $i$ -го элемента вектора параметров  $\mathbf{w}$ :

$$r_{ij} = \frac{v_{ij}^2/\lambda_j^2}{\sum_{j=1}^t v_{ij}^2/\lambda_j^2}.$$

- 1 Вычислим индексы обусловленности  $\eta_j$  и матрицу долевых коэффициентов  $\mathfrak{R} = [r_{gj}]$  для  $\mathbf{w}_{\mathcal{A}_q}$ .
- 2 Найдем индекс  $g^*$  максимального индекса обусловленности  $\eta_{\max}$ .
- 3 В матрице долевых коэффициентов  $\mathfrak{R}$  найдем индекс столбца  $j^*$ .

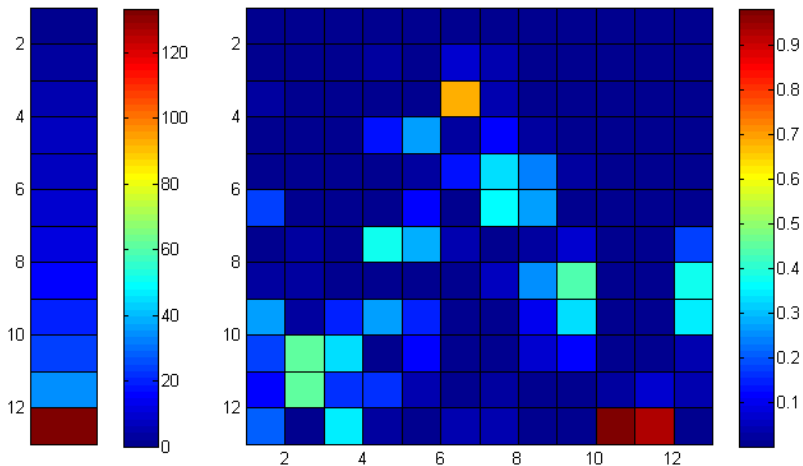
$$j^* = \operatorname{argmax}_{j \in \mathcal{A}_q} s_{g^*j}.$$

- 4 Удалим индекс  $j^*$  из множества  $\mathcal{A}_q$ :

$$\mathcal{A}_{q+1} = \mathcal{A}_q \setminus \{j^*\}$$

и будем повторять эту процедуру, пока

$$\left| S(\mathbf{w}_{\mathcal{A}_{q+1}}^* | \mathcal{D}, \mathbf{f}_{\mathcal{A}_{q+1}}) - S(\mathbf{w}_{\mathcal{A}_q}^* | \mathcal{D}, \mathbf{f}_{\mathcal{A}_q}) \right| \leq \Delta S_2.$$



**Рис.:** индексы обусловленностей и матрица дисперсионных долей.



В качестве критерия стабилизации структуры модели предлагается использовать энтропию:

$$H(\mathcal{A}, \mathcal{A}') = -\rho(\mathbf{z}, \mathbf{z}') \ln(\rho(\mathbf{z}, \mathbf{z}')),$$

где  $\rho(\cdot, \cdot)$  — расстояние Хэмминга между векторами  $\mathbf{z}$  и  $\mathbf{z}'$ , где

$$z_j = \begin{cases} 0, & \text{если } w_j + \Delta w_j = 0, \text{ т. е. } j \notin \mathcal{A}; \\ 1, & \text{иначе.} \end{cases}$$

Процесс считается стабильным, если изменение энтропии  $H(\mathcal{A}, \mathcal{A}')$  не превосходит заданного порога.

Рассмотрим последовательно порожденные модели  $\mathbf{f}_{\mathcal{A}_q}$  и  $\mathbf{f}_{\mathcal{A}_{q+1}}$ . Начальные приближения гиперпараметров для модели  $\mathbf{f}_{\mathcal{A}_{q+1}}$  будут выглядеть следующим образом:

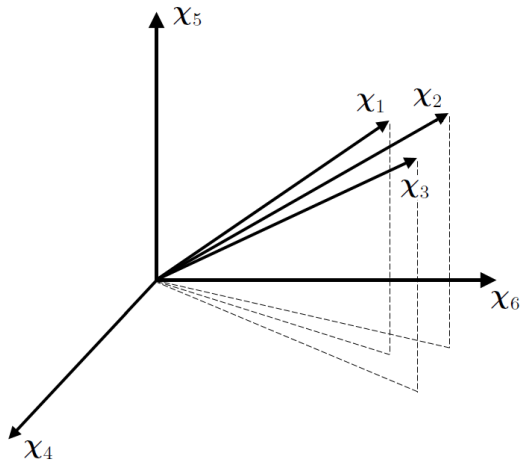
$$\tilde{\mathbf{w}}_{\mathcal{A}_q} = [\mathbf{w}_{\mathcal{A}_q}; 0]^\top; \quad \tilde{\mathbf{A}}_q = \begin{pmatrix} \mathbf{A}_q & 0 \\ 0 & 1 \end{pmatrix}.$$

Запишем функцию правдоподобия данных и априорное распределение параметров модели  $\mathbf{f}_{\mathcal{A}_{q+1}}$  следующим образом:

$$p(\mathcal{D} | \mathbf{w}_{\mathcal{A}_{q+1}}, \mathbf{B}_q) = \frac{1}{(2\pi)^{\frac{m}{2}} \det^{\frac{1}{2}}(\mathbf{B}_q^{-1})} \cdot \exp\left(-\frac{1}{2}(\mathbf{y} - \mathbf{f}_{\mathcal{A}_{q+1}})^\top \mathbf{B}_q (\mathbf{y} - \mathbf{f}_{\mathcal{A}_{q+1}})\right);$$

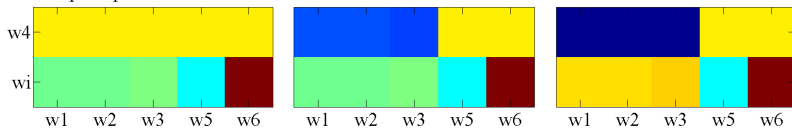
$$p(\mathbf{w}_{\mathcal{A}_{q+1}} | \tilde{\mathbf{A}}_q) = \frac{1}{(2\pi)^{\frac{t}{2}} \det^{\frac{1}{2}}(\tilde{\mathbf{A}}_q^{-1})} \cdot \exp\left(-\frac{1}{2}(\mathbf{w}_{\mathcal{A}_{q+1}} - \tilde{\mathbf{w}}_{\mathcal{A}_q})^\top \tilde{\mathbf{A}}_q (\mathbf{w}_{\mathcal{A}_{q+1}} - \tilde{\mathbf{w}}_{\mathcal{A}_q})\right).$$

$$y = 0.3\chi_5 + 0.3\chi_6 + 0.3\chi_3 + \mathcal{U}[0, 1].$$

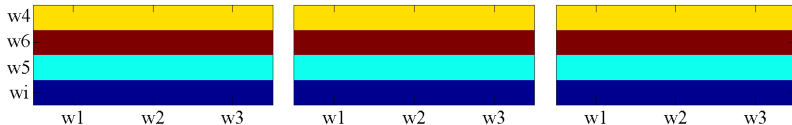
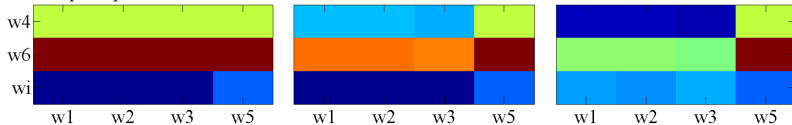


$$\mathbf{y} = 0.3\chi_5 + 0.3\chi_6 + 0.3\chi_3 + \mathcal{U}[0, 1], \quad \hat{\mathbf{y}} = 0.3\chi_5 + 0.8\chi_6 + 0.5\chi_4.$$

Выбран признак #6



Выбран признак #5



Алгоритм Левенберга-Марквардта предназначен для оптимизации параметров нелинейных регрессионных моделей.

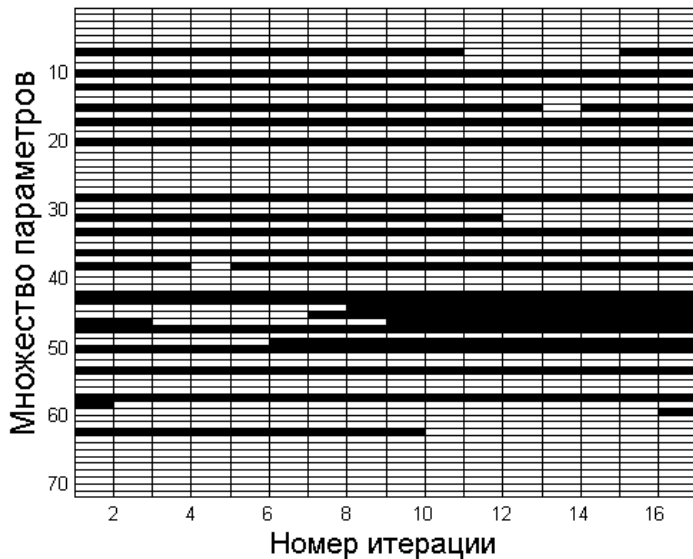
$\mathbf{f}(\mathbf{w} + \Delta\mathbf{w}, \mathbf{X}) \approx \mathbf{f}(\mathbf{w}, \mathbf{X}) + \mathbf{J}\Delta\mathbf{w}$ , где  $\mathbf{J}$  — якобиан функции  $\mathbf{f}$  в точке  $\mathbf{w}$ .

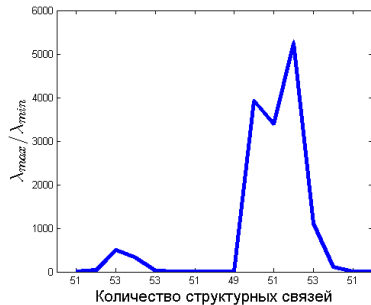
Функция ошибки:

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

Для нахождения экстремума приравняем вектор частных производных  $S$  по  $\mathbf{w}$  к нулю:  $\nabla S = \frac{1}{2}(\mathbf{w} + \Delta\mathbf{w})^T (\mathbf{A} + \mathbf{A}^T) + \frac{1}{2} [(\mathbf{J}(\mathbf{w} + \Delta\mathbf{w}) - \mathbf{y})^T \mathbf{B}^T \mathbf{J} + (\mathbf{J}(\mathbf{w} + \Delta\mathbf{w}) - \mathbf{y})^T \mathbf{B} \mathbf{J}] = 0$ .

Выразив приращение  $\Delta\mathbf{w}$ , получим следующую рекуррентную формулу:  $\Delta\mathbf{w} = [(\mathbf{A} + \mathbf{A}^T + \mathbf{J}^T (\mathbf{B}^T + \mathbf{B}) \mathbf{J})^{-1}]^T (-\mathbf{w}^T (\mathbf{A} + \mathbf{A}^T) + (\mathbf{y} - \mathbf{f}(\mathbf{w}, \mathbf{X}))^T (\mathbf{B}^T + \mathbf{B}) \mathbf{J})^T$ .





## Предложено

- стратегия пошаговой модификации структуры модели;
- модификация алгоритма Левенберга-Марквардта для градиентной оптимизации в случае нелинейных моделей.

## Публикации ВАК

- 1 Токмакова А.А., Стрижов В.В. *Оценивание гиперпараметров линейных регрессионных моделей при отборе шумовых и коррелирующих признаков.* — Информатика и её применения, 2012. Том 6(4), сс. 66-75.
- 2 Оценка гиперпараметров регрессионных моделей методом максимального правдоподобия. — Информационные технологии, 2013. Том 2, сс. 11-15.
- 3 Токмакова А.А. *Алгоритм стохастического отбора объектов и признаков в задаче банковского кредитного скоринга.* — Информационные технологии, 2014. Том 3, сс. 30-35.
- 4 Kuznetsov M.P., Tokmakova A.A., Strijov V.V. *Analytic and stochastic methods of structure parameter estimation.* — Machine Learning Journal, 2014.