

My first scientific paper

Week 7

Error analysis

Vadim Strijov

Moscow Institute of Physics and Technology

2021

Significant increase in complexity and modest increase in accuracy

	train	test	out-of-time	# parameters
Logistic regression	53,08%	55,18%	57,50%	= 12
Neural network	59,85%	57,04%	58,27%	~ 240
Regression forest	61,85%	57,01%	59,61%	> 1000
Gradient boosting	63,58%	58,31%	59,50%	> 10,000

Model selection is an important problem!

... it was a banking credit scoring model

Задача нахождения наиболее правдоподобных параметров

Задана выборка $D = \{(\mathbf{x}_i, y_i)\}$, $i \in \mathcal{I}$, функция ошибки модели S и модель — параметрическое семейство функций $f(\mathbf{w}, \mathbf{x})$. Требуется найти такие параметры \mathbf{w} модели, которые бы доставляли минимум функции ошибки

$$\mathbf{w}^* = \arg \min_{\mathbf{w} \in \mathbb{W}} S(\mathbf{w} | D, f). \quad (1)$$

Функция ошибки, определенная посредством логарифмической функции правдоподобия

$$S(\mathbf{w}) = -\ln(p(D | \mathbf{w}, f)),$$

обеспечивает максимизацию правдоподобия параметров. Параметры, найденные минимизацией этой функции ошибок, будут называться наиболее правдоподобными.

Примеры функции ошибки в регрессии и классификации

Регрессия

Гипотеза порождения данных: $\mathbf{y} \sim \mathcal{N}(\mathbf{f}, \mathbf{I})$.

Функция ошибки:

$$S(\mathbf{w}) = \|\mathbf{y} - \mathbf{f}\|_2^2.$$

Классификация

Гипотеза порождения данных: $\mathbf{y} \sim \mathcal{B}(f, 1 - f)$.

Функция ошибки:

$$S(\mathbf{w}) = \sum_{i \in \mathcal{I}} y_i \ln f(\mathbf{w}^T \mathbf{x}_i) + (1 - y_i) \ln(1 - f(\mathbf{w}^T \mathbf{x}_i)).$$

Задача выбора оптимального набора признаков

- ▶ Задана выборка $D = \{(\mathbf{x}_i, y_i)\}, i \in \mathcal{I}$.
- ▶ Задано случайное разбиение множество индексов элементов выборки $\mathcal{I} = \mathcal{L} \sqcup \mathcal{C}$.
- ▶ Множество независимых переменных $\mathbf{x} = [x_1, \dots, x_j, \dots, x_n]$ проиндексировано $j \in \mathcal{J} = \{1, \dots, n\}$.
- ▶ Задано множество моделей-претендентов $\mathfrak{F} = \{f(\mathbf{w}, \mathbf{x})\}$.
- ▶ Модель — параметрическое семейство функций $f(\mathbf{w}, \mathbf{x}) = \mu(\mathbf{w}^T \mathbf{x})$, где μ — функция связи (в случае регрессии $\mu = \text{id}$, в случае классификации $\mu = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$).
- ▶ Структура модели $f_{\mathcal{A}}$ задана множеством индексов $\mathcal{A} \subseteq \mathcal{J}$ и означает включение переменных $\mathbf{x}_{\mathcal{A}}$. Иначе, используются только признаки-столбцы матрицы \mathbf{X} с индексами из множества \mathcal{A} .
- ▶ Задана функция ошибки S .

Среднее значение и стандартное отклонение ошибки

Задана выборка $\mathcal{D} = \{\mathbf{x}_i, y_i\}$. Ее элементы проиндексированы:

$$i \in \mathcal{I} = \{1, \dots, m\}.$$

Разобьем выборку равномерно случайно, на две равномошные подвыборки, обучение и контроль, K раз:

$$\mathcal{I} \longrightarrow \mathcal{L}_k \sqcup \mathcal{C}_k, \quad k \in \{1, \dots, K\}.$$

Задана модель $f(\mathbf{w}, \mathbf{w})$ и функция ошибки $S(\mathbf{w}|\mathcal{D})$. Параметры модели оптимизированы на обучении $\mathcal{D}_{\mathcal{L}}$ как

$$\hat{\mathbf{w}} = \arg \min S(\mathbf{w}|f, \mathcal{D}_{\mathcal{L}}).$$

Для каждого из K разбиений вычисляем ошибку на обучении и на контроле. Получаем два набора ошибок:

$$\{S_k(\hat{\mathbf{w}}_k|f, \mathcal{D}_{\mathcal{L}k})\}, \quad \{S_k(\hat{\mathbf{w}}_k|f, \mathcal{D}_{\mathcal{C}k})\}, \quad k \in \{1, \dots, K\}.$$

Задача выбора оптимального набора признаков

Требуется найти такое подмножество индексов $\mathcal{A} \subseteq \mathcal{J}$, которое бы доставляло минимум функции

$$\mathcal{A}^* = \arg \min_{\mathcal{A} \subseteq \mathcal{J}} S(f_{\mathcal{A}} | \mathbf{w}^*, D_{\mathcal{C}})$$

на разбиении выборки D , определенном множеством индексов \mathcal{C} .

При этом параметры \mathbf{w}^* модели должны доставлять минимум функции

$$\mathbf{w}^* = \arg \min_{\mathbf{w} \in \mathbb{W}} S(\mathbf{w} | D_{\mathcal{L}}, f_{\mathcal{A}})$$

на разбиении выборки, определенном множеством \mathcal{L} .

Зависимость среднего значения ошибки от объема выборки

Для двух наборов ошибок вычислим среднее значение и поправленное стандартное отклонение:

$$\bar{S} = \frac{1}{K} \sum_{k=1}^K S_k, \quad \sigma = \frac{1}{K-1} \sqrt{\sum_{k=1}^K (\bar{S} - S_k)^2}.$$

Повторим процедуру на ограниченном объеме выборки, например:

$$m = \overline{1, M},$$

где M — наибольший объем доступной выборки.

Построим график зависимости ошибки и стандартного отклонения от объема выборки.

Домашнее задание 3 находится по адресу <http://bit.ly/16UIIQH>

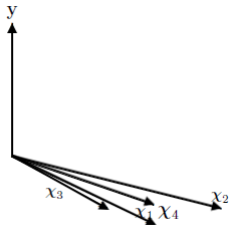
Некоторые задачи машинного обучения

- ▶ Задача оценки параметров модели,
- ▶ задача выбора признаков или объектов выборки,
- ▶ задача выбора модели оптимальной сложности,
- ▶ задача построения и выбора структуры модели,
- ▶ задача проверки гипотезы порождения данных.

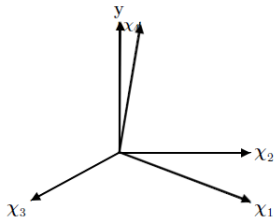
Предполагается, что функция ошибки $S(\mathbf{w}|D, f)$ задана исходя из

- ▶ гипотезы порождения данных,
- ▶ либо из практических соображений.

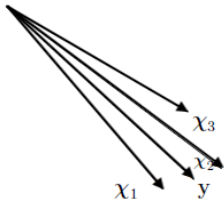
Configurations of feature space



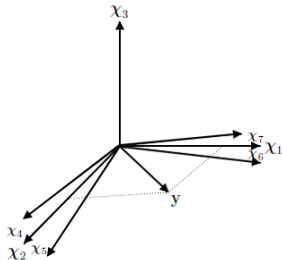
Non-adequate correlated



Adequate random



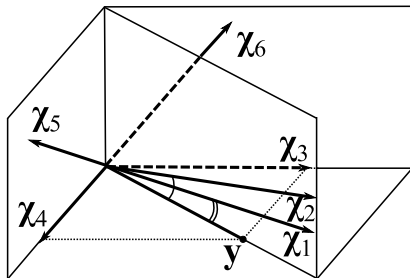
Adequate redundant



Adequate correlated

Выбор устойчивой и точной модели

Выборка содержит мультикоррелирующие χ_1, χ_2 и устойчивые χ_5, χ_6 признаки — столбцы матрицы «объект-признак» \mathbf{X} . Требуется выбрать два признака из шести.



Точность и устойчивость при заданной сложности

Решение: χ_3, χ_4 — набор ортогональных признаков с наименьшим значением функции ошибки.

Minimize number of similar and maximize number of relevant features

The model is defined by a vertex point in the n -dimensional cube.

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

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

Number of correlated features Sim \rightarrow min, number of correlated to the target Rel \rightarrow max, where matrix $\mathbf{Q} \in \mathbb{R}^{n \times n}$ of pairwise similarities of features χ_i and χ_j is

$$\mathbf{Q} = [q_{ij}] = \text{Sim}(\chi_i, \chi_j) = |\text{Corr}(\chi_i, \chi_j)|$$

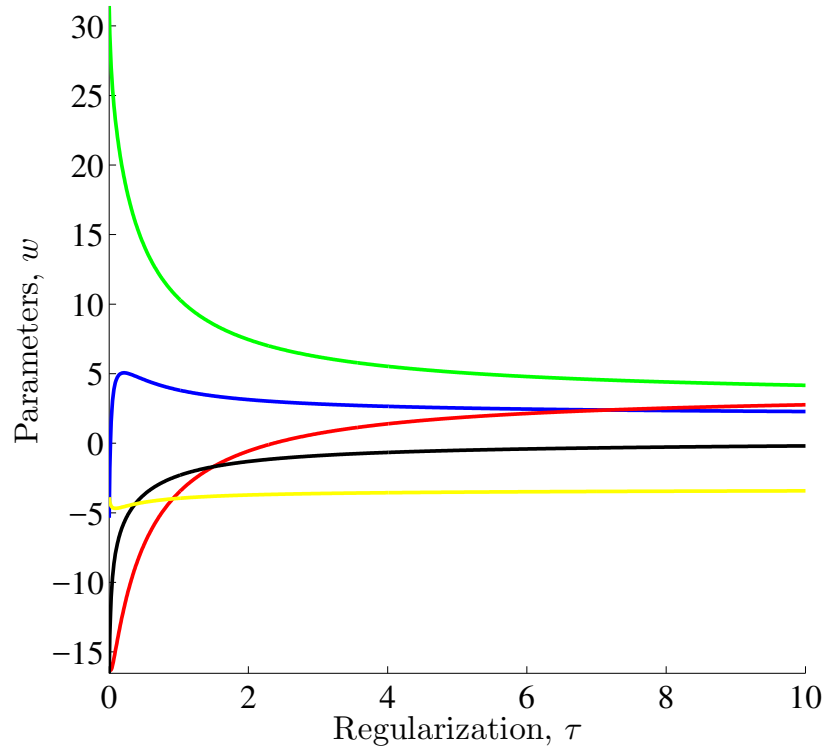
and vector $\mathbf{b} \in \mathbb{R}^n$ of feature relevancies to the target is

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

elements $b_i = \|\text{Corr}(\chi_i, \mathbf{y})\|$.

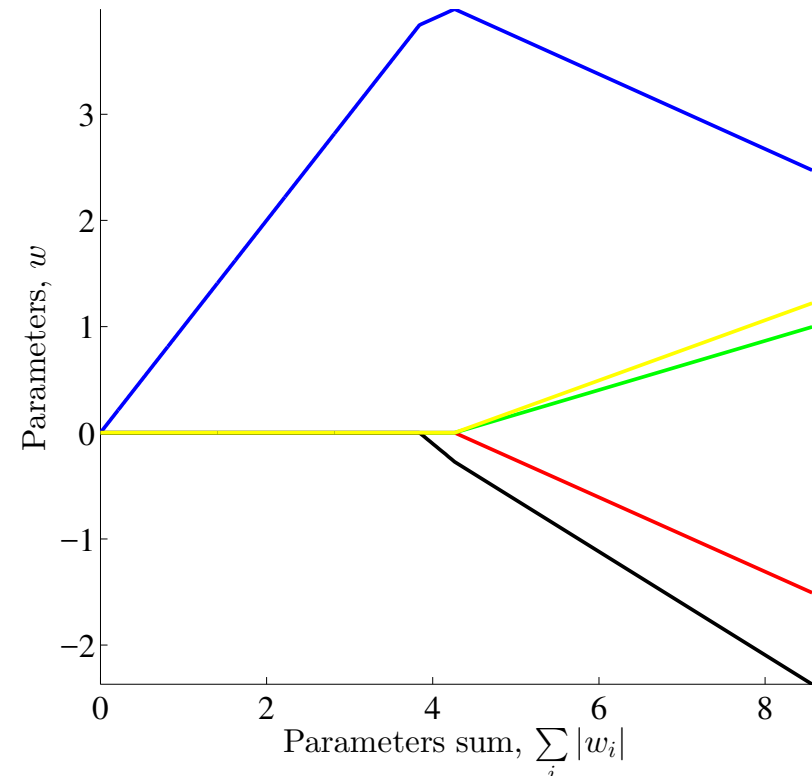
Model parameters with regularization

Ridge regression



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

Lasso



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

Discrete genetic algorithm for feature selection (simple ver.)

- 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:

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

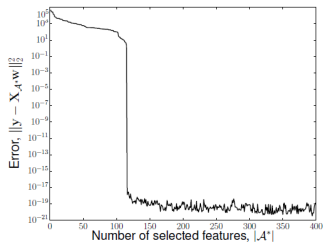
$$[a_{q,1}, \dots, a_{q,\nu}, a_{p,\nu+1}, \dots, a_{p,n}] \rightarrow \mathbf{a}'_q;$$

- 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.

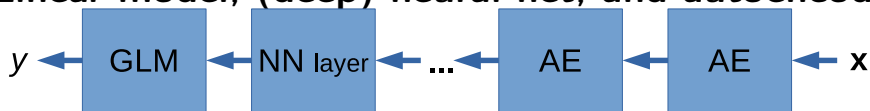
Evaluation criteria for the NIR spectra data set

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).

Linear model, (deep) neural net, and autoencoder



$$f = \sigma_k \circ \underbrace{\mathbf{w}_k^T}_{1 \times 1} \sigma_{k-1} \circ \mathbf{W}_{k-1} \sigma_{k-2} \circ \dots \circ \underbrace{\mathbf{W}_2 \sigma_1 \circ \mathbf{W}_1}_{\substack{n_2 \times 1 \\ n_1 \times n \quad n \times 1}} \mathbf{x} \in \mathcal{D}$$

$$E_x = \sum_{\mathbf{x}_i \in \mathcal{D}} \|\mathbf{x}_i - \mathbf{r}(\mathbf{x}_i)\|_2^2$$

$$E_D = \sum_{(\mathbf{x}_i, y_i) \in \mathcal{D}} (y_i - f(\mathbf{x}_i))^2$$

$$S = \lambda_1 E_D + \lambda_2 E_x + \lambda_3 E_w = \lambda^T \mathbf{s}$$

E_w is some regularisation error, for

principal component analysis: $\mathbf{W}^T \mathbf{W} = \mathbf{I}_n$,

skip block: $\mathbf{W} = \mathbf{I}_n$, $\sigma = \text{id}$,

classification: $\sigma \in \{\text{logistic}, \text{softmax}, \text{ReLU}, \dots\}$.

... including LM, LR, PCA, AE, SAE, 2NN, DLL, CNN, etc.

Probabilistic model selection

Bayesian inference delivers the error function $S(\mathbf{w})$

$$p(\mathbf{w}|\mathcal{D}, \mathbf{A}, \mathbf{B}, \mathbf{f}) = \frac{\overset{\text{Likelihood}}{p(\mathcal{D}|\mathbf{w}, \mathbf{B}, \mathbf{f})} \overset{\text{Prior}}{p(\mathbf{w}|\mathbf{A}, \mathbf{f})}}{\underset{\substack{\text{Evidence} \\ (\text{to select a model})}}{p(\mathcal{D}|\mathbf{A}, \mathbf{B}, \mathbf{f})}}.$$

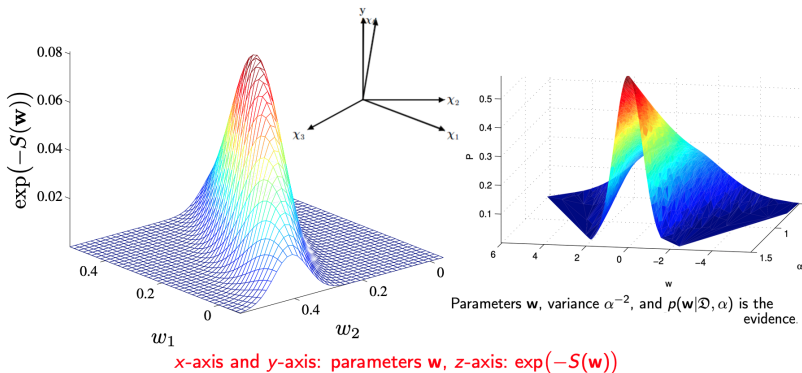
Write the error function given hyperparameters \mathbf{A}, \mathbf{B}

$$S(\mathbf{w}) = \underbrace{\frac{1}{2}(\mathbf{y} - \mathbf{f})^T \mathbf{B}(\mathbf{y} - \mathbf{f})}_{\text{approximation error}} + \underbrace{\frac{1}{2}(\mathbf{w} - \hat{\mathbf{w}})^T \mathbf{A}(\mathbf{w} - \hat{\mathbf{w}})}_{\text{regularisation error}},$$

$$S = E_D + E_w = \lambda^T s, \quad \text{metaparameters } \lambda = \frac{1}{2}.$$

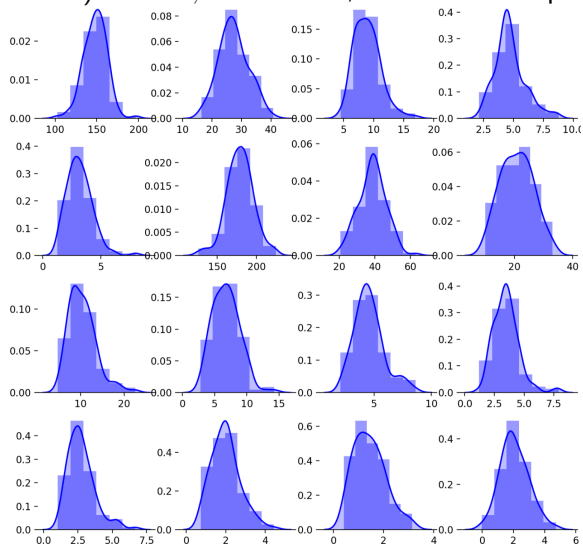
Empirical distribution of model parameters

The value of error function $S(\mathbf{w}|\mathcal{D}, f)$ depends on parameters.

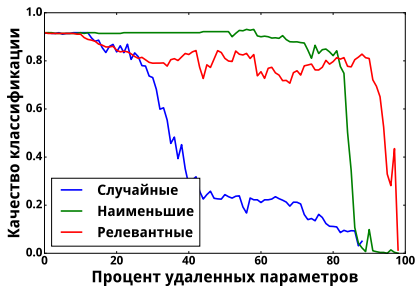


How to check the i.i.d hypothesis

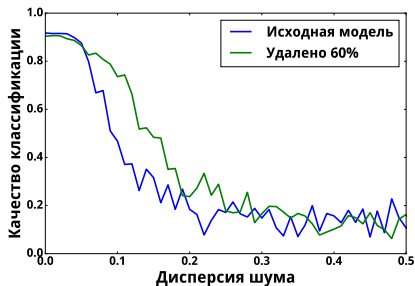
T-test) $E\varepsilon = 0, D\varepsilon = \text{const}$, as well as the spectrum analysis



Правдоподобие моделей с избыточным числом параметров не изменяется значительно при их удалении



Избыточность параметров модели



Устойчивость модели

Глубокое обучение предполагает оптимизацию моделей с заведомо избыточной сложностью.

Bakhteev, Strijov. 2019. Comprehensive analysis of gradient-based hyperparameter optimization algorithms // Annals of Operations Research

Двоичное представление структуры модели

Модель f выбирается из множества моделей-претендентов \mathfrak{F} путем оптимизации двоичного вектора $\mathbf{a} \in \mathbb{B}^n$,

$$\hat{y} = f(\mathbf{w}, \mathbf{x}) = a_1 w_1 x_1 + \dots + a_n w_n x_n$$

для линейной модели $f(\mathbf{w}, \mathbf{x}) = \mathbf{x}^T \mathbf{w}$
и для нейронной сети

$$\mathbf{f}(\mathbf{w}, \mathbf{x}) = \frac{\exp(\mathbf{h}(\mathbf{x}))}{\sum_k \exp(h_k(\mathbf{x}))}, \quad \mathbf{h}(\mathbf{x}) = \mathbf{W}_2^T \tanh(\mathbf{W}_1^T \mathbf{x}), \quad \mathbf{w} = \text{vec}(\mathbf{W}_1 : \mathbf{W}_2)$$

путем зануления соответствующего параметра

$$w_j = 0$$

или согласно методу оптимального прореживания

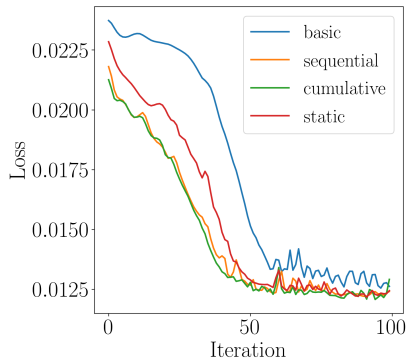
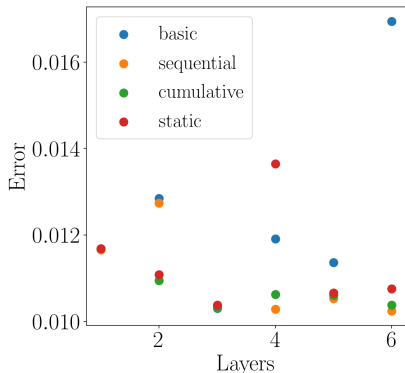
$$\mathbf{e}_j^T \Delta \mathbf{w} + w_j = 0$$

где j -й элемент вектора \mathbf{e} равен 1, прочие равны 0.

Модель задана вершиной двоичного n -мерного куба.

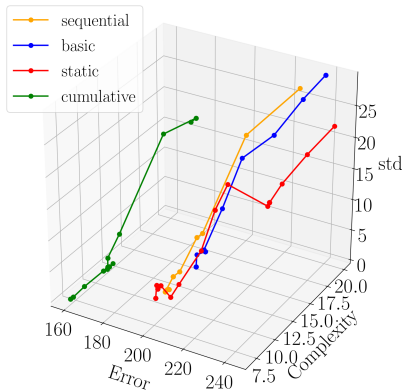
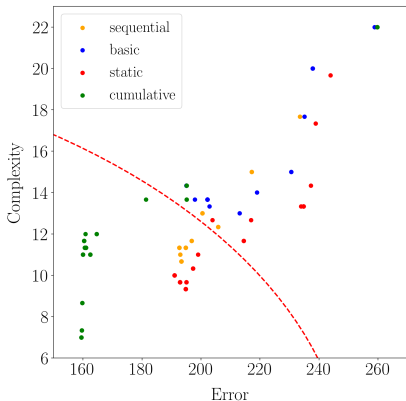
Процедура оптимизации параметров

Результаты работы алгоритма выбора расписания оптимизации. Слева — зависимость ошибки нейросети на тестовой выборке от используемого метода задания весов регуляризации и количества слоев. Справа — ошибка на валидационной выборке во время обучения.



Использование регуляризации приводит к более быстрой сходимости и меньшей ошибке.

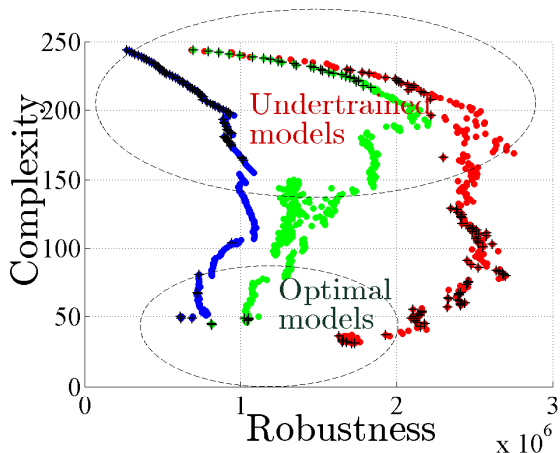
Процедура оптимизации структуры



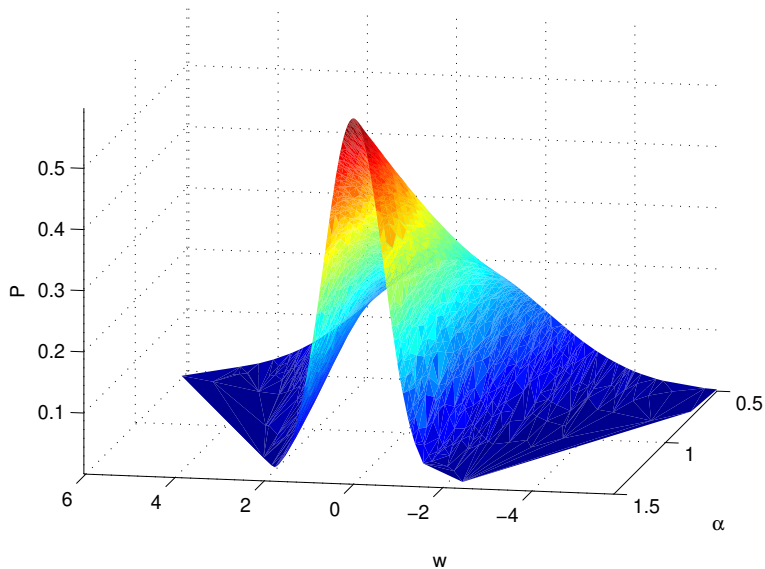
Изменение точности, сложности и устойчивости моделей при итерациях генетического алгоритма

Последовательный выбор моделей:

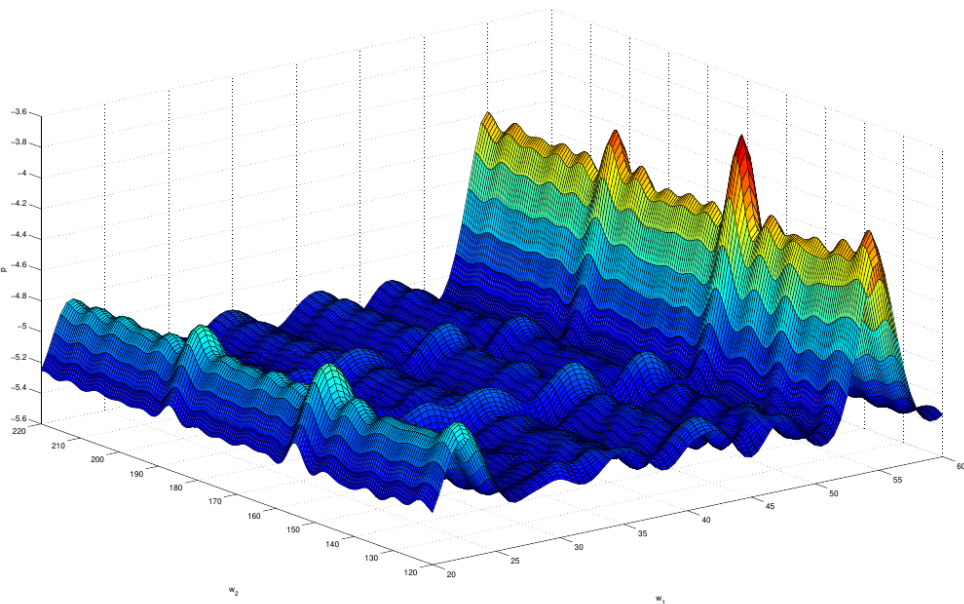
точность, сложность, устойчивость



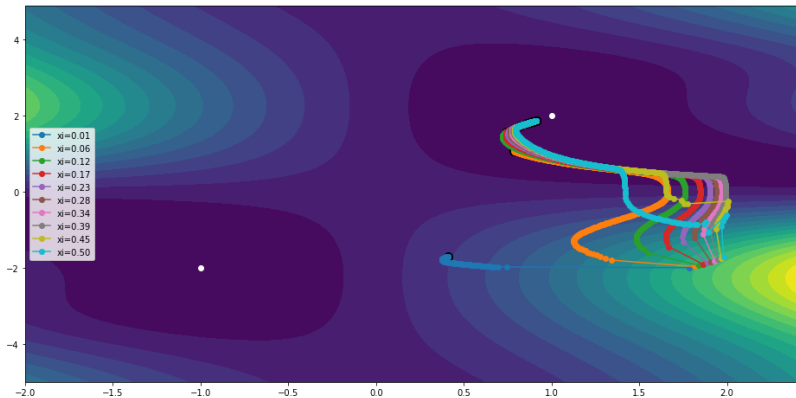
Точность или устойчивость



Многоэкстремальность функции ошибки

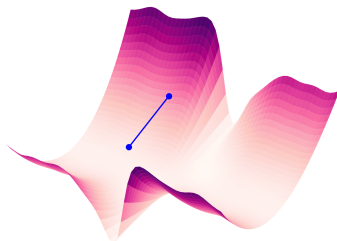
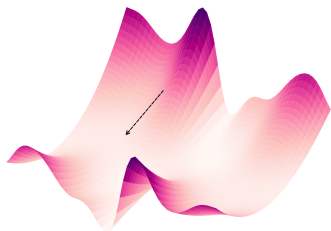


Многоэкстремальность, сходимость и мультистарт

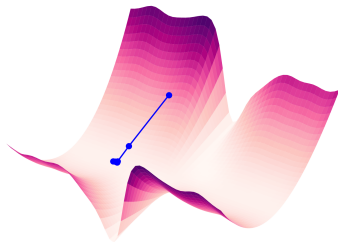
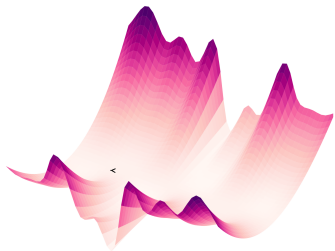


$$\mathbf{w} \in \mathbb{R}^2$$

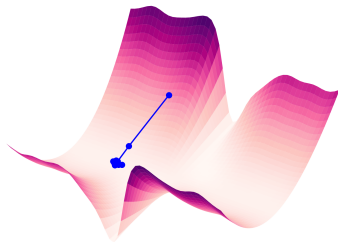
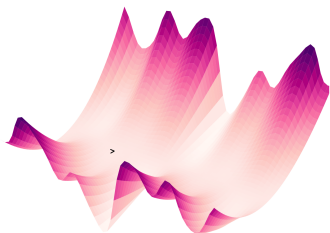
Стохастический градиент и ковариация параметров



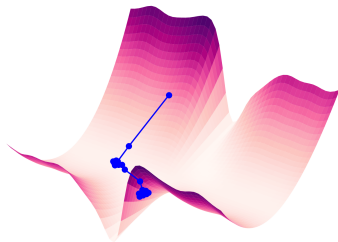
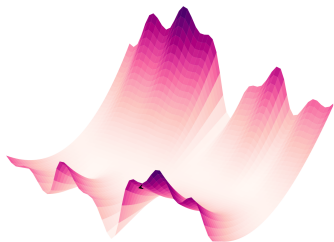
Стохастический градиент и ковариация параметров



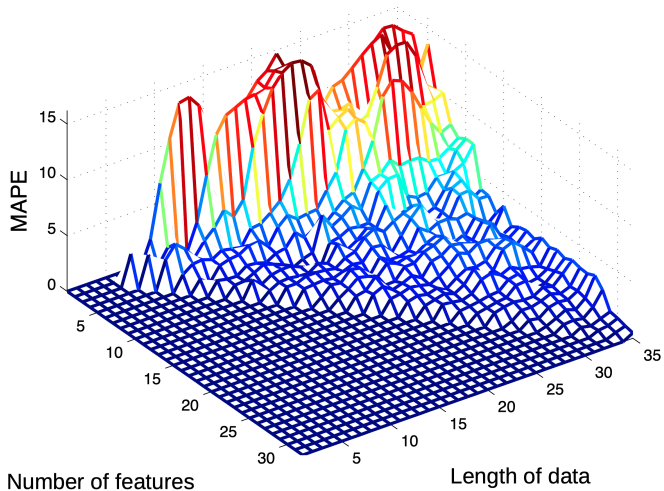
Стохастический градиент и ковариация параметров



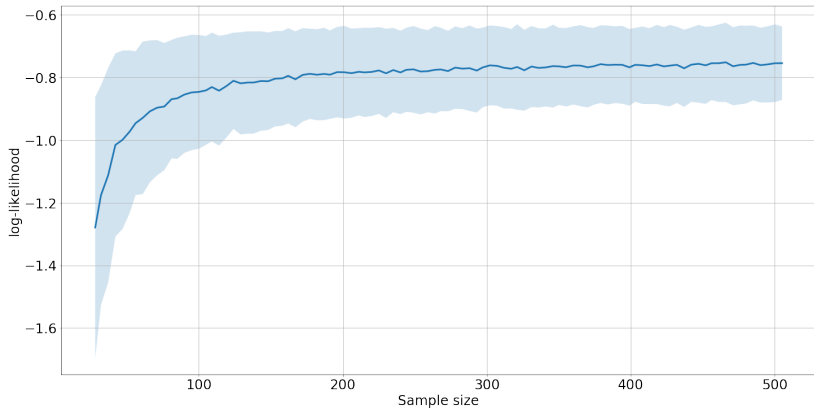
Стохастический градиент и ковариация параметров



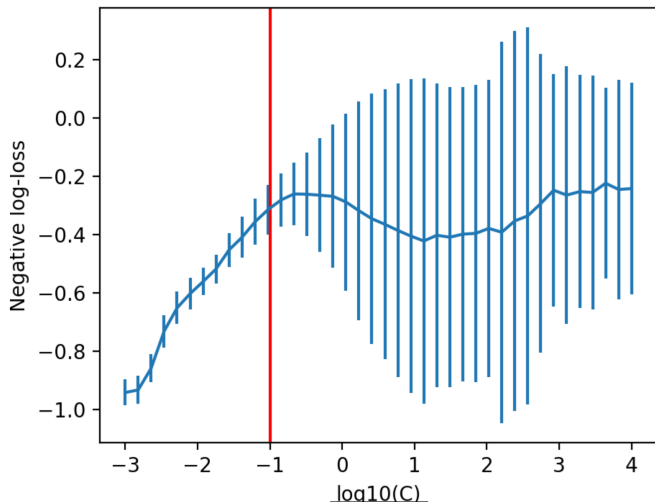
Ошибка (переобученной!) модели



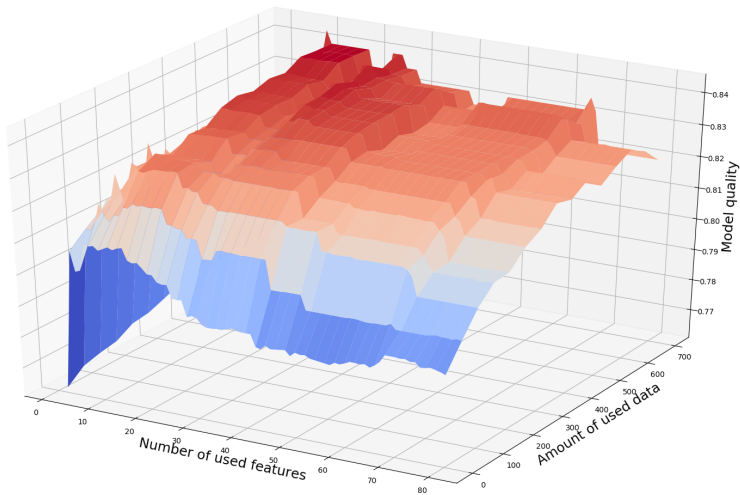
– Ошибка и её дисперсия при пополнении выборки



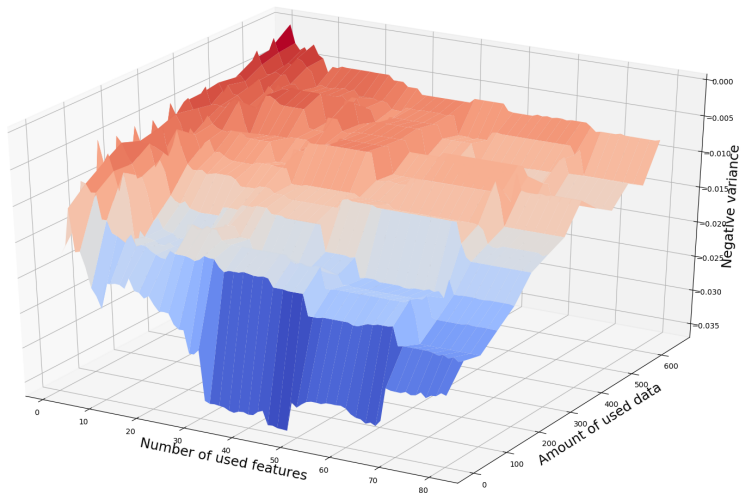
Дисперсия ошибки при повышении сложности модели



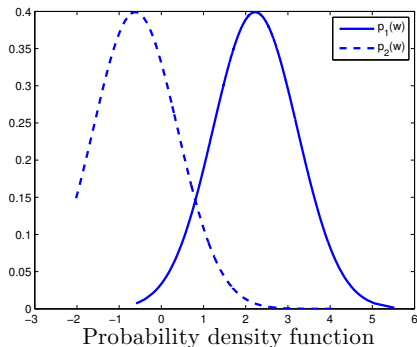
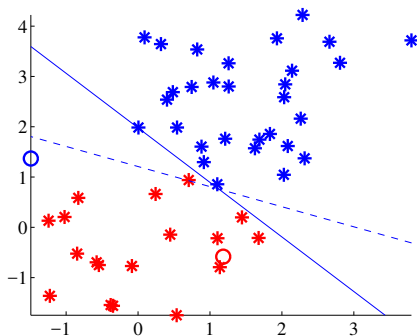
– Ошибка при различных объемах выборки



– Дисперсия ошибки при различных объемах выборки

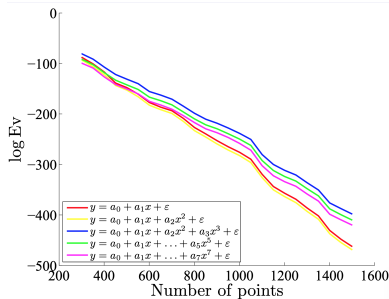
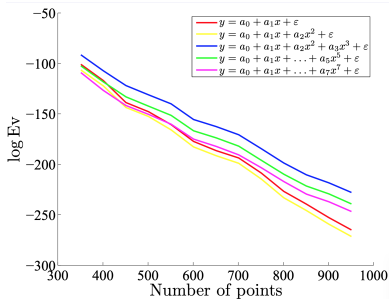
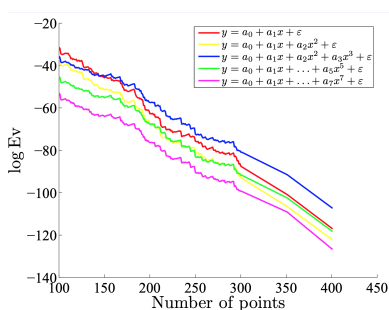
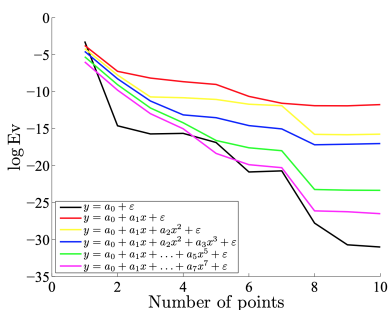


Изменение эмпирического распределения параметров

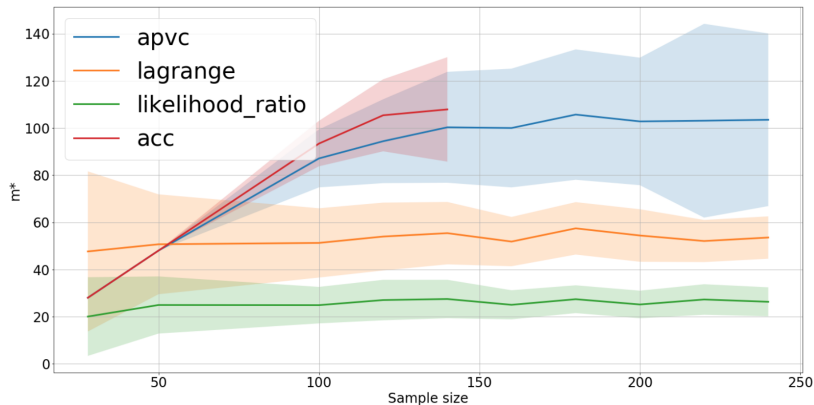


Объем выборки m^* из распределения P достаточен, если выборки X_1, X_2 размера $m > m^*$ из P схожи согласно функции сходства $D(\hat{P}_1, \hat{P}_2)$ между эмпирическими распределениями, полученными на этих выборках.

Правдоподобные модели при разных объемах выборки



Объем выборки, спрогнозированной на раннем этапе сбора данных



Имея выборку объема t требуется спрогнозировать оптимальный объем m^* .