

Significant increase in complexity and modest increase in accuracy

	train	test	out-of-time	Number of 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%	> 1 000
Gradient boosting	63,58%	58,31%	59,50%	> 10 000

... it was a banking credit scoring model

To start an applied project **an expert** and **an analyst** set

1. Project goal (**the expected result of development**)
main purpose of research
2. Project application (**how the project result will be applied**)
environment of measures and impacts
3. Historical data description (**data formats and timing**)
algebraic structures of data
4. Quality criteria (**how the project quality is measured**)
error function
5. Feasibility of the project (**how to prove the project feasibility, list possible risks**)
error analysis

How long the model lives after being put on operation? What replaces it after?

Problem statement for machine learning

Formal problem statement, **an analyst has to set**

- 1) an algebraic structure for the dataset from measurements
- 2) a data generation hypothesis from 1)
- 3) a model, or a mixture from 2)
- 4) an error function (quality criteria with restrictions) from 2)
- 5) an optimization algorithm from 3) and 4)

The result of the model construction is a Cartesian product

{models \times datasets \times quality criteria}.

Def: Big data rejects the i.i.d. (independent and identically distributed random variables) data generation hypothesis from 2). It requests a mixture model.

Analyst creates an optimal model for expert to put it to operation

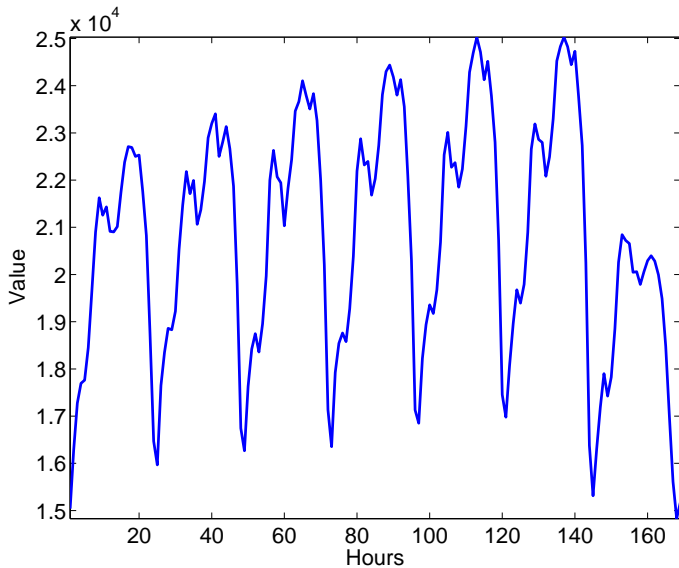
Origins of quality criteria

1. **Theory**: statistical hypotheses of data generation, algebraic structures of data, models of measurement
2. **Computations**: a criterion is useful to an optimisation procedure
3. **Deployment**: revenue, loss, failure rate

Quality criteria

- ▶ **Accuracy**: MAPE, AUC, F1 score
- ▶ **Stability**: forecasting variance, failure rate, parameter variance
- ▶ **Complexity**: number of parameters, Kolmogorov complexity

Source time series, one week



The autoregressive matrix to forecast periodic time series

- There given the time series $\{s_1, \dots, s_T, \dots, s_{T-1}\}$, the length of a period is κ .
- One must to forecast the next sample T .
- The autoregressive matrix:
 - its i -th row is a period of samples,
 - its j -th column is a phase of the period and
 - they map into the time series sample number such that $(i-1)\kappa \mapsto \tau$; let $\text{mod } \frac{T}{\kappa} = 0$;

$$X^*_{(m+1) \times (n+1)} = \begin{pmatrix} s_T & s_{T-1} & \dots & s_{T-\kappa+1} \\ s_{(m-1)\kappa} & s_{(m-1)\kappa-1} & \dots & s_{(m-2)\kappa+1} \\ \dots & \dots & \dots & \dots \\ s_{n\kappa} & s_{n\kappa-1} & \dots & s_{n(\kappa-1)+1} \\ \dots & \dots & \dots & \dots \\ s_\kappa & s_{\kappa-1} & \dots & s_1 \end{pmatrix}.$$

$$X^*_{(m+1) \times (n+1)} = \begin{pmatrix} S_T & S_{T-1} & \dots & S_{T-\kappa+1} \\ S_{(m-1)\kappa} & S_{(m-1)\kappa-1} & \dots & S_{(m-2)\kappa+1} \\ \dots & \dots & \dots & \dots \\ S_{n\kappa} & S_{n\kappa-1} & \dots & S_{n(\kappa-1)+1} \\ \dots & \dots & \dots & \dots \\ S_\kappa & S_{\kappa-1} & \dots & S_1 \end{pmatrix}.$$

In a nutshell,

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

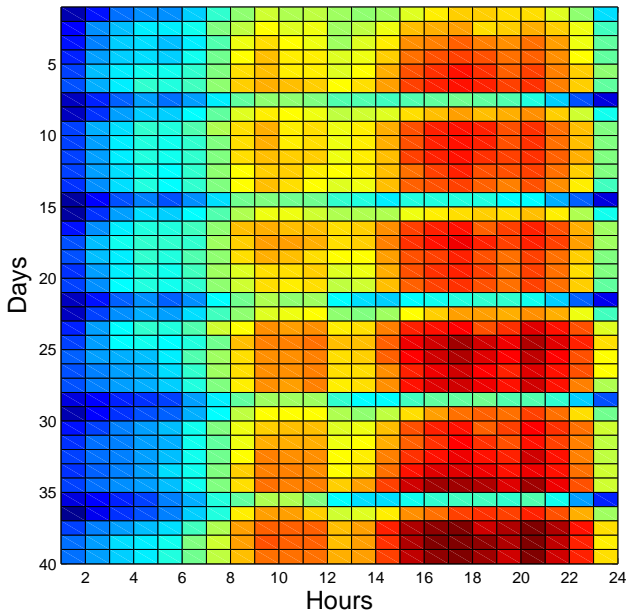
$\begin{matrix} 1 \times 1 & 1 \times n \\ m \times 1 & m \times n \end{matrix}$

In terms of linear regression:

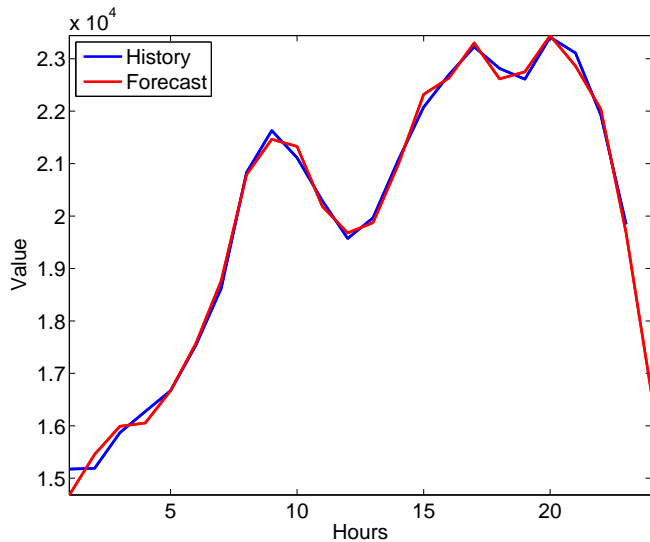
$$\mathbf{y} = X\mathbf{w},$$

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

The autoregressive matrix, five week-ends



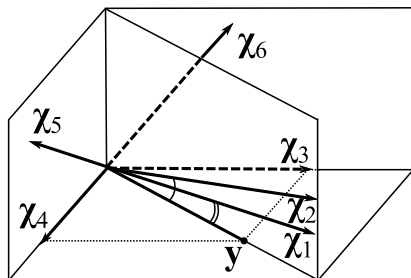
The one-day forecast, an example



Select an accurate and stable set of features

Features χ_1, \dots, χ_6 are columns of the design matrix \mathbf{X} .
 3×6

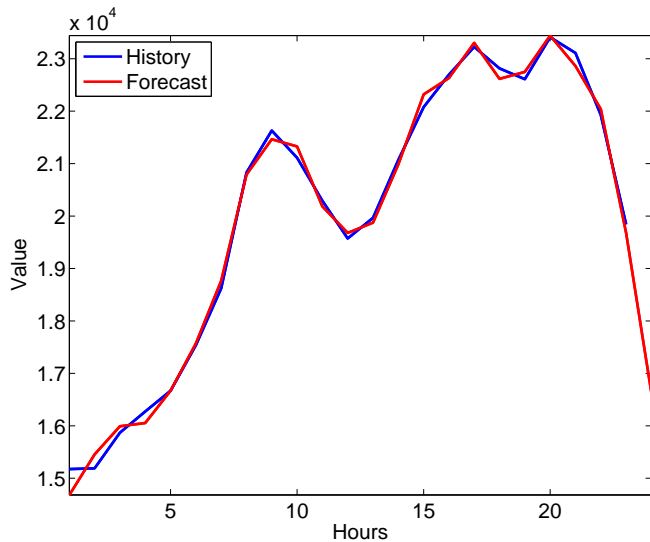
The sample contains multicollinear χ_1, χ_2 and noisy χ_5, χ_6 features, columns of the design matrix \mathbf{X} . One has to select two features from six.



Solution: χ_3, χ_4 are orthogonal; their linear combination fits \mathbf{y} .

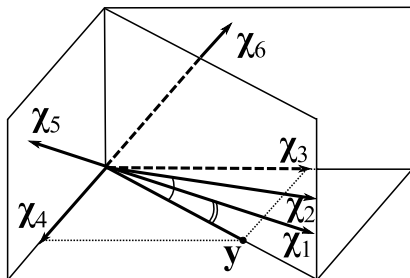
Katrutsa, Strijov. 2015. Stress-test procedure for feature selection // Chemometrics

The one-day forecast, an example



Выбор устойчивой и точной модели

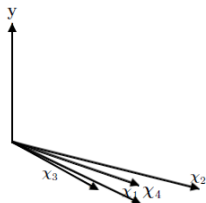
Выборка содержит мультикоррелирующие χ_1, χ_2 и устойчивые χ_5, χ_6 признаки — столбцы матрицы «объект-признак» \mathbf{X} . Требуется выбрать два признака из шести.



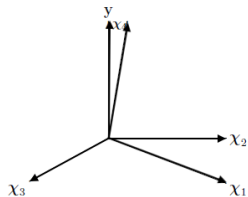
Точность и устойчивость при заданной сложности

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

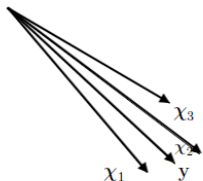
Configurations of design space



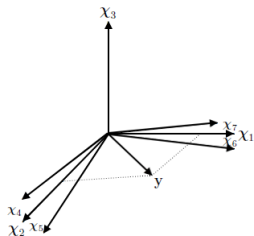
Non-adequate correlated



Adequate random



Adequate redundant



Adequate correlated

Katrutsa, Strijov. 2017. Comprehensive study of feature selection methods to solve multicollinearity problem // Expert Systems with Applications

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.

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 .

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} .