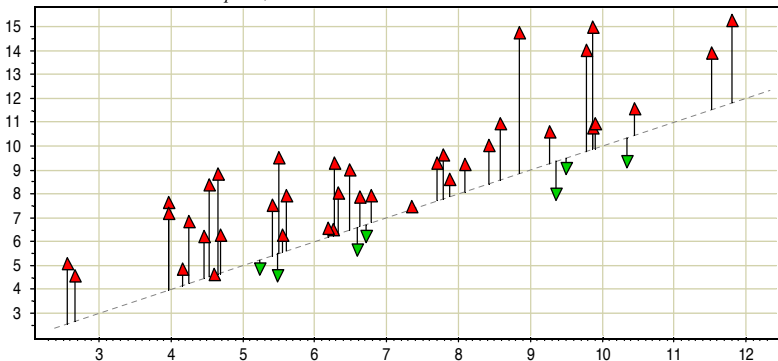


**Задача** предсказания отдалённого результата хирургического лечения атеросклероза. Точки — различные алгоритмы.

*Частота ошибок на контроле, %*



*Частота ошибок на обучении, %*

Зависимость  $y(x) = \frac{1}{1 + 25x^2}$  на отрезке  $x \in [-2, 2]$ .

Признаковое описание  $x \mapsto (1, x^1, x^2, \dots, x^n)$ .

Модель полиномиальной регрессии

$$a(x, \theta) = \theta_0 + \theta_1 x + \dots + \theta_n x^n \text{ — полином степени } n.$$

Обучение методом наименьших квадратов:

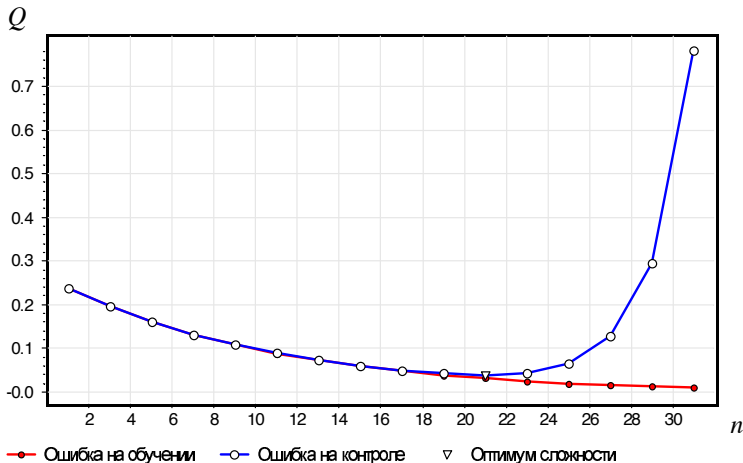
$$Q(a, X^\ell) = \sum_{i=1}^{\ell} (\theta_0 + \theta_1 x_i + \dots + \theta_n x_i^n - y_i)^2 \rightarrow \min_{\theta_0, \dots, \theta_n}.$$

Обучающая выборка:  $X^\ell = \{x_i = 4\frac{i-1}{\ell-1} - 2 \mid i = 1, \dots, \ell\}$ .

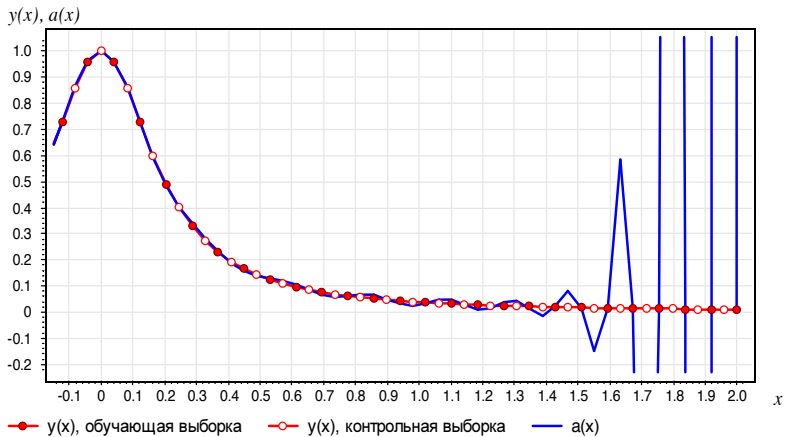
Контрольная выборка:  $X^k = \{x_i = 4\frac{i-0.5}{\ell-1} - 2 \mid i = 1, \dots, \ell - 1\}$ .

Что происходит с  $Q(a, X^\ell)$  и  $Q(a, X^k)$  при увеличении  $n$ ?

Переобучение — это когда  $Q(\mu(X^\ell), X^k) \gg Q(\mu(X^\ell), X^\ell)$ :



$$y(x) = \frac{1}{1 + 25x^2}; \quad a(x) \text{ — полином степени } n = 38$$



- Эмпирический риск на тестовых данных (hold-out):

$$\text{HO}(\mu, X^\ell, X^k) = Q(\mu(X^\ell), X^k) \rightarrow \min$$

- Скользящий контроль (leave-one-out),  $L = \ell + 1$ :

$$\text{LOO}(\mu, X^L) = \frac{1}{L} \sum_{i=1}^L \mathcal{L}(\mu(X^L \setminus \{x_i\}), x_i) \rightarrow \min$$

- Кросс-проверка (cross-validation) по  $N$  разбиениям,  $X^L = X_n^\ell \sqcup X_n^k$ ,  $L = \ell + k$ :

$$\text{CV}(\mu, X^L) = \frac{1}{N} \sum_{n=1}^N Q(\mu(X_n^\ell), X_n^k) \rightarrow \min$$

## Эксперименты на конкретной прикладной задаче:

- цель — решить задачу как можно лучше
- важно понимание задачи и данных
- важно придумывать информативные признаки
- конкурсы по анализу данных: <http://www.kaggle.com>

## Эксперименты на наборах прикладных задач:

- цель — протестировать метод в разнообразных условиях
- нет необходимости (и времени) разбираться в сути задач : (
- признаки, как правило, уже кем-то придуманы
- репозиторий UC Irvine Machine Learning Repository  
<http://archive.ics.uci.edu/ml> (308 задач, 09-02-2015)

Используются для тестирования новых методов обучения.  
Преимущество — мы знаем истинную  $y(x)$  (ground truth)

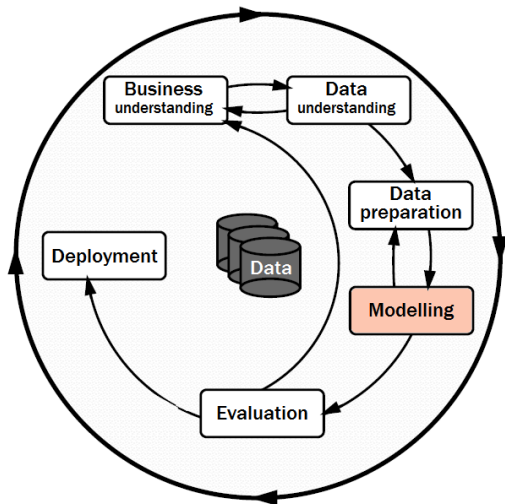
### Эксперименты на модельных (synthetic) данных:

- цель — отладить метод, выявить границы применимости
- объекты  $x_i$  из придуманного распределения (часто 2D)
- ответы  $y_i = y(x_i)$  для придуманной функции  $y(x)$
- двумерные данные + визуализация выборки

### Эксперименты на полумодельных (semi-synthetic) данных:

- цель — протестировать помехоустойчивость модели
- объекты  $x_i$  из реальной задачи (+ шум)
- ответы  $y_i = a(x_i)$  для полученного решения  $a(x)$  (+ шум)

CRISP-DM — межотраслевой стандарт решения задач интеллектуального анализа данных





## Этапы решения задач машинного обучения:

- понимание задачи и данных;
- предобработка данных и изобретение признаков;
- построение модели;
- сведение обучения к оптимизации;
- решение проблем оптимизации и переобучения;
- оценивание качества решения;
- внедрение и эксплуатация.