

# Нелинейная регрессия. Непараметрическая регрессия. Нестандартные функции потерь

К. В. Воронцов  
vokov@forecsys.ru

Этот курс доступен на странице вики-ресурса  
<http://www.MachineLearning.ru/wiki>  
«Машинное обучение (курс лекций, К.В.Воронцов)»  
Видеолекции: <http://shad.yandex.ru/lectures>

27 марта 2014

## Содержание

- 1 Нелинейная регрессия**
  - Нелинейная модель регрессии
  - Логистическая регрессия
  - Нелинейные преобразования признаков
- 2 Непараметрическая регрессия**
  - Формула Надарая–Ватсона
  - Выбор ядра  $K$  и ширины окна  $h$
  - Отсев выбросов
- 3 Неквадратичные функции потерь**
  - Квантильная регрессия
  - Робастная регрессия
  - SVM-регрессия

## Метод наименьших квадратов

- $X$  — объекты (часто  $\mathbb{R}^n$ );  $Y$  — ответы (часто  $\mathbb{R}$ , реже  $\mathbb{R}^m$ );  
 $X^\ell = (x_i, y_i)_{i=1}^\ell$  — обучающая выборка;  
 $y_i = y(x_i)$ ,  $y: X \rightarrow Y$  — неизвестная зависимость;
- $a(x) = f(x, \alpha)$  — модель зависимости,  
 $\alpha \in \mathbb{R}^p$  — вектор параметров модели.
- Метод наименьших квадратов (МНК):

$$Q(\alpha, X^\ell) = \sum_{i=1}^{\ell} w_i (f(x_i, \alpha) - y_i)^2 \rightarrow \min_{\alpha},$$

где  $w_i$  — вес, степень важности  $i$ -го объекта.

$Q(\alpha^*, X^\ell)$  — остаточная сумма квадратов  
(residual sum of squares, RSS).

## Нелинейная модель регрессии

Нелинейная модель регрессии  $f(x, \alpha)$ ,  $\alpha \in \mathbb{R}^p$ .

Функционал среднеквадратичного отклонения:

$$Q(\alpha, X^\ell) = \sum_{i=1}^{\ell} (f(x_i, \alpha) - y_i)^2 \rightarrow \min_{\alpha}.$$

**Метод Ньютона–Рафсона.**

1. Начальное приближение  $\alpha^0 = (\alpha_1^0, \dots, \alpha_p^0)$ .
2. Итерационный процесс

$$\alpha^{t+1} := \alpha^t - \eta_t (Q''(\alpha^t))^{-1} Q'(\alpha^t),$$

$Q'(\alpha^t)$  — градиент функционала  $Q$  в точке  $\alpha^t$ ,

$Q''(\alpha^t)$  — гессиан функционала  $Q$  в точке  $\alpha^t$ ,

$\eta_t$  — величина шага (можно полагать  $\eta_t = 1$ ).

## Метод Ньютона-Рафсона

Компоненты градиента:

$$\frac{\partial Q(\alpha)}{\partial \alpha_j} = 2 \sum_{i=1}^{\ell} (f(x_i, \alpha) - y_i) \frac{\partial f(x_i, \alpha)}{\partial \alpha_j}.$$

Компоненты гессиана:

$$\frac{\partial^2 Q(\alpha)}{\partial \alpha_j \partial \alpha_k} = 2 \sum_{i=1}^{\ell} \frac{\partial f(x_i, \alpha)}{\partial \alpha_j} \frac{\partial f(x_i, \alpha)}{\partial \alpha_k} - 2 \underbrace{\sum_{i=1}^{\ell} (f(x_i, \alpha) - y_i) \frac{\partial^2 f(x_i, \alpha)}{\partial \alpha_j \partial \alpha_k}}_{\text{при линейризации полагается} = 0}.$$

Не хотелось бы обращать гессиан на каждой итерации...

Линеаризация  $f(x_i, \alpha)$  в окрестности текущего  $\alpha^t$ :

$$f(x_i, \alpha) = f(x_i, \alpha^t) + \sum_{j=1}^p \frac{\partial f(x_i, \alpha_j)}{\partial \alpha_j} (\alpha_j - \alpha_j^t) + o(\alpha_j - \alpha_j^t).$$

## Метод Ньютона-Гаусса

Матричные обозначения:

$F_t = \left( \frac{\partial f}{\partial \alpha_j}(x_i, \alpha^t) \right)_{i=1, \ell}^{j=1, p}$  —  $\ell \times p$ -матрица первых производных;

$f_t = (f(x_i, \alpha^t))_{i=1, \ell}$  — вектор значений  $f$ .

Формула  $t$ -й итерации метода Ньютона-Гаусса:

$$\alpha^{t+1} := \alpha^t - h_t \underbrace{(F_t^T F_t)^{-1} F_t^T (f^t - y)}_{\beta}.$$

$\beta$  — это решение задачи многомерной линейной регрессии

$$\|F_t \beta - (f^t - y)\|^2 \rightarrow \min_{\beta}.$$

Нелинейная регрессия сведена к серии линейных регрессий.

Скорость сходимости — как и у метода Ньютона-Рафсона, но для вычислений можно применять стандартные методы.

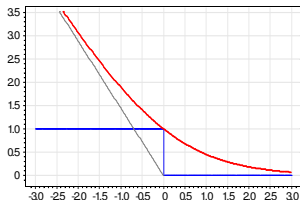
## Задача классификации. Логистическая регрессия

$Y = \{-1, +1\}$  — два класса,  $a(x, w) = \text{sign}(w^T x)$ ,  $x, w \in \mathbb{R}^n$ .

Функционал аппроксимированного эмпирического риска:

$$Q(w) = \sum_{i=1}^{\ell} [M_i(w) < 0] \leq \sum_{i=1}^{\ell} \mathcal{L}(w^T x_i y_i) \rightarrow \min_w,$$

где  $\mathcal{L}(M) = \log(1 + e^{-M})$  — логарифмическая функция потерь



$$M_i = w^T x_i y_i$$

## Метода Ньютона-Рафсона

Метода Ньютона-Рафсона для минимизации функционала  $Q(w)$ :

$$w^{t+1} := w^t - h_t(Q''(w^t))^{-1}Q'(w^t),$$

Элементы градиента — вектора первых производных  $Q'(w^t)$ :

$$\frac{\partial Q(w)}{\partial w_j} = - \sum_{i=1}^{\ell} (1 - \sigma_i) y_i f_j(x_i), \quad j = 1, \dots, n.$$

Элементы гессиана — матрицы вторых производных  $Q''(w^t)$ :

$$\frac{\partial^2 Q(w)}{\partial w_j \partial w_k} = \sum_{i=1}^{\ell} (1 - \sigma_i) \sigma_i f_j(x_i) f_k(x_i), \quad j, k = 1, \dots, n,$$

где  $\sigma_i = \sigma(y_i w^T x_i)$ ,  $\sigma(z) = \frac{1}{1+e^{-z}}$  — сигмоидная функция.

$\sigma_i = P(y_i|x_i)$  — вероятность правильной классификации  $x_i$ .



## Матричные обозначения

$F_{\ell \times n} = (f_j(x_i))$  — матрица «объекты–признаки»;

$\Gamma_{\ell \times \ell} = \text{diag}(\sqrt{(1 - \sigma_i)\sigma_i})$  — диагональная матрица;

$\tilde{F} = \Gamma F$  — взвешенная матрица «объекты–признаки»;

$\tilde{y}_i = y_i \sqrt{(1 - \sigma_i)/\sigma_i}$ ,  $\tilde{y} = (\tilde{y}_i)_{i=1}^{\ell}$  — взвешенный вектор ответов.

Тогда в методе Ньютона-Рафсона:

$$(Q''(w))^{-1} Q'(w) = -(F^T \Gamma^2 F)^{-1} F^T \Gamma \tilde{y} = -(\tilde{F}^T \tilde{F})^{-1} \tilde{F}^T \tilde{y} = -\tilde{F}^+ \tilde{y}.$$

Это совпадает с МНК-решением линейной задачи регрессии со взвешенными объектами и модифицированными ответами:

$$Q(w) = \|\tilde{F}w - \tilde{y}\|^2 = \sum_{i=1}^{\ell} (1 - \sigma_i)\sigma_i \left( w^T x - \frac{y_i}{\sigma_i} \right)^2 \rightarrow \min_w.$$

## МНК с итерационным перевзвешиванием объектов IRLS — Iteratively Reweighted Least Squares

**Вход:**  $F, y$  — матрица «объекты–признаки» и вектор ответов;

**Выход:**  $w$  — вектор коэффициентов линейной комбинации.

- 
- 1:  $w := (F^T F)^{-1} F^T y$  — нулевое приближение, обычный МНК;
  - 2: **для**  $t := 1, 2, 3, \dots$
  - 3:  $\sigma_i = \sigma(y_i w^T x_i)$  для всех  $i = 1, \dots, \ell$ ;
  - 4:  $\gamma_i := \sqrt{(1 - \sigma_i) \sigma_i}$  для всех  $i = 1, \dots, \ell$ ;
  - 5:  $\tilde{F} := \text{diag}(\gamma_1, \dots, \gamma_\ell) F$ ;
  - 6:  $\tilde{y}_i := y_i \sqrt{(1 - \sigma_i) / \sigma_i}$  для всех  $i = 1, \dots, \ell$ ;
  - 7: выбрать градиентный шаг  $h_t$ ;
  - 8:  $w := w + h_t (\tilde{F}^T \tilde{F})^{-1} \tilde{F}^T \tilde{y}$ ;
  - 9: **если**  $\{\sigma_i\}$  мало изменились **то** выйти из цикла;

## Обобщение линейной модели регрессии

Пусть  $\varphi_j: \mathbb{R} \rightarrow \mathbb{R}$  — некоторые нелинейные преобразования исходных признаков. Модель регрессии:

$$f(x, \alpha) = \sum_{j=1}^n \varphi_j(f_j(x)).$$

В частности, при  $\varphi_j(f_j(x)) = \alpha_j f_j(x)$  это линейная регрессия.

**ИДЕЯ:** будем по очереди уточнять функции  $\varphi_j$  по обучающей выборке  $(f_j(x_i), z_i)_{i=1}^{\ell}$ :

$$Q(\varphi_j, X^{\ell}) = \sum_{i=1}^{\ell} \left( \varphi_j(f_j(x_i)) - \underbrace{\left( y_i - \sum_{k=1, k \neq j}^n \varphi_k(f_k(x_i)) \right)}_{z_i = \text{const}(\varphi_j)} \right)^2 \rightarrow \min_{\varphi_j}.$$

## Метод backfitting [Хасты, Тибширани, 1986]

**Вход:**  $F, y$  — матрица «объекты–признаки» и вектор ответов;

**Выход:**  $\varphi_j(x)$  — все функции преобразования признаков.

1: нулевое приближение:

$\alpha :=$  решение задачи МЛР с признаками  $f_j(x)$ ;

$\varphi_j(x) := \alpha_j f_j(x), j = 1, \dots, n$ ;

2: **повторять**

3: **для**  $j = 1, \dots, n$

4:  $z_i := y_i - \sum_{k=1, k \neq j}^n \varphi_k(f_k(x_i)), i = 1, \dots, \ell$ ;

5:  $\varphi_j := \arg \min_{\varphi} \sum_{i=1}^{\ell} (\varphi(f_j(x_i)) - z_i)^2$ ; — одномерная регрессия

6:  $Q_j := \sum_{i=1}^{\ell} (\varphi_j(f_j(x_i)) - z_i)^2$ ;

7: **пока** значения  $Q_j$  не стабилизируются

## Формула Надарая–Ватсона

Приближение константой  $a(x) = \alpha$  в окрестности  $x \in X$ :

$$Q(\alpha; X^\ell) = \sum_{i=1}^{\ell} w_i(x) (\alpha - y_i)^2 \rightarrow \min_{\alpha \in \mathbb{R}}$$

где  $w_i(x) = K\left(\frac{\rho(x, x_i)}{h}\right)$  — веса объектов  $x_i$  относительно  $x$ ;  
 $K(r)$  — ядро, невозрастающее, ограниченное, гладкое;  
 $h$  — ширина окна сглаживания.

**Формула ядерного сглаживания Надарая–Ватсона:**

$$a_h(x; X^\ell) = \frac{\sum_{i=1}^{\ell} y_i w_i(x)}{\sum_{i=1}^{\ell} w_i(x)} = \frac{\sum_{i=1}^{\ell} y_i K\left(\frac{\rho(x, x_i)}{h}\right)}{\sum_{i=1}^{\ell} K\left(\frac{\rho(x, x_i)}{h}\right)}.$$

## Обоснование формулы Надарая–Ватсона

### Теорема

Пусть выполнены следующие условия:

- 1) выборка  $X^\ell = (x_i, y_i)_{i=1}^\ell$  простая, из распределения  $p(x, y)$ ;
- 2) ядро  $K(r)$  ограничено:  $\int_0^\infty K(r) dr < \infty$ ,  $\lim_{r \rightarrow \infty} rK(r) = 0$ ;
- 3) зависимость  $E(y|x)$  не имеет вертикальных асимптот:  
 $E(y^2|x) = \int_Y y^2 p(y|x) dy < \infty$  при любом  $x \in X$ ;
- 4) последовательность  $h_\ell$  убывает, но не слишком быстро:  
 $\lim_{\ell \rightarrow \infty} h_\ell = 0$ ,  $\lim_{\ell \rightarrow \infty} \ell h_\ell = \infty$ .

Тогда имеет место сходимость по вероятности:

$$a_{h_\ell}(x; X^\ell) \xrightarrow{P} E(y|x) \text{ в любой точке } x \in X,$$

в которой  $E(y|x)$ ,  $p(x)$  и  $D(y|x)$  непрерывны и  $p(x) > 0$ .

- Ядро  $K(r)$ 
  - существенно влияет на гладкость функции  $a_h(x)$ ,
  - слабо влияет на качество аппроксимации.
- Ширина окна  $h$ 
  - существенно влияет на качество аппроксимации.
- При неравномерной сетке  $\{x_i\}$  — переменная ширина окна:

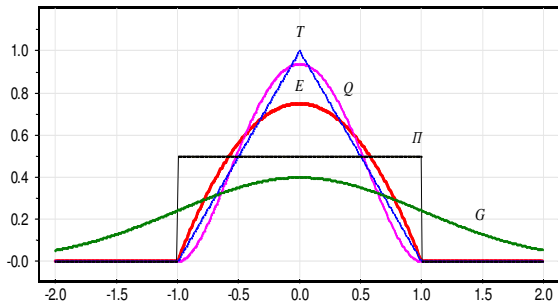
$$w_i(x) = K\left(\frac{\rho(x, x_i)}{h(x)}\right),$$

где  $h(x) = \rho(x, x^{(k+1)})$ ,  $x^{(k+1)}$  —  $k$ -й сосед объекта  $x$ .

- Оптимизация ширины окна по скользящему контролю:

$$\text{LOO}(h, X^\ell) = \sum_{i=1}^{\ell} \left( a_h(x_i; X^\ell \setminus \{x_i\}) - y_i \right)^2 \rightarrow \min_h.$$

## Часто используемые ядра $K(r)$



$E(r) = (1 - r^2) [ |r| \leq 1 ]$  — оптимальное (Епанечникова);

$Q(r) = (1 - r^2)^2 [ |r| \leq 1 ]$  — четвертое;

$T(r) = (1 - |r|) [ |r| \leq 1 ]$  — треугольное;

$G(r) = \exp(-\frac{1}{2}r^2)$  — гауссовское;

$\Pi(r) = [ |r| \leq 1 ]$  — прямоугольное.



## Локально взвешенное сглаживание (LOWESS — LOcally WEighted Scatter plot Smoothing)

### Основная идея:

чем больше величина ошибки  $\varepsilon_i = |a_h(x_i; X^\ell \setminus \{x_i\}) - y_i|$ , тем в большей степени прецедент  $(x_i, y_i)$  является выбросом, и тем меньше должен быть его вес  $w_i(x)$ .

### Эвристика:

домножить веса  $w_i(x)$  на коэффициенты  $\gamma_i = \tilde{K}(\varepsilon_i)$ ,  
где  $\tilde{K}$  — ещё одно ядро, вообще говоря, отличное от  $K(r)$ .

### Рекомендация:

квартичное ядро  $\tilde{K}(\varepsilon) = K_Q\left(\frac{\varepsilon}{6 \operatorname{med}\{\varepsilon_i\}}\right)$ ,  
где  $\operatorname{med}\{\varepsilon_i\}$  — медиана вариационного ряда ошибок.

## Алгоритм LOWESS

**Вход:**  $X^\ell$  — обучающая выборка;

**Выход:** коэффициенты  $\gamma_i$ ,  $i = 1, \dots, \ell$ ;

1: инициализация:  $\gamma_i := 1$ ,  $i = 1, \dots, \ell$ ;

2: **повторять**

3: **для всех** объектов  $i = 1, \dots, \ell$

4: вычислить оценки скользящего контроля:

$$a_i := a_h(x_i; X^\ell \setminus \{x_i\}) = \frac{\sum_{j=1, j \neq i}^{\ell} y_j \gamma_j K\left(\frac{\rho(x_i, x_j)}{h(x_i)}\right)}{\sum_{j=1, j \neq i}^{\ell} \gamma_j K\left(\frac{\rho(x_i, x_j)}{h(x_i)}\right)}$$

5: **для всех** объектов  $i = 1, \dots, \ell$

6:  $\gamma_i := \tilde{K}(|a_i - y_i|)$ ;

7: **пока** коэффициенты  $\gamma_i$  не стабилизируются;

## Метод наименьших модулей

$\varepsilon_i = (a(x_i) - y_i)$  — ошибка

$\mathcal{L}(\varepsilon_i)$  — функция потерь

$Q = \sum_{i=1}^{\ell} \mathcal{L}(\varepsilon_i) \rightarrow \min_a$  — критерий обучения модели по выборке

Метод наименьших квадратов,  $\mathcal{L}(\varepsilon) = \varepsilon^2$ :

$$\sum_{i=1}^{\ell} (a - y_i)^2 \rightarrow \min_a, \quad a = \frac{1}{\ell} \sum_{i=1}^{\ell} y_i.$$

Метод наименьших модулей,  $\mathcal{L}(\varepsilon) = |\varepsilon|$ :

$$\sum_{i=1}^{\ell} |a - y_i| \rightarrow \min_a, \quad a = \text{median}\{y_1, \dots, y_{\ell}\} = y^{(\ell/2)},$$

где  $y^{(1)}, \dots, y^{(\ell)}$  — вариационный ряд значений  $y_i$

## Квантильная регрессия

$$\text{Квантильная регрессия, } \mathcal{L}(\varepsilon) = \begin{cases} C_+|\varepsilon|, & \varepsilon > 0 \\ C_-|\varepsilon|, & \varepsilon < 0; \end{cases}$$

$$\sum_{i=1}^{\ell} \mathcal{L}(a - y_i) \rightarrow \min_a, \quad a = y^{(q)}, \quad q = \frac{\ell C_-}{C_- + C_+}$$

где  $y^{(1)}, \dots, y^{(\ell)}$  — вариационный ряд значений  $y_i$

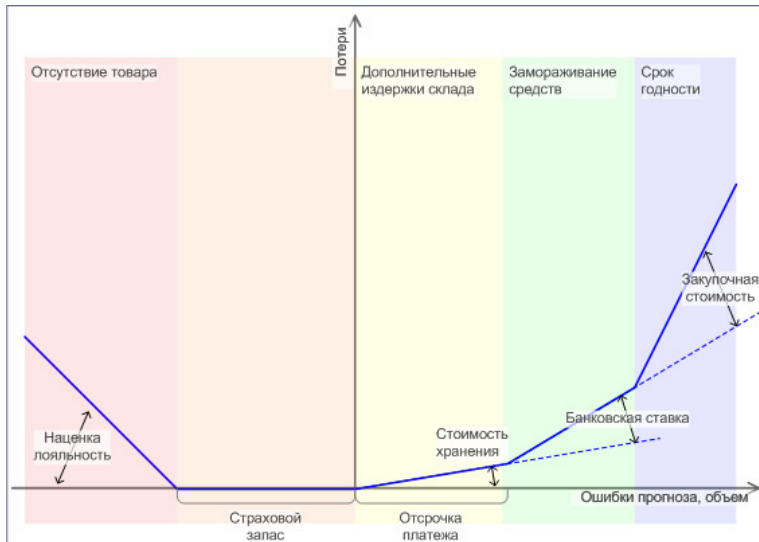
Линейная модель регрессии:  $a(x_i) = \langle x_i, w \rangle$ .

**Сведение к задаче линейного программирования:**

замена переменных  $\varepsilon_i^+ = (a(x_i) - y_i)_+$ ,  $\varepsilon_i^- = (y_i - a(x_i))_+$ ;

$$\begin{cases} Q = \sum_{i=1}^{\ell} C_+ \varepsilon_i^+ + C_- \varepsilon_i^- \rightarrow \min_w \\ \langle x_i, w \rangle - y_i = \varepsilon_i^+ - \varepsilon_i^-; \\ \varepsilon_i^+ \geq 0; \quad \varepsilon_i^- \geq 0. \end{cases}$$

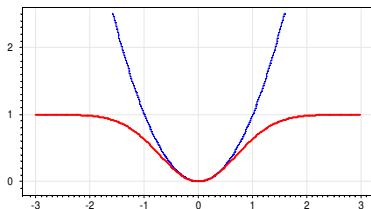
## Задача прогнозирования объёмов продаж



## Робастная регрессия

Модель регрессии:  $a(x) = f(x, \alpha)$

Функция Мешалкина:  $\mathcal{L}(\varepsilon) = 1 - \exp(\varepsilon^2)$



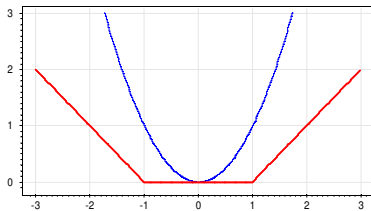
$$\sum_{i=1}^{\ell} \exp\left(-\frac{1}{2\sigma^2}(f(x_i, \alpha) - y_i)^2\right) \rightarrow \max_{\alpha}.$$

Задача решается методом Ньютона-Рафсона.

## SVM-регрессия

Модель регрессии:  $a(x) = \langle x, w \rangle + w_0$ ,  $w \in \mathbb{R}^n$ ,  $w_0 \in \mathbb{R}$ .

Функция потерь:  $\mathcal{L}(\varepsilon) = (|\varepsilon| - \delta)_+$



$$\sum_{i=1}^{\ell} (|\langle w, x_i \rangle - w_0 - y_i| - \delta)_+ + \frac{1}{2C} \|w\|^2 \rightarrow \min_{w, w_0}.$$

## SVM-регрессия

Замена переменных:

$$\begin{aligned}\xi_i^+ &= (\langle w, x_i \rangle - w_0 - y_i - \delta)_+; \\ \xi_i^- &= (-\langle w, x_i \rangle + w_0 + y_i - \delta)_+;\end{aligned}$$

Постановка задачи SVM-регрессии:

$$\begin{cases} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^{\ell} (\xi_i^+ + \xi_i^-) \rightarrow \min_{w, w_0, \xi^+, \xi^-}; \\ y_i - \delta - \xi_i^- \leq \langle w, x_i \rangle - w_0 \leq y_i + \delta + \xi_i^+, \quad i = 1, \dots, \ell; \\ \xi_i^- \geq 0, \quad \xi_i^+ \geq 0, \quad i = 1, \dots, \ell. \end{cases}$$

Это задача квадратичного программирования с линейными ограничениями-неравенствами



## Резюме в конце лекции

- **Н**елинейная регрессия
  - сводится к последовательности линейных
- **Н**епараметрическая регрессия (сглаживание)
  - очень просто, но главное — подобрать ширину окна
- **Н**еквадратичные функции потерь
  - зависят от задачи