

# Методы машинного обучения. Обобщения линейных моделей регрессии и классификации

Воронцов Константин Вячеславович

[www.MachineLearning.ru/wiki?title=User:Vokov](http://www.MachineLearning.ru/wiki?title=User:Vokov)

вопросы к лектору: [k.vorontsov@iai.msu.ru](mailto:k.vorontsov@iai.msu.ru)

материалы курса:

[github.com/MSU-ML-COURSE/ML-COURSE-25-26](https://github.com/MSU-ML-COURSE/ML-COURSE-25-26)

орг.вопросы по курсу: [ml.cmc@mail.ru](mailto:ml.cmc@mail.ru)

## 1 Нелинейная регрессия

- Нелинейная модель регрессии
- Логистическая регрессия
- Обобщённая аддитивная модель

## 2 Обобщённая линейная модель

- Экспоненциальное семейство распределений
- Максимизация правдоподобия для GLM
- Логистическая регрессия как частный случай GLM

## 3 Неквадратичные функции потерь

- Квантильная регрессия
- Робастная регрессия
- SVM-регрессия

## Нелинейная модель регрессии

**Дано:** обучающая выборка  $X^\ell = (x_i, y_i)_{i=1}^\ell$ ,  $x_i \in \mathbb{R}^n$ ,  $y_i \in \mathbb{R}$   
 $y_i = y(x_i)$ ,  $y: X \rightarrow Y$  — неизвестная регрессионная зависимость

**Найти:** параметры  $\alpha \in \mathbb{R}^p$  модели регрессии  $f(x, \alpha)$

**Критерий:** метод наименьших квадратов (МНК)

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

**Метод Ньютона–Рафсона:** итерационный процесс стартует из начального приближения  $\alpha^0 = (\alpha_1^0, \dots, \alpha_p^0)$ :

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

$\nabla Q(\alpha^t)$  — градиент функционала  $Q$  в точке  $\alpha^t$ , вектор из  $\mathbb{R}^p$   
 $Q''(\alpha^t)$  — гессиан функционала  $Q$  в точке  $\alpha^t$ , матрица из  $\mathbb{R}^{p \times p}$   
 $h_t$  — величина шага (можно полагать  $h_t = 1$ ).

## 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:

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

Не хотелось бы обращаться гессиан на каждой итерации...

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

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

## Метод Ньютона-Гаусса

Матричные обозначения:

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

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

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

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

$\beta$  — 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 линейные методы.

## Напоминание. Двухклассовая логистическая регрессия

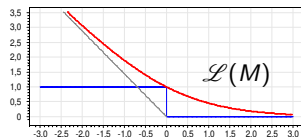
**Дано:** обучающая выборка  $(x_i, y_i)_{i=1}^{\ell}$ ,  $x_i \in \mathbb{R}^n$ ,  $y_i \in \{-1, +1\}$

**Найти:** параметр  $w \in \mathbb{R}^n$  линейной модели  $a(x, w) = \text{sign}(w^T x)$

$M = (w^T x)y$  — отступ (margin)

Логарифмическая функция потерь:

$$\mathcal{L}(M) = \ln(1 + e^{-M})$$

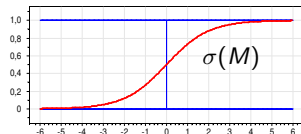


$M$

Модель условной вероятности:

$$P(y|x, w) = \sigma(M) = \frac{1}{1 + e^{-M}},$$

$\sigma(M)$  — сигмоидная функция,



$M$

**Критерий** максимума log правдоподобия, без регуляризации:

$$Q(w) = \sum_{i=1}^{\ell} \ln P(y_i|x_i, w) = \sum_{i=1}^{\ell} \ln(1 + \exp(-w^T x_i y_i)) \rightarrow \min_w$$

## Метода Ньютона-Рафсона

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

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

Элементы градиента — вектора первых производных  $\nabla Q(w^t)$ :

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

$\sigma_i = \sigma(y_i w^T x_i) = P(y_i | x_i, w)$  — вероятность верной классификации

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

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

## Снова сведение к многомерной линейной регрессии

Матричные обозначения:  $F = (f_j(x_i))_{\ell \times n}$ ,  $D = \text{diag}(\sqrt{(1 - \sigma_i)\sigma_i})$

$$-(Q''(w))^{-1} \nabla Q(w) = \underbrace{(F^T D D F)^{-1} F^T D}_{\tilde{F}^+ = (\tilde{F}^T \tilde{F})^{-1} \tilde{F}^T} \underbrace{\left( y_i \sqrt{\frac{1 - \sigma_i}{\sigma_i}} \right)}_{\tilde{y}_i}_{\ell \times 1}$$

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

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

Интерпретация:

- чем выше вероятность ошибки, тем больше  $\frac{1}{\sigma_i}$
- чем ближе  $x_i$  к границе, тем больше вес  $(1 - \sigma_i)\sigma_i$

Таким образом, на каждой итерации происходит более точная настройка на «наиболее трудных» объектах.



## МНК с итерационным перевзвешиванием объектов

### Метод 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:    $\tilde{F} := \text{diag}(\sqrt{(1 - \sigma_i)\sigma_i}) F$ ;
  - 5:    $\tilde{y}_i := y_i \sqrt{\frac{1 - \sigma_i}{\sigma_i}}$  для всех  $i = 1, \dots, \ell$ ;
  - 6:   выбрать градиентный шаг  $h_t$ ;
  - 7:    $w := w + h_t (\tilde{F}^T \tilde{F})^{-1} \tilde{F}^T \tilde{y}$ ;
  - 8:   **если**  $\{\sigma_i\}$  мало изменились **то** выйти из цикла;

## Обобщённая аддитивная модель (Generalized Additive Model)

**Дано:** обучающая выборка  $(x_i, y_i)_{i=1}^{\ell}$ ,  $x_i \in \mathbb{R}^n$ ,  $y_i \in \mathbb{R}$

**Найти:** нелинейные преобразования признаков  $\varphi_j(f_j, \alpha_j)$

(в частности, при  $\varphi_j(f_j(x), \alpha_j) = \alpha_j f_j(x)$  это линейная модель):

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

**Критерий:** метод наименьших квадратов

**Идея 1:** поочерёдно уточнять  $\varphi_j$  по выборке  $(f_j(x_i), z_i)_{i=1}^{\ell}$ :

$$\sum_{i=1}^{\ell} \left( \underbrace{\varphi_j(f_j(x_i), \alpha_j) - \left( y_i - \sum_{k \neq j} \varphi_k(f_k(x_i), \alpha_k) \right)}_{z_i} \right)^2 + \tau R(\alpha_j) \rightarrow \min_{\alpha_j}$$

**Идея 2:** используя в качестве  $\varphi_j$  ядерное сглаживание или сплайны, постепенно уменьшать  $\tau$  у регуляризатора гладкости

$$R(\alpha_j) = \int (\varphi_j''(\zeta, \alpha_j))^2 d\zeta$$

## Метод backfitting [Хасты, Тибширани, 1986]

Многомерная задача сводится к серии одномерных задач.

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

**Выход:**  $\varphi_j(f_j, \alpha_j)$  — обучаемые преобразования признаков.

- 
- 1: начальное приближение:  $\alpha := (F^T F)^{-1} F^T y$ ;  
 $\varphi_j(f_j, \alpha_j) := \alpha_j f_j(x), \quad j = 1, \dots, n$ ;
  - 2: **повторять**
  - 3:   **для**  $j = 1, \dots, n$
  - 4:        $z_i := y_i - \sum_{k=1, k \neq j}^n \varphi_k(f_k(x_i), \alpha_k), \quad i = 1, \dots, \ell$ ;
  - 5:        $\alpha_j := \arg \min_{\alpha} \sum_{i=1}^{\ell} (\varphi(f_j(x_i), \alpha) - z_i)^2 + \tau R(\alpha)$ ;
  - 6:   уменьшить коэффициент регуляризации  $\tau$ ;
  - 7: **пока**  $Q(\alpha, X^{\ell})$  и/или  $Q(\alpha, X^k)$  заметно уменьшаются;

---

*T.J.Hastie, R.J.Tibshirani. Generalized Additive Models. 1990.*

## Напоминание. Вероятностная постановка задачи регрессии

**Дано:** выборка  $(x_i, y_i)_{i=1}^{\ell}$ ,  $x_i \in X$ ,  $y_i \in \mathbb{R}$

**Найти:** параметр  $w$  модели регрессии  $y_i = a(x_i, w) + \varepsilon_i$ , где шум  $\varepsilon_i \sim \mathcal{N}(0, \sigma_i^2)$  гауссовский, некоррелированный  $E\varepsilon_i \varepsilon_k = 0$ :

$$y_i \sim \mathcal{N}(\mu_i, \sigma_i^2), \quad \mu_i = E y_i = a(x_i, w), \quad i = 1, \dots, \ell$$

**Критерий** максимума правдоподобия эквивалентен МНК:

$$p(\varepsilon_1, \dots, \varepsilon_{\ell} | w) = \prod_{i=1}^{\ell} \frac{1}{\sigma_i \sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma_i^2} \varepsilon_i^2\right) \rightarrow \max_w$$

$$-\ln p(\varepsilon_1, \dots, \varepsilon_{\ell} | w) = \text{const} + \frac{1}{2} \sum_{i=1}^{\ell} \frac{1}{\sigma_i^2} (a(x_i, w) - y_i)^2 \rightarrow \min_w$$

Что использовать вместо метода наименьших квадратов, если  $y_i$  не гауссовские, в частности, если  $y_i$  дискретнозначные?

## Обобщение: экспоненциальное распределение шума

**Дано:** выборка  $(x_i, y_i)_{i=1}^{\ell}$ ,  $x_i \in X$ ,  $y_i \in \mathbb{R}$

**Найти:** параметр  $w$  при более общем предположении о шуме:

$$y_i \sim \text{Exp}(\theta_i, \phi_i), \quad \theta_i = g(Ey_i) = a(x_i, w), \quad i = 1, \dots, \ell$$

**Exp** — экспоненциальное семейство распределений

с параметрами  $\theta_i$ ,  $\phi_i$  и параметрами-функциями  $c(\theta)$ ,  $h(y, \phi)$ :

$$p(y_i | \theta_i, \phi_i) = \exp\left(\frac{y_i \theta_i - c(\theta_i)}{\phi_i} + h(y_i, \phi_i)\right)$$

Математическое ожидание и дисперсия с.в.  $y_i \sim \text{Exp}(\theta_i, \phi_i)$ :

$$\mu_i = Ey_i = c'(\theta_i) \Rightarrow \theta_i = [c']^{-1}(\mu_i) = g(Ey_i)$$

$$Dy_i = \phi_i c''(\theta_i)$$

$g(\mu) = [c']^{-1}(\mu)$  — монотонная функция связи (link function)

## Примеры распределений из экспоненциального семейства

Нормальное (гауссовское) распределение,  $y_i \in \mathbb{R}$ :

$$\begin{aligned} p(y_i | \mu_i, \sigma_i^2) &= \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left(-\frac{1}{2\sigma_i^2}(y_i - \mu_i)^2\right) = \\ &= \exp\left(\frac{y_i\mu_i - \frac{1}{2}\mu_i^2}{\sigma_i^2} - \frac{y_i^2}{2\sigma_i^2} - \frac{1}{2}\ln(2\pi\sigma_i^2)\right); \end{aligned}$$

$$\theta_i = g(\mu_i) = \mu_i, \quad c(\theta_i) = \frac{1}{2}\mu_i^2 = \frac{1}{2}\theta_i^2, \quad \phi_i = \sigma_i^2.$$

Распределение Бернулли,  $y_i \in \{0, 1\}$ :

$$p(y_i | \mu_i) = \mu_i^{y_i} (1 - \mu_i)^{1-y_i} = \exp\left(y_i \ln \frac{\mu_i}{1-\mu_i} + \ln(1 - \mu_i)\right);$$

$$\theta_i = g(\mu_i) = \ln \frac{\mu_i}{1-\mu_i}, \quad c(\theta_i) = -\ln(1 - \mu_i) = \ln(1 + e^{\theta_i}).$$

## Примеры распределений из экспоненциального семейства

Биномиальное распределение,  $y_i \in \{0, 1, \dots, n_i\}$ :

$$p(y_i | \mu_i, n_i) = C_{n_i}^{y_i} \left( \frac{\mu_i}{n_i} \right)^{y_i} \left( 1 - \frac{\mu_i}{n_i} \right)^{n_i - y_i} =$$

$$= \exp \left( y_i \ln \frac{\mu_i}{n_i - \mu_i} + n_i \ln(n_i - \mu_i) + \ln C_{n_i}^{y_i} - n_i \ln n_i \right);$$

$$\theta_i = g(\mu_i) = \ln \frac{\mu_i}{n_i - \mu_i}, \quad c(\theta_i) = -n_i \ln(n_i - \mu_i) = n_i \ln \frac{1 + e^{\theta_i}}{n_i}.$$

Пуассоновское распределение,  $y_i \in \{0, 1, 2, \dots\}$ :

$$p(y_i | \mu_i) = \frac{e^{-\mu_i} \mu_i^{y_i}}{y_i!} = \exp \left( \frac{y_i \ln(\mu_i) - \mu_i}{1} - \ln y_i! \right);$$

$$\theta_i = g(\mu_i) = \ln(\mu_i), \quad c(\theta_i) = \mu_i = e^{\theta_i}, \quad \phi_i = 1.$$

## Примеры распределений из экспоненциального семейства

- нормальное (гауссовское)
- распределение Пуассона
- биномиальное и мультиномиальное
- геометрическое
- $\chi^2$ -распределение
- бета-распределение
- гамма-распределение
- распределение Дирихле
- распределение Лапласа с фиксированным матожиданием

**Контр-примеры** не экспоненциальных распределений:

- $t$ -распределение Стьюдента, Коши, гипергеометрическое



## Обобщённая линейная модель (Generalized Linear Model, GLM)

**Дано:** выборка  $(x_i, y_i)_{i=1}^{\ell}$ ,  $x_i \in \mathbb{R}^n$ ,  $y_i \in \mathbb{R}$

**Найти:**  $w \in \mathbb{R}^n$  линейной модели  $\theta_i = \langle x_i, w \rangle = \sum_{j=1}^n w_j f_j(x_i)$

**Критерий** максимума правдоподобия для оценивания  $w$ :

$$Q(w) = \ln \prod_{i=1}^{\ell} p(y_i | \theta_i, \phi_i) = \sum_{i=1}^{\ell} \frac{y_i \theta_i - c(\theta_i)}{\phi_i} \rightarrow \max_w,$$

Метод Ньютона-Рафсона:  $w^{t+1} := w^t + h_t (Q''(w^t))^{-1} \nabla Q(w^t)$

Компоненты вектора градиента  $\nabla Q(w)$ :

$$\frac{\partial Q(w)}{\partial w_j} = \sum_{i=1}^{\ell} \frac{y_i - c'(\theta_i)}{\phi_i} f_j(x_i).$$

Компоненты матрицы Гессе  $Q''(w)$ :

$$\frac{\partial^2 Q(w)}{\partial w_j \partial w_k} = - \sum_{i=1}^{\ell} \frac{c''(\theta_i)}{\phi_i} f_j(x_i) f_k(x_i).$$

## Матричные обозначения

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

$\tilde{F} = DF$ ,  $D = \text{diag}\left(\sqrt{\frac{1}{\phi_i} c''(\theta_i)}\right)$  — веса объектов,  $\theta_i = \langle x_i, w^t \rangle$

$\tilde{y} = (\tilde{y}_i)_{\ell \times 1}$ ,  $\tilde{y}_i = \frac{y_i - c'(\theta_i)}{\sqrt{\phi_i c''(\theta_i)}}$  — модифицированный вектор ответов

Тогда метод Ньютона-Рафсона снова приводит к IRLS:

$$w^{t+1} := w^t - h_t \underbrace{(F^\top D D F)^{-1} F^\top D}_{\tilde{F}^+ = (\tilde{F}^\top \tilde{F})^{-1} \tilde{F}^\top} \underbrace{\left( \sqrt{\frac{\phi_i}{c''(\theta_i)}} \frac{y_i - c'(\theta_i)}{\phi_i} \right)}_{\tilde{y}_i}_{\ell \times 1}$$

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

$$\|\tilde{F}w - \tilde{y}\|^2 \rightarrow \min_w$$

## МНК с итерационным перевзвешиванием объектов

### Метод IRLS — Iteratively Reweighted Least Squares

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

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

---

- 1: начальное приближение:  $w := (F^T F)^{-1} F^T y$ ;
- 2: **для**  $t := 1, 2, 3, \dots$
- 3:    $\theta_i = \langle x_i, w^t \rangle$  для всех  $i = 1, \dots, \ell$ ;
- 4:    $\tilde{F} := \text{diag}\left(\sqrt{\frac{1}{\phi_i} c''(\theta_i)}\right) F$ ;
- 5:    $\tilde{y}_i := \frac{y_i - c'(\theta_i)}{\sqrt{\phi_i c''(\theta_i)}}$  для всех  $i = 1, \dots, \ell$ ;
- 6:   выбрать градиентный шаг  $h_t$ ;
- 7:    $w := w + h_t (\tilde{F}^T \tilde{F})^{-1} \tilde{F}^T \tilde{y}$ ;
- 8:   **если**  $\{\theta_i\}$  мало изменились **то** выйти из цикла;

## Двухклассовая логистическая регрессия

Распределение Бернулли,  $y_i \in \{0, 1\}$ :  $p(y_i|\mu_i) = \mu_i^{y_i}(1 - \mu_i)^{1-y_i}$   
 $\theta_i = g(\mu_i) = \ln \frac{\mu_i}{1-\mu_i} = \langle x_i, w \rangle$ ,  $\mu_i = g^{-1}(\theta_i) = \frac{1}{1+e^{-\theta_i}} \equiv \sigma(\theta_i) = \mathbb{E} y_i$

**Дано:** выборка  $(x_i, y_i)_{i=1}^{\ell}$ ,  $x_i \in \mathbb{R}^n$ ,  $y_i \in \{0, 1\} \sim p(y_i|\mu_i)$

**Найти:** вероятностную модель  $\mathbb{E}(y|x) = P(y=1|x) = \sigma(\langle x, w \rangle)$

**Критерий:** максимум log-правдоподобия (log-loss)

$$Q(w) = \sum_{i=1}^{\ell} \ln p(y_i|\mu_i) = \sum_{i=1}^{\ell} y_i \ln \mu_i + (1 - y_i) \ln(1 - \mu_i) \rightarrow \max_w$$

Альтернативная кодировка:  $y_i \in \{0, 1\} \rightarrow \tilde{y}_i = 2y_i - 1 \in \{\pm 1\}$

$$-\sum_{i=1}^{\ell} \ln p(\tilde{y}_i|x_i) = \sum_{i=1}^{\ell} \ln(1 + \exp(-\underbrace{\langle w, x_i \rangle \tilde{y}_i}_{\text{margin}})) \rightarrow \min_w$$

## Логистическая регрессия как частный случай GLM

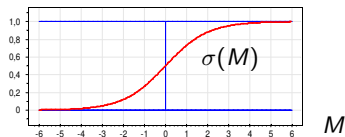
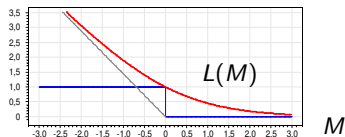
Всего лишь из двух предположений:

- $y_i$  — бернуллиевские случайные величины с  $Ey_i = \mu_i$
- параметр связан с линейной моделью:  $\theta_i = g(\mu_i) = \langle x_i, w \rangle$

следуют важнейшие свойства логистической регрессии:

- логарифмическая функция потерь  $\ln(1 + \exp(-\langle x_i, w \rangle \tilde{y}_i))$ ;
- сигмоидная функция связи  $P(y_i|x_i) = \sigma(\langle x_i, w \rangle \tilde{y}_i)$ ;
- связь линейной модели с *отношением шансов* (odds ratio):

$$\langle x_i, w \rangle = \ln \frac{\mu_i}{1 - \mu_i} = \ln \frac{P(y_i=1|x_i)}{P(y_i=0|x_i)}$$



## Многоклассовая логистическая регрессия

Категориальное (дискретное) распределение,  $y_i \in Y$ ,  $|Y| < \infty$ :

$$p(y_i|\mu_i) = \prod_{y \in Y} \mu_{yi}^{[y=y_i]} = \exp\left(\sum_{y \in Y} [y=y_i] \ln \mu_{yi}\right), \quad \mu_i = (\mu_{yi})_{y \in Y}$$

$$\theta_{yi} = g(\mu_{yi}) = \ln \mu_{yi} = \langle x_i, w_y \rangle, \quad \sum_y \mu_{yi} = 1, \quad \mu_{yi} > 0, \quad \varphi_i = 1$$

$$\mu_{yi} = g^{-1}(\theta_{yi}) = \frac{\exp(\theta_{yi})}{\sum_{z \in Y} \exp(\theta_{zi})} = \text{SoftMax}_{y \in Y} \theta_{yi} = P(y|x_i)$$

**Дано:** выборка  $(x_i, y_i)_{i=1}^{\ell}$ ,  $x_i \in \mathbb{R}^n$ ,  $y_i \in Y \sim p(y_i|\mu_i)$

**Найти:** линейную вероятностную модель классификации

$$P(y|x, w) = \text{SoftMax}_{y \in Y} \langle w_y, x \rangle, \quad a(x) = \arg \max_{y \in Y} \langle w_y, x \rangle, \quad w = (w_y)_{y \in Y}$$

**Критерий:** максимум log-правдоподобия (log-loss)

$$Q(w) = \sum_{i=1}^{\ell} \ln P(y_i|x_i, w) \rightarrow \max_w$$

## Метод наименьших модулей (Least Absolute Deviation Regression)

$\mathcal{L}(\varepsilon_i)$  — функция потерь;  $\varepsilon_i = (a(x_i, w) - y_i)$  — ошибка;

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

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

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

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

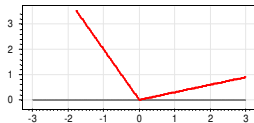
$$\sum_{i=1}^{\ell} |a - y_i| \rightarrow \min_a \Rightarrow a = \text{median}\{y_1, \dots, y_{\ell}\} = y^{(\ell/2)},$$

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

Медиана более устойчива к редким большим выбросам  $y_i$ .

## Квантильная регрессия (Quantile Regression)

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



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

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

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

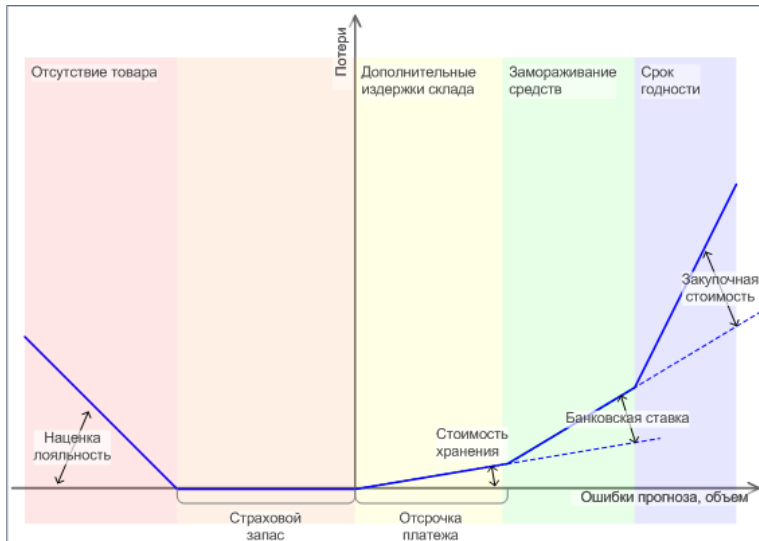
**Сведение к задаче линейного программирования:**

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

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



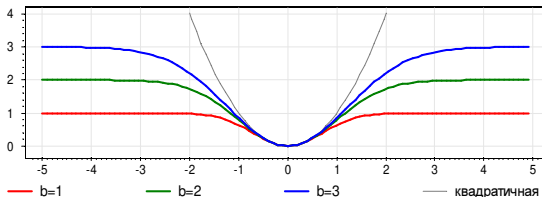
## Пример. Задача прогнозирования объёмов продаж



## Робастная регрессия (Robust Regression)

$a(x, w)$  — модель регрессии;  $\varepsilon_i = (a(x_i, w) - y_i)$  — ошибка;  
 $\mathcal{L}(\varepsilon)$  — функция потерь, устойчивая к большим выбросам  $\varepsilon$

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



Постановка задачи:

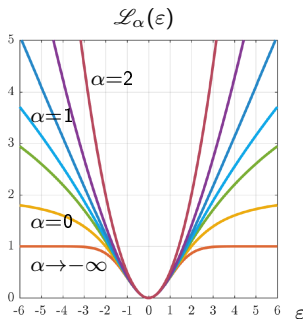
$$\sum_{i=1}^{\ell} \exp\left(-\frac{1}{b}(a(x_i, w) - y_i)^2\right) \rightarrow \max_w.$$

Эта задача также решается методом Ньютона-Рафсона.

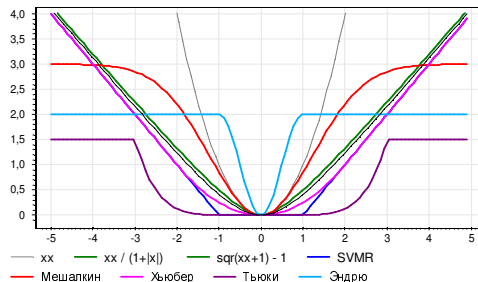
## Функции потерь для робастной регрессии

Семейство функций потерь Баррона с параметром  $\alpha$ :

$$\mathcal{L}_\alpha(\varepsilon) = \frac{|\alpha - 2|}{\alpha} \left( \left( \frac{\varepsilon^2}{|\alpha - 2|} + 1 \right)^{\alpha/2} - 1 \right)$$



другие примеры робастных  $\mathcal{L}(\varepsilon)$

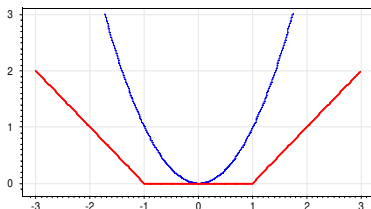


Jonathan T. Barron. A General and Adaptive Robust Loss Function. 2019.

## Напоминание: SVM-регрессия. Тоже робастная регрессия

$a(x, w, w_0) = \langle x, w \rangle - w_0$  — модель регрессии,  $w \in \mathbb{R}^n$ ,  $w_0 \in \mathbb{R}$

$\mathcal{L}(\varepsilon) = (|\varepsilon| - \delta)_+$  — кусочно-линейная функция потерь



Постановка задачи:

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

Задача решается путём замены переменных  
и сведения к задаче квадратичного программирования

- Нелинейная регрессия
  - сводится к последовательности линейных регрессий
- Логистическая регрессия
  - не регрессия, а классификация
  - метод Ньютона-Рафсона приводит к IRLS
- Обобщённая линейная модель (GLM)
  - мощно обобщает обычную и логистическую регрессию
  - метод Ньютона-Рафсона приводит к IRLS
- Обобщённая аддитивная регрессия (GAM, backfitting)
  - сводится к серии одномерных сглаживаний
- Неквадратичные функции потерь
  - проблемно-ориентированные (зависят от задачи)
  - в том числе робастная регрессия
  - приводят к разным методам, отличным от МНК
  - в некоторых случаях к методу Ньютона-Рафсона и IRLS