

Методы машинного обучения. Нелинейная регрессия, обобщённые линейные модели, нестандартные функции потерь

Воронцов Константин Вячеславович

www.MachineLearning.ru/wiki?title=User:Vokov

вопросы к лектору: k.vorontsov@iai.msu.ru

материалы курса:

github.com/MSU-ML-COURSE/ML-COURSE-24-25

орг.вопросы по курсу: ml.cmc@mail.ru

1 Нелинейная регрессия

- Нелинейная модель регрессии
- Логистическая регрессия
- Обобщённая аддитивная модель

2 Обобщённая линейная модель

- В каких случаях нельзя использовать МНК
- Экспоненциальное семейство распределений
- Максимизация правдоподобия для GLM

3 Неквадратичные функции потерь

- Квантильная регрессия
- Робастная регрессия
- SVM-регрессия

Нелинейная модель регрессии

$X^\ell = (x_i, y_i)_{i=1}^\ell$ — обучающая выборка, $x_i \in \mathbb{R}^n$, $y_i \in \mathbb{R}$
 $y_i = y(x_i)$, $y: X \rightarrow Y$ — неизвестная регрессионная зависимость
 $a(x, w)$ — нелинейная модель регрессии, $w \in \mathbb{R}^p$

Метод наименьших квадратов (МНК):

$$Q(w, X^\ell) = \sum_{i=1}^{\ell} (a(x_i, w) - y_i)^2 \rightarrow \min_w.$$

Метод Ньютона–Рафсона:

1. Начальное приближение $w^0 = (w_1^0, \dots, w_p^0)$.
2. Итерационный процесс

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

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

Метод Ньютона-Рафсона

Компоненты градиента:

$$\frac{\partial Q(w)}{\partial w_j} = 2 \sum_{i=1}^{\ell} (a(x_i, w) - y_i) \frac{\partial a(x_i, w)}{\partial w_j}$$

Компоненты гессиана:

$$\frac{\partial^2 Q(w)}{\partial w_j \partial w_k} = 2 \sum_{i=1}^{\ell} \frac{\partial a(x_i, w)}{\partial w_j} \frac{\partial a(x_i, w)}{\partial w_k} - 2 \underbrace{\sum_{i=1}^{\ell} (a(x_i, w) - y_i) \frac{\partial^2 a(x_i, w)}{\partial w_j \partial w_k}}_{\text{при линеаризации полагается} = 0}$$

Не хотелось бы обращаться к гессиану на каждой итерации...

Линеаризация $a(x_i, w)$ в окрестности текущего w^t :

$$a(x_i, w) = a(x_i, w^t) + \sum_{j=1}^p \frac{\partial a(x_i, w_j^t)}{\partial w_j} (w_j - w_j^t) + o(w_j - w_j^t)$$

Метод Ньютона-Гаусса

Матричные обозначения:

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

$a_t = (a(x_i, w^t))_{\ell \times 1}$ — вектор значений a .

Формула t -й итерации метода Ньютона-Гаусса:

$$w^{t+1} := w^t - h_t \underbrace{(F_t^\top F_t)^{-1} F_t^\top (a_t - y)}_u.$$

u — это решение задачи многомерной линейной регрессии

$$\|F_t u - (a_t - y)\|^2 \rightarrow \min_u.$$

Нелинейная регрессия сведена к серии линейных регрессий.

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

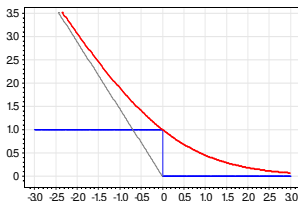
Задача классификации. Логистическая регрессия

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

Снова сведение к задаче линейной регрессии

В матричных обозначениях $F = (f_j(x_i))_{\ell \times n}$, $D = \text{diag}((1 - \sigma_i)\sigma_i)$

$$(Q''(w))^{-1} Q'(w) = -(F^T D F)^{-1} F^T \left(\frac{y_i}{\sigma_i} \right).$$

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

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

Интерпретация (как будут доказано далее, слайды 20–21):

- $\sigma_i = P(y_i | x_i)$ — вероятность правильной классификации x_i
- чем ближе x_i к границе, тем больше вес $(1 - \sigma_i)\sigma_i$
- чем выше вероятность ошибки, тем больше $\frac{1}{\sigma_i}$

ВЫВОД: 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 модели.

-
- 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;

Обобщённая аддитивная модель (Generalized Additive Model)

Регрессия с нелинейными преобразованиями признаков φ_j :

$$a(x, w) = \sum_{j=1}^n \varphi_j(f_j(x), w_j)$$

В частности, при $\varphi_j(f_j(x), w_j) = w_j f_j(x)$ это линейная модель

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

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

Идея 2: постепенно уменьшать τ у регуляризатора гладкости

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

В качестве φ_j использовать сплайны или ядерное сглаживание

Метод backfitting [Хасты, Тибширани, 1986]

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

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

1: начальное приближение:

$w := (F^T F)^{-1} F^T y$ — линейная регрессия;

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

2: **повторять**

3: **для** $j = 1, \dots, n$

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

5: $w_j := \arg \min_{w_j} \sum_{i=1}^{\ell} (\varphi(f_j(x_i), w_j) - z_i)^2 + \tau R(w_j);$

6: уменьшить коэффициент регуляризации τ ;

7: **пока** $Q(w, X^{\ell})$ и/или $Q(w, X^k)$ заметно уменьшаются;

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

Связь МНК с методом максимума правдоподобия

Модель данных с некоррелированным гауссовским шумом:

$$y_i = a(x_i, w) + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, \sigma_i^2), \quad i = 1, \dots, \ell.$$

Эквивалентная запись: $y_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$, $\mu_i = \mathbb{E}y_i = a(x_i, w)$.

МНК эквивалентен методу максимума правдоподобия (ММП):

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

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

Как использовать линейные модели, если y_i не гауссовские, в частности, если y_i дискретнозначные?

Обобщённая линейная модель (Generalized Linear Model, GLM)

Нормальная линейная модель для математического ожидания:

$$y_i \sim \mathcal{N}(\mu_i, \sigma_i^2), \quad \mu_i = x_i^\top w = \mathbb{E}y_i$$

Обобщённая линейная модель для математического ожидания:

$$y_i \sim \text{Exp}(\theta_i, \phi_i), \quad \theta_i = x_i^\top w = g(\mathbb{E}y_i) \text{ — почему?}$$

Exp — экспоненциальное семейство распределений

с параметрами θ_i , ϕ_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. $y_i \sim \text{Exp}(\theta_i, \phi_i)$:

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

$$\text{D}y_i = \phi_i c''(\theta_i)$$

$g(\mu) = [c']^{-1}(\mu)$ — 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, $y_i \in \{0, 1, \dots, n_i\}$:

$$\begin{aligned} p(y_i | \mu_i, n_i) &= C_{n_i}^{y_i} \left(\frac{\mu_i}{n_i} \right)^{y_i} \left(1 - \frac{\mu_i}{n_i} \right)^{n_i - y_i} = \\ &= \exp \left(y_i \ln \frac{\mu_i}{n_i - \mu_i} + n_i \ln(n_i - \mu_i) + \ln C_{n_i}^{y_i} - n_i \ln n_i \right); \end{aligned}$$

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

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
- χ^2 -распределение
- бета-распределение
- гамма-распределение
- распределение Дирихле
- распределение Лапласа с фиксированным матожиданием

Контр-примеры не экспоненциальных распределений:

- t -распределение Стьюдента, Коши, гипергеометрическое

Максимизация правдоподобия для GLM

Принцип максимума правдоподобия:

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

где θ_i линейно зависит от w : $\theta_i = x_i^T w = \sum_{j=1}^n w_j f_j(x_i)$.

Метод Ньютона-Рафсона: $w^{t+1} := w^t + h_t (L''(w^t))^{-1} L'(w^t)$.

Компоненты вектора градиента $L'(w)$:

$$\frac{\partial L(w)}{\partial w_j} = \sum_{i=1}^{\ell} \frac{y_i - c'(x_i^T w)}{\phi_i} f_j(x_i).$$

Компоненты матрицы Гессе $L''(w)$:

$$\frac{\partial^2 L(w)}{\partial w_j \partial w_k} = - \sum_{i=1}^{\ell} \frac{c''(x_i^T w)}{\phi_i} f_j(x_i) f_k(x_i).$$

Матричные обозначения

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

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

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

Тогда метод Ньютона-Рафсона снова приводит к IRLS:

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

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

$$Q(w) = \|\tilde{F}w - \tilde{y}\|^2 \rightarrow \min_w.$$

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 модели.

-
- 1: $w := (F^T F)^{-1} F^T y$ — нулевое приближение, обычный МНК;
 - 2: **для** $t := 1, 2, 3, \dots$
 - 3: $\theta_i = x_i^T w^t$ для всех $i = 1, \dots, \ell$;
 - 4: $\gamma_i := \sqrt{\frac{1}{\phi_i} c''(\theta_i)}$ для всех $i = 1, \dots, \ell$;
 - 5: $\tilde{F} := \text{diag}(\gamma_1, \dots, \gamma_\ell) F$;
 - 6: $\tilde{y}_i := \frac{y_i - c'(\theta_i)}{\phi_i \gamma_i}$ 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** $\{\theta_i\}$ мало изменились **то** выйти из цикла;

Логистическая регрессия как частный случай GLM

Распределение Бернулли, $y_i \in \{0, 1\}$: $p(y_i|\mu_i) = \mu_i^{y_i}(1 - \mu_i)^{1-y_i}$

$$\theta_i = g(\mu_i) = \ln \frac{\mu_i}{1-\mu_i} \quad \mu_i = g^{-1}(\theta_i) = \frac{1}{1+\exp(-\theta_i)} \equiv \sigma(\theta_i)$$

Апостериорная вероятность классов, $\tilde{y}_i = 2y_i - 1 \in \{-1, +1\}$:

$$\left. \begin{aligned} P(y_i=1|x_i) &= E y_i = \mu_i = \sigma(\theta_i) \\ P(y_i=0|x_i) &= 1 - \mu_i = \sigma(-\theta_i) \end{aligned} \right\} \quad p(y_i|x_i) = \sigma(\underbrace{\tilde{y}_i x_i^T w}_{\text{margin}})$$

Максимум правдоподобия \iff минимум критерия log-loss:

$$\begin{aligned} \sum_{i=1}^{\ell} \ln p(y_i|x_i) &= \sum_{i=1}^{\ell} y_i \ln \mu_i + (1 - y_i) \ln(1 - \mu_i) \rightarrow \max_w \\ &\iff \sum_{i=1}^{\ell} \ln(1 + \exp(-\tilde{y}_i x_i^T w)) \rightarrow \min_w \end{aligned}$$

Логистическая регрессия как частный случай GLM

Распределение Бернулли, $y_i \in \{0, 1\}$: $p(y_i|\mu_i) = \mu_i^{y_i}(1 - \mu_i)^{1-y_i}$

$$\theta_i = g(\mu_i) = \ln \frac{\mu_i}{1-\mu_i} \quad \mu_i = g^{-1}(\theta_i) = \frac{1}{1+\exp(-\theta_i)} \equiv \sigma(\theta_i)$$

Всего лишь из двух предположений:

- y_i — бернуллиевские случайные величины с $E y_i = \mu_i$
- модель линейна и μ_i монотонно зависит от $\theta_i = x_i^T w$

вытекают все основные свойства логистической регрессии:

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

$$x_i^T w = \theta_i = \ln \frac{\mu_i}{1 - \mu_i} = \ln \frac{P(y_i = 1|x_i)}{P(y_i = 0|x_i)}.$$

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

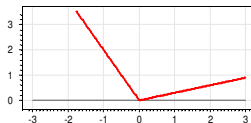
$$\sum_{i=1}^{\ell} |a - y_i| \rightarrow \min_a \Rightarrow a = \text{median}\{y_1, \dots, y_{\ell}\} = y^{(\ell/2)},$$

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 ;

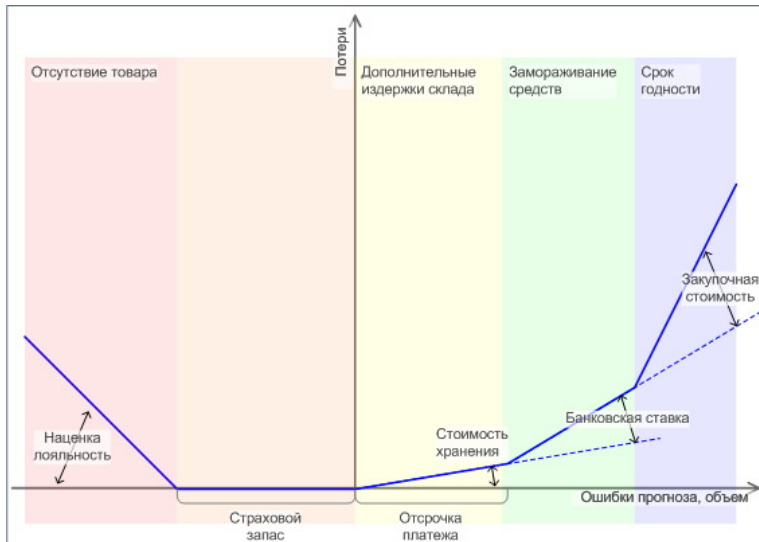
Линейная модель регрессии: $a(x_i, w) = \langle x_i, w \rangle$.

Сведение к задаче линейного программирования:

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

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

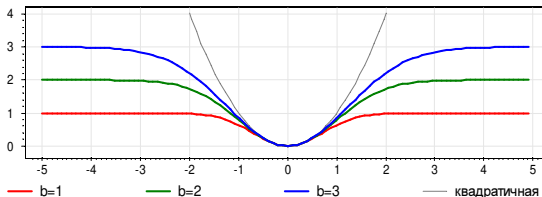
Пример. Задача прогнозирования объёмов продаж



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 ε

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



Постановка задачи:

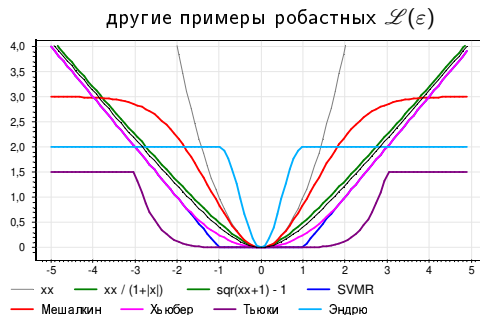
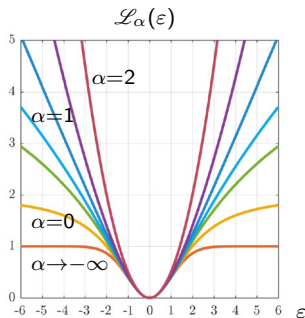
$$\sum_{i=1}^{\ell} \exp\left(-\frac{1}{b}(a(x_i, w) - y_i)^2\right) \rightarrow \max_w.$$

Эта задача также решается методом Ньютона-Рафсона.

Функции потерь для робастной регрессии

Семейство функций потерь Баррона с параметром α :

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

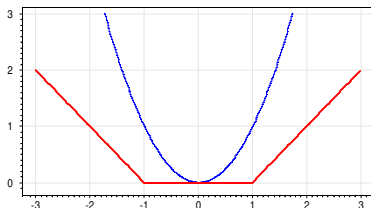


Jonathan T. Barron. A General and Adaptive Robust Loss Function. 2019.

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
 - метод Ньютона-Рафсона приводит к IRLS
- 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