

Линейная и нелинейная регрессия

Воронцов Константин Вячеславович

vokov@forecsys.ru

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

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>

19 марта 2018

1 Методы линейной регрессии

- Многомерная линейная регрессия
- Сингулярное разложение
- Гребневая регрессия

2 Методы нелинейной регрессии

- Нелинейная модель регрессии
- Обобщённая аддитивная модель
- Неквадратичные функции потерь

3 Обобщённая линейная модель

- Обобщённая линейная модель
- Максимизация правдоподобия для GLM
- Логистическая 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):

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

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

Многомерная линейная регрессия

$f_1(x), \dots, f_n(x)$ — числовые признаки;

Модель многомерной линейной регрессии:

$$f(x, \alpha) = \sum_{j=1}^n \alpha_j f_j(x), \quad \alpha \in \mathbb{R}^n.$$

Матричные обозначения:

$$F_{\ell \times n} = \begin{pmatrix} f_1(x_1) & \dots & f_n(x_1) \\ \dots & \dots & \dots \\ f_1(x_\ell) & \dots & f_n(x_\ell) \end{pmatrix}, \quad y_{\ell \times 1} = \begin{pmatrix} y_1 \\ \dots \\ y_\ell \end{pmatrix}, \quad \alpha_{n \times 1} = \begin{pmatrix} \alpha_1 \\ \dots \\ \alpha_n \end{pmatrix}.$$

Функционал квадрата ошибки:

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

Нормальная система уравнений

Необходимое условие минимума в матричном виде:

$$\frac{\partial Q}{\partial \alpha}(\alpha) = 2F^T(F\alpha - y) = 0,$$

откуда следует *нормальная система* задачи МНК:

$$F^T F \alpha = F^T y,$$

где $F^T F$ — матрица размера $n \times n$.

Решение системы: $\alpha^* = (F^T F)^{-1} F^T y = F^+ y$.

Значение функционала: $Q(\alpha^*) = \|P_F y - y\|^2$,

где $P_F = F F^+ = F(F^T F)^{-1} F^T$ — *проекционная матрица*.

Геометрическая интерпретация МНК

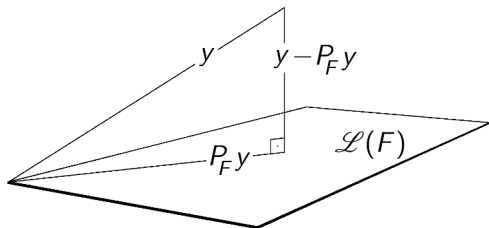
Линейная оболочка столбцов матрицы $F = (f_1, \dots, f_n)$, $f_j \in \mathbb{R}^\ell$:

$$\mathcal{L}(F) = \left\{ \sum_{j=1}^n \alpha_j f_j \mid \alpha \in \mathbb{R}^n \right\}$$

$P_F = F(F^T F)^{-1} F^T$ — проекционная матрица

$P_F y$ — проекция вектора $y \in \mathbb{R}^\ell$ на подпространство $\mathcal{L}(F)$

$(I_\ell - P_F)y$ — проекция y на его ортогональное дополнение



МНК — это опускание перпендикуляра в \mathbb{R}^ℓ из y на $\mathcal{L}(F)$

Сингулярное разложение

Произвольная $\ell \times n$ -матрица представима в виде *сингулярного разложения* (singular value decomposition, SVD):

$$F = VDU^T.$$

Основные свойства сингулярного разложения:

- 1 $\ell \times n$ -матрица $V = (v_1, \dots, v_n)$ ортогональна, $V^T V = I_n$, столбцы v_j — собственные векторы матрицы FF^T ;
- 2 $n \times n$ -матрица $U = (u_1, \dots, u_n)$ ортогональна, $U^T U = I_n$, столбцы u_j — собственные векторы матрицы $F^T F$;
- 3 $n \times n$ -матрица D диагональна, $D = \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n})$, $\lambda_j \geq 0$ — собственные значения матриц $F^T F$ и FF^T .

Решение МНК через сингулярное разложение

Псевдообратная F^+ , вектор МНК-решения α^* ,
 МНК-аппроксимация целевого вектора $F\alpha^*$:

$$F^+ = (UDV^T VDU^T)^{-1}UDV^T = UD^{-1}V^T = \sum_{j=1}^n \frac{1}{\sqrt{\lambda_j}} u_j v_j^T;$$

$$\alpha^* = F^+ y = UD^{-1}V^T y = \sum_{j=1}^n \frac{1}{\sqrt{\lambda_j}} u_j (v_j^T y);$$

$$F\alpha^* = P_F y = (VDU^T)UD^{-1}V^T y = VV^T y = \sum_{j=1}^n v_j (v_j^T y);$$

$$\|\alpha^*\|^2 = \|UD^{-1}V^T y\|^2 = \|D^{-1}V^T y\|^2 = \sum_{j=1}^n \frac{1}{\lambda_j} (v_j^T y)^2.$$

Регуляризация L_2 (гребневая регрессия)

Штраф за увеличение нормы вектора весов $\|\alpha\|$:

$$Q_\tau(\alpha) = \|F\alpha - y\|^2 + \frac{1}{\sigma}\|\alpha\|^2,$$

где $\tau = \frac{1}{\sigma}$ — неотрицательный *параметр регуляризации*.

Вероятностная интерпретация: априорное распределение вектора α — гауссовское с ковариационной матрицей σI_n .

Модифицированное МНК-решение (τI_n — «гребень»):

$$\alpha_\tau^* = (F^T F + \tau I_n)^{-1} F^T y.$$

Преимущество сингулярного разложения:

можно подбирать параметр τ , вычислив SVD только один раз.

Регуляризованный МНК через сингулярное разложение

Вектор регуляризованного МНК-решения α_τ^*
 и МНК-аппроксимация целевого вектора $F\alpha_\tau^*$:

$$\alpha_\tau^* = U(D^2 + \tau I_n)^{-1} D V^T y = \sum_{j=1}^n \frac{\sqrt{\lambda_j}}{\lambda_j + \tau} u_j (v_j^T y);$$

$$F\alpha_\tau^* = V D U^T \alpha_\tau^* = V \operatorname{diag}\left(\frac{\lambda_j}{\lambda_j + \tau}\right) V^T y = \sum_{j=1}^n \frac{\lambda_j}{\lambda_j + \tau} v_j (v_j^T y);$$

$$\|\alpha_\tau^*\|^2 = \|(D^2 + \tau I_n)^{-1} D V^T y\|^2 = \sum_{j=1}^n \frac{\lambda_j}{(\lambda_j + \tau)^2} (v_j^T y)^2.$$

$F\alpha_\tau^* \neq F\alpha^*$, но зато решение становится гораздо устойчивее.

Выбор параметра регуляризации τ

Контрольная выборка: $X^k = (x'_i, y'_i)_{i=1}^k$;

$$F'_{k \times n} = \begin{pmatrix} f_1(x'_1) & \dots & f_n(x'_1) \\ \dots & \dots & \dots \\ f_1(x'_k) & \dots & f_n(x'_k) \end{pmatrix}, \quad y'_{k \times 1} = \begin{pmatrix} y'_1 \\ \dots \\ y'_k \end{pmatrix}.$$

Вычисление функционала Q на контрольных данных T раз потребует $O(kn^2 + knT)$ операций:

$$Q(\alpha_\tau^*, X^k) = \|F' \alpha_\tau^* - y'\|^2 = \left\| \underbrace{F' U}_{k \times n} \operatorname{diag} \left(\frac{\sqrt{\lambda_j}}{\lambda_j + \tau} \right) \underbrace{V^T y}_{n \times 1} - y' \right\|^2.$$

Зависимость $Q(\tau)$ обычно имеет характерный минимум.

Регуляризация сокращает «эффективную размерность»

Сжатие (shrinkage) или сокращение весов (weight decay):

$$\|\alpha_\tau^*\|^2 = \sum_{j=1}^n \frac{\lambda_j}{(\lambda_j + \tau)^2} (v_j^T y)^2 < \|\alpha^*\|^2 = \sum_{j=1}^n \frac{1}{\lambda_j} (v_j^T y)^2.$$

Почему говорят о сокращении эффективной размерности?

Роль размерности играет след проекционной матрицы:

$$\text{tr } F(F^T F)^{-1} F^T = \text{tr}(F^T F)^{-1} F^T F = \text{tr } I_n = n.$$

При использовании регуляризации:

$$\text{tr } F(F^T F + \tau I_n)^{-1} F^T = \text{tr } \text{diag} \left(\frac{\lambda_j}{\lambda_j + \tau} \right) = \sum_{j=1}^n \frac{\lambda_j}{\lambda_j + \tau} < n.$$

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

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

$Q'(\alpha^t)$ — градиент функционала Q в точке α^t , вектор из \mathbb{R}^p

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

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

Не хотелось бы обращать гессиан на каждой итерации...

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

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

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

Формула t -й итерации метода Ньютона-Гаусса:

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

β — 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.

Обобщённая аддитивная модель (Generalized Additive Model)

Регрессия с нелинейными функциями признаков $\varphi_j: \mathbb{R} \rightarrow \mathbb{R}$:

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

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

ИДЕЯ: поочерёдно уточнять φ_j по выборке $(f_j(x_i), z_i)_{i=1}^{\ell}$, постепенно ослабляя регуляризатор гладкости $R(\alpha_j)$:

$$Q(\alpha_j) + \tau R(\alpha_j) \rightarrow \min_{\alpha_j}$$

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

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

Метод backfitting [Хасты, Тибширани, 1986]

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

Выход: $\varphi_j(f_j, \alpha_j)$ — все функции преобразования признаков;

нулевое приближение:

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

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

повторять

для всех признаков $j = 1, \dots, n$

$$\left[\begin{array}{l} z_i := y_i - \sum_{k=1, k \neq j}^n \varphi_k(f_k(x_i), \alpha_k), \quad i = 1, \dots, \ell; \\ \alpha_j := \arg \min_{\alpha} \sum_{i=1}^{\ell} (\varphi(f_j(x_i), \alpha) - z_i)^2 + \tau R(\alpha); \end{array} \right.$$

уменьшить коэффициент регуляризации τ ;

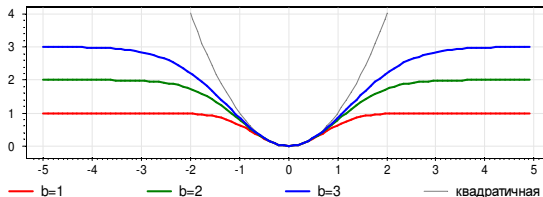
пока $Q(\alpha, X^{\ell})$ и/или $Q(\alpha, X^k)$ заметно уменьшаются;

Робастная регрессия

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

Неквадратичная функция потерь:

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



Постановка задачи:

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

Задача также решается методом Ньютона-Рафсона.

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

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 .

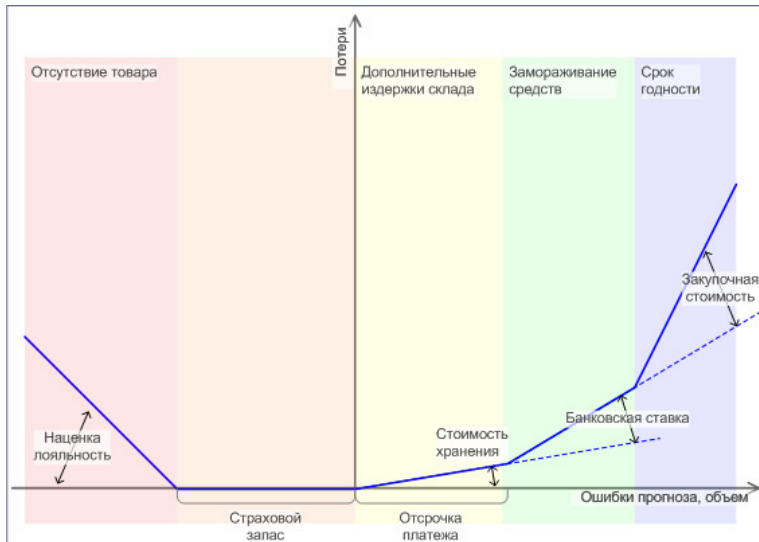
Линейная модель регрессии: $a(x_i) = \langle x_i, w \rangle$.

Сведение к задаче линейного программирования:

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

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

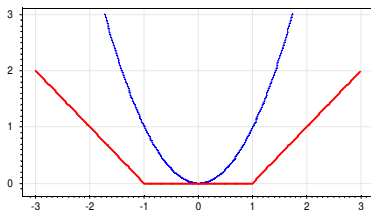
Пример. Задача прогнозирования объёмов продаж



SVM-регрессия (напоминание)

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

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



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 также решается путём замены переменных и сведения к задаче математического (квадратичного) программирования

Связь МНК с методом максимума правдоподобия

Модель данных с некоррелированным гауссовским шумом:

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

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

МНК эквивалентен методу максимума правдоподобия (ММП):

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

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

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:

$$y_i \sim \mathcal{N}(\mu_i, \sigma_i^2), \quad \mu_i = \mathbb{E}y_i = \mathbf{x}_i^\top \boldsymbol{\alpha},$$

Обобщённая линейная модель для математического ожидания:

$$y_i \sim \text{Exp}(\mu_i, \varphi_i), \quad \mu_i = \mathbb{E}y_i, \quad \mathbf{g}(\mu_i) = \theta_i = \mathbf{x}_i^\top \boldsymbol{\alpha},$$

$\mathbf{g}(\mu)$ — монотонная функция связи (link function),

Exp — экспоненциальное семейство распределений

с параметрами θ_i , φ_i и параметрами-функциями $c(\theta)$, $h(y, \varphi)$:

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

Замечательные свойства экспоненциального семейства:

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

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

Примеры распределений из экспоненциального семейства

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

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

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

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

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

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

Примеры распределений из экспоненциального семейства

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

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

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

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

- 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:

$$L(\alpha; X^\ell) = \sum_{i=1}^{\ell} \frac{y_i \theta_i - c(\theta_i)}{\varphi_i} \rightarrow \max_{\alpha}, \quad \theta_i = x_i^\top \alpha = \sum_{j=1}^n \alpha_j f_j(x_i)$$

Метод Ньютона-Рафсона:

$$\alpha^{t+1} := \alpha^t - h_t(L''(\alpha^t))^{-1} L'(\alpha^t).$$

Компоненты вектора градиента $L'(\alpha)$:

$$\frac{\partial L(\alpha)}{\partial \alpha_j} = \sum_{i=1}^{\ell} \frac{y_i - c'(x_i^\top \alpha)}{\varphi_i} f_j(x_i).$$

Компоненты матрицы Гессе $L''(\alpha)$:

$$\frac{\partial^2 L(\alpha)}{\partial \alpha_j \partial \alpha_k} = - \sum_{i=1}^{\ell} \frac{c''(x_i^\top \alpha)}{\varphi_i} f_j(x_i) f_k(x_i).$$

Матричные обозначения

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

$y = (y_i)_{\ell \times 1}$ — вектор ответов.

$\mu^t = (c'(\theta_i))_{\ell \times 1}$ — вектор матожиданий, $\theta_i = x_i^T \alpha^t$.

$W^t = \text{diag}(\frac{1}{\varphi_i} c''(\theta_i))$ — диагональная матрица.

Тогда метод Ньютона-Рафсона приводит к МНК с итеративным перевзвешиванием объектов (IRLS, Iteratively Reweighted Least Squares):

$$\alpha^{t+1} := \alpha^t - h_t (F^T W^t F)^{-1} F^T W^t \text{diag}(\frac{1}{c''(\theta_i)}) (y - \mu^t).$$

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

$$Q(\alpha) = \|\tilde{F}\alpha - (\tilde{y} - \tilde{\mu})\|^2 \rightarrow \min_{\alpha}.$$

Логистическая регрессия как частный случай GLM

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

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

Принцип максимума правдоподобия приводит к log-loss:

$$\sum_{i=1}^{\ell} \ln \mu_i^{y_i} (1 - \mu_i)^{1-y_i} = \sum_{i=1}^{\ell} y_i \ln \mu_i + (1 - y_i) \ln(1 - \mu_i)$$

Выражение для апостериорной вероятности класса +1:

$$P(y_i=1|x_i) = E y_i = \mu_i = \frac{1}{1 + \exp(-\theta_i)} = \sigma(\theta_i) = \sigma(x_i^T \alpha)$$

Линейный классификатор и *отношение шансов* (odds ratio):

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

- Метод наименьших квадратов
 - нормальный некоррелированный шум
- Многомерная линейная регрессия
 - через сингулярное разложение
- Гребневая регрессия
 - тоже через сингулярное разложение
- Нелинейная регрессия
 - сводится к последовательности линейных регрессий
 - используется метод Ньютона-Рафсона
- Логистическая регрессия
 - не регрессия, а классификация
 - используется метод Ньютона-Рафсона
- Обобщённая линейная регрессия
 - обобщает обычную и логистическую регрессию
 - используется метод Ньютона-Рафсона
- Неквадратичные функции потерь
 - проблемно-ориентированные (зависят от задачи)
 - приводят к разным методам, отличным от МНК