

# Семинары по байесовским методам

Евгений Соколов  
[sokolov.evg@gmail.com](mailto:sokolov.evg@gmail.com)

5 декабря 2014 г.

## 2 Нормальный дискриминантный анализ

Нормальный дискриминантный анализ — это частный случай байесовской классификации, когда предполагается, что функции правдоподобия классов  $p(x | y)$  являются нормальными.

### §2.1 Векторное дифференцирование

Выведем формулы векторного дифференцирования, которые пригодятся нам при работе с плотностями нормальных распределений.

**Задача 2.1.** Покажите, что

$$\nabla_X a^T X b = ab^T,$$

где  $a \in \mathbb{R}^m$ ,  $b \in \mathbb{R}^n$ ,  $X \in \mathbb{R}^{m \times n}$ .

**Решение.** Вспомним, что производная по матрице — это матрица частных производных по компонентам этой матрицы. Найдем их:

$$\frac{\partial}{\partial x_{ij}} a^T X b = \frac{\partial}{\partial x_{ij}} \sum_{k=1}^m \sum_{s=1}^n x_{ks} a_k b_s = a_i b_j.$$

Таким образом,

$$\nabla_X a^T X b = ab^T = (a_i b_j)_{i,j} = ab^T.$$

■

**Задача 2.2.** Покажите, что

$$\nabla_X \log \det X = X^{-T},$$

где  $X \in \mathbb{R}^{n \times n}$  — положительно определенная матрица<sup>1</sup>.

---

<sup>1</sup> Если матрица не положительно определена, то ее определитель может быть отрицательным или равным нулю, и логарифм от него будет неопределен.

**Решение.** Запишем производную по  $x_{ij}$ :

$$\frac{\partial}{\partial x_{ij}} \log \det X = \frac{1}{\det X} \frac{\partial \det X}{\partial x_{ij}}.$$

Вспомним *теорему Лапласа* из линейной алгебры и несколько связанных с ней определений. *Минором*  $M_{ij}$  матрицы  $X$  называется определитель матрицы, полученной из  $X$  вычеркиванием  $i$ -й строки и  $j$ -го столбца<sup>2</sup>. *Алгебраическим дополнением*  $C_{ij}$  матрицы  $X$  называется величина  $(-1)^{i+j} M_{ij}$ . Теорема Лапласа гласит, что определитель матрицы  $X$  можно выразить через ее алгебраические дополнения:

$$\det X = \sum_{i=1}^n x_{ij} C_{ij}. \quad (2.1)$$

Вернемся к вычислению производной  $\partial \det X / \partial x_{ij}$ . Заметим, что в разложении (2.1) все алгебраические дополнения вычисляются по матрицам, в которых отсутствует элемент  $x_{ij}$ , и поэтому они могут быть вынесены за знак производной. Получаем, что

$$\frac{\partial \det X}{\partial x_{ij}} = \frac{\partial}{\partial x_{ij}} \sum_{k=1}^n x_{kj} C_{kj} = \sum_{k=1}^n C_{kj} \frac{\partial}{\partial x_{ij}} x_{kj} = C_{ij}.$$

Отсюда следует, что

$$\nabla_X \log \det X = \frac{1}{\det X} (C_{ij})_{i,j=1}^n = \frac{1}{\det X} (X^*)^T.$$

Матрица  $X^* = (C_{ji})$ , составленная из алгебраических дополнений к матрице  $X$ , называется *союзной* или *присоединенной*. Из линейной алгебры известно, что союзная матрица пропорциональна обратной:

$$X^{-1} = \frac{1}{\det X} X^*.$$

Учитывая это, получаем:

$$\nabla_X \log \det X = \frac{1}{\det X} (X^*)^T = \frac{1}{\det X} (X^{-1} \det X)^T = \frac{\det X}{\det X} X^{-T} = X^{-T}.$$

■

**Задача 2.3.** Покажите, что

$$\nabla_X \log \det X^{-1} = -X^{-T},$$

где  $X \in \mathbb{R}^{n \times n}$  — положительно определенная матрица.

**Решение.**

$$\nabla_X \log \det X^{-1} = \nabla_X \log(\det X)^{-1} = -\nabla_X \log \det X = -X^{-T}.$$

■

---

<sup>2</sup> Стого говоря, минор — это определитель произвольной подматрицы, но здесь нам понадобятся миноры именно такого вида

## §2.2 Нормальное распределение

**Одномерное нормальное распределение.** Случайная величина  $x$  имеет нормальное распределение, если ее плотность имеет вид

$$p(x | \mu, \sigma^2) = \mathcal{N}(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right).$$

Вычисляя соответствующие интегралы, можно показать, что параметры  $\mu$  и  $\sigma^2$  соответствуют матожиданию и дисперсии:

$$\begin{aligned} \mathbb{E}x &= \mu; \\ \mathbb{D}x &= \sigma^2. \end{aligned}$$

*Центральная предельная теорема* гласит, что среднее арифметическое независимых одинаково распределенных случайных величин стремится к нормальному распределению:

$$\sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n X_i - \mu \right) \xrightarrow{d} \mathcal{N}(x | 0, \sigma^2),$$

где  $\mu$  и  $\sigma^2$  — матожидание и дисперсия данных случайных величин.

Известно, что нормальное распределение имеет легкие хвосты — вероятность того, что нормальная случайная величина отклонится от своего среднего больше, чем на  $3\sigma$ , не превышает 0.3%. Этот факт называют «правилом трех сигм».

**Многомерное нормальное распределение.** Случайный вектор  $x = (x_1, \dots, x_d)$  имеет многомерное нормальное распределение, если его плотность имеет вид

$$p(x | \mu, \Sigma) = \mathcal{N}(x | \mu, \Sigma) = \frac{1}{(2\pi)^{d/2}(\det \Sigma)^{1/2}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right).$$

Матрица  $\Sigma$  должна быть симметричной и положительно определенной.

Вычисляя соответствующие интегралы, можно показать, что параметры  $\mu$  и  $\Sigma$  соответствуют матожиданию и ковариационной матрице:

$$\begin{aligned} \mathbb{E}x &= \mu; \\ \mathbb{E}(x - \mu)(x - \mu)^T &= \Sigma; \\ \mathbb{D}x_i &= \Sigma_{ii}; \\ \text{Cov}(x_i, x_j) &= \mathbb{E}(x_i - \mu_i)(x_j - \mu_j) = \Sigma_{ij}. \end{aligned}$$

Можно показать, что все моменты многомерной случайной величины выражаются через среднее  $\mu$  и ковариационную матрицу  $\Sigma$ .

Существует обобщение центральной предельной теоремы на многомерный случай, которое гласит, что среднее арифметическое независимых одинаково распределенных случайных векторов стремится к многомерному нормальному распределению:

$$\sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n X_i - \mu \right) \xrightarrow{d} \mathcal{N}(x | 0, \Sigma),$$

где  $\mu$  и  $\Sigma$  — матожидание и ковариационная матрица случайных величин.

Линии уровня плотности нормального распределения соответствуют линиям уровня квадратичной формы  $(x - \mu)^T \Sigma^{-1} (x - \mu)$  и представляют собой эллипсы. Ранее мы подробно выводили вид линий уровня таких квадратичных форм, когда сталкивались с расстоянием Махалонобиса (см. семинары по метрическим методам).

## §2.3 Нормальный дискриминантный анализ

Оптимальный байесовский классификатор при бинарной функции потерь имеет вид

$$a(x) = \arg \max_{y \in Y} p(y)p(x | y).$$

В нормальном дискриминантном анализе предполагается, что распределения объектов внутри классов  $p(x | y)$  — нормальные:

$$p(x | y) = \mathcal{N}(x | \mu_y, \Sigma_y).$$

Параметрами алгоритма являются средние  $\mu_i$  и ковариационные матрицы классов  $\Sigma_y$ , которые оцениваются по выборке методом максимального правдоподобия.

**Задача 2.4.** Выведите оценку максимального правдоподобия на вектор матожиданий  $\mu_y$ , если к классу  $y$  относятся объекты выборки  $X_y = \{x_1, \dots, x_m\}$ .

**Решение.** Для краткости будем обозначать вектор матожиданий и ковариационную матрицу для класса  $y$  через  $\mu$  и  $\Sigma$ . Нам нужно решить задачу

$$p(X_y | \mu, \Sigma) = \prod_{i=1}^m \mathcal{N}(x_i | \mu, \Sigma) \rightarrow \max_{\mu}.$$

Перейдем к логарифму:

$$\log p(X_y | \mu, \Sigma) = -\frac{m}{2} \log \det \Sigma - \frac{1}{2} \sum_{i=1}^m (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) + \text{const.}$$

Найдем производную по  $\mu$  и приравняем ее к нулю:

$$\begin{aligned} \nabla_{\mu} \log p(X_y | \mu, \Sigma) &= -\frac{1}{2} \nabla_{\mu} \left( \sum_{i=1}^m x_i^T \Sigma^{-1} x_i - 2 \sum_{i=1}^m x_i^T \Sigma^{-1} \mu + \sum_{i=1}^m \mu^T \Sigma^{-1} \mu \right) = \\ &= -\frac{1}{2} \left( -2 \sum_{i=1}^m \underbrace{\Sigma^{-T}}_{=\Sigma^{-1}} x_i + \sum_{i=1}^m 2 \Sigma^{-1} \mu \right) = \\ &= \Sigma^{-1} \left( m\mu - \sum_{i=1}^m x_i \right) = \\ &= 0. \end{aligned}$$

Домножая слева на матрицу  $\Sigma$ , получаем

$$\mu = \frac{1}{m} \sum_{i=1}^m x_i.$$

■

**Задача 2.5.** Выведите оценку максимального правдоподобия на ковариационную матрицу  $\Sigma$ , если к классу  $y$  относятся объекты выборки  $X_y = \{x_1, \dots, x_m\}$ .

**Решение.** Как и в предыдущей задаче, будем обозначать вектор матожиданий и ковариационную матрицу для класса  $y$  через  $\mu$  и  $\Sigma$ .

Для удобства перейдем в правдоподобии к матрице точности  $\Lambda = \Sigma^{-1}$ :

$$\log p(X_y | \mu, \Lambda) = -\frac{m}{2} \log \det \Lambda^{-1} - \frac{1}{2} \sum_{i=1}^m (x_i - \mu)^T \Lambda (x_i - \mu) + \text{const.}$$

Найдем производную по  $\Lambda$  и приравняем ее к нулю:

$$\begin{aligned} \nabla_\Lambda \log p(X_y | \mu, \Lambda) &= -\frac{m}{2} \nabla_\Sigma \log \det \Lambda^{-1} - \frac{1}{2} \sum_{i=1}^m \nabla_\Lambda (x_i - \mu)^T \Lambda (x_i - \mu) = \\ &= \frac{m}{2} \underbrace{\Lambda^{-T}}_{=\Lambda^{-1}} - \frac{1}{2} \sum_{i=1}^m (x_i - \mu)(x_i - \mu)^T = \\ &= 0 \end{aligned}$$

Отсюда

$$\Lambda = \frac{1}{m} \left( \sum_{i=1}^m (x_i - \mu)(x_i - \mu)^T \right)^{-1}.$$

Переходя обратно к ковариационной матрице  $\Sigma$ , получаем

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x_i - \mu)(x_i - \mu)^T.$$

■

### 2.3.1 Линейный дискриминант Фишера

Если предположить, что ковариационные матрицы классов равны, и оценивать их по всей выборке, то мы получим алгоритм, называемый *линейным дискриминантом Фишера*. Можно показать, что он является линейным:

$$a(x) = \arg \max_{y \in Y} (\langle w_y, x \rangle + w_{0y}),$$

причем  $w_y = \Sigma^{-1} \mu_y$ . В случае двух классов ( $Y = \{-1, +1\}$ ) классификатор принимает вид

$$a(x) = \text{sign}(\langle w, x \rangle + b) \quad w = \Sigma^{-1}(\mu_2 - \mu_1). \quad (2.2)$$

Разберем другую интерпретацию линейного дискриминанта Фишера. Будем классифицировать объекты следующим образом: выберем прямую с направляющим вектором  $w$  и спроектируем объект на нее; если значение проекции окажется больше порога  $-b$ , то отнесем объект к классу  $+1$ , иначе к классу  $-1$ . Таким образом, классификатор будет иметь вид  $a(x) = \text{sign}(\langle w, x \rangle + b)$ . Обучение классификатора сводится

к поиску проекционной прямой. Будем выбирать ее так, чтобы после проецирования разброс точек из одного класса был как можно меньше, а расстояние между центрами классов было как можно больше. Формализуем эти требования. Обозначим через  $m_k$  центр  $k$ -го класса,  $k \in Y$ :

$$m_k = \frac{1}{N_k} \sum_{i:y_i=k} x_i.$$

Пусть  $s_k^2$  — внутриклассовая дисперсия класса  $k$ :

$$s_k^2 = \sum_{i:y_i=k} (w^T x_i - w^T m_k)^2.$$

В качестве меры «сгруппированности» точек внутри своих классов возьмем сумму внутриклассовых дисперсий  $s_{-1}^2 + s_{+1}^2$ . В качестве меры расстояния между центрами проекций классов («межклассовой дисперсии») возьмем квадрат расстояния между этими центрами:  $(w^T m_{-1} - w^T m_{+1})^2$ . Чтобы совместить минимизацию первой величины и максимизацию второй, возьмем в качестве функционала их отношение. Получим следующую оптимизационную задачу:

$$J(w) = \frac{(w^T m_{-1} - w^T m_{+1})^2}{s_{-1}^2 + s_{+1}^2} \rightarrow \min_w.$$

Распишем данный функционал:

$$\begin{aligned} J(w) &= \frac{(w^T m_{-1} - w^T m_{+1})^2}{s_{-1}^2 + s_{+1}^2} = \\ &= \frac{(w^T (m_{-1} - m_{+1}))^2}{\sum_{i:y_i=-1} (w^T (x_i - m_{-1}))^2 + \sum_{i:y_i=+1} (w^T (x_i - m_{+1}))^2} = \\ &= \frac{w^T (m_{-1} - m_{+1})(m_{-1} - m_{+1})^T w}{\sum_{i:y_i=-1} w^T (x_i - m_{-1})(x_i - m_{-1})^T w + \sum_{i:y_i=+1} w^T (x_i - m_{+1})(x_i - m_{+1})^T w} = \\ &= \frac{w^T (m_{-1} - m_{+1})(m_{-1} - m_{+1})^T w}{w^T \left( \sum_{i:y_i=-1} (x_i - m_{-1})(x_i - m_{-1})^T + \sum_{i:y_i=+1} (x_i - m_{+1})(x_i - m_{+1})^T \right) w}. \end{aligned}$$

Введем обозначения для ковариационных матриц:

$$\begin{aligned} S_b &= (m_{-1} - m_{+1})(m_{-1} - m_{+1})^T; \\ S_w &= \sum_{i:y_i=-1} (x_i - m_{-1})(x_i - m_{-1})^T + \sum_{i:y_i=+1} (x_i - m_{+1})(x_i - m_{+1})^T. \end{aligned}$$

Тогда функционал примет вид

$$J(w) = \frac{w^T S_b w}{w^T S_w w} \rightarrow \min_w.$$

Нам понадобится следующее правило векторного дифференцирования.

**Задача 2.6.** Покажите, что если  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  и  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  — вещественные функции, то

$$\nabla_x \frac{f(x)}{g(x)} = \frac{g(x)\nabla_x f(x) - f(x)\nabla_x g(x)}{g^2(x)}.$$

Воспользуемся полученным правилом, чтобы вычислить градиент функционала  $J(w)$  и приравнять его нулю:

$$\begin{aligned}\nabla_w J(w) &= \frac{(S_b + S_b^T)w(w^T S_w w) - (S_w + S_w^T)w(w^T S_b w)}{(w^T S_w w)^2} = \\ &= 2 \frac{S_b w(w^T S_w w) - S_w w(w^T S_b w)}{(w^T S_w w)^2} = \\ &= 0.\end{aligned}$$

Приходим к уравнению

$$S_b w(w^T S_w w) = S_w w(w^T S_b w). \quad (2.3)$$

Пусть минимум функционала  $J(w)$  достигается на векторе  $w_*$ . Тогда этот вектор удовлетворяет уравнению (2.3). Поскольку классификатор (2.2) зависит только от направления вектора  $w$  и не зависит от его длины, мы можем проигнорировать скалярные множители. Получаем:

$$\begin{aligned}S_w w_* &= \\ &= \underbrace{\frac{w_*^T S_w w_*}{w_*^T S_b w_*}}_{\in \mathbb{R}} S_b w_* \propto \\ &\propto S_b w_* = \\ &= (m_{-1} - m_{+1}) \underbrace{(m_{-1} - m_{+1})^T w_*}_{\in \mathbb{R}} \propto \\ &\propto (m_{-1} - m_{+1}).\end{aligned}$$

Значит,

$$w_* = S_w^{-1}(m_{-1} - m_{+1}).$$

Мы пришли к такому же вектору весов  $w$ , который может быть получен при нормальном дискриминантном анализе в предположении о равенстве ковариационных матриц классов.