

# Байесовская регрессия и классификация

Зинина Анастасия

- Задача: найти классификатор  $a : X \rightarrow Y$  с минимальной вероятностью ошибки
- Пусть известна совместная плотность  $p(x, y) = p(x)P(y | x) = p(y)P(x | y)$
- Принцип максимума апостериорной вероятности:

$$a(x) = \underset{y \in Y}{\operatorname{argmax}} P(y | x)$$

- $L(y, a(x))$ -функция потерь
- Функционал риска

$$R(a, x) = E(L(y, a(x)) | x)$$

- Минимизируем ожидаемые потери:

- задача классификации:

$$a(x) = \underset{s}{\operatorname{argmin}} R(s, x) = \underset{s}{\operatorname{argmin}} \sum_{y \in Y} L(y, s) P(y | x) = \underset{s}{\operatorname{argmin}} \sum_{y \in Y} L(y, s) P(y) P(x | y)$$

- задача регрессии:

$$a(x) = \underset{s}{\operatorname{argmin}} R(s, x) = \underset{s}{\operatorname{argmin}} \int_{y \in Y} L(y, s) P(y | x) dy$$

- Оценим не риск для конкретного  $x$ , а в среднем

$$R(a) = E_x R(a(x), x)$$

- Для задачи классификации с дискретными признаками:

$$R(a) = \sum_{x \in X} R(a(x), x) P(x) \leq \sum_{x \in X} \min_s R(s, x) P(x)$$

- Оптимальный байесовский классификатор минимизирует и средний риск

- Знаем, как построить классификатор, минимизирующий ожидаемые потери, если известны априорные вероятности классов  $P(y)$  и функции правдоподобия  $p(x | y)$
- Найдём эмпирические оценки  $\hat{P}(y)$  и  $\hat{p}(x | y)$

- Оценка априорных вероятностей частотами:

$$\hat{P}(y) = \frac{l_y}{l}, l_y - \text{количество объектов с данным значением } y$$

- Оценка функций правдоподобия:

- параметрическое оценивание  $\hat{p}(x) = \phi(x, \theta)$ ;
- восстановление смеси распределений  $\hat{p}(x) = \sum_{j=1}^k w_j \phi(x, \theta)$ ;
- непараметрическое оценивание  $\hat{p}(x) = \sum_{i=1}^m \frac{1}{mV(h)} K\left(\frac{\rho(x, x_i)}{h}\right)$ ;

- "Наивная" гипотеза: признаки  $f_j : X \rightarrow D_j$  - независимые случайные величины с плотностями распределения  $p_{y,j}(\xi)$

$$p_y(x) = p_{y,1}(\xi_1) \dots p_{y,n}(\xi_n), x = (\xi_1, \dots, \xi_n)$$

- Получим классификатор:

$$a(x) = \max_{y \in Y} (\ln \lambda_y \hat{P}_y + \sum_{j=1}^n \ln \hat{p}_{yj}(\xi_j))$$

- Предположим, что распределение имеет определенный вид:

$$p(x) = \varphi(x; \theta)$$

- Определим параметр  $\theta$  согласно принципу максимума правдоподобия:

$$L(\theta; X^m; G^m) = \sum_{j=1}^n g_j \ln \varphi(x_j; \theta) \rightarrow \max_{\theta},$$

где  $(g_1, \dots, g_m)$ -вектор весов объектов

- Необходимое условие экстремума:

$$\frac{\partial}{\partial \theta} L(\theta; X^m, G^m) = \sum_{i=1}^m g_i \frac{\partial}{\partial \theta} \ln \varphi(x_i; \theta) = 0$$



- Пусть классы имеют  $n$ -мерные гауссовские плотности:

$$p_y(x) = N(x; \mu_y, \Sigma_y) = \frac{e^{-0.5(x-\mu_y)^T \Sigma_y^{-1} (x-\mu_y)}}{\sqrt{(2\pi)^n \det \Sigma_y}}$$

- Разделяющая поверхность

$$\{x \in X \mid \lambda_t P_t p_t(x) = \lambda_s P_s p_s(x)\}$$

квадратична для всех  $y, s \in Y, y \neq s$

- Если  $\Sigma_y = \Sigma_s$ , то она вырождается в линейную

- ОМП в нашем случае:

$$\hat{\mu}_y = \frac{1}{G_y} \sum_{i:y_i=y} g_i x_i;$$

$$\hat{\Sigma}_y = \frac{1}{G_y} \sum_{i:y_i=y} g_i (x_i - \hat{\mu}_y)(x_i - \hat{\mu}_y)^T,$$

где  $G_y = \sum_{i:y_i=y} g_i$

- Получаем алгоритм-квадратичный дискриминант:

$$a(x) = \underset{y \in Y}{\operatorname{argmax}} (\ln \lambda_y P_y - \frac{1}{2} (x - \hat{\mu}_y)^T \hat{\Sigma}_y^{-1} (x - \hat{\mu}_y) - 0.5 \ln \det \Sigma_y)$$

- Допустим, что ковариационные матрицы классов равны  $\Sigma_y = \Sigma$
- Получаем алгоритм-линейный дискриминант:

$$a(x) = \underset{y \in Y}{\operatorname{argmax}} (\ln \lambda_y P_y - \frac{1}{2} \hat{\mu}_y^T \hat{\Sigma}_y^{-1} \hat{\mu}_y + x^T \hat{\Sigma}^{-1} \hat{\mu}_y) = \underset{y \in Y}{\operatorname{argmax}} (x^T \alpha_y + \beta_y)$$

- Модель плотности:

$$p(x) = \sum_{j=1}^k w_j p_j(x), \quad \sum_{j=1}^k w_j = 1, \quad w_j \geq 0, \quad p_j(x) = \varphi(x; \theta_j)$$






- Задачи:
  - оценить параметры  $w_j$  и  $\theta_j$
  - оценить  $k$

- Латентные переменные нельзя выделить из данных напрямую, однако они сильно влияют вид наших данных.
- Два типа скрытых переменных:
  - дополнительные: нужны в дополнение к уже доступным данным для объяснения определённого вида данных в рамках выбранной модели. Примеры: пропущенные данные, смеси распределений
  - скрытые: можно извлечь из доступных данных и понять причины, влияющие на наблюдаемые данные. Пример: понижение размерности.
- Модели с латентными переменными часто имеют меньше параметров, чем модели, в которых корреляции представлены в "видимом" пространстве признаков

## Идея ЕМ-алгоритма на примере

Простой эксперимент. Две монеты (нет гарантии, что честные). Нужно оценить их смещение, повторяя 5 раз следующую процедуру: случайно выбираем одну из двух монет (с равной вероятностью) и делаем 10 независимых подбрасываний выбранной монеты.

### Maximum likelihood approach

	Coin A	Coin B
 H T T T H H T H T H		5 H, 5 T
 H H H H T H H H H H	9 H, 1 T	
 H T H H H H H T H H	8 H, 2 T	
 H T H T T T H H T T		4 H, 6 T
 T H H H T H H H T H	7 H, 3 T	
5 sets, 10 tosses per set	24 H, 6 T	9 H, 11 T

$$\hat{\theta}_A = \frac{24}{24 + 6} = 0.80$$

$$\hat{\theta}_B = \frac{9}{9 + 11} = 0.45$$

$$\text{Bin}(k|n, \theta) = C_n^k \theta^k (1 - \theta)^{n-k}$$

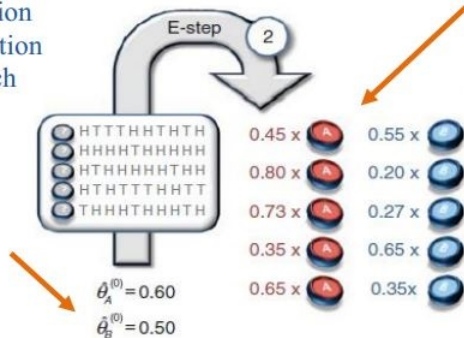
# Идея EM-алгоритма на примере

Теперь предположим, что мы не знаем, какая из монет подбрасывалась.

$$p(z = A | \mathbf{X}, \theta) = \frac{p(z = A | \theta) p(\mathbf{X} | z = A, \theta)}{p(z = A | \theta) p(\mathbf{X} | z = A, \theta) + p(z = B | \theta) p(\mathbf{X} | z = B, \theta)}$$

Expectation  
maximization  
approach

Initial  
guess

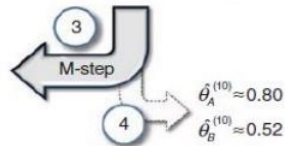


Coin A	Coin B
≈ 2.2 H, 2.2 T	≈ 2.8 H, 2.8 T
≈ 7.2 H, 0.8 T	≈ 1.8 H, 0.2 T
≈ 5.9 H, 1.5 T	≈ 2.1 H, 0.5 T
≈ 1.4 H, 2.1 T	≈ 2.6 H, 3.9 T
≈ 4.5 H, 1.9 T	≈ 2.5 H, 1.1 T
≈ 21.3 H, 8.6 T	≈ 11.7 H, 8.4 T



$$\hat{\theta}_A^{(1)} \approx \frac{21.3}{21.3 + 8.6} \approx 0.71$$

$$\hat{\theta}_B^{(1)} \approx \frac{11.7}{11.7 + 8.4} \approx 0.58$$



$$\hat{\theta}_A^{(10)} \approx 0.80$$

$$\hat{\theta}_B^{(10)} \approx 0.52$$

$$\text{Bin}(k | n, \theta) = C_n^k \theta^k (1 - \theta)^{n-k}$$

- Не применяем принцип максимума правдоподобия к модели смеси распределений, а вводим скрытые переменные.
- 2 повторяющихся до стабилизации параметров шага - Е и М
- **Е-шаг** - вычисляем вероятность принадлежности  $x_i$  к определенному классу с параметрами  $w_j, \theta_j$ , вычисленными на предыдущем шаге.
  - Плотность вероятности того, что объект получен из  $j$ -ой компоненты смеси:

$$p(x, \theta_j) = p(x)P(\theta_j | x) = w_j p_j(x)$$

- Обозначим

$$g_{ij} \equiv P(\theta_j | x_i) - \text{скрытые переменные}$$

- $\sum_{j=1}^k g_{ij} = 1$  для всех  $i = 1, \dots, l$

- Зная параметры, можем вычислить скрытые переменные:

$$g_{ij} = \frac{w_j p_j(x)}{p(x_i)} = \frac{w_j p_j(x)}{\sum_{s=1}^k w_s p_s(x_i)}$$



- **М-шаг** - обновляем параметры, максимизируя логарифм правдоподобия:

$$Q(\theta) = \ln \prod_{i=1}^m p(x_i) = \sum_{i=1}^m \ln \sum_{j=1}^k w_j p_j(x_i) \rightarrow \max_{\theta}, \quad \sum_{j=1}^k w_j = 1, w_j \geq 0$$

- Запишем функцию Лагранжа:

$$L(\Theta, X) = \sum_{i=1}^m \ln \left( \sum_{j=1}^k w_j p_j(x_i) \right) - \lambda \left( \sum_{j=1}^k w_j - 1 \right)$$

- Производная равна нулю:

$$\frac{\partial L}{\partial w_j} = \sum_{i=1}^m \frac{p_j(x_i)}{\sum_{s=1}^k w_s p_s(x_i)} - \lambda = 0, \quad j = 1, \dots, k. \quad (1)$$

$$\sum_{i=1}^m \sum_{j=1}^k \frac{w_j p_j(x_i)}{\sum_{s=1}^k w_s p_s(x_i)} = \lambda \sum_{j=1}^k w_j$$

$$\lambda = m$$

- Умножим (1) на  $w_j$ , подставим  $\lambda = m$ :

$$w_j = \frac{1}{m} \sum_{i=1}^m \frac{w_j p_j(x_i)}{\sum_{s=1}^k w_s p_s(x_i)} = \frac{1}{m} \sum_{i=1}^m g_{ij}, \quad j = 1, \dots, k$$

- Теперь приравняем нулю производную функции Лагранжа по  $\theta_j$ :

$$\begin{aligned}\frac{\partial L}{\partial \theta_j} &= \sum_{i=1}^m \frac{w_j}{\sum_{s=1}^k w_s p_s(x_i)} \frac{\partial}{\partial \theta_j} p_j(x_i) = \sum_{i=1}^m \frac{w_j p_j}{\sum_{s=1}^k w_s p_s(x_i)} \frac{\partial}{\partial \theta_j} \ln p_j(x_i) \\ &= \sum_{i=1}^m g_{ij} \frac{\partial}{\partial \theta_j} \ln p_j(x_i) = \frac{\partial}{\partial \theta_j} \sum_{i=1}^m g_{ij} \ln p_j(x_i) = 0, \quad j = 1, \dots, k.\end{aligned}$$

Допущения:

- функции правдоподобия классов представимы в виде смесей  $k_y$  компонент,  $y \in Y = \{1, \dots, M\}$
- компоненты имеют  $n$ -мерные гауссовские плотности с некоррелированными признаками:

$$\mu_{yj} = (\mu_{yj1}, \dots, \mu_{yjn}), \quad \Sigma_{yj} = \text{diag}(\sigma_{yj1}^2, \dots, \sigma_{yjn}^2), \quad j = 1, \dots, k_y :$$

$$p_y(x) = \sum_{j=1}^{k_y} w_{yj} p_{yj}(x), \quad p_{yj}(x) = N(x; \mu_{yj}, \Sigma_{yj}), \quad \sum_{j=1}^{k_y} w_{yj} = 1, \quad w_{yj} \geq 0$$

- $\phi(x; \theta_j) = \prod_{d=1}^n \frac{1}{\sigma_{jd} \sqrt{2\pi}} \exp(-0.5(\frac{x_{id} - \mu_{jd}}{\sigma_{jd}})^2)$

- M-шаг:

$$\hat{\mu}_{jd} = \frac{1}{mw_j} \sum_{i=1}^m g_{ij} x_{id},$$

$$\hat{\sigma}_{jd}^2 = \frac{1}{mw_j} \sum_{i=1}^m g_{ij} (x_{id} - \hat{\mu}_{jd})^2$$

- Подставим гауссовскую смесь в байесовский классификатор:

$$a(x) = \underset{y \in Y}{\operatorname{argmax}} \lambda_y P_y \sum_{j=1}^{k_y} w_{yj} N_{yj} \exp(-0.5 \rho_{yj}^2(x, \mu_{yj})),$$

$N_{yj} = (2\pi)^{-\frac{n}{2}} (\sigma_{yj1} \dots \sigma_{yjn})^{-1}$  — нормировочные множители;

$\rho_{yj}(x, \mu_{yj})$  — взвешенная евклидова метрика в  $X = R^n$  :

Определим плотность и её оценку:

- Дискретный случай:

$$\hat{p}(x) = \frac{1}{m} \sum_{i=1}^m [x_i = x]$$

- Одномерный непрерывный случай:

$$p(x) = \lim_{h \rightarrow 0} \frac{1}{2h} P[x - h, x + h]$$

Эмпирическая оценка плотности по окну ширины  $h$ :

$$\hat{p}_h(x) = \frac{1}{2mh} \sum_{i=1}^m [|x - x_i| < h]$$

- $\hat{p}_h(x) = \frac{1}{mh} \sum_{i=1}^m \frac{1}{2} \left[ \left| \frac{x-x_i}{h} \right| < 1 \right]$

- Обобщение: оценка Парзена-Розенблатта по окну ширины  $h$ :

$$\hat{p}_h(x) = \frac{1}{mh} \sum_{i=1}^m K\left(\frac{x-x_i}{h}\right),$$

где  $K(r)$  -ядро, удовлетворяющее требованиям:

- чётная функция;
- нормированная функция;
- невозрастающая, неотрицательная функция

## Теорема

При выполнении следующих условий:

1.  $X^m$ -простая выборка из распределения  $p(x)$ ;
2. ядро  $K(z)$  непрерывно и ограничено:  $\int_X K^2(z) dz < \infty$
3.  $\lim_{m \rightarrow \infty} h_m = 0$  и  $\lim_{m \rightarrow \infty} mh_m = \infty$

имеет место утверждение:

$$\hat{p}_{h_m} \rightarrow p(x) \quad \text{при } m \rightarrow \infty \text{ для почти всех } x \in X$$



1. Если объекты описываются  $n$  числовыми признаками  $f_j : X \rightarrow R, j = 1, \dots, n$

$$\hat{p}_h(x) = \frac{1}{m} \sum_{j=1}^m \prod_{i=1}^n \frac{1}{h_j} K\left(\frac{f_j(x) - f_j(x_i)}{h_j}\right)$$

2. Если на  $X$  задана функция расстояния  $\rho(x, x')$ :

$$\hat{p}_h(x) = \frac{1}{mV(h)} \sum_{j=1}^m \frac{1}{h_j} K\left(\frac{\rho(x, x_i)}{h}\right),$$

где  $V(h) = \int_X K\left(\frac{\rho(x, x_i)}{h}\right) dx$  -нормирующий множитель, не зависящий от  $x_i$

$$a(x; X^l, h) = \underset{y \in Y}{\operatorname{argmax}} \Gamma_y(x),$$

$$\Gamma_y(x) = \lambda_y \frac{P_y}{l_y} \sum_{j=1}^m \frac{1}{h_j} K\left(\frac{\rho(x, x_i)}{h}\right)$$

Варианты ядер:

- $Q(r) = \frac{15}{16}(1 - r^2)^2 [ |r| \leq 1 ]$  — кватрическое;
- $T(r) = (1 - |r|) [ |r| \leq 1 ]$  — треугольное;
- $G(r) = (2\pi)^{-1/2} \exp(-\frac{1}{2}r^2)$  — гауссовское;
- $(r) = \frac{1}{2} [ |r| \leq 1 ]$  — прямоугольное