

Лекция 3. Классификация

Основы интеллектуального анализа данных

Полузёров Т. Д.

БГУ ФПМИ

1 Байесовские методы

- Оптимальный классификатор
- Наивный подход
- Непараметрическое восстановление плотности
- Параметрическое восстановление плотности

2 Линейные методы

- Оценка сверху эмперического риска
- Логистическая регрессия

Вероятностная постановка задачи

\mathbb{X} - множество объектов, \mathbb{Y} - множество классов.

$(\mathbb{X} \times \mathbb{Y})$ - вероятностное пространство с совместной плотностью $p(x, y) = P(y)p(x|y)$

$P_y := P(y)$ - **априорные вероятности** классов (prior)

$p_y(x) := p(x|y)$ - **функции правдоподобия** классов (likelihood)

Задачи:

- 1 По выборке $(X, Y) \in (\mathbb{X}, \mathbb{Y})$ построить оценки распределений \hat{P}_y и $\hat{p}_y(x)$
- 2 По известным распределениям $p_y(x)$ и P_y построить алгоритм $a : \mathbb{X} \rightarrow \mathbb{Y}$ минимизирующий вероятность ошибочной классификации

Функционал среднего риска

Алгоритм $a(x)$ разбивает \mathbb{X} на непересекающиеся области $A_y = \{x \in \mathbb{X} | a(x) = y\}$

Каждой паре $(y, s) \in (\mathbb{Y} \times \mathbb{Y})$ соответствует величина потери λ_{ys} при классификации объекта класса y к классу s , $\lambda_{yy} = 0$ и $\lambda_{ys} > 0$ при $y \neq s$

Функционал среднего риска:

$$R(a) = \sum_{y \in \mathbb{Y}} \sum_{s \in \mathbb{Y}} \lambda_{ys} P_y P(A_s | y)$$

где $P(A_s | y) = \int_{A_s} p_y(x) dx$ - вероятность отнесения к классу s объекта класса y .

Теорема о минимуме среднего риска

Если известны априорные вероятности P_y и функции правдоподобия $p_y(x)$, то минимум среднего риска

$$R(a) = \sum_{y \in \mathbb{Y}} \sum_{s \in \mathbb{Y}} \lambda_{ys} P_y P(A_s | y)$$

достигается алгоритмом

$$a(x) = \arg \min_{s \in \mathbb{Y}} \sum_{y \in \mathbb{Y}} \lambda_{ys} P_y p_y(x)$$

Байесовское решающее правило

Если предположить что потери от ошибочной классификации зависят только от истинного класса объекта, т.е. $\lambda_{ys} = \lambda_y$, то алгоритм называется **Байесовским решающим правилом**:

$$a(x) = \arg \max_{y \in \mathbb{Y}} \lambda_y P_y p_y(x)$$

Апостериорные вероятности

Вероятность $P(y|x)$ - называется **апостериорной вероятностью** (posterior).

Зная $p_y(x)$ и P_y , то по формуле Байеса:

$$P(y|x) = \frac{p(x, y)}{p(x)} = \frac{p_y(x)P_y}{\sum_{s \in \mathbb{Y}} p_s(x)P_s}$$

Величина ожидаемых потерь на объекте x :

$$R(x) = \sum_{y \in \mathbb{Y}} \lambda_y P(y|x)$$

Принцип максимума апостериорной вероятности

Оптимальный байесовский классификатор через апостериорные вероятности:

$$a(x) = \arg \max_{y \in \mathbb{Y}} \lambda_y P(y|x)$$

- Если классы равнозначны ($\lambda_y = \lambda_s \forall y, s \in \mathbb{Y}$), то байесовское решающее правило называют **принципом максимума апостериорной вероятности**.
- В случае равновероятных (сбалансированных) классов ($P_y = \frac{1}{|\mathbb{Y}|}$), объект x просто относится к классу s наибольшим значением плотности $p_y(x)$.

Получен оптимальный байесовский классификатор

$$a(x) = \arg \max_{y \in \mathbb{Y}} \lambda_y P(y|x)$$

Но плотности $P(y|x)$ неизвестны. Чтобы построить итоговый классификатор, ставится задача их оценить.

Основные подходы:

- 1 непараметрическое
- 2 параметрическое оценивание плотности
- 3 наивный подход

Наивный подход

Суть наивного подхода - предположение о независимости признаков между собой. Это позволяет упростить

$$P(x|y) = \prod_{i=1}^n P(x_i|y)$$

Для построения итоговой плотности $P(x|y)$ нужно оценить все индивидуальные распределения признаков $P(x_i), i = 1, \dots, n$.

Наивный байесовский классификатор

Оценив априорные плотности и индивидуальные функции правдоподобия, получим **наивный байесовский классификатор**

$$a(x) = \arg \max_{y \in \mathbb{Y}} \ln \lambda_y \hat{P}(y) + \sum_{i=1}^n \ln \hat{P}(x_i | y)$$

Предположение о независимости признаков является очень сильным, но на практике почти никогда не выполняется.

Одномерный непрерывный случай

Пусть $\mathbb{X} = \mathbb{R}$. Эмпирической оценкой плотности есть доля элементов выборки внутри окна шириной h

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

Результат есть кусочно-постоянная функция. Это приводит к появлению зон неуверенности оптимальном классификаторе. Идея состоит в применении ядра

Ядерная оценка плотности

Функция $K(z)$ называется ядром, если она чётная и нормированная $\int K(z)dz = 1$

Тогда **ядерная оценка плотности** имеет вид:

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

Примеры ядер:

- $K(z) = \frac{1}{2}[|z| < 1]$ - прямоугольное
- $K(z) = \frac{1}{\sqrt{2\pi}}e^{-\frac{1}{2}z^2}$ -гауссово

Многомерный случай

Ядерная оценка плотности для многомерной величины $X \in \mathbb{R}^n$

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

Метод парзеновского окна

Применяя ядерную оценку плотности в байесовском решающем правиле

$$a(x) = \arg \max_{y \in \mathbb{Y}} \lambda_y \sum_{i=1}^{\ell} [y_i = y] K\left(\frac{\rho(x, x_i)}{h}\right)$$

где $\rho(z_1, z_2)$ некоторая функция расстояния между объектами

Параметрический подход

Имеется выборка $X = (x_1, \dots, x_\ell) \in \mathbb{X}$. Предполагается, что плотность, порождающая данные, известна с **точностью до параметра**, $p(x) = \phi(x; \theta)$. Подбор параметров θ приводится по выборке X с помощью **метода максимального правдоподобия**.

Нормальный дискриминантный анализ - случай байесовской классификации в предположении о нормальном распределении всех классов, $p_y(x) \sim N(\mu_y, \sigma_y^2)$, $y \in \mathbb{Y}$.

Теорема о разделяющей поверхности

Если классы имеют n -мерные нормальные плотности распределения

$$p_y(x) = N(x; \mu_y, \Sigma_y), y \in \mathbb{Y}$$

то баейсовский классификатор задаёт квадратичную разделяющую поверхность. Она вырождается в линейную, если ковариационные матрица классов равны $\Sigma_y = \Sigma, y \in \mathbb{Y}$

Байесовский нормальный классификатор

Оценим параметры $\hat{\mu}_y$ и $\hat{\Sigma}_y$ n -мерной плотности по имеющимся данным для каждого класса $y \in \mathbb{Y}$.

И воспользуемся идеей оптимального байесовского классификатора

$$a(x) = \arg \max_{y \in \mathbb{Y}} \lambda_y P(y) P(x|y)$$

Такой классификатор называется **байесовский нормальный классификатор** или **подстановочный**

Линейный дискриминант Фишера

Предположив, что ковариационные матрицы классов равны $\Sigma_y = \Sigma, y \in \mathbb{Y}$ и применяя подстановочный алгоритм, получим метод линейного дискриминанта Фишера.

В таком случае разделяющая поверхность является линейной (в случае нескольких классов - кусочно линейная) и обладает определенными свойствами устойчивости.

Бинарная классификация

Рассмотрим задачу бинарной классификации,
 $X = (x_1, \dots, x_\ell) \in \mathbb{X}$, $Y \in \mathbb{Y} = \{-1, 1\}$

Функцию $a(x, \theta) = \text{sign } f(x, \theta)$ будем называть
дискриминантной функцией.

Если $f(x, \theta) > 0$, то x относится к классу $+1$, $f(x, \theta) < 0$ то -1 .

А множество точек $\{x | f(x, \theta) = 0\}$ - **разделяющая поверхность**.

Величина $M_i(\theta) = y_i f(x_i, \theta)$ - **отступ (margin)** классификатора
 $a(x, \theta) = \text{sign } f(x, \theta)$ относительно объекта x_i .

Если $M_i(\theta) < 0$ то алгоритм a допускает ошибку на объекте x_i .
Чем больше $M_i(\theta)$ тем правильнее и надежнее классификация.

Функция потерь

Требуется подобрать параметры θ при которых классификатор a допускает как можно меньше ошибок:

$$Q(a, X) = \frac{1}{\ell} \sum_{i=1}^{\ell} [M_i(\theta) < 0] \rightarrow \min_{\theta}$$

Однако в таком виде Q - кусочно постоянная функция

Идея - мажорирование (оценка сверху) индикатора ошибки $[M_i(\theta) < 0]$ с помощью "удобной" функцией потерь $\mathcal{L}(M_i)$:

$$Q(a, X) = \frac{1}{\ell} \sum_{i=1}^{\ell} [M_i(\theta) < 0] \leq \tilde{Q}(a, X) = \frac{1}{\ell} \sum_{i=1}^{\ell} \mathcal{L}(M_i(\theta))$$

Популярные функции потерь

- ❶ $[M < 0]$ - индикатор ошибки
- ❷ $(1 - M)^2$ - квадратичная
- ❸ $(1 - M)_+$ - кусочно линейная
- ❹ $\frac{2}{1+e^M}$ - сигмоидная
- ❺ $\log_2(1 + e^{-M})$ - логистическая
- ❻ e^{-M} - экспоненциальная

Функция потерь и совместное распределение

Пусть множество $(\mathbb{X} \times \mathbb{Y})$ - вероятностное пространство. Имея выборку (X, Y) и предполагаемый вид совместной плотности $p(x, y; \theta)$, применим метод максимального правдоподобия

$$L(\theta) = \prod_{i=1}^{\ell} p(x_i, y_i; \theta) \rightarrow \max_{\theta}$$

$$\ln L = \sum_{i=1}^{\ell} \ln p(x_i, y_i; \theta) \rightarrow \max_{\theta}$$

$$-\ln p(x_i, y_i; \theta) = \mathcal{L}(y_i f(x_i, \theta))$$

По виду плотности $p(x, y; \theta)$ восстанавливается f и \mathcal{L} . И наоборот, используя некоторые разделяющую поверхность и функцию потерь - предполагаем определенное распределение в данных.

Линейная модель

Случай $f(x, \omega) = \langle x, \omega \rangle$ - класс линейных моделей классификации.

$$a(x, \omega) = \text{sign} \langle x, \omega \rangle$$

Разделяющая поверхность $\text{sign} \langle x, \omega \rangle = 0$ является гиперплоткостью в \mathbb{R}^n . Причем объекты по одну сторону от гиперплоткости относятся к одному классу, по другую - к другому.

Метод обучения

Метод минимизации мажорированного эмперического риска

$$\tilde{Q}(a, X) = \sum_{i=1}^{\ell} \mathcal{L}(\langle x_i, \omega \rangle y_i) \rightarrow \min_{\omega}$$

Необходимое условие минимума:

$$\frac{\partial Q}{\partial \omega} = \sum_{i=1}^{\ell} x_i y_i \mathcal{L}'(\langle x_i, \omega \rangle y_i) = 0$$

Логистическая регрессия

Логистическая регрессия - линейный алгоритм бинарной классификации.

При достаточно сильных свойствах обладает свойствами:

- оптимальный байесовский классификатор
- однозначно определена функция потерь
- возможность оценивать вероятности классов

Пусть $(\mathbb{X} \times \mathbb{Y}) = (\mathbb{R}^n \times \{-1, 1\})$ - вероятностное пространство с плотностью $p(x, y)$. Выборка (X, Y) получена из этого распределения.

Экспонентный класс распределений

Плотность $p(x)$, $x \in \mathbb{R}^n$ называется экспонентной, если

$$p(x) = \exp(c(\delta)\langle\theta, x\rangle + b(\delta, \theta) + d(x, \delta))$$

где θ - параметр сдвига, δ - масштаба, b, c, d - произвольные числовые функции.

Принадлежат к классу экспонентных:

- Равномерное, Нормальное, Гамма
- Гипергеометрическое, Пуассоновское, Биномиальное
- и другие

Обоснование линейной модели

$$a(x) = \text{sign}(\lambda_+ P(+1|x) - \lambda_- P(-1|x)) = \text{sign}\left(\frac{P(+1|x)}{P(-1|x)} - \frac{\lambda_-}{\lambda_+}\right)$$

Если функции правдоподобия $p(x|y)$ принадлежат экспонентному классу, причем параметры d и δ одинаковы, а отличаются только параметры сдвига θ ,

то:

- ❶ байесовский классификатор является линейным:

$$a(x) = \text{sign}\langle \omega, x \rangle$$

- ❷ апостериорная вероятность $p(y|x) = \sigma(\langle \omega, x \rangle y)$

где $\sigma(z) = \frac{1}{1+e^{-z}}$ - сигмоидная функция, $\sigma(-z) = 1 - \sigma(z)$

Модель оценки вероятностей

Построим модель которая оценивает не сами метки классов, а **вероятности** принадлежности к ним.

$$a(x, \omega) = P(+1|x) = \sigma(\langle w, x \rangle)$$

Другими словами, в каждой точке x величина $y \sim \text{Bernoulli}(\sigma(\langle w, x \rangle))$

Задача классификации решается путем выбора порога $t \in [0, 1]$. И тогда итоговый классификатор имеет вид:

$$b(x, t) = \text{sign}(a(x, \omega) - t)$$

Метод максимального правдоподобия

С учетом вероятностной постановки задачи, воспользуемся методом максимального правдоподобия:

$$L(\omega) = \prod_{i=1}^{\ell} p(y_i|x_i) = \prod_{i=1}^{\ell} \sigma(\langle w, x_i \rangle y_i) \rightarrow \max_{\omega}$$

$$\ln L(\omega) = \sum_{i=1}^{\ell} \ln \sigma(\langle w, x_i \rangle y_i) =$$

$$= - \sum_{i=1}^{\ell} \ln(1 + e^{-y_i \langle w, x_i \rangle}) \rightarrow \max_{\omega}$$

-совпадает с логистической функцией потерь

$$Q(\omega) = \sum_{i=1}^{\ell} \ln(1 + e^{-M_i}) \rightarrow \min_{\omega}$$

Решение оптимизационной задачи

Имеем

$$Q(\omega) = \sum_{i=1}^{\ell} \ln(1 + e^{-\langle \omega, x_i \rangle}) \rightarrow \min_{\omega}$$

Аналитического решения нет, поэтому применяются градиентные методы

$$\nabla Q(\omega) = \sum_{i=1}^{\ell} x_i y_i \sigma(-\langle \omega, x_i \rangle)$$