

Обучение с учителем. Классификация. Дискриминантный анализ.

Е. Ларин, Ф. Ежов, И. Кононыхин

1 Обучение с учителем

Рассмотрим задачу обучения с учителем, частным случаем которой являются задачи классификации и регрессии.

Алгоритм в общем виде имеет вид:

- *Вход:* \mathbf{X} — выборка ξ , случайной величины признаков, \mathbf{y} — выборка η , случайной величины «ответов» (принадлежность к классу для классификации, либо значение функции для регрессии). Предполагаем, что существует неизвестное отображение $y^* : \xi \rightarrow \eta$ (гипотеза непрерывности или компактности)
- *Задача:* По \mathbf{X} и \mathbf{y} найти такое отображение $\hat{y}^* : \xi \rightarrow \eta$, которое приблизит отображение y^* .
- *Оценка:* Функция потерь $\mathcal{L}(y^*(x), \hat{y}^*(x))$. Здесь x — реализация ξ

2 Классификация

Перейдём к задаче классификации. Как и в задаче регрессии, данные должны происходить из некоторой генеральной совокупности.

Будем рассматривать выборку признаков и ответов 1

$$\mathbf{X} \in \mathbb{R}^{n \times p}, \mathbf{y} \in \mathbb{A}^n. \quad (1)$$

Отметим, что множество \mathbb{A}^n не является непрерывным. Размерность этого множества $k \times n$, где k — количество возможных классов.

Для обоснования применения методов классификации используется **гипотеза компактности:**

«Близкие» объекты, как правило, принадлежат одному классу.

		Actual Values	
		Positive (1)	Negative (0)
Predicted Values	Positive (1)	TP	FP
	Negative (0)	FN	TN

2.1 Классификация: вероятностная постановка

Поставим задачу классификации в терминах генеральных случайных величин.

Дано:

- $\xi \in \mathbb{R}^p$ — вектор признаков
- $\eta \in \mathbb{A}$ — классовая принадлежность

Предположение об их зависимости можно записать в виде 2.

$$\eta = \Phi(\xi) \quad (2)$$

Задача: найти Φ

При переходе к выборкам, случайная величина признаков ξ заменяется на матрицу наблюдений \mathbf{X} , а случайная величина ответов η — на вектор классовой принадлежности \mathbf{y} .

Предположение принимает вид 3.

$$y_i = \Phi(x_i), \quad i = 1, \dots, n \quad (3)$$

2.2 Классификация: оценка качества

На основе матрицы ошибок 2.2 есть большое количество разных метрик. Приведём некоторые из них:

- $accuracy = \frac{TP+TN}{TP+TN+FP+FN}$,
- $recall = \frac{TP}{TP+FN}$,
- $precision = \frac{TP}{TP+FP}$,
- $F_\beta = (1 - \beta^2) \frac{precision \times recall}{(\beta^2 \times precision) + recall}$,
- $ROC-AUC$

2.3 Классификация: этапы обучения модели

- Выбор модели (класс рассматриваемых Φ из 3). Здесь будут рассмотрены модели LDA и QDA.
- Выбор функции потерь. Чаще всего это $\sum_{i=1}^n (y_i \neq \hat{y}_i)$.
- Выбор метода обучения. Выбор способа подбора параметров для минимизации функции потерь на обучающем множестве.
- Выбор метода проверки. Выбор оценки качества модели, например, с помощью метрик из раздела 2.2.

2.4 Классификация: общий подход к решению

Как построить функционал Φ ?

Общий подход — построить набор классифицирующих функций f_i , $i = 1, \dots, K$. Каждая функция $f_i(\mathbf{x})$ показывает меру принадлежности \mathbf{x} классу i .

Таким образом, решение о принадлежности классу принимается при обнаружении классифицирующей функции с наибольшим значением:

$$\Phi(\mathbf{x}) = \arg \max_i (f_i(\mathbf{x})). \quad (4)$$

3 Дискриминантный анализ

Примем за функции f_i из 4 оценку вероятности принадлежности к i -му классу.

$$\Phi(\mathbf{x}) = \arg \max_i (P(C_i|\mathbf{x})).$$

C_i — класс, состоящий из одного события: \mathbf{x} принадлежит i -му классу.

Если известны априорные вероятности получения i -го класса (π_i), применим формулу Байеса

$$P(C_i|\mathbf{x}) = \frac{\pi_i P(\mathbf{x}|C_i)}{\sum_{j=1}^K \pi_j P(\mathbf{x}|C_j)}.$$

Отбросим знаменатель

$$f_i = P(C_i|\mathbf{x}) = \pi_i P(\mathbf{x}|C_i).$$

3.1 LDA

Предположим, что искомые классы имеют многомерное нормальное распределение с равными дисперсиями.

Запишем это в виде формулы:

$$P(\xi|\eta = A_i) = N(\mu_i, \Sigma)$$

Построим классифицирующие функции:

$$f_i(x) = \frac{\pi_i}{(2\pi)^{p/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_i)\Sigma^{-1}(x - \mu_i)^T\right)$$

Немного упростим (подробнее было изложено в одном из предыдущих курсов):

$$h_i(x) = -0.5\mu_i\Sigma^{-1}\mu_i^T + \mu_i\Sigma^{-1}x + \log \pi_i \quad (5)$$

Функции 5 применяются при классификации данным методом.

3.2 QDA

Предположим, что искомые классы имеют многомерное нормальное распределение с различными дисперсиями.

Запишем это в виде формулы:

$$P(\xi|\eta = A_i) = N(\mu_i, \Sigma_i)$$

Построим классифицирующие функции:

$$f_i(x) = \frac{\pi_i}{(2\pi)^{p/2}|\Sigma_i|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_i)\Sigma_i^{-1}(x - \mu_i)^T\right)$$

Немного упростим (подробнее было изложено в одном из предыдущих курсов):

$$g_i(x) = -0.5(x - \mu_i)\Sigma_i^{-1}(x - \mu_i)^T - 0.5 \log |\Sigma_i| + \log \pi_i \quad (6)$$

Функции 6 применяются при классификации данным методом.

4 Классификация и регрессия

4.1 Регрессия

Обучающая выборка: $X \in \mathbb{R}^{n \times p}$, $y \in \mathbb{R}^n$.

1. Модель регрессии:

$$\hat{y} = \Phi(\mathbf{x}, \beta) = \langle \mathbf{x}, \beta \rangle = \sum_{j=1}^p \beta_j x_j, \beta \in \mathbb{R}^p$$

2. Функция потерь:

$$\mathfrak{L}(\hat{y}, y) = (\hat{y} - y)^2$$

3. Метод обучения — метод наименьших квадратов:

$$Q(\beta) = \sum_{i=1}^n (\Phi(\mathbf{x}_i, \beta) - y_i)^2 \rightarrow \min_{\beta}$$

4.2 Классификация

Обучающая выборка: $\mathbf{X} \in \mathbb{R}^{n \times p}$, $y \in \{-1, 1\}$.

1. Модель классификации:

$$\hat{y} = \Phi(\mathbf{x}, \beta) = \text{sign} \langle \mathbf{x}, \beta \rangle = \text{sign} \sum_{j=1}^p \beta_j x_j, \beta \in \mathbb{R}^p$$

2. Функция потерь:

$$\mathfrak{L}(\hat{y}, y) = [\hat{y}y < 0] = [\langle \mathbf{x}, \beta \rangle y < 0] \leq \hat{\mathfrak{L}}(\langle \mathbf{x}, \beta \rangle y)$$

3. Метод обучения — минимизация эмпирического риска:

$$Q(\beta) = \sum_{i=1}^n [\langle \mathbf{x}_i, \beta \rangle y_i < 0] \leq \sum_{i=1}^n \hat{\mathfrak{L}}(\langle \mathbf{x}_i, \beta \rangle y_i) \rightarrow \min_{\beta}$$

5 Отступы

$\Phi(\mathbf{x}, \beta) = \text{sign}(g(\mathbf{x}, \beta))$ — разделяющий классификатор,

$g(\mathbf{x}, \beta)$ — разделяющая функция,

$g(\mathbf{x}, \beta) = 0$ — уравнение разделяющей поверхности.

$M_i(\beta) = g(\mathbf{x}_i, \beta) y_i$ — отступ объекта \mathbf{x}_i .

Если $M_i(\beta) < 0$, тогда алгоритм ошибается на \mathbf{x}_i .

6 Отступы

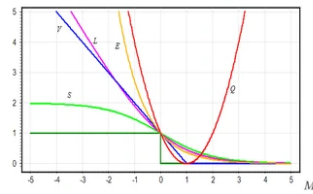


Рис. 7. Непрерывные аппроксимации пороговой функции потерь $[M < 0]$.

$Q(M) = (1 - M)^2$ — квадратичная;
 $V(M) = (1 - M)_+$ — кусочно-линейная;
 $S(M) = 2(1 + e^M)^{-1}$ — сигмоидная;
 $L(M) = \log_2(1 + e^{-M})$ — логистическая;
 $E(M) = e^{-M}$ — экспоненциальная.