

Распознавание образов

12.02.14

Классификация или отнесение объектов к одному из классов

Объекты, кот. классиф-ся на образы

Хар-ки объекта — это признаки

Несколько признаков отнесет объект к одному из классов на основе алгоритма — этот процесс на классификации

Литература:

1. Местечкий "Классификация"
2. Золотулов → МАЛ (школа анализа данных) → лекции (есть сессия на конспекты)
3. Филалетко (ПОМИ) — кафедра мат. логики уделяет внимание алгоритмической стороне вопроса — см курс лекций.

Признаки бывают:

1. Двоичные (0 или 1) (да или нет)
2. Номинальные — конечное число неповторяющихся значений.
3. Порядковые (ординальные) т.е. их можно упорядочить, но они не числового характера.
4. Числовые

Классиф-ся → с учетом или без учета \rightarrow дано несколько объектов проклассиф-ых заранее на основе тех же признаков даны результаты классиф- заранее классиф-кация

Байесовские методы классификации

$(x_1, y_1) \dots (x_n, y_n) \leftarrow \text{сл. вел.}$
 $\sim F(x, y)$ где y — принимает конечное число значений.

x_i — признаки
В действ-ти мы наблюдаем x_1, \dots, x_n

1

Прогнозируем Y .

Бинарная классификация

$Y = 0$ или 1 .

Классификатор

$$g(x) = \begin{cases} 0 & \text{или} \\ 1 \end{cases}$$

$$\eta(x) = P(Y=1|X=x)$$

$Y = \begin{cases} 0, & \text{с вероятностью} \\ 1, & \text{с вероятностью} \end{cases}$

$$\begin{cases} 1 - \eta(x) \\ \eta(x) \end{cases}$$

$$0 \leq \eta(x) \leq 1$$

(при фикс. x , и задан. Y)

Ошибка классификации: $P(g(x) \neq Y) = P(g(X) \neq Y)$

Байесовский классификатор

$$g(x) = \begin{cases} 1, & \eta(x) \geq \frac{1}{2} \\ 0, & \eta(x) < \frac{1}{2} \end{cases}$$

$$g = \eta$$

$$\boxed{\Gamma} \quad P(g^*(X) \neq Y) \leq P(g(X) \neq Y)$$

Доказ-во: $P(g(X) \neq Y) = E_x E(\mathbb{1}_{g(X) \neq Y} | X=x) \Rightarrow$

\Rightarrow Дост-но минимизировать ошибку классификации

$$E(\mathbb{1}_{g(X) \neq Y} | X=x) =$$

$$= 1 - E(\mathbb{1}_{g(X)=Y} | X=x) =$$

$$= 1 - (P(g(X)=1, Y=1 | X=x) + P(g(X)=0, Y=0 | X=x)) =$$

$$= 1 - (\mathbb{1}_{g(X)=0} P(Y=0 | X=x) + \mathbb{1}_{g(X)=1} P(Y=1 | X=x)) =$$

$$= 1 - (\mathbb{1}_{g(X)=0} (1 - \eta(x)) + \mathbb{1}_{g(X)=1} \eta(x)) =$$

$$= P(g^*(X) \neq Y | X=x) - P(g(X) \neq Y | X=x) =$$

$$= (\mathbb{1}_{g(X)=0} - \mathbb{1}_{g^*(X)=0}) (1 - \eta(x)) +$$

$$+ (\mathbb{1}_{g(X)=1} - \mathbb{1}_{g^*(X)=1}) \eta(x)$$

замечание: $\mathbb{1}_{g(X)=1} = 1 - \mathbb{1}_{g(X)=0}$

$$= (\mathbb{1}_{g^*(X)=1} - \mathbb{1}_{g(X)=1}) (1 - \eta(x)) + (\mathbb{1}_{g(X)=1} -$$

$$- \mathbb{1}_{g^*(X)=1}) \eta(x) =$$

$$= (\mathbb{1}_{g^*(X)=1} - \mathbb{1}_{g(X)=1}) (1 - 2\eta(x)) \leq 0$$

2

$$1 - 2\eta(x) > 0 \Rightarrow \mathbb{1}_{g^*(X)=1} = 0$$

□

Обучающий пример:

$$(x_1, y_1), \dots, (x_k, y_k) \quad y_i = \begin{cases} 1 \\ \vdots \\ k \end{cases}$$

$$P(Y_i = c | X_i = x) = \eta_c(x) \quad \sum_c \eta_c(x) = 1$$

а) $\min P(g(X) = Y)$

Оптимальный классификатор: $\hat{c} = \arg \max_c \eta_c(x)$

Предположение смеси вероятностных распределений:

$$p(x) = \sum_{c=0}^k \eta_c \cdot p_c(x)$$

Ответ: $P(Y = c | X = x) = \frac{\eta_c p_c(x)}{\sum \eta_c p_c(x)}$

$$\hat{c} = \arg \max_c \eta_c p_c(x)$$

Классификация для смеси нормальных распределений. Дискриминантный анализ.

$p(x) = \sum \eta_c N(x | \mu_c, \Sigma_c)$ берем по максимуму того или иного класса
 параметры норм. расп-ия

Найд. класс:

$$\arg \max_c \eta_c N(x | \mu_c, \Sigma_c)$$

Если есть обучающая выборка, то оцениваем μ_c и Σ_c . Если нет, то есть эмпирические

$$\hat{c} = \arg \max_c \ln \eta_c + \ln N_c(x, \mu_c, \Sigma_c)$$

(*) $\ln \eta_c - \frac{d}{2} \ln |2\pi| - \frac{1}{2} \ln |\Sigma_c| - \frac{1}{2} (x - \mu_c)^T \Sigma_c^{-1} (x - \mu_c)$
 чиреем
 оценка Σ_c
 оценка μ_c
 расстояние Махаланобиса

Дискриминант Фишера

$$J \quad \Sigma_1 = \Sigma_2 = \dots = \Sigma_n$$

Тогда и построим общую оценку где две переменные S .

$$(*1): \ln \eta_c - \frac{d}{2} \ln |2\pi| - \frac{1}{2} \ln |\Sigma| - \frac{1}{2} (x - \bar{x}_c)' S_c^{-1} (x - \bar{x}_c)$$

Тогда $x S^{-1} x$ сокращается т.к не зависит от x
 argmax $\underbrace{C_c + \bar{x}_c' S_c^{-1} \bar{x}_c - \frac{1}{2} \bar{x}_c' S_c^{-1} \bar{x}_c}_{= C_c}$ — не зависит от x .
 Дискриминант Фишера
 // Некоторые параметры не зависят от x поэтому
 все max и там остаются классификация

02.14.

$$p(x) = \sum_{i=1}^k \omega_i p_i(x, \theta_i) \quad \text{— модель где } k \text{ — число классов}$$

распределение внутри класса.

Предположим, что существуют выборки нет. \Rightarrow придумаем EM-алгоритм (expectation, maximization) — состоит из этих 2х шагов.

Задача моделирования смеси распределений:

1. Шаги Z_1 : $P(Z_i = i) = \omega_i$

2. Моделируем x_i с плотностью распределения

4 // Разбираем номер, после $p_i(x)$ даем с.в.
 $(z_1, x_1), \dots, (z_n, x_n)$

$$\hat{w}_i = \frac{1}{n} \sum_{s=1}^n \mathbb{1}_{\{i_2^s = i\}} \quad \text{оценка } w_i(x)$$

$\exists x_1, \dots, x_n$ — см. величины x_1, \dots, x_n не модально зависят
 как построить
 лотинг примеров. см. в. i_1, \dots, i_n
 $P(i_1 = i | X = x) = \text{по г. Байеса} = \frac{w_i p_i(x)}{\sum_{j=1}^m w_j p_j(x)} = g_i(x)$ w_i как оценка

После того как i_1, \dots, i_n \rightarrow оценки i_1, \dots, i_n
 и оценки $\hat{w}_i(x)$

\forall ранжирующие \uparrow дисперсию оценки \Rightarrow
 индикаторы $\delta(x)$ замены их мат. ожид. $\delta(x)$

$$E(\hat{w}_i | x_1, \dots, x_n) = \frac{1}{n} \sum_{s=1}^n g_{is}(x_s) =: \hat{w}_i \quad \text{необязательно}$$

$\approx \frac{1}{n} \sum_{s=1}^n E(\mathbb{1}_{\{i_2^s = i\}} | x_s) = g_i(x_s)$
 // То научились
 просто оценивать
 $\hat{w}_i \rightarrow$ и напи-
 сать алг. модаль.

Алгоритм:

1° Берем нач. приближения $w_{i0}, p_{i0}, \ell=0$

2° $\hat{w}_{i\ell} = \frac{1}{n} \sum_{s=1}^n \frac{\hat{w}_{i\ell-1} p_i(x_s; \hat{\theta}_{i\ell-1})}{\sum_{j=1}^m \hat{w}_{j\ell-1} p_j(x_s; \hat{\theta}_{j\ell-1})}$
 делаем оценку некой
 мат. ожидания

3° $\hat{\theta}_{i\ell}$
 находим оценку \forall max правдоподобия
 для $g_{i\ell}$ при условии,
 что $w_i = \hat{w}_{i\ell}$ — известны.

p.s. Уч. max правд. написать не совсем
 в очевидной форме, обходим ее
 позже, это уравнение не зависит от θ_i ,

$$\sum g_{i\ell}(x_s) \frac{p_i(x_s; \theta_i)}{p_i(x_s; \theta_i)} = 0 \quad \leftarrow \text{решаем отдельно}$$

$$g_{i\ell}(x) = \frac{w_{i\ell} p_i(x; \theta_{i\ell-1})}{\sum_{j=1}^m \hat{w}_{j\ell} p_j(x; \theta_{j\ell-1})} \quad \leftarrow \text{при фикс. } i$$

— реш. св-во метода

Целью это нам нужно сейчас обосн.
Общее обоснование (*):

$$\max \sum_{s=1}^n \ln \left(\sum_{j=1}^K \hat{w}_j p_j(x, \theta_j) \right) = \text{приэфф. по } \theta_i$$

$$\text{Получаем: } \sum_{s=1}^n \left(\hat{w}_{i,s-1} / \sum_{j=1}^K \hat{w}_{j,s-1} p_j(x_s, \theta_j) \right) \cdot p_{i,\theta}(x_s, \theta_i) = 0$$

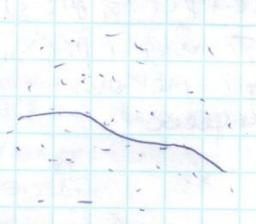
// Потом выведем это для $N(a, \sigma^2)$ //

Умножим на $p_i(x_s, \theta_i) / p_i(x_s, \theta_i)$, тогда:

$$\sum \frac{\hat{w}_{i,s-1} \cdot p_i(x_s, \theta_i)}{\sum_{j=1}^K \hat{w}_{j,s-1} p_j(x_s, \theta_j)} \cdot \frac{p_{i,\theta}(x_s, \theta_i)}{p_i(x_s, \theta_i)}$$

$$0 = \sum_{s=1}^n g_{i,s-1}(x_s) - \frac{p_{i,\theta}(x_s, \hat{\theta}_{i,s})}{p_i(x_s, \hat{\theta}_{i,s})} = 0$$

Замечание: и показать, что оценки E-
класс явл. су. макс правд-ые w_i при
факте, что θ_i - фиксер. Таким образом
мы можем ОМП поочередно параметр θ_i и фиксировать
и все хорошо алгоритм? один из них



Можно взять для
каждого класса смесь
нормальных распреде-
лений. - это повлияет
на границу

Вставка (см. после Каштанова)

Смесь нормальных распределений

$$p(x) = \sum w_i N(\mu_i, R_i)$$

↑ плотность норм. распр.

Шаг 1: Нам предложены w_{i0}, μ_{i0}, R_{i0} .

6

$$\text{Шаг 2: } g_{i,s}(x) = \frac{\hat{w}_{i,s} N(\hat{\mu}_{i,s}, \hat{R}_{i,s})}{\sum_{j=1}^K \hat{w}_{j,s} N(\hat{\mu}_{j,s}, \hat{R}_{j,s})}$$

$$w_{i,s} = \frac{1}{n} \sum g_{i,s-1}(x_s)$$

Шаг 3: $\hat{\mu}_{ce} = \sum g_{ce}(x_s) \cdot \vec{x}_s$ // μ_{ce} и x_s - вектора

$$\hat{K}_{ce} = \sum_{s=1}^n g_{ce}(x_s) (\vec{x}_s - \hat{\mu}_{ce}) (\vec{x}_s - \hat{\mu}_{ce})^T$$

$\hat{\mu}_{ce} = \frac{1}{n} \sum_{s=1}^n g_{ce}(x_s) \vec{x}_s$ - берется в среднем
 $\hat{K}_{ce} = \frac{1}{n} \sum_{s=1}^n g_{ce}(x_s) (\vec{x}_s - \hat{\mu}_{ce}) (\vec{x}_s - \hat{\mu}_{ce})^T$ - берется в среднем
 = 0

Непараметрические методы оценки плотности классификации
 Применение методов непараметрической оценки (классификация с обучением) для задач классификации

$$p(x) = \sum w_j p_j(x)$$

Если есть обучающая выборка $(x_1, y_1), \dots, (x_n, y_n)$, то мы берем относительные частоты $\hat{w}_j = \frac{1}{n} \sum_{s=1}^n \mathbb{1}_{\{y_s=j\}}$

$\hat{w}_j = \frac{1}{n} \sum_{s=1}^n \mathbb{1}_{\{y_s=j\}}$ как оценки w_j после введения подвыборки $x_{s_1}, \dots, x_{s_{r_j}}$

$y_{s_{r_j}} = j$ и для них непараметрически оценить плотность $p_j(x)$ и задача будет решена

P.S. основная проблема - вектор признаков X многомерный, а мы его в задаче разбиваем на кусочки.

Выход:
 Близкая оценка: $\hat{f}_h(x) = \frac{1}{nh} \sum_{s=1}^n K\left(\frac{p(x_s, x)}{h}\right)$

где $p(x_s, x)$ - расстояние между x_s и x .

Возникает вопрос выбора h .

Идея: выбрать h т.ч. что $\#\{x_s: p(x_s, x) < h\} = k$ - фикс.

Если $k(x) = \#\{x_s: p(x_s, x) < h\}$, то мы получаем в качестве оценки:

$$\hat{f}_h(x) = \frac{k}{nh} p(x, x^{(k)})$$

где $x^{(k)}$ - это k -ое удавленное на-
 деление от x .

$$p(x, x^{(1)}) \leq p(x, x^{(2)}) \leq \dots \leq p(x, x^{(n)})$$

Поэтому в реальных задачах обучение
выборка разбивается на 2 части: ~~на~~ первой
строим $f(x_i, w)$, но 2ой проверим, не прои-
зойдет ли эффект переподготовки.

Т.о. в задачах классификации не только
строим классификатор, но и проверяем
его.

Фун: Веса $f(x_i, w) y_i$ наз. остатками.
 $M_i := f(x_i, w) y_i$

Как бороться с \perp в формуле (*)?
(?): \perp заменяем на $\mathcal{L}(f(x_i, w) y_i)$ — непре-
рывная функция. $\perp \cdot \mathcal{L}(M)$

Пр: $\mathcal{L}(M) = \frac{1}{1+e^{-M}}$ — логистическая регрессия ✓

$(1-M)_+ = \max(1-M, 0)$
↑ SVM — support vector machines

e^{-M} — ADA boosting

(*) заменяем на

это остаток M

R^d (***) $Q = \sum \mathcal{L}(M_i) = \sum_{i=1}^n \mathcal{L}(y_i f(x_i, w))$ — мин
// хотим получить Q, мин //

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

p — плотность распределения x_i, y_i
З: Тогда задача (***) эквивалентна задаче на-
хождения мин где $p(x, y, w)$

(*)
Т.е. $Q = -\sum_{i=1}^n \ln p(x_i, y_i, w)$

е — интуитивно функция аналогична функции распределения
и имеет регрессию допускают также
вектор интерпретации

Классификация на основе логистической регрессии

Есть наблюдения $(x_1, y_1), \dots, (x_n, y_n)$ $y_i \in \{-1, 1\}$
 и написать совместную плотность распредел.

$$p(x, y) = p(y|x) \cdot p(x) \quad , \quad y = \{-1, 1\}$$

Предположим, что $p(1|x) = \exp\{\sigma(\delta) \langle w, x \rangle + h(\delta, x) + \tau(\delta, w)\}$

$$p(-1|x) = \exp\{\sigma(\delta) \langle w_{-1}, x \rangle + h(\delta, x) + \tau(\delta, w_{-1})\}$$

Они δ имеют сред. вид $\prod p(y_i, x_i, w_i)$

$$\sum \ln p(y_i, x_i, w_i) = \sum \ln p_w(y_i | x_i) + \sum \ln p(x_i) \rightarrow \max$$

не зависит от $w \Rightarrow$

\Rightarrow В результате остается найти $\sum \ln p_w(y_i | x_i)$ —
 — зависимость на $p(y_i | x_i)$

и значение не зависит от параметра w
 может быть взято из $p(x)$

$$\textcircled{4} \sum \ln \frac{p(y_i | x_i)}{p(1|x_i) + p(-1|x_i)} \rightarrow \max$$

$$y=1: \frac{\exp\{\sigma(\delta) \langle w_1, x \rangle + \tau(\delta, w_1)\}}{\exp\{\sigma(\delta) \langle w_1, x \rangle + \tau(\delta, w_1)\} + \exp\{\sigma(\delta) \langle w_{-1}, x \rangle + \tau(\delta, w_{-1})\}}$$

$$\frac{1}{1 + \exp\{\sigma(\delta) \langle w_1 - w_{-1}, x \rangle + \tau(\delta, w_1) - \tau(\delta, w_{-1})\}}$$

$$= \frac{1}{1 + \exp\{\sigma(\delta) \langle w_1 - w_{-1}, x \rangle + \tau(\delta, w_1) - \tau(\delta, w_{-1})\}} =$$

Добавим к x еще одну координату $(x, 1)$

$$= \frac{1}{1 + \exp\{\sigma(\delta) \langle w_1 - w_{-1}, x \rangle\}} \quad (*)$$

В случае $y = -1$ $\frac{p(-1|x)}{p(1|x) + p(-1|x)} - \text{дугерн} = (*)$

$\sigma(t) = \frac{1}{1+e^{-t}}$ - сигмоидальная функция

$Q = \sum \ln \sigma(y_i < w, x_i >)$ → минимизация $w = w_1, \dots, w_n$
/ метод, найти экстремум //

$\sigma(t)$ присутствует во всех нейронных сетях.
Она гарантирует то, что $\sigma'(t) = \sigma(t)(1-\sigma(t)) = \sigma(t)\sigma(-t)$

Градиентный метод:

$w_k = w_{k-1} - \eta \nabla Q$ // метод оценки w //

$$\nabla Q = \sum_{i=1}^n \sigma(-y_i < w, x_i >) y_i x_i$$

Метод стохастического градиента

1. Берем нач. приближение w_1, \dots, w_d
2. Вычисляем Q
3. Выбираем случайное $i = 1, \dots, n$
 $p(i = \tau) = 1/n$, $1 \leq \tau \leq n$ (т.е. равные)
4. Обновляем $w := w - \eta \sigma(-y_i < w, x_i >) y_i x_i$
5. Вычисляем $\varepsilon = \ln \sigma(-y_i < w, x_i >)$
6. Берем $Q = \frac{n-1}{n} Q + \varepsilon$ // посчитать Q //
7. Останавливаемся, если знаем w и Q стабильно, иначе до 3?

Проблемы, связанные со сходимость метода;
(Куда возникают?)

1. Число параметров (w) м.б. >, т.е. не сопоставимо с объемом обучающей выборки
2. x_i м.б. сильно коррелированы между собой
3. Нест. приращки и. Нест. мало инф-ции
⇒ коэф-ты w и b оцениваются плохо