

Оценка объема выборки в задачах классификации и прогнозирования

А. П. Мотренко, В. В. Стрижов

Московский физико-технический институт,
ВЦ ФИЦ ИУ РАН

Математические методы распознавания образов 17
г. Светлогорск, 23 сентября 2015.

Цель: разработать метод оценки объема выборки, необходимого для получения статистически достоверных результатов классификации.

Задача: выбрать оптимальную модель к решению задачи классификации — порождающую, разделяющую, либо комбинированную.

Предлагается для выбора оптимальной модели при решении задач классификации оценить объем выборки в рамках каждой из моделей.

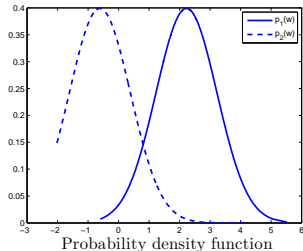
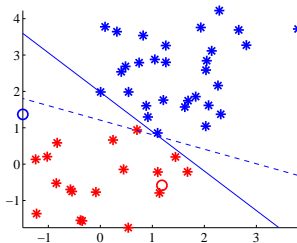
Существующие методы оценки объема выборки

Подход	Формула
<p>Метод доверительных интервалов $Z = \frac{\bar{X} - \mu}{\sigma} \sqrt{m} \rightarrow \mathcal{N}(0, 1)$ при $H_0 : EX = \mu$</p>	$m = \left(\frac{z_{\alpha/2} \sigma}{\bar{X} - \mu} \right)^2$
<p>Тест на равенство: $Z = \frac{\hat{p} - p_0}{\sqrt{\hat{p}(1-\hat{p})}} \sqrt{m} \rightarrow \mathcal{N}(0, 1)$ при $H_0 : p = p_0$ против $H_1 : p \neq p_0$</p>	$m = \frac{(z_{\text{Pow}} + z_{\alpha/2})^2 p(1-p)}{(\hat{p} - p_0)^2}$
<p>Тест отношения правдоподобия: $\gamma_m : \chi_{p, 1-\text{Pow}}^2(\gamma_m) = \chi_{p, \alpha}^2$</p>	$m = \frac{\gamma_m}{\Delta^*}, \text{ где}$ $\Delta^* = E_{\mathcal{X}} \left[\frac{-X(\beta - \beta^*)}{1 + e^{-X\beta}} - \log \left(\frac{1 + e^{-X\beta}}{1 + e^{-X\beta^*}} \right) \right]$
<p>Статистика Вальда: $Z = \frac{\hat{\beta} - \beta^0}{\sqrt{\hat{V}}} \sqrt{m} \rightarrow \mathcal{N}(0, 1)$ при $H_0 : \beta = \beta^0$</p>	$\hat{m} = \frac{(\sqrt{V_1} z_{\text{Pow}} - \sqrt{V_0} z_{\alpha/2})^2}{(\beta^1 - \beta^0)^2}$
<p>Заданная точность регрессии: $\hat{\beta}_j = t_{1-\alpha/2}(m - n - 1) \sqrt{\frac{1 - R^2}{(1 - R_j^2)(m - n - 1)}}$</p>	$m^* = \frac{z_{\alpha/2}^2}{\delta^2} \left(\frac{1 - R^2}{1 - R_j^2} \right) \left(\frac{\chi_{1-\gamma}^2(m-1)}{m-n-1} \right) + n + 1,$ <p>где $\mathbf{R} = \rho'_{yx} \mathbf{R}_{xx}^{-1} \rho_{yx}$</p>
<p>С помощью метода Bootstrap</p>	$m = \left(\frac{z_{\alpha/2} \sigma}{\bar{X} - \mu} \right)^2 \text{ и}$ $m = \frac{z_{\alpha/2}^2}{(\bar{X} - \mu)^2} \left(\frac{1 - R^2}{1 - R_j^2} \right) + n$

Пример выборки недостаточного объема

Дана выборка $Z = \{\mathbf{x}_i, y_i\}, i = 1, \dots, m$, где $y_i \in \{0, 1\}$ — метка класса объекта, $\mathbf{x}_i \in \mathbb{R}^n$ — признакововое описание объекта.

При изменении состава выборки существенно изменяются параметры модели $\mathbf{x}^T \mathbf{w} + c = 0$ и оценка априорного распределения $p(\mathbf{w})$.



Назовем объем m^* выборки из распределения P достаточным, если для всех выборок X_1, X_2 объема $m > m^*$ из P распределения $\hat{p}_1(\mathbf{w})$ и $\hat{p}_2(\mathbf{w})$ близки согласно некоторой функции расстояния $D(\hat{p}_1 || \hat{p}_2)$.

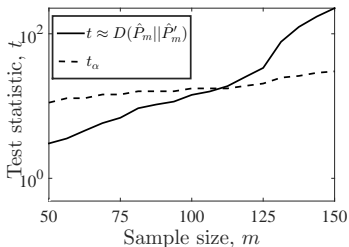
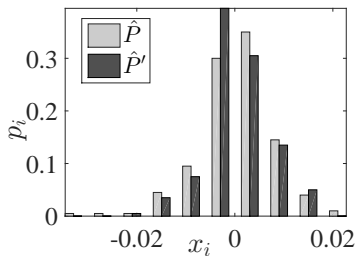
Непараметрический случай: сравнение гистограмм

Рассмотрим выборку Z объема m и разбиение множества значений Z на N интервалов (a_i, a_{i+1}) . Набор оценок

$$\hat{P}_m(a_i < z \leq a_{i+1}) = \frac{n_i}{m} = \hat{p}_i, \quad i = 1, \dots, N$$

вероятности $p_i = P(a_i < z \leq a_{i+1})$ назовем гистограммой \hat{P}_m .

Требуется ввести меру сходства гистограмм $D(\hat{P}_m || \hat{P}'_m)$ и найти m^* : $D(\hat{P}_m || \hat{P}'_m) > D_0$:



f -дивергенция^a между распределениями Q и P :

$$D_f(Q||P) = \sum P \cdot f\left(\frac{Q}{P}\right).$$

Дивергенция Кульбака-Лейблера: $f(t) = t \ln t$, $D_{\text{KL}}(Q||P) = \sum Q \ln\left(\frac{Q}{P}\right)$.
Свойства D_f (для $f(t)$, строго выпуклых и дважды дифференцируемых в $t = 1$):

- 1 Дивергенция $D_f(Q, P)$ определена на всех парах распределений с одинаковым носителем.
- 2 $D_f(Q, P)$ достигает минимума при $P = Q$.
- 3 При $m \rightarrow \infty$ выполнено

$$\frac{2m}{f''(1)} \cdot D_f(\hat{P}_m||P) \rightarrow \chi_N^2,$$

где \hat{P}_m — гистограмма, построенная по выборке X_m из распределения P .

^aS. M. Ali, S. D. Silvey. 1966. A general class of coefficients of divergence of a distribution from another. Journal of Royal Statistical Society. Series B (Methodological). 1(28):131-142.

Пусть

- (a) функция f строго выпукла и дважды дифференцируема в единице,
- (b) распределение $P : \mathbb{R} \rightarrow [0, 1]$ таково, что для всех $x \in \mathbb{R}$ $P(x) \neq 0$,
- (c) выборки Z и Z' выбираются независимо из распределения P , причем $\frac{1}{\rho} \leq m'/m \leq \rho$ при всех m, m' для некоторого $0 < \rho \leq 1$.

Тогда случайная величина $\frac{2m}{f''(1)} \cdot D_f(\hat{P}_m || \hat{P}_{m'})$ в пределе ограничена сверху и снизу случайными величинами из распределения χ_N^2

$$C_1 \chi_N^2 \leq \frac{2m}{f''(1)} \cdot D_f(\hat{P}_m || \hat{P}_{m'}) \leq C_2 \chi_N^2, \text{ при } m, m' \rightarrow \infty.$$

Аппроксимация D_{KL} с помощью χ^2

Схема доказательства: при $P \rightarrow Q : P(x)f\left(\frac{Q(x)}{P(x)}\right) =$

$$P(x)f\left(\frac{Q(x)}{P(x)}\right) + \rightarrow 0, \text{ т.к. } f(1) = 0$$

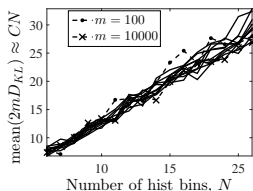
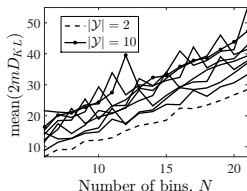
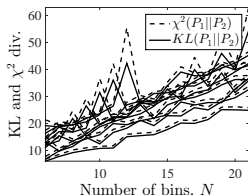
$$f'(1)(Q(x) - P(x)) + \text{ суммируется по } x \text{ к } 0,$$

$$\frac{f''(1)}{2} \frac{(Q(x) - P(x))^2}{P(x)} + P(x) o\left(\left(\frac{Q(x)}{P(x)} - 1\right)^3\right),$$

тогда

$$D_f(\hat{P}_m || P) = \sum_{i=1}^N p_i f\left(\frac{\hat{p}_i}{p_i}\right) \approx \frac{f''(1)}{2m} \sum_i \frac{(n_i - mp_i)^2}{mp_i}.$$

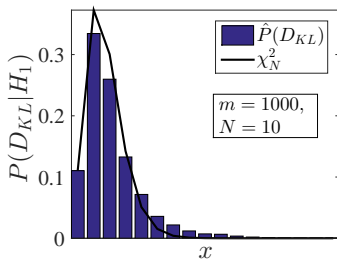
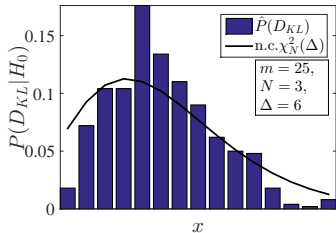
Зависимость параметров аппроксимации C_1, C_2 от размера выборки m и классов $|\mathcal{Y}|$.



Для оценки объема выборки сформулируем задачу двух выборок:

$$H_0 : P(x) \neq P'(x) \text{ при альтернативе } H_1 : P(x) = P'(x).$$

- При H_1 с.в. $C_1' \chi_N^2(\Delta) \leq D_{\text{KL}}(\hat{P}_m, \hat{P}'_m) \leq C_2' \chi_N^2(\Delta)$, где $\chi^2(\Delta)$ — нецентральное хи-квадрат распределение с параметром нецентральности $\Delta = \frac{1}{\mathbb{D}p_i} \sum_{i=1}^N [\mathbb{E}(p_i - p'_i)]^2$.
- При H_1 с.в. $C_1 \chi_N^2 \leq D_{\text{KL}}(\hat{P}_m, \hat{P}'_m) \leq C_2 \chi_N^2$.



Для оценки объема выборки m^* : $2mD_{\text{KL}}(\hat{P}_m || \hat{P}_{m'}) = t_{m^*} \in U(\alpha)$
приблизим критическую область

$$U(\alpha) = \{t : t < t_\alpha\}, \text{ где } P(t < t_\alpha | H_0) = \alpha,$$

с помощью $U^{\chi^2}(\alpha) = \{t : t < n.c.\chi_\alpha^2\}$, где $P(t > \chi_\alpha^2 | t \sim \chi_N^2(\Delta)) = \alpha$.

Для соблюдения условия $P(H_0 | H_1) = P(t_m > t_\beta | H_1) \geq 1 - \beta$
приблизим t_β квантилью $\chi_{N,\beta}^2$.

Из предельных неравенств следует:

$$t_\alpha < \chi_{N,\alpha}^2(\Delta) \Rightarrow U(\alpha) \subseteq U^{\chi^2}(\alpha), \text{ т.е. } \hat{\alpha} < \alpha,$$

$$t_\beta < \chi_{N,\beta}^2 \Rightarrow \hat{\beta} < \beta.$$

Тогда

$$m^* : \chi_{N,\beta}^2 < t_{m^*} < \chi_{N,\alpha}^2(\Delta).$$

Задача классификации:

$$\hat{y} = \operatorname{argmax}_{y \in \{0,1\}} p(y|\mathbf{x}, \mathbf{w}).$$

- Разделяющая модель: параметры \mathbf{w} максимизируют условное правдоподобие

$$\mathbf{w} = \mathbf{w}_D = \operatorname{argmax}_{\mathbf{w} \in \mathcal{W}} L_D(\mathbf{w}; \mathbf{y}, X), \quad L_D(\mathbf{w}) = \ln p(\mathbf{y}|X, \mathbf{w}).$$

- Порождающая модель: параметры \mathbf{w} максимизируют совместное правдоподобие

$$\mathbf{w} = \mathbf{w}_G = \operatorname{argmax}_{\mathbf{w} \in \mathcal{W}} L_G(\mathbf{w}; \mathbf{y}, X), \quad L_G(\mathbf{w}) = \ln p(\mathbf{y}, X|\mathbf{w}).$$

- Разделяющий классификатор имеет вид $\hat{y} = a_D(\mathbf{x}) = [\mathbf{w}_D^T \mathbf{x} > 0]$.
 $\mathbf{y} = [y_1, \dots, y_m]^T$ — случайный вектор с независимыми компонентами y_i :

$$p(y_i = 1 | \mathbf{x}_i, \mathbf{w}) = \frac{1}{1 + \exp(-\mathbf{x}_i^T \mathbf{w})},$$

- Порождающий классификатор: $\hat{y} = a_G(\mathbf{x}) = [\mathbf{w}_G^T \mathbf{x} > 0]$, где $\mathbf{w} = [c, \beta]$.

Параметры $\beta = \Sigma^{-1}(\mu_1 - \mu_0)$, $c = \ln \frac{P}{1-P} - \frac{1}{2} \beta^T (\mu_1 + \mu_0)$ максимизируют совместное правдоподобие

$$\ln p(\mathbf{y}, X | \mathbf{w}) = L_G(\mu_1, \mu_0, \Sigma, P) \rightarrow \max_{\mu_1, \mu_0, \Sigma, P},$$

Предполагается, что $\mathbf{x} | y \sim \mathcal{N}(\mu_y, \Sigma)$, т.е.

$$p(\mathbf{x} | y) = \frac{1}{\sqrt{(2\pi)^n |\Sigma^{-1}|}} \exp\left(-\frac{1}{2} (\mathbf{x} - \mu_y)^T \Sigma^{-1} (\mathbf{x} - \mu_y)\right).$$

тогда имеет место

$$\frac{p(1 | \mathbf{x})}{p(0 | \mathbf{x})} = \frac{P p(\mathbf{x} | 1)}{(1 - P) p(\mathbf{x} | 0)} = \frac{P}{1 - P} \exp(\beta^T \mathbf{x} + \tilde{c}) = \exp(\mathbf{w}^T \mathbf{x}),$$

Оценим вероятность ошибки ε каждого классификаторов a_D и a_G частотой ошибок на выборке $Z = \{\mathbf{x}_i, y_i\}_{i=1}^m$ *:

$$\hat{\varepsilon}_m(a) = \frac{1}{m} \sum_{i=1}^m [a(\mathbf{x}_i) \neq y_i].$$

• При $m \rightarrow \infty$ выполняется $\varepsilon_G > \varepsilon_D$.

Учитывая соотношение $\ln L_G(\mathbf{w}) = \ln L_D(\mathbf{w}) + \ln \prod_{i=1}^m p(\mathbf{x}_i)$, получаем:

$$\prod_{i=1}^m p(y_i | \mathbf{x}_i, \mathbf{w}_D) > \prod_{i=1}^m p(y_i | \mathbf{x}_i, \mathbf{w}_G),$$

поэтому при $m \rightarrow \infty$ с высокой вероятностью выполняется

$$a_D(\mathbf{x}) > a_G(\mathbf{x}) \text{ для } y = 1, \quad a_D(\mathbf{x}) < a_G(\mathbf{x}) \text{ для } y = 0.$$

• При $m \approx 1$ за счет регуляризатора $\ln \prod_{i=1}^m p(\mathbf{x}_i)$ в L_G выполняется $\varepsilon_G < \varepsilon_D$.

* A. Y. Ng, M. I. Jordan. On Discriminative vs. Generative Classifiers: A comparison of logistic regression and naive Bayes, 2002.

Достоинства моделей:

- Разделяющие модели более точны $\varepsilon_D < \varepsilon_G$ в случае если априорное распределение $p(X)$ плохо описывает истинное распределение данных.
- При малых m , $\varepsilon_D > \varepsilon_G$; порождающие модели также позволяют учитывать неразмеченные данные.

Линейная комбинация:

$$\mathbf{w}_\lambda = \underset{\mathbf{w}}{\operatorname{argmax}} L_\lambda(\mathbf{w}), \quad L_\lambda(\mathbf{w}) = \lambda L_D(\mathbf{w}) + (1 - \lambda) L_G(\mathbf{w})$$

- $\lambda \rightarrow 1 \Leftrightarrow \mathbf{w} \rightarrow \mathbf{w}_D$;
- $\lambda \rightarrow 0 \Leftrightarrow \mathbf{w} \rightarrow \mathbf{w}_G$.

Построим комбинированную модель с параметрами

$\mathbf{w}_\lambda = (\mathbf{w}_D, \mathbf{w}_G)$, где

$$\mathbf{w}_\lambda = \operatorname{argmax}_{\mathbf{w}} \ln p(X, \mathbf{y}, \mathbf{w}_\lambda),$$

$$p(X, \mathbf{y}, \mathbf{w}_\lambda) = p(\mathbf{w}_D, \mathbf{w}_G) \left[\prod_{i \in 1}^n p(y_i | \mathbf{x}_i, \mathbf{w}_D) \right] \left[\prod_{i \in 1}^n p(\mathbf{x}_i | \mathbf{w}_G) \right].$$

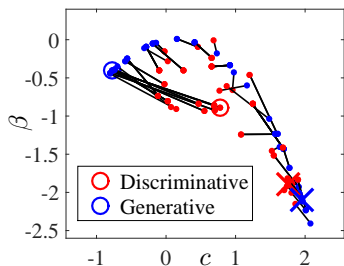
Задание нормального апостериорного распределения параметров \mathbf{w}_λ

$$p(\mathbf{w}_D, \mathbf{w}_G) \propto p(\mathbf{w}_D)p(\mathbf{w}_G) \frac{1}{\sigma} e^{\left(-\frac{1}{2\sigma^2} \|\mathbf{w}_D - \mathbf{w}_G\|^2\right)}, \quad \sigma(\lambda) = \left(\frac{\lambda}{1-\lambda}\right)^2.$$

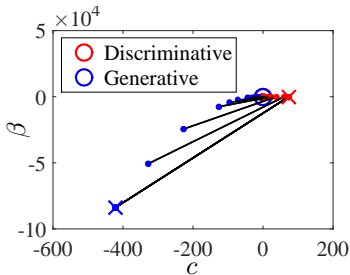
позволяет параметризовать $\lambda \in [0, 1]$ переход от разделяющей к порождающей модели:

- $\lambda \rightarrow 0 \Rightarrow \sigma(\lambda) \rightarrow 0 \Rightarrow p(\mathbf{w}_D, \mathbf{w}_G) \approx p(\mathbf{w}_D)\delta(\mathbf{w}_D - \mathbf{w}_G)$.
- $\lambda \rightarrow 1 \Rightarrow \sigma(\lambda) \rightarrow \infty \Rightarrow p(\mathbf{w}_D, \mathbf{w}_G) \approx p(\mathbf{w}_D)p(\mathbf{w}_G)$.

Итерации настройки параметров $\mathbf{w} = [\beta, c]$ при экстремальных значениях параметра λ без дополнительной регуляризации (бесконечная дисперсия априорных распределений $p(\mathbf{w})$).

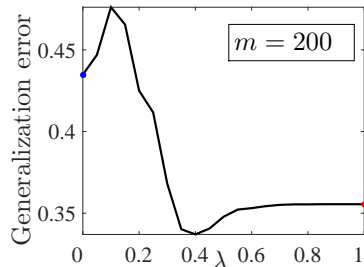
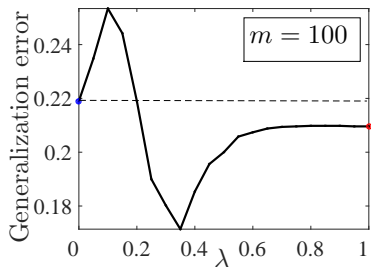
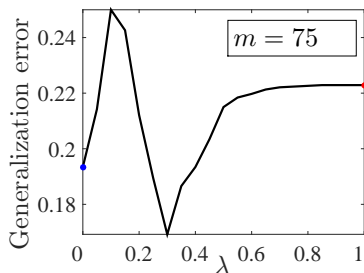
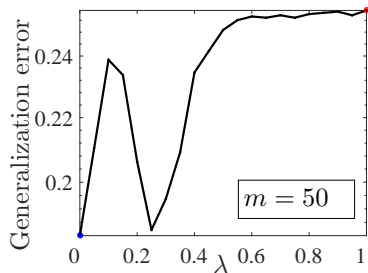


$\lambda = 0$



$\lambda = 1$

Пример: ошибка классификации $\hat{\varepsilon}_m(\lambda)$, $x \sim \mathcal{N}(\mu_y, \sigma_y^2)$



$$D_{\text{KL}}(\lambda) = \int_{\mathbf{w}} p(\mathbf{w}|D_1, \alpha, \lambda) \ln \frac{p(\mathbf{w}|D_1, \alpha, \lambda)}{p(\mathbf{w}|D_2, \alpha, \lambda)} d\mathbf{w}.$$

Выражение для $p(\mathbf{w}|D, \alpha, \lambda)$ найдем с помощью формулы Байеса

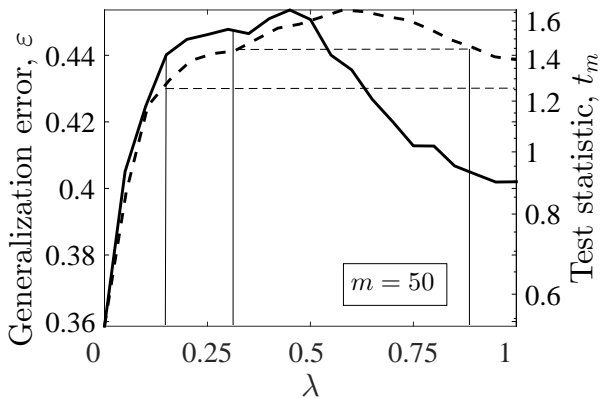
$$p(\mathbf{w}_\lambda|D, \alpha, \lambda) = \frac{p(D|\mathbf{w}_\lambda)p(\mathbf{w}_\lambda|\alpha, \lambda)}{p(D|\alpha, \lambda)},$$

где $p(D|\mathbf{w}_\lambda) = p(\mathbf{y}|X, \mathbf{w}_D)p(X|\mathbf{w}_G)$ — правдоподобие данных,
 $p(\mathbf{w}_\lambda|\alpha, \lambda) = p(\mathbf{w}_D, \mathbf{w}_G|\alpha, \lambda)$ — плотность распределения параметров модели, и

$$p(D|\alpha, \lambda) = \int p(D|\mathbf{w}_\lambda)p(\mathbf{w}_\lambda|\alpha, \lambda)d\mathbf{w}.$$

Для определения m^*

- в разделяющем случае положим
 $p(D|\mathbf{w}_\lambda) = p(\mathbf{y}|X, \mathbf{w}_D), p(\mathbf{w}_\lambda|\alpha, \lambda) = p(\mathbf{w}_D|\alpha);$
- в порождающем случае:
 $p(D|\mathbf{w}_\lambda) = p(\mathbf{y}, X|\mathbf{w}_G), p(\mathbf{w}_\lambda|\alpha, \lambda) = p(\mathbf{w}_G|\alpha).$



m	50	75	100	200	500
$\cos(D_{\text{KL}}(\lambda), \varepsilon_m(\lambda))$	0.916	0.933	0.926	0.969	0.8622

- Рассмотрена задача оценки объема выборки в задаче классификации с учетом используемой модели.
- Предложен способ оценки объема выборки, основанный на анализе распределения параметров модели.
- Предложенный способ проиллюстрирован задачей выбора между порождающим и разделяющим подходами к решению задачи классификации.

А. П. Мотренко, В. В. Стрижов. Построение агрегированных прогнозов объемов железнодорожных грузоперевозок с использованием расстояния Кульбака–Лейблера. Информатика и ее применения, 2014, 8(2), 86-97.

A. Motrenko, V. Strijov and G.-W. Weber. Sample Size Determination for Logistic Regression. Journal of Computational and Applied Mathematics, 2014, 255, 743-752.