

# *Regresyjne metody łączenia klasyfikatorów*

Tomasz Górecki, Mirosław Krzyśko

Wydział Matematyki i Informatyki  
Uniwersytet im. Adama Mickiewicza

XXXV Konferencja Statystyka Matematyczna  
Wiśła 7-11.12.2009

Założmy, że dysponujemy  $c$  różnymi klasyfikatorami  $\hat{d}_1, \hat{d}_2, \dots, \hat{d}_c$  skonstruowanymi z próby uczącej

$$\mathcal{L}_n = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\},$$

gdzie  $\mathbf{x}_i$  jest wartością wektora  $\mathbf{X} = (X_1, \dots, X_p)'$  obserwowanych  $p$  cech obiektu, a  $y_i \in \{1, 2, \dots, K\}$  etykietą skalarną jednej z ustalonych  $K$  klas klasyfikowanych obiektów. Ponieważ każdy z nich może być oparty na innej koncepcji, to może się zdarzyć, że żaden z nich nie jest jednoznacznie lepszy od pozostałych oraz, że poszczególne klasyfikatory mogą błędnie klasyfikować inne obserwacje.

LeBlanc i Tibshirani (1996) zaproponowali, by ostateczną decyzję klasyfikacyjną podejmować za pomocą metaklasyfikatora (klasyfikatora łączonego), wykorzystującego informację płynącą ze wszystkich  $c$  klasyfikatorów. Klasyfikator łączony otrzymany tą metodą jest kombinacją liniową ocen prawdopodobieństw a posteriori  $K$  klas uzyskanych przez  $c$  różnych klasyfikatorów. Współczynniki tej kombinacji liniowej znajdowane są metodą najmniejszych kwadratów w modelu liniowej regresji wielokrotnej. Metoda ta została nazwana regresją stosową.

W prezentowanym referacie proponujemy w charakterze metaklasyfikatorów wykorzystać, obok regresji stosowej, następujące modele regresyjne:

- regresję grzbietową,
- regresję składowych głównych,
- regresją logistyczną,
- regresję nieparametryczną.

Załóżmy, że każdej wartości  $\mathbf{x}$  wektora obserwowanych cech  $\mathbf{X}$ , obok etykiety skalarnej  $y \in \{1, 2, \dots, K\}$ , przypisana jest etykieta wektorowa  $\mathbf{Y} = (Y_1, Y_2, \dots, Y_K)'$ , gdzie  $Y_k = 1$ , jeżeli obserwacja  $\mathbf{x}$  pochodzi z klasy  $k$ ,  $k = 1, 2, \dots, K$ . Niech rozkład prawdopodobieństwa pary losowej  $(\mathbf{X}, \mathbf{Y})$  opisuje para  $(\mu, r)$ , gdzie  $\mu$  jest miarą probabilistyczną wektora losowego  $\mathbf{X}$ , a  $r$  jest funkcją regresji  $\mathbf{Y}$  względem  $\mathbf{X}$ . Dokładniej, dla dowolnego zbioru borelowskiego  $\mathfrak{B} \subseteq \mathbb{R}^p$ ,  $\mu(\mathfrak{B}) = P(\mathbf{X} \in \mathfrak{B})$  oraz dla każdego  $\mathbf{x} \in \mathbb{R}^p$

$$\begin{aligned} r(\mathbf{x}) &= (r_1(\mathbf{x}), \dots, r_K(\mathbf{x}))' = E(\mathbf{Y} | \mathbf{X} = \mathbf{x}) = \\ &= (E(Y_1 | \mathbf{X} = \mathbf{x}), \dots, E(Y_K | \mathbf{X} = \mathbf{x}))'. \end{aligned}$$

Ponieważ  $Y_k$  przyjmuje tylko dwie wartości 1 lub 0, to

$$r_k(\mathbf{x}) = E(Y_k | \mathbf{X} = \mathbf{x}) = P(Y_k = 1 | \mathbf{X} = \mathbf{x}), \quad k = 1, 2, \dots, K.$$

Jest to zatem prawdopodobieństwo a posteriori przynależności obserwacji  $\mathbf{x}$  do klasy o etykiecie  $k$ , gdzie  $k = 1, 2, \dots, K$ .

Uzyskaliśmy zatem klasyfikator postaci

$$\hat{d}(\mathbf{x}) = \arg \max_{1 \leq k \leq K} \hat{r}_k(\mathbf{x}),$$

gdzie  $\hat{r}_k(\mathbf{x})$  jest oceną z próby uczącej  $\mathfrak{L}_n = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  funkcji regresji  $r_k(\mathbf{x})$ ,  $k = 1, 2, \dots, K$ .

Funkcje regresji  $E(Y_k | \mathbf{X} = \mathbf{x})$  możemy estymować na wiele sposobów.

Dysponujemy  $c$  różnymi klasyfikatorami oraz  $K$  klasami. Mamy zatem  $cK$  ocen prawdopodobieństw a posteriori odpowiadających obserwacji  $\mathbf{x}_i$ ,  $i = 1, 2, \dots, n$ . Oceny te zapisujemy w postaci wektorów wierszowych

$$\hat{\mathbf{p}}'(\mathbf{x}_i) = (\hat{p}_1(1|\mathbf{x}_i), \dots, \hat{p}_1(K|\mathbf{x}_i), \dots, \hat{p}_c(1|\mathbf{x}_i), \dots, \hat{p}_c(K|\mathbf{x}_i))$$

i układamy w stos jako wiersze macierzy

$$\mathbf{P} = \begin{pmatrix} \hat{\mathbf{p}}'(\mathbf{x}_1) \\ \vdots \\ \hat{\mathbf{p}}'(\mathbf{x}_n) \end{pmatrix}$$

rozmiaru  $n \times cK$ .

Wielowymiarowy model regresji liniowej ma postać

$$\mathbf{Y}^* = \mathbf{PB} + \mathbf{E},$$

gdzie

$$\mathbf{Y}^* = \begin{pmatrix} Y_{11} & Y_{12} & \dots & Y_{1K} \\ Y_{21} & Y_{22} & \dots & Y_{2K} \\ \vdots & \vdots & \vdots & \vdots \\ Y_{n1} & Y_{n2} & \dots & Y_{nK} \end{pmatrix}$$

jest macierzą rozmiaru  $n \times K$  wartości etykiety wektorowej

$\mathbf{Y} = (Y_1, Y_2, \dots, Y_K)'$  przypisanej obserwacjom  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ ,  $\mathbf{B}$  jest macierzą rozmiaru  $cK \times K$  współczynników regresji, natomiast  $\mathbf{E}$  jest macierzą rozmiaru  $n \times cK$  błędów losowych.



## Regresja stosowa (3)

Ponieważ kolumny macierzy  $\mathbf{P}$  podlegają  $c$  ograniczeniom liniowym, to  $\mathbf{P}$  nie jest pełnego rzędu kolumnowego i  $\mathbf{P}'\mathbf{P}$  jest macierzą osobliwą. Stąd za estymator macierzy  $\mathbf{B}$ , uzyskany metodą najmniejszych kwadratów, przyjmujemy

$$\hat{\mathbf{B}} = (\mathbf{P}'\mathbf{P})^+ \mathbf{P}'\mathbf{Y}^*,$$

gdzie  $(\mathbf{P}'\mathbf{P})^+$  jest uogólnioną odwrotnością Moore'a–Penrose'a macierzy  $\mathbf{P}'\mathbf{P}$ .

Niech

$$\hat{\mathbf{u}}(\mathbf{x}) = (\hat{u}_1(\mathbf{x}), \hat{u}_2(\mathbf{x}), \dots, \hat{u}_K(\mathbf{x})) = \hat{\mathbf{p}}'(\mathbf{x})\hat{\mathbf{B}}.$$

Wówczas

$$\hat{d}_{\text{stos}}(\mathbf{x}) = \arg \max_{1 \leq k \leq K} \hat{u}_k(\mathbf{x}).$$

Ponieważ problemy w metodzie regresji stosowej pojawiają się w związku z niemożnością odwrócenia macierzy  $\mathbf{P}'\mathbf{P}$ , to do jej przekątnej można dodać pewną stałą  $\lambda \geq 0$ . Dla takiego zagadnienia otrzymuje się następujące rozwiązanie:

$$\hat{\mathbf{B}} = (\mathbf{P}'\mathbf{P} + \lambda\mathbf{I})^{-1}\mathbf{P}'\mathbf{Y}^*,$$

które nazywane jest regresją grzbietową. Optymalna wartość parametru  $\lambda$  wybierana jest za pomocą metody sprawdzania krzyżowego.

Próba uniknięcia problemu zależności zmiennych objaśniających (a co za tym idzie problemów z odwracaniem macierzy) jest regresja składowych głównych (regresja PCR). Zamiast oryginalnych zmiennych objaśniających (prawdopodobieństw) używamy składowych głównych, które są nieskorelowane. W praktyce używamy jedynie kilku pierwszych składowych, które w zadowalający sposób odzwierciedlają zmienność oryginalnych danych.

Funkcję regresji

$$p_k = r_k(\mathbf{x}) = E(Y_k | \mathbf{X} = \mathbf{x}) = P(Y_k = 1 | \mathbf{X} = \mathbf{x})$$

modelujemy za pomocą funkcji logistycznej

$$p_k = \frac{\exp(\beta_{0k} + \beta'_k \hat{\mathbf{p}}(\mathbf{x}))}{1 + \sum_{k=1}^{K-1} \exp(\beta_{0k} + \beta'_k \hat{\mathbf{p}}(\mathbf{x}))},$$
$$p_K = \frac{1}{1 + \sum_{k=1}^{K-1} \exp(\beta_{0k} + \beta'_k \hat{\mathbf{p}}(\mathbf{x}))}$$

gdzie

$$\beta_k = \left( \beta_{1k}^{(1)}, \dots, \beta_{Kk}^{(1)}, \dots, \beta_{1k}^{(c)}, \dots, \beta_{Kk}^{(c)} \right)', \quad k = 1, 2, \dots, K-1$$

oraz

$$p_1 + p_2 + \dots + p_K = 1.$$

Parametry  $\beta_{0k}$  i  $\beta_k$  estymujemy metodą największej wiarygodności.

Nieparametryczny estymator Nadarayi–Watsona funkcji regresji  $r_k(\mathbf{x})$  ma postać:

$$\hat{r}_k(\mathbf{x}) = \frac{\sum_{i=1}^n K\left(\frac{\hat{\mathbf{p}}(\mathbf{x}) - \hat{\mathbf{p}}(\mathbf{x}_i)}{h}\right) I(y_i = k)}{\sum_{i=1}^n K\left(\frac{\hat{\mathbf{p}}(\mathbf{x}) - \hat{\mathbf{p}}(\mathbf{x}_i)}{h}\right)}, \quad k = 1, 2, \dots, K,$$

gdzie  $K(\cdot)$  jest funkcją jądrową zależną od parametru gładkości  $h$ .  
Optymalną wartość parametru gładkości  $h$  dobieramy metodą sprawdzania krzyżowego.

Jedna z najpopularniejszych metod klasyfikacji. Idea jej działania jest prosta i intuicyjna. Nowy obiekt otrzymuje etykietę klasy obiektu, który jest najbliżej (w uogólnieniu tej metody, czyli metodzie  $k$  najbliższych sąsiadów, etykietę klasy, która występuje najczęściej pośród jego  $k$  sąsiadów). W tym miejscu należy podkreślić, że do oceny odległości może zostać wykorzystana bardzo szeroka klasa funkcji. Jest to metoda nieparametryczna, nie wymaga zatem żadnych założeń co do rozkładów danych w klasach.

Jest to prosty klasyfikator probabilistyczny oparty na założeniu, że cechy opisujące obiekty są wzajemnie niezależne. Założenie to nie ma raczej nic wspólnego z rzeczywistością i właśnie z tego powodu metoda nazywana jest naiwną. Pomimo tego klasyfikator ten często działa lepiej niż można się po nim było spodziewać (zwłaszcza jeśli jest dużo cech). W praktyce estymuje się gęstość każdej cechy w każdej klasie, a następnie bada iloczyn takich gęstości dla każdej klasy. Obserwacja klasyfikowana jest do klasy, dla której ten iloczyn jest największy. Najczęściej zakłada się, że cechy mają rozkłady normalne.

Jest to klasyfikator oparty na sekwencyjnym dzieleniu podzbiorów przestrzeni próby. W każdym kroku podział dokonywany jest tak, aby uzyskane części były jak najbardziej jednorodne. Cały proces przedstawiany jest za pomocą drzewa binarnego (stąd nazwa). Drzewo składa się z korzenia oraz gałęzi prowadzących z korzenia do kolejnych węzłów. W każdym węźle sprawdzany jest pewien warunek dotyczący danej obserwacji, i na jego podstawie wybierana jest jedna z gałęzi prowadząca do kolejnego węzła poniżej. Na dole znajdują się liście, z których odczytujemy do której z klas należy przypisać daną obserwację.






	Liczba przypadków	Liczba cech		Liczba klas
		C	D	
Balance	625	4	0	3
Breast-W	699	0	9	2
Car	1728	0	6	4
Diabets	768	8	0	2
Echo	132	5	1	2
German	1000	7	13	2
Glass	214	9	0	6
Heart-C	303	6	7	5
Heart-H	294	6	7	5
Heart	270	7	6	2
Hepatitis	155	6	13	2
Ionosphere	351	34	0	2
Iris	150	4	0	3
Sonar	208	60	0	2
Wine	178	13	0	3

	$k$	$h$	$\lambda$
Balance	10	0.6	0.09
Breast-W	5	0.6	0.01
Car	7	0.4	0.01
Diabets	20	0.6	0.01
Echo	9	1.2	0.01
German	14	0.8	0.01
Glass	1	0.2	0.02
Heart-C	35	0.4	0.21
Heart-H	17	0.6	0.01
Heart	24	0.8	0.01
Hepatitis	9	0.4	0.01
Ionosphere	2	0.8	0.01
Iris	10	0.4	0.01
Sonar	2	0.6	0.01
Wine	1	0.4	0.09

## Porównanie regresji (% błędnej klasyfikacji)

	R	R Grzb.	R PCR	R Logist.	R NParam.
Balance	22.26	22.26	22.66	22.26	<b>14.93</b>
Breast-W	4.95	4.95	4.95	4.93	<b>3.97</b>
Car	<b>5.21</b>	<b>5.21</b>	5.46	5.39	5.43
Diabets	34.90	34.90	34.90	34.90	<b>33.22</b>
Echo	37.61	37.61	37.61	37.61	<b>34.99</b>
German	32.66	32.66	32.66	32.66	<b>31.94</b>
Glass	36.29	<b>36.24</b>	36.48	38.15	37.14
Heart-C	53.28	51.83	51.08	53.82	<b>41.54</b>
Heart-H	53.96	54.00	47.32	54.00	<b>35.59</b>
Heart	36.63	36.63	36.63	36.63	<b>28.59</b>
Hepatitis	19.78	19.78	19.78	19.78	<b>18.47</b>
Ionosphere	12.94	12.94	12.94	12.91	<b>12.19</b>
Iris	7.67	7.67	7.73	9.27	<b>6.13</b>
Sonar	30.60	30.60	30.60	29.55	<b>25.35</b>
Wine	4.94	4.82	4.88	6.46	<b>3.65</b>
Średnia	19.78	19.72	19.64	20.47	<b>16.80</b>

Jak widać najlepiej wypada metoda oparta o regresję nieparametryczną, która okazuje się najskuteczniejsza dla wszystkich zbiorów z wyjątkiem dwóch. Na tych dwóch zbiorach ustępuje jednak innym metodom minimalnie. Najgorsza wydaje się metoda oparta o regresję logistyczną. Natomiast klasyczna regresja stosowa oraz jej wersja grzbietowa oraz regresja składowych głównych wypadają bardzo podobnie.

-  Krzyśko M., Wołyński W., Górecki T., Skorzybut M. (2008). *Systemy uczące się*. WNT.
-  LeBlanc M., Tibshirani R. (1996), Combining Estimates in Regression and Classification. *JASA* 91:641–1650.
-  Merz C.J., Murphy P.M. (1998). UCI repository of machine learning databases. Machine readable data repository <http://www.ics.uci.edu/~mlern/MLRepository.html>. Irvine, CA: University of California, Department of Information and Computer Science.