

Zadanie uczenia.

Podział metod uczenia

- Uczenie pod nadzorem – w przypadku gdy zbiór uczący składa się z wektorów cech (atrybutów) oraz wektorów odpowiedzi.
 - wektor wejściowy – wektor cech (zmiennych, atrybutów) opisujących
 - wektor wyjściowy – wektor zmiennych opisywanych
 - cel uczenia – nauczenie systemu odpowiedzi na wektory wejściowe
 - Przykład – zadanie klasyfikacji. Wyjściem etykieta klasy

A. Brückner

Podstawy sztucznej inteligencji

- 191 -

Zadanie uczenia.

Podział metod uczenia

- Uczenie nienadzorowane – zbiór uczący składający się tylko z wejściowych wektorów cech
 - cel uczenia – opisanie, objaśnienie zbioru danych wejściowych tylko na podstawie ich samych. Wykrycie wewnętrznej struktury danych, wykrycie współzależności
 - Przykład – zadanie grupowania danych w rozłączne klasy (klasyfikacja nienadzorowana)

A. Brückner

Podstawy sztucznej inteligencji

- 192 -

Zadanie uczenia klasyfikacji...

Założenia:

- Dysponujemy N niezależnymi obserwacjami (próbkami) pochodzącymi z k populacji (klas, grup)
- Wszystkie obserwacje są wektorami losowymi (wektory cech) o tym samym skończonym wymiarze p

Zadanie klasyfikacji polega na przypisaniu obiektowi numeru klasy na podstawie jego wektora cech.

A. Brückner

Podstawy sztucznej inteligencji

- 193 -

Zadanie uczenia klasyfikacji...

Definicje

- Zbiór obserwacji (próbek) nazywamy *zbiorem wektorów uczących* - oznaczamy symbolem V
- Każdy element zbioru V opisywany jest przez zestaw atrybutów (cech) a_1, \dots, a_p ze zbioru A oraz przez klasę, którą reprezentuje
- Atrybuty tworzą tak zwaną *przestrzeń cech*

A. Brückner

Podstawy sztucznej inteligencji

- 194 -

Zadanie uczenia klasyfikacji...

- Zakładamy, że przestrzeń cech jest podzbiorem R^p , a każdy obiekt reprezentowany jest p wymiarowym wektorem cech
- Zakładamy, że klasę identyfikujemy poprzez etykietę oznaczającą jej numer:

- Przy powyższych założeniach

$$V = \{(x_1, y_1), \dots, (x_N, y_N)\}, \text{ gdzie } x_i = [x_{i1}, \dots, x_{im}]^T,$$

$$y_i \in \{1, \dots, k\} \text{ dla } i = 1, \dots, N,$$

A. Brückner

Podstawy sztucznej inteligencji

- 195 -

Zadanie uczenia klasyfikacji...

Definicja:

Algorytmem klasyfikacji (regułą decyzyjną) nazywamy funkcję odwzorowującą przestrzeń cech w zbiór numerów klas, to znaczy (przy przyjętych założeniach) funkcję:

$$\psi: R^p \rightarrow \{1, \dots, K\}$$

taką, że $\psi(x) = i$ gdy wektor x jest elementem klasy i .

- Każda próbka może należeć tylko i wyłącznie do jednej klasy

A. Brückner

Podstawy sztucznej inteligencji

- 196 -

Zadanie uczenia klasyfikacji...

- Reguła decyzyjna generuje rozkład przestrzeni cech na tak zwane obszary decyzyjne $D_x^i = \{x \in R^p : \psi(x) = i\}$ o następujących własnościach:
 - $D_x^i \cap D_x^j = \Phi$ dla $i, j \in \{1, \dots, k\}$, $i \neq j$, co znaczy, że żaden obiekt nie może należeć do dwóch klas jednocześnie
 - $\bigcup_{i \in \{1, \dots, k\}} D_x^i = R^p$, co znaczy, że każdy obiekt musi należeć do jednej z klas.

Zadanie uczenia klasyfikacji...

Definicja:

Funkcja klasyfikującą (dyskryminującą) i -tej klasy nazywamy funkcję odwzorowującą przestrzeń cech w zbiór liczb rzeczywistych $g_i : R^p \rightarrow R$, taką, że:

$$\forall_{i \in \{1, \dots, k\}} \forall_{x \in D_x^i} g_i(x) = \max_{j \in \{1, \dots, k\}} g_j(x)$$

- Definicja oznacza, że, na i -tym obszarze decyzyjnym wartość i -tej funkcji klasyfikującej jest największa.

Zadanie uczenia klasyfikacji...

- Zadanie klasyfikacji sprowadza się w ten sposób do przydzielenia wektora cech x do klasy, dla której odpowiednia funkcja klasyfikująca na wektorze x osiąga największą wartość:
$$\psi(x) = i, \text{ gdy } g_i(x) = \max_{j \in \{1, \dots, k\}} g_j(x)$$
- wystarczy więc wyliczyć wartości wszystkich funkcji klasyfikujących, znaleźć wartość największą, a jako wynik działania algorytmu zwrócić indeks funkcji dla której ta wartość została osiągnięta.
- Uczenie klasyfikacji** – polega na znalezieniu postaci funkcji klasyfikujących

Przykłady zadania klasyfikacji

- Wspomaganie diagnostyki medycznej
 - klasyfikacja pacjenta ze względu na jednostkę chorobową
 - określenie przynależności pacjenta do grupy o podwyższonym ryzyku
- Rozpoznawanie cyfr / pisma
- Określenie zdolności kredytowej klienta banku
- Rozpoznawanie niechcianych wiadomości elektronicznych (spamu)

Optymalny klasyfikator statystyczny

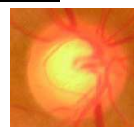
- W praktyce zadanie klasyfikacji na ogół nie jest jednoznaczne, ponieważ zdarza się, że obiekty należące do dwóch różnych klas reprezentowane są przez ten sam wektor cech.
 - Klasyfikowany obiekt przydziela się do klasy w której prawdopodobieństwo jego wystąpienia jest największe
- optymalny klasyfikator statystyczny**

Optymalny klasyfikator statystyczny

Model matematyczny zadania klasyfikacji:

Zakładamy, że wektor $x = (x_1, \dots, x_p)^T \in R^p$ oraz numer klasy j do której należy jest realizacją pary zmiennych losowych (X, J) , $X : \Omega \rightarrow R^p$, gdzie Ω przestrzeń klasyfikowanych obserwacji, J - dyskretna zmienna losowa o wartościach w zbiorze $\{1, \dots, K\}$

Przykład:



klasyfikowany obiekt

$$x = \begin{bmatrix} 0,7875 \\ 1,561 \\ -2,432 \end{bmatrix}$$

wektor cech

$$P(J = j | X = x) = ??$$

prawdopodobieństwa przynależności do klas

<p>Prawdopodobieństwo całkowite, wzór Bayesa.</p> <p>Twierdzenie (o prawdopodobieństwie całkowitym)</p> <p>Niech $A, B_1, \dots, B_n \in \mathfrak{S} \subset \Omega$, $P(B_i) > 0$ dla $i = 1, \dots, n$,</p> $B_1 \cup \dots \cup B_n = \Omega \text{ oraz } B_i \cap B_j = \emptyset \text{ dla } i \neq j. \text{ Wtedy:}$ $P(A) = P(A B_1)P(B_1) + \dots + P(A B_n)P(B_n).$ <p>Prawdopodobieństwo $P(A)$ nazywamy całkowitym (zupełnym) prawdopodobieństwem zajścia zdarzenia A.</p> <div> <div>A. Brückner</div> <div>Podstawy sztucznej inteligencji</div> <div>- 203 -</div> </div>	<p>Prawdopodobieństwo całkowite, wzór Bayesa.</p> <p>Twierdzenie (Bayesa)</p> <p>Niech $A, B_1, \dots, B_n \in \mathfrak{S} \subset \Omega$, $P(B_i) > 0$ dla $i = 1, \dots, n$,</p> $B_1 \cup \dots \cup B_n = \Omega \text{ oraz } B_i \cap B_j = \emptyset \text{ dla } i \neq j. \text{ Wtedy:}$ $P(B_i A) = \frac{P(A B_i)P(B_i)}{P(A)}.$ <p><u>Uwaga:</u> Na mocy twierdzenia o prawdopodobieństwie całkowitym:</p> $P(B_i A) = \frac{P(A B_i)P(B_i)}{P(A B_1)P(B_1) + \dots + P(A B_n)P(B_n)}$ <div> <div>A. Brückner</div> <div>Podstawy sztucznej inteligencji</div> <div>- 204 -</div> </div>
<p>Prawdopodobieństwo całkowite, wzór Bayesa.</p> <p><u>Przykład:</u></p> <p>Jeden z testów na obecność w organizmie wirusa HCV daje pozytywny rezultat w 97% przypadków osób zarażonych tą chorobą i mylnie wskazuje na jego obecność u osób zdrowych w 3% przypadków. Ile wynosi prawdopodobieństwo, że osoba, u której wykryto tym testem obecność wirusa HCV, jest faktycznie chora, jeśli wiadomo, że u 2% populacji stwierdza się tego wirusa.</p> <div> <div>A. Brückner</div> <div>Podstawy sztucznej inteligencji</div> <div>- 205 -</div> </div>	<p>Zmienna losowa</p> <p>Definicja</p> <p>Niech $(\Omega, \mathfrak{S}, P)$ będzie dowolną przestrzenią probabilistyczną.</p> <p>Zmienną losową nazywamy każdą funkcję X określoną na przestrzeni zdarzeń elementarnych Ω o wartościach rzeczywistych, taką że każdemu przedziałowi wartości zmiennej X postaci $(-\infty, x)$ odpowiada zdarzenie losowe, to znaczy</p> $\{\omega : X(\omega) < x\} \in \mathfrak{S} \text{ dla każdego } x \in R$ <p>Gdy przestrzeń zdarzeń elementarnych jest skończona, a każdy jej podzbiór jest zdarzeniem elementarnym to każda funkcja $X : \Omega \rightarrow R$ jest zmienną losową.</p> <div> <div>A. Brückner</div> <div>Podstawy sztucznej inteligencji</div> <div>- 206 -</div> </div>
<p>Zmienna losowa</p> <p>Oznaczenia</p> <ul style="list-style-type: none"> $P(X = x)$ - prawdopodobieństwo, że zmienna losowa X przyjmuje wartość x, czyli $P(X = x) = P\{\omega \in \Omega : X(\omega) = x\}$ $P(X < x)$ - prawdopodobieństwo, że zmienna losowa X przyjmuje wartości mniejsze od x. $P(X \geq x)$ - prawdopodobieństwo, że zmienna losowa X przyjmuje wartości większe lub równe x. $P(a \leq X \leq b)$ - prawdopodobieństwo, że zmienna losowa X przyjmuje wartości z przedziału $[a, b]$ <div> <div>A. Brückner</div> <div>Podstawy sztucznej inteligencji</div> <div>- 207 -</div> </div>	<p>Zmienna losowa</p> <p>Uwaga</p> <ul style="list-style-type: none"> Zmienną losową, która przyjmuje skończoną bądź przeliczalną liczbę wartości nazywamy zmienną losową dyskretną (typu dyskretnego) Zmienną losową, która przyjmuje wartości z pewnego przedziału nazywamy zmienną losową ciągłą (typu ciągłego) <div> <div>A. Brückner</div> <div>Podstawy sztucznej inteligencji</div> <div>- 208 -</div> </div>

Zmienna losowa

Definicja

Niech $X : \Omega \rightarrow R$ będzie zmienną losową. Funkcję $F : R \rightarrow [0,1]$ określoną wzorem $F(x) = P(X < x)$ nazywamy dystrybuantą zmiennej losowej X .

Własności:

1. $0 \leq F(x) \leq 1$ dla każdego $x \in R$
2. $\lim_{x \rightarrow -\infty} F(x) = 0$ oraz $\lim_{x \rightarrow \infty} F(x) = 1$
3. $F(x)$ jest funkcją niemalejącą
4. $F(x)$ jest funkcją lewostronnie ciągłą, to znaczy
$$\lim_{x \rightarrow x_0^-} F(x) = F(x_0)$$
5. $P(a \leq X < b) = P(X \in [a, b)) = F(b) - F(a)$
6. $P(X = x_0) = \lim_{x \rightarrow x_0^-} F(x) - F(x_0)$
7. Każda funkcja spełniająca własności 2-4 jest dystrybuantą pewnej zmiennej losowej.

Zmienna losowa

Zmienna losowa dyskretna

- Zmienna losowa, która przyjmuje skończoną bądź przeliczalną liczbę wartości x_1, \dots, x_n, \dots
- $P(X = x_i) = p_i > 0$ - to znaczy p_i oznacza prawdopodobieństwo przyjęcia przez zmienną losową X wartości x_i
- $\sum_i p_i = 1$

Zmienna losowa

Definicja

Niech X będzie zmienną losową dyskretną. Funkcję p określoną na zbiorze wartości zmiennej losowej określoną równością $p(x_i) = P(X = x_i) = p_i$ nazywamy **funkcją rozkładu prawdopodobieństwa** zmiennej losowej X lub krócej **rozkładem zmiennej losowej X** .
Równoważnie rozkład prawdopodobieństwa dyskretnej zmiennej losowej podaje się najczęściej w postaci tablicy

x_i	x_1	x_2	\dots	x_n	\dots
p_i	p_1	p_2	\dots	p_n	\dots

Oczywiście dla dowolnego rozkładu $\sum_i x_i$

Zmienna losowa

Definicja

Zmienną losową X przyjmującą wszystkie wartości z pewnego przedziału liczbowego, bądź przedziałów, dla której istnieje nieujemna funkcja $f : R \rightarrow R^+ \cup \{0\}$ taka, że dystrybuantę zmiennej losowej X przedstawić można w postaci
$$F(x) = \int_{-\infty}^x f(t) dt$$
 nazywamy **zmienną losową typu ciągłego**, funkcję f natomiast nazywamy **gęstością rozkładu** zmiennej losowej X .

Zmienna losowa

Zmienna losowa typu ciągłego. Definicja.

Mówimy że dany jest rozkład zmiennej losowej X typu ciągłego gdy dana jest jej funkcja gęstości bądź dystrybuanta.

Własności

1. $F'(x) = f(x)$
2. $\int_{-\infty}^{\infty} f(x) = 1$
3. $P(a \leq X \leq b) = F(b) - F(a) = \int_a^b f(x) dx$
4. zarówno funkcja gęstości oraz dystrybuanta są funkcjami ciągłymi

<div data-bbox="52 313 266 344" data-label="Section-Header"> <h2>Zmienna losowa</h2> </div> <div data-bbox="52 371 525 400" data-label="Section-Header"> <h3>Charakterystyki liczbowe zmiennej losowej</h3> </div> <div data-bbox="52 427 668 678" data-label="List-Group"> <ul style="list-style-type: none"> Opisują takie własności zmiennej losowej jak wartość najbardziej prawdopodobna, rozrzut wartości, kształt histogramu czy krzywej gęstości Podstawowe charakterystyki liczbowe zmiennej losowej <ul style="list-style-type: none"> Wartość oczekiwana (przeciętna) Wariancja Odchylenie standardowe </div> <div data-bbox="52 795 761 813" data-label="Page-Footer"> <div>A. Brückner</div> <div>Podstawy sztucznej inteligencji</div> <div>- 215 -</div> </div>	<div data-bbox="850 313 1064 344" data-label="Section-Header"> <h2>Zmienna losowa</h2> </div> <div data-bbox="850 371 1323 400" data-label="Section-Header"> <h3>Charakterystyki liczbowe zmiennej losowej</h3> </div> <div data-bbox="850 407 1540 790" data-label="List-Group"> <ol style="list-style-type: none"> Wartość oczekiwana zmiennej losowej X, oznaczana $EX = m$, określa wartość wokół której skupiają się realizacje zmiennej losowej uzyskiwane w wyniku wielokrotnego powtarzania tego samego eksperymentu. Niech X będzie zmienną losową dyskretną. Wówczas jej wartość oczekiwana jest równa $EX = \sum_i x_i p_i$. Niech X będzie zmienną losową ciągłą. Wówczas jej wartość oczekiwana wyraża się wzorem $EX = \int_R x f(x) dx$, gdzie $f(x)$ jest gęstością. </div> <div data-bbox="850 795 1559 813" data-label="Page-Footer"> <div>A. Brückner</div> <div>Podstawy sztucznej inteligencji</div> <div>- 216 -</div> </div>
<div data-bbox="52 880 266 911" data-label="Section-Header"> <h2>Zmienna losowa</h2> </div> <div data-bbox="52 918 400 947" data-label="Section-Header"> <h3>Własności wartości oczekiwanej</h3> </div> <div data-bbox="52 956 721 1137" data-label="List-Group"> <ol style="list-style-type: none"> $E(C) = C$, gdzie C jest stałą $E(CX) = CE(X)$ $E(X \pm Y) = E(X) \pm E(Y)$ $E(XY) = E(X) \cdot E(Y)$ gdy X, Y są zmiennymi niezależnymi. </div> <div data-bbox="52 1153 710 1299" data-label="Text"> <p>Definicja. Zmienne losowe X, Y określone na tej samej przestrzeni zdarzeń elementarnych nazywamy niezależnymi gdy dla dowolnych x, y niezależne są zdarzenia $\{X < x\}$ oraz $\{Y < y\}$ tzn. gdy:</p> </div> <div data-bbox="52 1314 469 1348" data-label="Equation-Block"> $P(X < x \wedge Y < y) = P(X < x) \cdot P(Y < y)$ </div> <div data-bbox="52 1361 761 1379" data-label="Page-Footer"> <div>A. Brückner</div> <div>Podstawy sztucznej inteligencji</div> <div>- 217 -</div> </div>	<div data-bbox="850 880 1064 911" data-label="Section-Header"> <h2>Zmienna losowa</h2> </div> <div data-bbox="850 918 1323 947" data-label="Section-Header"> <h3>Charakterystyki liczbowe zmiennej losowej</h3> </div> <div data-bbox="850 956 1522 1122" data-label="List-Group"> <ol style="list-style-type: none"> Wariancja zmiennej losowej oznaczana D^2X jest miarą rozproszenia wokół wartości średniej, wzrasta wraz ze wzrostem rozproszenia kolejnych realizacji zmiennej losowej. Wariancję zmiennej losowej określamy za pomocą następujących wzorów: </div> <div data-bbox="890 1128 1345 1220" data-label="Equation-Block"> $D^2X = E(X - EX)^2$ $D^2X = E(X^2) - (EX)^2 \text{ (Dowód. ćwiczenie)}$ </div> <div data-bbox="850 1232 1412 1359" data-label="List-Group"> <ul style="list-style-type: none"> W przypadku dyskretnym: $D^2X = \sum_i x_i^2 p_i - (EX)^2$ W przypadku ciągłym: $D^2X = \int_R x_i^2 f(x) dx - (EX)^2$ </div> <div data-bbox="850 1361 1559 1379" data-label="Page-Footer"> <div>A. Brückner</div> <div>Podstawy sztucznej inteligencji</div> <div>- 218 -</div> </div>
<div data-bbox="52 1447 266 1478" data-label="Section-Header"> <h2>Zmienna losowa</h2> </div> <div data-bbox="52 1505 272 1534" data-label="Section-Header"> <h3>Własności wariancji</h3> </div> <div data-bbox="52 1559 647 1787" data-label="List-Group"> <ol style="list-style-type: none"> $D^2(C) = 0$, gdzie C jest stałą $D^2(CX) = C^2 D^2(X)$ $D^2(X \pm C) = D^2(X)$ $D^2(X \pm Y) = D^2(X) + D^2(Y)$ gdy X, Y są zmiennymi niezależnymi. </div> <div data-bbox="52 1926 761 1944" data-label="Page-Footer"> <div>A. Brückner</div> <div>Podstawy sztucznej inteligencji</div> <div>- 219 -</div> </div>	<div data-bbox="850 1447 1064 1478" data-label="Section-Header"> <h2>Zmienna losowa</h2> </div> <div data-bbox="850 1505 1121 1534" data-label="Section-Header"> <h3>Wybrane rozkłady ciągłe</h3> </div> <div data-bbox="850 1559 1522 1879" data-label="List-Group"> <ol style="list-style-type: none"> Rozkład normalny: Mówimy, że zmienna losowa X ma rozkład normalny w wartością oczekiwaną μ i odchyleniem standardowym σ, co oznaczamy $X \sim N(\mu, \sigma)$ jeśli jej funkcja gęstości określona jest wzorem: $f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$ Uwaga: Dystrybuanta rozkładu normalnego nie wyraża się za pomocą funkcji elementarnych – jest stabilizowana </div> <div data-bbox="850 1926 1559 1944" data-label="Page-Footer"> <div>A. Brückner</div> <div>Podstawy sztucznej inteligencji</div> <div>- 220 -</div> </div>

<div data-bbox="52 313 268 347" data-label="Section-Header"> <h2>Zmienna losowa</h2> </div> <div data-bbox="52 371 403 400" data-label="Section-Header"> <h3>Własności rozkładu normalnego</h3> </div> <div data-bbox="52 409 743 788" data-label="List-Group"> <ul style="list-style-type: none"> • Symetryczny względem prostej $x = \mu$ co oznacza, że $P(X < \mu) = P(X > \mu) = 0,5$ • Funkcja gęstości osiąga maksimum w punkcie $x = \mu$ wynoszące $f(\mu) = \frac{1}{\sqrt{2\pi}\sigma}$ • Prawdopodobieństwo, że zmienna losowa X przyjmuje wartości z przedziału $[m - 3\sigma, m + 3\sigma]$ jest w przybliżeniu równe 1 • Rozkład normalny $N(0,1)$ nazywamy rozkładem normalnym standardowym </div> <div data-bbox="52 797 761 813" data-label="Page-Footer"> <div>A. Brückner</div> <div>Podstawy sztucznej inteligencji</div> <div>- 221 -</div> </div>	<div data-bbox="850 313 1345 347" data-label="Section-Header"> <h2>Optymalny klasyfikator statystyczny...</h2> </div> <div data-bbox="850 398 1380 465" data-label="Text"> <p>Rozkład zmiennej losowej I scharakteryzowany jest prawdopodobieństwami a priori klas:</p> </div> <div data-bbox="1061 472 1332 506" data-label="Equation-Block"> $P(J = j) = p_j \text{ dla } j=1,...,k.$ </div> <div data-bbox="850 560 1516 672" data-label="Text"> <p>Wektor losowy X dla każdego $j \in \{1,...,k\}$ ma natomiast rozkład prawdopodobieństwa wyrażony gęstością, zwaną gęstością warunkową cech w klasie.</p> </div> <div data-bbox="1054 678 1340 712" data-label="Equation-Block"> $f(x j) = f_j(x), \text{ dla } x \in R^p.$ </div> <div data-bbox="850 797 1559 813" data-label="Page-Footer"> <div>A. Brückner</div> <div>Podstawy sztucznej inteligencji</div> <div>- 222 -</div> </div>
<div data-bbox="52 880 549 913" data-label="Section-Header"> <h2>Optymalny klasyfikator statystyczny...</h2> </div> <div data-bbox="52 965 625 994" data-label="Text"> <p>Gęstość bezwarunkowa zmiennej X wyraża się wzorem:</p> </div> <div data-bbox="234 1001 568 1070" data-label="Equation-Block"> $f(x) = \sum_{j \in \Theta} p_j f_j(x) = \sum_{j=1}^M p_j f_j(x).$ </div> <div data-bbox="52 1122 537 1153" data-label="Text"> <p>Przyjmujemy, że $f(x) > 0$ dla każdego $x \in R^p$.</p> </div> <div data-bbox="52 1202 707 1267" data-label="Text"> <p>Obliczmy prawdopodobieństwo warunkowe, że obiekt, któremu odpowiada wektor cech x należy do klasy j, to znaczy</p> </div> <div data-bbox="186 1274 614 1308" data-label="Equation-Block"> $P(J = j X = x) \text{ dla } j \in \{1,...,k\}, x \in R^p.$ </div> <div data-bbox="52 1361 761 1377" data-label="Page-Footer"> <div>A. Brückner</div> <div>Podstawy sztucznej inteligencji</div> <div>- 223 -</div> </div>	<div data-bbox="850 880 1345 913" data-label="Section-Header"> <h2>Optymalny klasyfikator statystyczny...</h2> </div> <div data-bbox="850 965 1153 994" data-label="Text"> <p>Stosując wzór Bayesa, mamy:</p> </div> <div data-bbox="924 1001 1473 1093" data-label="Equation-Block"> $P(J = j X = x) = \frac{P(X = x J = j) \cdot P(J = j)}{\sum_{j \in \{1,...,k\}} P(X = x J = j) \cdot P(J = j)},$ </div> <div data-bbox="850 1144 973 1173" data-label="Text"> <p>Oznaczamy:</p> </div> <div data-bbox="928 1182 1072 1216" data-label="Equation-Block"> $p_j = P(J = j)$ </div> <div data-bbox="928 1227 1189 1263" data-label="Equation-Block"> $p_j(x) = P(J = j X = x),$ </div> <div data-bbox="928 1274 1347 1332" data-label="Equation-Block"> $f(x) = \sum_{j \in \{1,...,k\}} P(X = x J = j) \cdot P(J = j)$ </div> <div data-bbox="850 1361 1559 1377" data-label="Page-Footer"> <div>A. Brückner</div> <div>Podstawy sztucznej inteligencji</div> <div>- 224 -</div> </div>
<div data-bbox="52 1447 549 1480" data-label="Section-Header"> <h2>Optymalny klasyfikator statystyczny...</h2> </div> <div data-bbox="52 1494 108 1523" data-label="Text"> <p>Skąd:</p> </div> <div data-bbox="312 1529 489 1599" data-label="Equation-Block"> $p_j(x) = \frac{p_j f_j(x)}{f(x)}.$ </div> <div data-bbox="57 1615 700 1648" data-label="Text"> <p>$p_j(x)$ - nazywamy prawdopodobieństwem a posteriori klasy j.</p> </div> <div data-bbox="52 1700 730 1877" data-label="List-Group"> <ul style="list-style-type: none"> • Optymalny klasyfikator statystyczny (klasyfikator Bayesa) obiekt reprezentowany przez wektor cech x przyporządkowuje do klasy, dla której wartość prawdopodobieństwa a posteriori jest największa, czyli do klasy w której prawdopodobieństwo jego wystąpienia jest największe </div> <div data-bbox="52 1928 761 1944" data-label="Page-Footer"> <div>A. Brückner</div> <div>Podstawy sztucznej inteligencji</div> <div>- 225 -</div> </div>	<div data-bbox="850 1447 1345 1480" data-label="Section-Header"> <h2>Optymalny klasyfikator statystyczny...</h2> </div> <div data-bbox="850 1532 1442 1603" data-label="Text"> <p>Ponieważ dla dowolnego x wartość $f(x)$ jest stała, reguła decyzyjna optymalnego klasyfikatora Bayesa jest postaci:</p> </div> <div data-bbox="1007 1610 1386 1657" data-label="Equation-Block"> $\psi^*(x) = i, \text{ gdy } p_i(x) = \max_{j \in \{1,...,k\}} p_j(x),$ </div> <div data-bbox="850 1671 903 1700" data-label="Text"> <p>czyli:</p> </div> <div data-bbox="995 1706 1398 1751" data-label="Equation-Block"> $\psi^*(x) = i, \text{ gdy } p_i f_i(x) = \max_{1 \leq j \leq k} p_j f_j(x).$ </div> <div data-bbox="850 1765 1543 1830" data-label="Text"> <p>Ponieważ logarytm jest funkcją rosnącą regułę powyższą można zapisać w równoważnej postaci:</p> </div> <div data-bbox="952 1836 1441 1883" data-label="Equation-Block"> $\psi^*(x) = i, \text{ gdy } \ln(p_i f_i(x)) = \max_{1 \leq j \leq k} \ln(p_j f_j(x)).$ </div> <div data-bbox="850 1928 1559 1944" data-label="Page-Footer"> <div>A. Brückner</div> <div>Podstawy sztucznej inteligencji</div> <div>- 226 -</div> </div>

Optimalny klasyfikator statystyczny...

Przypadek normalności rozkładu cech w klasach:

Często spotykamy się z sytuacją gdy rozkłady wektorów cech obiektów w klasach pochodzą z wielowymiarowego rozkładu normalnego. Wtedy funkcje gęstości rozkładów cech w klasach wyrażają się wzorami:

$$f_j(x) = (2\pi)^{-\frac{p}{2}} |\Sigma_j|^{-1/2} \exp\left[-\frac{1}{2}(x-m_j)^T \Sigma_j^{-1}(x-m_j)\right], \text{ dla } j=1,2,\dots,M,$$

gdzie: m_j p -wymiarowy wektor wartości oczekiwanych,

Σ_j macierz kowariancji rozkładu w klasie j .

Optimalny klasyfikator statystyczny...

Podstawiając (patrz: 195) otrzymujemy:

$$\begin{aligned} \ln(g_j(x)) &= \ln\left(p_j (2\pi)^{-\frac{p}{2}} |\Sigma_j|^{-1/2} \exp\left[-\frac{1}{2}(x-m_j)^T \Sigma_j^{-1}(x-m_j)\right]\right) = \\ &= \ln(p_j) + \ln\left((2\pi)^{-\frac{p}{2}}\right) + \ln(|\Sigma_j|^{-1/2}) - \frac{1}{2}(x-m_j)^T \Sigma_j^{-1}(x-m_j) = \\ &= \ln(p_j) - \frac{p}{2} \ln(2\pi) + \frac{1}{2} \ln|\Sigma_j| - \frac{1}{2}(x-m_j)^T \Sigma_j^{-1}(x-m_j) \end{aligned}$$

pomijając stałą $-\frac{p}{2} \ln(2\pi)$ otrzymujemy funkcje klasyfikujące:

$$g_j(x) = \ln(p_j) + \frac{1}{2} \ln|\Sigma_j| - \frac{1}{2}(x-m_j)^T \Sigma_j^{-1}(x-m_j), \quad j=1,2,\dots,M,$$

Optimalny klasyfikator statystyczny...

- W szczególnych przypadkach funkcje klasyfikujące mogą dalej się redukować:

- gdy prawdopodobieństwa a priori wystąpienia obiektu są równe:

$$g_j(x) = \frac{1}{2} \ln|\Sigma_j| - \frac{1}{2}(x-m_j)^T \Sigma_j^{-1}(x-m_j)$$

- w przypadku, gdy macierze kowariancji rozkładów w każdej z klas są identyczne:

$$g_j(x) = \ln(p_j) - \frac{1}{2}(x-m_j)^T \Sigma^{-1}(x-m_j)$$

Optimalny klasyfikator statystyczny...

Przykład

Rozważmy problem klasyfikacji do dwóch klas, wiadomo, że rozkłady cech w klasach są dwuwymiarowymi rozkładami normalnymi z wektorami wartości oczekiwanych dla klas 1, 2 odpowiednio:

$$m_1 = \begin{bmatrix} -3 \\ -3 \end{bmatrix}, \quad m_2 = \begin{bmatrix} 3 \\ 3 \end{bmatrix}$$

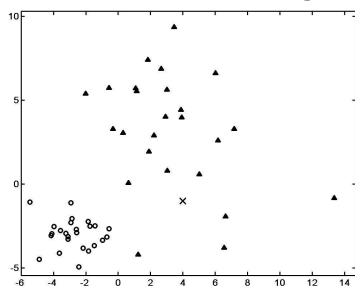
i macierzami kowariancji

$$\Sigma_1 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \Sigma_2 = \begin{bmatrix} 9 & 0 \\ 0 & 9 \end{bmatrix}$$

Prawdopodobieństwo wystąpienia każdej z klas wynosi 0,5.

Przykład...

Po 25 próbek z wymienionych rozkładów ilustruje rysunek. Krzyżykiem zaznaczono przykładową próbkę podlegającą klasyfikacji opisaną wektorem cech równym $[4 \quad -1]^T$.



Przykład...

Rozkłady cech w klasach są dwuwymiarowymi rozkładami normalnymi. Uczenie klasyfikatora polega na wyznaczeniu funkcji klasyfikujących, które dla rozkładu normalnego wyrażają się

wzorem: $g_j(x) = \ln(p_j) - \frac{1}{2} \ln|\Sigma_j| - \frac{1}{2}(x-m_j)^T \Sigma_j^{-1}(x-m_j)$.

Wstawiając odpowiednie wartości za m oraz Σ otrzymujemy:

$$\begin{aligned} g_1(x) &= \ln(0,5) - \frac{1}{2} \ln|I| - \frac{1}{2} \begin{bmatrix} x_1+3 & x_2+3 \end{bmatrix} \cdot \begin{bmatrix} x_1+3 \\ x_2+3 \end{bmatrix} = \\ &= \ln(0,5) - \frac{1}{2} ((x_1+3)^2 + (x_2+3)^2) = \\ &= \ln(0,5) - \frac{1}{2} (x_1^2 + x_2^2 + 6x_1 + 6x_2 + 18) \end{aligned}$$

Przykład...

$$\begin{aligned} g_2(x) &= \ln(0,5) - \frac{1}{2} \ln|81| - \frac{1}{2} \left(\begin{bmatrix} x_1 - 3 & x_2 - 3 \end{bmatrix} \cdot \begin{bmatrix} \frac{1}{9} & 0 \\ 0 & \frac{1}{9} \end{bmatrix} \cdot \begin{bmatrix} x_1 - 3 \\ x_2 - 3 \end{bmatrix} \right) = \\ &= \ln(0,5) - \frac{1}{2} \ln|81| - \frac{1}{2} \left(\begin{bmatrix} \frac{x_1 - 3}{9} & \frac{x_2 - 3}{9} \end{bmatrix} \cdot \begin{bmatrix} x_1 - 3 \\ x_2 - 3 \end{bmatrix} \right) = \\ &= \ln(0,5) - \frac{1}{2} \ln|81| - \frac{1}{2} \left(\frac{(x_1 - 3)^2}{9} + \frac{(x_2 - 3)^2}{9} \right) = \\ &= \ln(0,5) - \frac{1}{2} \ln|81| - \frac{1}{2} \left(\frac{x_1^2}{9} + \frac{x_2^2}{9} - \frac{6}{9}x_1 - \frac{6}{9}x_2 + 2 \right) \end{aligned}$$

A. Brückner

Podstawy sztucznej inteligencji

- 233 -

Przykład...

Uwzględniając fakt, że prawdopodobieństwa a priori klas są równe otrzymujemy następujące funkcje klasyfikujące:

$$g_1(x) = -\frac{1}{2}(x_1^2 + x_2^2 + 6x_1 + 6x_2 + 18),$$

$$g_2(x) = -2\ln|3| - \frac{1}{2} \left(\frac{x_1^2}{9} + \frac{x_2^2}{9} - \frac{6}{9}x_1 - \frac{6}{9}x_2 + 2 \right).$$

Powierzchnia decyzyjna jest postaci $g_1(x) = g_2(x)$, a zatem wstawiając wyliczone funkcje klasyfikujące otrzymujemy:

$$(x_1^2 + x_2^2 + 6x_1 + 6x_2 + 18) = \frac{1}{2} \ln|81| + \left(\frac{x_1^2}{9} + \frac{x_2^2}{9} - \frac{6}{9}x_1 - \frac{6}{9}x_2 + 2 \right)$$

$$8x_1^2 + 8x_2^2 + 60x_1 + 60x_2 + 144 - 9\ln|81| = 0.$$

$$2x_1^2 + 2x_2^2 + 15x_1 + 15x_2 + 36 - 9\ln|3| = 0$$

A. Brückner

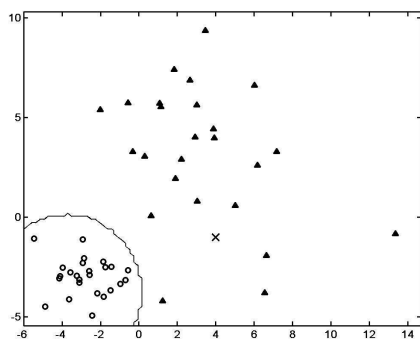
Podstawy sztucznej inteligencji

- 234 -

Przykład...

Powierzchnię decyzyjną $2x_1^2 + 2x_2^2 + 15x_1 + 15x_2 + 36 - 9\ln|3| = 0$

ilustruje rysunek:



A. Brückner

Podstawy sztucznej inteligencji

- 235 -

Przykład...

Zaklassyfikujmy punkt $\begin{bmatrix} 4 & -1 \end{bmatrix}^T$. Mamy

$$g_1(x) = -\frac{1}{2}(x_1^2 + x_2^2 + 6x_1 + 6x_2 + 18),$$

$$g_1\left(\begin{bmatrix} 4 \\ -1 \end{bmatrix}\right) = \frac{16 + 1 + 24 - 6 + 18}{-2} = -\frac{53}{2}$$

oraz $g_2(x) = -2\ln|3| - \frac{1}{2} \left(\frac{x_1^2}{9} + \frac{x_2^2}{9} - \frac{6}{9}x_1 - \frac{6}{9}x_2 + 2 \right).$

$$g_2\left(\begin{bmatrix} 4 \\ -1 \end{bmatrix}\right) = -2\ln|3| - \frac{1}{2} \left(\frac{-1}{9} + 2 \right) \approx -3,142$$

Klasyfikowany obiekt należy więc do klasy 2.

A. Brückner

Podstawy sztucznej inteligencji

- 236 -

Parametryczny klasyfikator Bayesa

- Optymalna klasyfikacja bayesowska wymaga znajomości, rozkładu zmiennej losowej J i warunkowej gęstości rozkładu prawdopodobieństwa cech w klasach.
- W praktycznych zadaniach klasyfikacji tej wiedzy na ogół nie mamy
 - Znana jest postać rozkładu lecz nieznane jego parametry
 - Nie znana jest ani postać rozkładu
- W przypadku, gdy znamy postać warunkowego rozkładu cech w klasach, natomiast nie znamy jego parametrów, brakującą wiedzę rekompensujemy informacjami zawartymi w zbiorze uczącym.

A. Brückner

Podstawy sztucznej inteligencji

- 237 -

Parametryczny klasyfikator Bayesa...

- Zbiór uczący $V = \{(x_1, y_1), \dots, (x_N, y_N)\}$ - zbiór N niezależnych realizacji pary zmiennych losowych (X, J)
- Niech $V_i = \{x_j : y_j = i\}$ oznacza zbiór tych wektorów cech, które są elementami klasy i , to znaczy pochodzą z populacji o warunkowej gęstości $f_j(x)$. Niech $|V_j| = N_j$.
- Zadanie polega na estymacji brakujących parametrów i wstawieniu uzyskanych wartości do funkcji klasyfikujących optymalnego algorytmu Bayesa.
- Algorytm uzyskany w ten sposób nazywa się parametrycznym klasyfikatorem Bayesa.

A. Brückner

Podstawy sztucznej inteligencji

- 238 -

Parametryczny klasyfikator Bayesa...

- Prawdopodobieństwa a priori klas przybliżamy częstościami występowania poszczególnych klas w zbiorze uczącym:

p_j = N_j / N

- Estymacja parametrów gęstości warunkowych cech w klasach zależy od samej postaci tych gęstości i może odbywać się na przykład metodą największej wiarygodności.

Przykład:

Założmy, że rozkład cech w klasach jest dwuwymiarowym rozkładem normalnym o nieznanych parametrach oraz dany jest zbiór uczący:

x ₁	x ₂	kl.	x ₁	x ₂	kl.	x ₁	x ₂	kl.	x ₁	x ₂	kl.
-4,887	-4,489	1	-1,755	-2,504	1	1,241	-4,208	2	-0,56	5,748	2
-2,928	-1,104	1	-2,586	-2,703	1	-0,318	3,301	2	6,636	-1,92	2
-0,707	-3,145	1	-4,112	-2,949	1	6,026	6,620	2	3,015	5,635	2
-2,566	-2,883	1	-2,852	-2,046	1	3,895	4,423	2	1,826	7,418	2
-2,442	-4,940	1	-1,481	-3,658	1	2,210	2,911	2	6,180	2,612	2
-1,861	-2,228	1	-3,633	-4,108	1	2,652	6,866	2	3,049	0,808	2
-3,582	-2,765	1	-1,434	-2,475	1	13,373	-0,834	2	3,464	9,363	2
-3,976	-2,535	1	-5,462	-1,062	1	1,895	1,942	2	1,152	5,554	2
-3,247	-2,940	1	-4,157	-3,055	1	0,636	0,080	2	0,294	3,051	2
-1,844	-3,978	1	-3,100	-3,281	1	1,072	5,711	2	7,186	3,293	2
-2,176	-3,810	1	-0,565	-2,653	1	6,560	-3,787	2	5,028	0,581	2
-2,928	-2,286	1	-0,974	-3,332	1	-2,024	5,388	2	3,926	3,987	2
-3,099	-3,135	1				2,936	4,015	2			

Przykład:

Funkcje klasyfikujące są postaci:

g_j(x) = ln(p_j) - 1/2 * (x - m_j)^T * Sigma^-1 * (x - m_j)

Estymatory wartości średniej i macierzy kowariancji wyliczamy na podstawie ze zbioru uczącego. Estymator największej wiarygodności dla wektora wartości średnich dany jest wzorem:

m_j = 1/|V_j| * sum_{i in V_j} x_i

natomiast dla macierzy kowariancji:

S_j = 1/|V_j| * sum_{i in V_j} (x_i - m_j) * (x_i - m_j)^T

Przykład:

Otrzymujemy odpowiednio:

m_1 = [-2,734; -2,963], m_2 = [3,254; 3,1420]

S_1 = [1,575 -0,134; -0,134 0,827], S_2 = [10,421 -3,763; -3,762 11,804]

Wyznaczniki i macierze odwrotne do macierzy kowariancji:

|S_1| = 1,284, |S_2| = 108,86,

S_1^-1 = [0,644 0,107; 0,107 1,227], S_2^-1 = [0,108 0,035; 0,035 0,096]

Przykład:

Gęstości rozkładu cech w klasach:

f_1([x_1; x_2]) = 1/(2*pi*sqrt(1,284)) * exp[-1/2 * [x_1 + 2,734; x_2 + 2,963]^T * [0,644 0,107; 0,107 1,227] * [x_1 + 2,734; x_2 + 2,963]] = 1/(2*pi*sqrt(1,284)) * exp[-0,322*x_1^2 - 0,614*x_2^2 - 0,107*x_1*x_2 - 2,078*x_1 - 3,928*x_2 - 8,66]

oraz

f_2([x_1; x_2]) = 1/(2*pi*sqrt(108,86)) * exp[-1/2 * [x_1 - 3,254; x_2 - 3,142]^T * [0,108 0,035; 0,035 0,096] * [x_1 - 3,254; x_2 - 3,142]] = 1/(2*pi*sqrt(108,86)) * exp[-0,054*x_1^2 - 0,048*x_2^2 - 0,035*x_1*x_2 + 0,461*x_1 + 0,415*x_2 - 1,403]

Przykład:

Uwzględniając fakt, że prawdopodobieństwa a priori klas są równe oraz opuszczając stałą 1/(2*pi) otrzymujemy funkcje klasyfikujące

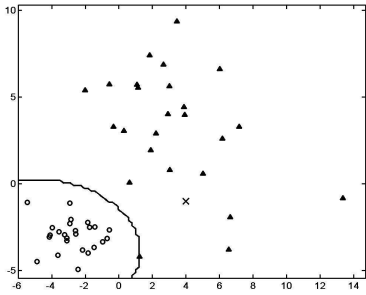
g_1([x_1; x_2]) = ln(1/sqrt(1,284)) - 0,322*x_1^2 - 0,614*x_2^2 - 0,107*x_1*x_2 - 2,078*x_1 - 3,928*x_2 - 8,66

g_2([x_1; x_2]) = ln(1/sqrt(108,86)) - 0,054*x_1^2 - 0,048*x_2^2 - 0,035*x_1*x_2 + 0,461*x_1 + 0,415*x_2 - 1,4

Przykład:

Powierzchnia decyzyjna w przypadku dwóch klas spełnia zależność $g_1(x)=g_2(x)$, a zatem także $\ln g_1(x)=\ln g_2(x)$, skąd wstawiając wyliczone funkcje klasyfikujące otrzymujemy:

$$\ln \sqrt{84,782}-0,268x_1^2-0,566x_2^2-0,072x_1x_2-2,539x_1-4,343x_2-7,257=0$$



Przykład:

- Zaklasyfikujmy wektor: $x=\begin{bmatrix} 4 & -1 \end{bmatrix}^T$. Mamy

$$g_1\left(\begin{bmatrix} 4 \\ -1 \end{bmatrix}\right)=\ln \frac{1}{\sqrt{1,284}}-0,322 \cdot 16-0,614 \cdot 1-0,107 \cdot 4 \cdot(-1)-2,078 \cdot 4+\\-3,928 \cdot(-1)-8,66=\ln \frac{1}{\sqrt{1,284}}-5,152-0,614-0,428-8,312+3,928-8,66=\\=-19,363$$

$$g_2\left(\begin{bmatrix} 4 \\ -1 \end{bmatrix}\right)=\ln \frac{1}{\sqrt{108,86}}-0,054 \cdot 16-0,048 \cdot 1-0,035 \cdot 4 \cdot(-1)+0,461 \cdot 4+\\+0,415 \cdot(-1)-1,403=\ln \frac{1}{\sqrt{108,86}}-0,864-0,048+0,14+1,844-0,415-1,403=\\=-7,059$$

Zgodnie z regułą klasyfikatora klasa 2