

Rachunek prawdopodobo- bieństwa i statystyka

Wykład 6

Warunkowe wartości oczekiwane

Niech będzie dany rozkład p-twa P_{km} dyskretnych zmiennych losowych k i m lub gęstość p-twa $f(x,y)$ ciągłych zmiennych losowych x i y .

Definicja: Wartością oczekiwaną dyskretnej zmiennej losowej k pod warunkiem, że zmienna losowa m przyjmuje wartość m nazywamy wielkość:

$$\mathcal{E}[k | m = m] \equiv \mu_k(m) \equiv \sum_k k P(k = k | m = m) = \frac{1}{P_{\bullet, m}} \sum_k k P_{km}$$

Definicja: Wartością oczekiwaną ciągłej zmiennej losowej x pod warunkiem, że zmienna losowa y przyjmuje wartość y nazywamy wielkość:

$$\mathcal{E}[x | y = y] \equiv \mu_x(y) \equiv \int_{-\infty}^{+\infty} x \cdot f(x | y) dx = \frac{1}{f_2(y)} \int_{-\infty}^{+\infty} x \cdot f(x, y) dx$$

Definicja: Podobnie określamy wartość oczekiwaną dyskretnej zm. losowej m pod warunkiem, że zmienna losowa k przyjmuje wartość k :

$$\mathcal{E}[m | k = k] \equiv \mu_m(k) \equiv \sum_m m P(m = m | k = k) = \frac{1}{P_{k, \bullet}} \sum_m m P_{km}$$

oraz wartość oczekiwaną ciągłej zmiennej losowej y pod warunkiem, że zmienna losowa x przyjmuje wartość x :

$$\mathcal{E}[y | x = x] \equiv \mu_y(x) \equiv \int_{-\infty}^{+\infty} y \cdot f(y | x) dy = \frac{1}{f_1(x)} \int_{-\infty}^{+\infty} y \cdot f(x, y) dy$$

Warunkowe wartości oczekiwane

Twierdzenie: Niech x, y, y_1 i y_2 będą zmiennymi losowymi, $g()$ funkcją oraz c stałą.

Zachodzą następujące związki:

$$\mathcal{E}[c | \mathbf{x} = x] = c$$

$$\mathcal{E}[y_1 + y_2 | \mathbf{x} = x] = \mathcal{E}[y_1 | \mathbf{x} = x] + \mathcal{E}[y_2 | \mathbf{x} = x]$$

$$\mathcal{E}[c y | \mathbf{x} = x] = c \cdot \mathcal{E}[y | \mathbf{x} = x]$$

$$\mathcal{E}[g(\mathbf{x}, y) | \mathbf{x} = x] = \mathcal{E}[g(x, y) | \mathbf{x} = x]$$

$$\mathcal{E}[g(x) \cdot y | \mathbf{x} = x] = g(x) \cdot \mathcal{E}[y | \mathbf{x} = x]$$

$$\mathcal{E}[y | \mathbf{x} = x] = \mathcal{E}[y] \quad \text{jeśli zmienne losowe } x \text{ i } y \text{ są niezależne}$$

Twierdzenie: Niech $\mathcal{E}|y| < \infty$, wówczas $\mathcal{E}[\mathcal{E}[y | \mathbf{x} = x]] = \mathcal{E}[y]$

Dowód:

$$\begin{aligned} \mathcal{E}[\mathcal{E}[y | \mathbf{x}]] &= \int_{-\infty}^{+\infty} \mathcal{E}[y | \mathbf{x} = x] f_x(x) dx = \int_{-\infty}^{+\infty} \left(\int_{-\infty}^{+\infty} y \cdot f(y | x) dy \right) f_x(x) dx = \\ &= \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} y \cdot f(x, y) dy dx = \int_{-\infty}^{+\infty} y \cdot f_y(y) dy = \mathcal{E}[y] \end{aligned}$$

Warunkowa wariancja

Definicja: Warunkową wariancją zmiennej losowej y pod warunkiem, że zmienna losowa x przyjmuje wartość x nazywamy wielkość:

$$\mathcal{V}[y | x = x] = \mathcal{E}[(y - \mathcal{E}[y | x = x])^2 | x = x]$$

Twierdzenie: Niech x i y będą zmiennymi losowymi, a $g(x)$ funkcją, oraz niech $\mathcal{E}[y^2] < \infty$ i $\mathcal{E}[g(x)^2] < \infty$, wówczas

$$\mathcal{E}[(y - g(x))^2] = \mathcal{E}[\mathcal{V}[y | x]] + \mathcal{E}[(\mathcal{E}[y | x] - g(x))^2]$$

Dowód:

$$\begin{aligned} \mathcal{E}[(y - g(x))^2] &= \mathcal{E}[(y - \mathcal{E}[y | x] + \mathcal{E}[y | x] - g(x))^2] = \\ &= \mathcal{E}[(y - \mathcal{E}[y | x])^2] + 2\mathcal{E}[(y - \mathcal{E}[y | x])(\mathcal{E}[y | x] - g(x))] + \mathcal{E}[(\mathcal{E}[y | x] - g(x))^2] = \\ &= \mathcal{E}[\mathcal{E}[(y - \mathcal{E}[y | x])^2 | x]] + 2\mathcal{E}[\mathcal{E}[(y - \mathcal{E}[y | x])(\mathcal{E}[y | x] - g(x)) | x]] + \mathcal{E}[(\mathcal{E}[y | x] - g(x))^2] = \\ &= \mathcal{E}[\mathcal{V}[y | x]] + 2\mathcal{E}[(\mathcal{E}[y | x] - g(x)) \underbrace{\mathcal{E}[(y - \mathcal{E}[y | x]) | x]}_{\parallel 0}] + \mathcal{E}[(\mathcal{E}[y | x] - g(x))^2] = \\ &= \mathcal{E}[\mathcal{V}[y | x]] + \mathcal{E}[(\mathcal{E}[y | x] - g(x))^2] \end{aligned}$$

W szczególności dla $g(x) = \mathcal{E}[y]$ mamy: $\mathcal{V}[y] = \mathcal{E}[\mathcal{V}[y | x]] + \mathcal{V}[\mathcal{E}[y | x]]$

Linie regresji I-go rodzaju

Linia regresji I-go rodzaju zmiennej losowej $k(x)$ względem zmiennej losowej $m(y)$ nazywamy zbiór punktów o współrzędnych (k,m) w przypadku zmiennych dyskretnych i (x,y) w przypadku zmiennych ciągłych, spełniających równania:

$$k = \mu_k(m) \quad x = \mu_x(y)$$

Analogicznie, **linia regresji I-go rodzaju** zmiennej losowej $m(y)$ względem zmiennej $k(x)$ nazywamy, odpowiednio dla zmiennych dyskretnych i ciągłych, zbiory punktów (k,m) lub (x,y) , spełniające równania: $m = \mu_m(k) \quad y = \mu_y(x)$

Własności:

Średnie odchylenie kwadratowe zmiennej losowej x od pewnej funkcji $g(y)$ jest najmniejsze, gdy ta funkcja z p-twem 1 jest równa $\mu_x(y)$:

$$\mathcal{E}[(x - g(y))^2] = \mathcal{E}[\mathcal{V}[x|y]] + \mathcal{E}[(\mu_x(y) - g(y))^2] \geq \mathcal{E}[\mathcal{V}[x|y]] = \mathcal{E}[(x - \mu_x(y))^2]$$

Średnie odchylenie kwadratowe zmiennej losowej y od pewnej funkcji $f(x)$ jest najmniejsze, gdy ta funkcja z p-twem 1 jest równa $\mu_y(x)$:

$$\mathcal{E}[(y - g(x))^2] = \mathcal{E}[\mathcal{V}[y|x]] + \mathcal{E}[(\mu_y(x) - g(x))^2] \geq \mathcal{E}[\mathcal{V}[y|x]] = \mathcal{E}[(y - \mu_y(x))^2]$$

Linie regresji II-go rodzaju

Linia regresji II-go rodzaju nazywamy daną krzywą $y = h(x; a, b, \dots)$ gdy spełnia ona warunek:

$$\mathcal{E} \left[(y - h(x; a, b, \dots))^2 \right] = \min(a, b, \dots)$$

W szczególności prostą regresji II-go rodzaju nazywamy linię prostą $y = ax + b$ spełniającą warunek:

$$\begin{aligned} \mathcal{E} \left[(y - ax - b)^2 \right] &= \mathcal{E} \left[((y - \mu_y) - a(x - \mu_x) + \mu_y - a\mu_x - b)^2 \right] = \\ &= \sigma_y^2 + a^2 \sigma_x^2 - 2a \operatorname{cov}[x, y] + (\mu_y - a\mu_x - b)^2 = \sigma_y^2 + a^2 \sigma_x^2 - 2a \sigma_x \sigma_y \rho + (\mu_y - a\mu_x - b)^2 = \min(a, b) \end{aligned}$$

Szukamy minimum ze względu na parametry a i b :

$$\begin{cases} \frac{\partial}{\partial a} \mathcal{E} \left[(y - ax - b)^2 \right] = 2a\sigma_x^2 - 2\sigma_x\sigma_y\rho - 2\mu_x(\mu_y - a\mu_x - b) = 0 \\ \frac{\partial}{\partial b} \mathcal{E} \left[(y - ax - b)^2 \right] = -2(\mu_y - a\mu_x - b) = 0 \end{cases} \Rightarrow \begin{cases} a = \rho \frac{\sigma_y}{\sigma_x} \\ b = \mu_y - a\mu_x \end{cases}$$

Prosta regresji II-go rodzaju zmiennej y względem zmiennej x ma postać:

$$y = \rho \frac{\sigma_y}{\sigma_x} x - \rho \frac{\sigma_y}{\sigma_x} \mu_x + \mu_y$$

Podobnie znajdujemy prostą regresji II-go rodzaju zmiennej x względem y :

$$x = \rho \frac{\sigma_x}{\sigma_y} y - \rho \frac{\sigma_x}{\sigma_y} \mu_y + \mu_x \quad \Rightarrow \quad y = \frac{1}{\rho} \frac{\sigma_y}{\sigma_x} x - \frac{1}{\rho} \frac{\sigma_y}{\sigma_x} \mu_x + \mu_y$$

Linie regresji I-go rodzaju

Przykład: Linie regresji I-go rodzaju dla dwuwymiarowego rozkładu normalnego.

$$\mathcal{N}(x, y; \mu_x, \mu_y, \sigma_x, \sigma_y) = \frac{1}{2\pi\sigma_x\sigma_y\sqrt{1-\rho^2}} \exp \left\{ -\frac{1}{2(1-\rho^2)} \left[\left(\frac{x-\mu_x}{\sigma_x} \right)^2 + \left(\frac{y-\mu_y}{\sigma_y} \right)^2 - 2\rho \left(\frac{x-\mu_x}{\sigma_x} \right) \left(\frac{y-\mu_y}{\sigma_y} \right) \right] \right\}$$

Rozkłady brzegowe to jednowymiarowe rozkłady normalne:

$$\mathcal{N}_1(x; \mu_x, \sigma_x) = \frac{1}{\sqrt{2\pi}\sigma_x} \exp \left(-\frac{(x-\mu_x)^2}{2\sigma_x^2} \right) \quad \mathcal{E}[x] = \mu_x \quad \mathcal{V}[x] = \sigma_x^2$$

$$\mathcal{N}_2(y; \mu_y, \sigma_y) = \frac{1}{\sqrt{2\pi}\sigma_y} \exp \left(-\frac{(y-\mu_y)^2}{2\sigma_y^2} \right) \quad \mathcal{E}[y] = \mu_y \quad \mathcal{V}[y] = \sigma_y^2$$

Rozkłady warunkowe i warunkowe wartości oczekiwane:

$$g(y|x) = \frac{\mathcal{N}(x, y; \mu_x, \mu_y, \sigma_x, \sigma_y)}{\mathcal{N}_1(x; \mu_x, \sigma_x)} = \frac{1}{\sqrt{2\pi}\sigma_y\sqrt{1-\rho^2}} \exp \left(-\frac{1}{2\sigma_y^2(1-\rho^2)} \left[y - \left(\mu_y + \rho \frac{\sigma_y}{\sigma_x} (x - \mu_x) \right) \right]^2 \right)$$

$$g(x|y) = \frac{\mathcal{N}(x, y; \mu_x, \mu_y, \sigma_x, \sigma_y)}{\mathcal{N}_2(y; \mu_y, \sigma_y)} = \frac{1}{\sqrt{2\pi}\sigma_x\sqrt{1-\rho^2}} \exp \left(-\frac{1}{2\sigma_x^2(1-\rho^2)} \left[x - \left(\mu_x + \rho \frac{\sigma_x}{\sigma_y} (y - \mu_y) \right) \right]^2 \right)$$

Linie regresji I-go rodzaju

Linie regresji I-go rodzaju to proste o równaniach:

$$y = \mathcal{E}[y | x] = \mu_y + \rho \frac{\sigma_y}{\sigma_x} (x - \mu_x)$$

$$x = \mathcal{E}[x | y] = \mu_x + \rho \frac{\sigma_x}{\sigma_y} (y - \mu_y) \quad \Rightarrow \quad y = \mu_y + \frac{1}{\rho} \frac{\sigma_y}{\sigma_x} (x - \mu_x)$$

Uwaga: W przypadku dwuwymiarowego rozkładu normalnego, linie regresji I-go rodzaju są jednocześnie prostymi regresji II-go rodzaju.

Linie regresji I-go rodzaju

Przykład: Chcemy umieć przewidzieć wzrost pewnego gatunku żyta na podstawie koncentracji fosforu w glebie. W tym celu badany cztery koncentracje 2, 4, 8, 16 ppm fosforu, i obserwujemy wzrost roślin przy każdej z nich począwszy od nasiona aż do momentu kwitnienia. Następnie dokonujemy pomiaru suchej masy każdej z roślin.

Wyniki pomiarów przedstawia tabela.

Roślina	Fosfor [ppm]	Masa [g]
1	2	4.1
2	2	3.8
3	2	4.0
4	2	3.9
5	4	5.2
6	4	4.9
7	4	5.0
8	4	4.8
9	8	5.7
10	8	5.9
11	8	6.0
12	8	6.2
13	16	11.7
14	16	8.9
15	16	10.1
16	16	10.3

Zmienne losowe: x - koncentrację fosforu, y - masa roślin.

$$\mu_x \equiv \mathcal{E}[x] = \frac{1}{16}(4 \times 2 + 4 \times 4 + 4 \times 8 + 4 \times 16) = 7.5$$

$$\mu_y \equiv \mathcal{E}[y] = \frac{1}{16}(4.1 + 3.8 + 4.0 + 3.9 + 5.2 + 4.9 + 5.0 + 4.8 + 5.7$$

$$+ 5.9 + 6.0 + 6.2 + 11.7 + 8.9 + 10.1 + 10.3) = \frac{100.5}{16} \cong 6.28$$

$$\mathcal{E}[x^2] = \frac{1360}{16} = 85 \quad \mathcal{E}[y^2] = \frac{727.49}{16} = 45.468 \quad \mathcal{E}[xy] = \frac{957.6}{16} = 59.85$$

$$\text{cov}[x, y] = \mathcal{E}[xy] - \mathcal{E}[x]\mathcal{E}[y] = 59.85 - 7.5 \times 6.28 = 12.75$$

$$\mathcal{D}[x] = \sqrt{\mathcal{E}[x^2] - (\mathcal{E}[x])^2} = \sqrt{85 - 7.5^2} = 5.36$$

$$\mathcal{D}[y] = \sqrt{\mathcal{E}[y^2] - (\mathcal{E}[y])^2} = \sqrt{45.468 - 6.28^2} = 2.46$$

$$\rho = \frac{\text{cov}[x, y]}{\mathcal{D}[x]\mathcal{D}[y]} = 0.967$$

$$y = 0.44x + 2.96$$

Randomizacja – rozkład dyskretny

Przykład: Zmienna losowa X będąca liczbą jaj znoszonych przez kurę ma rozkład Poissona. P-two wyklucia się pisklęcia z każdego jaja wynosi p i jest niezależne od p -twa wyklucia z innych jaj. Jaki jest rozkład zmiennej losowej Y będącej liczbą wyklutych piskląt.

$$P(X = n) = \mathcal{P}_n(\lambda) = \frac{\lambda^n}{n!} e^{-\lambda} \quad \text{gdzie } n = 0, 1, 2, \dots$$

$$P(Y = j | X = n) = \mathcal{B}_j(n, p) = \binom{n}{j} p^j (1-p)^{n-j} \quad \text{gdzie } j = 0, 1, \dots, n$$

$$P(Y = j, X = n) = P(Y = j | X = n) P(X = n)$$

$$\begin{aligned} P(Y = j) &= \sum_{n=j}^{\infty} P(X = n, Y = j) = \sum_{n=j}^{\infty} \frac{\lambda^n}{n!} e^{-\lambda} \frac{n!}{j!(n-j)!} p^j (1-p)^{n-j} = \\ &= \frac{(\lambda p)^j}{j!} e^{-\lambda} \sum_{n=j}^{\infty} \frac{[\lambda(1-p)]^{n-j}}{(n-j)!} = \frac{(\lambda p)^j}{j!} e^{-\lambda} \sum_{k=0}^{\infty} \frac{[\lambda(1-p)]^k}{k!} = \\ &= \frac{(\lambda p)^j}{j!} e^{-\lambda} e^{\lambda(1-p)} = \frac{(\lambda p)^j}{j!} e^{-\lambda p} = \mathcal{P}_j(\lambda p) \end{aligned}$$

Randomizacja – rozkład ciągły

Przykład: Losujemy punkt X z rozkładu jednostajnego na przedziale (a,b) . Następnie losujemy punkt Y z rozkładu jednostajnego na przedziale (X,b) . Znajdź rozkład zmiennej Y .

$$f_1(x) = \begin{cases} \frac{1}{b-a} & \text{dla } a \leq x \leq b \\ 0 & \text{dla pozostałych } x \end{cases} \quad g_{21}(y|x) = \begin{cases} \frac{1}{b-x} & \text{dla } x \leq y \leq b \\ 0 & \text{dla pozostałych } y \end{cases}$$

Znajdujemy łączny rozkład p-twa:

$$f(x, y) = g_{21}(y|x) f_1(x) = \begin{cases} \frac{1}{b-a} \cdot \frac{1}{b-x} & \text{dla } a \leq x \leq y \leq b \\ 0 & \text{dla pozostałych } x, y \end{cases}$$

Oraz interesujący nas rozkład brzegowy:

$$f_2(y) = \int_a^y \frac{1}{b-a} \cdot \frac{1}{b-x} dx = \frac{1}{b-a} \ln \frac{b-a}{b-y}$$

Ogólnie proces randomizacji dla rozkładów odpowiednio dyskretnych i ciągłych przebiega według schematów:

$$P(Y = j) = \sum_n P(Y = j | X = n) P(X = n)$$

$$f(y) = \int_{-\infty}^{\infty} g(y|x) h(x) dx$$

Randomizacja – rozkład mieszany

Przykład: Załóżmy, że zmienna losowa Y przyjmuje $n > 1$ wartości całkowitych z p -twami:

$$P(Y = i) = p_i \quad \text{przy czym} \quad p_1 + p_2 + \dots + p_n = 1$$

Założmy następnie, że dla ustalonej wartości $Y = i$ zmienna losowa X ma rozkład ciągły o gęstości $\varphi_i(x)$.

P-two że zmienna X przyjmuje wartości z przedziału (a, b) natomiast zmienna $Y = i$ wynosi:

$$P(a \leq X \leq b, Y = i) = p_i \int_a^b \varphi_i(x) dx$$

Korzystając z twierdzenia o p -twie całkowitym znajdujemy rozkład brzegowy zmiennej X :

$$\begin{aligned} P(a \leq X \leq b) &= \sum_{i=1}^n P(a \leq X \leq b | Y = i) P(Y = i) = \\ &= \sum_{i=1}^n p_i \int_a^b \varphi_i(x) dx = \int_a^b \sum_{i=1}^n p_i \varphi_i(x) dx \end{aligned}$$

A więc zmienna losowa X jest zmienną ciągłą o funkcji gęstości danej przez:

$$f_X(x) = \sum_{i=1}^n p_i \varphi_i(x) \quad \text{(mieszanina rozkładów)}$$