

1 Podstawowe pojęcia

1.1 Średnia

Niech $X = (X_1, X_2, \dots, X_n)$ będzie krotką n obserwacji (lub próbą statystyczną). Wówczas

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

oznacza *średnią* z X . Zapis z pionową linią na górze jest wygodny również do oznaczania średnich z innych wyrażeń, na przykład

$$\overline{X^2} = \frac{1}{n} \sum_{i=1}^n X_i^2.$$

Łatwo pokazać, że dla równolicznych krotek X, Y i liczb rzeczywistych a, b :

$$\overline{aX + bY} = a\bar{X} + b\bar{Y}.$$

1.2 Wariancja

Wariancję (bez uwzględniania, że X może być próbą – wtedy warto byłoby zastosować „estymator nieobciążony wariancji” lub „wariancję z próby”) obliczamy jako:

$$\text{Var } X = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2.$$

Dodatkowo można zauważyć, że:

$$\begin{aligned} \text{Var } X &= \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 = \overline{(X - \bar{X})^2} = \overline{X^2 - 2\bar{X}X + \bar{X}^2} = \\ &= \overline{X^2} - 2\overline{\bar{X}X} + \overline{\bar{X}^2} = \overline{X^2} - 2\bar{X}^2 + \bar{X}^2 = \overline{X^2} - \bar{X}^2. \end{aligned}$$

Stąd szczególnie łatwo pokazać, że dla dowolnej liczby rzeczywistej a :

$$\text{Var}(aX) = a^2 \text{Var } X$$

oraz

$$\text{Var}(X + a) = \text{Var } X.$$

Z ostatniego spostrzeżenia wynika także, że wariancja ze stałej zawsze wynosi 0. Wariancję można traktować jako miarę zróżnicowania elementów w próbie – czyli miarę odchylenia od sytuacji, w której wszystkie elementy są jednakowe.

1.3 Odchylenie standardowe

Odchylenie standardowe definiujemy jako:

$$\text{sd } X = \sqrt{\text{Var } X}.$$

Podobnie jak w przypadku wariancji, dodanie stałej do wszystkich elementów krotki nie zmienia wyniku, a wynik dla krotki jednakowych liczb wynosi 0.

Odchylenie standardowe jest przydatną, dodatnio liniową miarą rozrzutu obserwacji, gdyż dla dowolnej liczby rzeczywistej a :

$$\text{sd}(aX) = |a| \text{sd } X.$$

2 Kowariancja, korelacja liniowa

2.1 Kowariancja

Jeśli $X = (X_1, X_2, \dots, X_n)$ oraz $Y = (Y_1, Y_2, \dots, Y_n)$, to *kowariancję* definiujemy jako:

$$\text{Cov}(X, Y) = \overline{(X - \bar{X})(Y - \bar{Y})}.$$

Jest to miara „wspólnego” zróżnicowania X i Y od średniej. Jeśli X_i są powyżej średniej dla tych samych $i = 1, \dots, n$, co Y_i , to kowariancja jest dodatnia. Jeśli jest dokładnie przeciwnie (gdy X_i są powyżej średniej, to Y_i są poniżej), otrzymywane są wartości ujemne. Wartości bliskie zeru sugerują brak występowania takiej prostej, liniowej zależności.

Ważne własności (gdzie X, Y, Z – równoliczne krotki obserwacji, a, b – dowolne stałe, rzeczywiste współczynniki):

$$\text{Cov}(X, X) = \text{Var } X,$$

$$\text{Cov}(X, Y) = \text{Cov}(Y, X),$$

$$\text{Cov}(X + a, Y) = \text{Cov}(Y, X),$$

$$\text{Cov}(X, a) = 0,$$

$$\text{Cov}(aX, bY) = ab \text{Cov}(X, Y),$$

$$\text{Cov}(X, Y + Z) = \text{Cov}(X, Y) + \text{Cov}(X, Z).$$

2.2 Współczynnik korelacji liniowej Pearsona

To, ile wynosi kowariancja, zależy od zależności między próbami, ale również od charakterystyki samych prób – pomnożenie wszystkich elementów jednej z prób przez stały współczynnik powoduje również, że pomnożony przez niego zostaje wynik kowariancji. Prostszy w interpretacji wynikiem badania zależności między dwiema próbami jest *współczynnik korelacji liniowej Pearsona*, czyli:

$$r(X, Y) = \frac{\text{Cov}(X, Y)}{\text{sd } X \text{ sd } Y}.$$

Współczynnik ten można obliczać dla krotek X, Y , które mają odchylenia standardowe (lub – równoważnie – wariancje) różne od zera. Potrzeba i wystarcza zatem, aby nie były to stałe.

Dla tak obliczanego współczynnika otrzymujemy, że:

$$r(X, X) = 1,$$

$$r(X, -X) = -1,$$

$$r(X, aX) = 1 \quad \text{gdy } a > 0,$$

$$r(X, aX) = -1 \quad \text{gdy } a < 0,$$

$$-1 \leq r(X, Y) \leq 1,$$

$$r(X, Y) = r(Y, X),$$

$$r(X + a, Y) = r(X, Y).$$

Taki współczynnik łatwo jest interpretować – wartość 1 oznacza niezaburzoną, idealną, rosnącą zależność liniową. Podobnie wartość -1 oznacza idealną, malejącą zależność liniową. Z otrzymanej wartości 0 wynikałoby, że kowariancja wynosi 0, czyli brak jest tego rodzaju zależności.

Nawet w przypadku „idealnych zależności liniowych” nie wskazaliśmy jeszcze odpowiedzi, jakie to są zależności (na przykład – jakim wzorem można je opisać). Problem ten można rozwiązać stosując metodę ogólniejszą, która służy do szukania „najlepszych zależności liniowych” – najprostsza wersja takiej metody zostanie opisana w kolejnej sekcji.

3 Prosta regresja liniowa jednej zmiennej

Niech X i Y będą równolicznymi krotkami obserwacji (po n obserwacji). Przedstawiona metoda będzie służyła wskazywaniu współczynników „najlepszej” funkcji afinicznej, która pozwala obliczać Y na podstawie X . Zatem jeśli oznaczymy:

$$Z = f(X) = aX + b,$$

to chcemy, aby Z było „najbardziej zbliżone” do Y . Wzajemne podobieństwo dwóch sekwencji liczb można definiować w różny sposób, jednak w proponowanej metodzie zostanie zastosowana *metoda najmniejszych kwadratów*, lub minimalizacja średniego błędu kwadratowego. Najprościej to zapisać jako:

$$\overline{(Z - Y)^2} \rightarrow \min,$$

lub

$$\text{cost}(a, b) = \overline{(aX + b - Y)^2} \rightarrow \min.$$

Jeśli X jest stałą, to trudno rozwiązać problem w jakikolwiek sprytny lub przydatny sposób. Można wtedy co najwyżej przyjąć $a = 0$ oraz $b = \bar{Y}$ (łatwo pokazać, że takie rozwiązanie jest wówczas optymalne). Lub dla dowolnego a przyjąć $b = \bar{Y} - a\bar{X}$. Wtedy rozwiązanie nie jest nawet jednoznaczne.

Skupmy zatem uwagę na przypadku, gdy X stałą nie jest. Wtedy dążenie współczynnikiem a lub b do wartości ∞ lub $-\infty$ powoduje, że również koszt jest rozbieżny do nieskończoności. Funkcja $\text{cost}(a, b)$ jest zatem koercywną funkcją o dziedzinie \mathbb{R}^2 . Intuicyjnie lub z wykorzystaniem analizy matematycznej można stwierdzić, że taka funkcja na pewno osiąga gdzieś swoje minimum globalne.

Z twierdzenia Fermata o zerowaniu się pochodnej wynika, że kandydatami na takie minimum będą wyłącznie te punkty, w których pochodne cząstkowe $\text{cost}(a, b)$ po a oraz b zerują się, czyli:

$$\begin{cases} \frac{\partial \text{cost}(a, b)}{\partial a} = 0, \\ \frac{\partial \text{cost}(a, b)}{\partial b} = 0. \end{cases}$$

Zacznijmy od prostszego z przekształceń, które szybko da nam przydatne wnioski:

$$\begin{aligned} \frac{\partial \text{cost}(a, b)}{\partial b} &= \frac{\partial \overline{(aX + b - Y)^2}}{\partial b} = \frac{1}{n} \frac{\partial (\sum_{i=1}^n (aX_i + b - Y_i)^2)}{\partial b} = \\ &= \frac{1}{n} \sum_{i=1}^n \frac{\partial ((aX_i + b - Y_i)^2)}{\partial b} = \frac{1}{n} \sum_{i=1}^n \frac{\partial (b^2 + 2(aX_i - Y_i)b + (aX_i - Y_i)^2)}{\partial b} = \\ &= \frac{1}{n} \sum_{i=1}^n (2b + 2(aX_i - Y_i)) = 2(b + a\bar{X} - \bar{Y}). \end{aligned}$$

Skoro $\frac{\partial \text{cost}(a,b)}{\partial b} = 0$, to:

$$\begin{aligned} 2(b + a\bar{X} - \bar{Y}) &= 0 \\ b &= \bar{Y} - a\bar{X}. \end{aligned}$$

Badając pochodną cząstkową po a otrzymamy:

$$\begin{aligned} \frac{\partial \text{cost}(a,b)}{\partial a} &= \frac{\partial (a\bar{X} + b - \bar{Y})^2}{\partial a} = \frac{1}{n} \frac{\partial (\sum_{i=1}^n (aX_i + b - Y_i)^2)}{\partial a} = \\ &= \frac{1}{n} \sum_{i=1}^n \frac{\partial ((aX_i + b - Y_i)^2)}{\partial a} = \frac{1}{n} \sum_{i=1}^n \frac{\partial (a^2 X_i^2 + 2aX_i(b - Y_i) + (b - Y_i)^2)}{\partial a} = \\ &= \frac{1}{n} \sum_{i=1}^n (2aX_i^2 + 2X_i(b - Y_i)) \end{aligned}$$

Po podstawieniu $b = \bar{Y} - a\bar{X}$:

$$\begin{aligned} \frac{\partial \text{cost}(a,b)}{\partial a} &= \frac{1}{n} \sum_{i=1}^n (2aX_i^2 + 2X_i(b - Y_i)) = \\ &= \frac{1}{n} \sum_{i=1}^n (2aX_i^2 + 2X_i(\bar{Y} - a\bar{X} - Y_i)) = \overline{2aX^2 + 2X(\bar{Y} - a\bar{X} - Y)} = \\ &= \overline{2aX^2 + 2X\bar{Y} - 2aX\bar{X} - 2XY} = 2a(\overline{X^2 - X\bar{X}}) - 2(\overline{X(Y - \bar{Y})}) = \\ &= 2a(\overline{X^2} - \bar{X}^2) - 2(\overline{X(Y - \bar{Y})}) = 2a \text{Var } X - 2 \text{Cov}(X - \bar{X}, Y) = \\ &= 2a \text{Var } X - 2 \text{Cov}(X, Y). \end{aligned}$$

Skoro $\frac{\partial \text{cost}(a,b)}{\partial a} = 0$, to:

$$\begin{aligned} 2a \text{Var } X - 2 \text{Cov}(X, Y) &= 0 \\ a &= \frac{\text{Cov}(X, Y)}{\text{Var } X}. \end{aligned}$$

Współczynnik b możemy natychmiast obliczyć ze wskazanej wcześniej zależności:

$$b = \bar{Y} - a\bar{X}.$$

Warto zauważyć, że jedyne założenie tutaj dotyczyło $\text{Var } X \neq 0$. Gdy jest ono spełnione, to wskazany wynik stanowi jedyne miejsce zerowania się pochodnych cząstkowych. Jest to zatem globalne minimum kwadratowej funkcji kosztu.

4 Ocena przydatności wyniku regresji

Jeśli $Z = f(X)$ jest sekwencją proponowanych wyników (przy oznaczeniach jak w poprzedniej sekcji), to krotkę błędów można oznaczyć jako:

$$Err = Y - Z.$$

Dla optymalnego modelu na pewno $\text{Var } Err \leq \text{Var } Y$ – równość występowałaby dla modelu stałego $Z = 0$. Iloraz

$$FUV = \frac{\text{Var } Err}{\text{Var } E}$$

określa się zatem jako „współczynnik niewyjaśnionej wariancji”. Jest to liczba z przedziału $[0; 1]$ – im mniejsza, tym lepiej.

Jeszcze popularniejszym w literaturze jest „współczynnik determinacji” lub po prostu „R kwadrat”, czyli

$$R^2 = 1 - FUV.$$

W przypadku R^2 interpretacja jest nieco podobna, co w przypadku współczynnika korelacji – 0 oznacza zupełny brak zależności, która pasowałaby do modelu, a 1 oznacza zależność bezbłędną (bez uwzględniania takich szczegółów jak monotoniczność zależności liniowej). Wartość tę często podaje się procentowo, jako „odsetek wyjaśnionej wariancji”. W zależności od dziedziny zastosowań, wyboru modelu i wyboru źródła bibliograficznego, różne wartości mogą być interpretowane w różny sposób. Wartości poniżej 70% sugerują, że zależność występuje, ale model może nie być najlepszym z możliwych dla wybranych danych. Wartości powyżej 85% często są oceniane jako sukces w dopasowywaniu modelu do danych z pewnym szumem.

Satysfakcjonująca wartość R^2 może zostać wskazana jako warunek konieczny uznania regresji za udaną, ale nie powinien to być warunek wystarczający. W tym celu warto zbadać rozkład wartości Err . W przypadku regresji liniowej zawsze $\overline{Err} = 0$, ale to jeszcze nawet nie znaczy, że rozkład jest symetryczny. Jeśli błędy nie przystają do rozkładu normalnego, to trudno je nazwać przypadkowym szumem, a wskazane jest poszukiwanie modelu, który wyjaśniłby dane lepiej. Ocena, czy dane „przystają do rozkładu normalnego” może obejmować narysowanie histogramu i ocenę „na oko”, ale powinna w sobie zawierać również formalne podejście takie jak test Shapiro-Wilka albo test zgodności χ^2 Pearsona.

5 Regresja liniowa wielokolumnowa

Przyjmijmy, że mamy wiele objaśniających kolumn X_1, X_2, \dots, X_m , które dotyczą n opisanych w wierszach obserwacji. Wówczas model liniowy wyjaśniający zmienną Y polega na wskazaniu takiego

$$Z = f(X_1, X_2, \dots) = A_0 + A_1 X_1 + A_2 X_2 + \dots$$

gdzie A to wektor o odpowiedniej długości, dla którego minimalizowany jest średni błąd kwadratowy, czyli

$$\text{cost}(A) = \overline{(Z - Y)^2} \rightarrow \min.$$

Niech X będzie macierzą podsumowującą wszystkie „kolumny z danymi”, które są brane pod uwagę we wzorze na Z . Jeśli chcemy, żeby w tym wzorze występował składnik A_0 (nie jest to wymagane), potrzebna będzie odpowiednia kolumna, z której brane są liczby mnożone przez A_0 – czyli kolumna samych jedynek (oznaczymy ją również jako X_0). Taka macierz X ma wymiary $n \times (m+1)$.

$$X = \begin{bmatrix} 1 & (X_1)_1 & (X_2)_1 & \dots & (X_m)_1 \\ 1 & (X_1)_2 & (X_2)_2 & \dots & (X_m)_2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & (X_1)_n & (X_2)_n & \dots & (X_m)_n \end{bmatrix}.$$

Przyjmując, że A , Z i Y to macierze kolumnowe, w szczególności:

$$A = \begin{bmatrix} A_0 \\ A_1 \\ A_2 \\ \vdots \\ A_m \end{bmatrix},$$

otrzymamy, że

$$Z = XA$$

jest macierzą $n \times 1$, czyli da się porównywać z kolumną Y .

Zauważmy, że:

$$\begin{aligned} \text{cost}(A) &= \overline{(XA - Y)^2} = \frac{1}{n}(XA - Y)^T(XA - Y) = \\ &= \frac{1}{n}(A^T X^T - Y^T)(XA - Y) = \frac{1}{n}(A^T X^T XA - Y^T XA - A^T X^T Y + Y^T Y) = \\ &= \frac{1}{n}(Y^T Y - 2A^T X^T Y + A^T X^T XA). \end{aligned}$$

Obliczając gradient po A (wektor pochodnych cząstkowych po wszystkich elementach) otrzymujemy

$$\begin{aligned} \nabla_A(\text{cost}(A)) &= \\ &= \frac{1}{n}(\nabla_A(Y^T Y) - 2\nabla_A(A^T X^T Y) + \nabla_A(A^T X^T X A)) = \\ &= \frac{2}{n}(-X^T Y + X^T X A). \end{aligned}$$

Gradient ten się wyzeruje, gdy:

$$\begin{aligned} X^T X A - X^T Y &= 0 \\ X^T X A &= X^T Y \\ (X^T X)^{-1}(X^T X)A &= (X^T X)^{-1}X^T Y \\ A &= (X^T X)^{-1}X^T Y. \end{aligned}$$

Pod warunkiem, że macierz $X^T X$ jest odwracalna. Na pewno jest to macierz kwadratowa (taką macierzą nie jest raczej samo X ani równoważnie X^T , więc takiego skracania nie można było wykonać). Macierz $X^T X$ na pewno nie będzie odwracalna, jeśli układ kolumn X jest liniowo zależny. Taka sytuacja wydarzy się między innymi wtedy, gdy pewne dwie kolumny X będą jednakowe (lub proporcjonalne) lub gdy X ma więcej kolumn, niż wierszy. Warto zatem sprawdzić odwracalność macierzy $X^T X$, a nie zakładać, że wynik zawsze uda się wskazać – istnienie jednoznacznego zestawu współczynników jest równoważne tej odwracalności.

Ważna uwaga: w przypadku modeli realizujących kombinacje liniowe, na przykład $f(X_1, X_2) = a + bX_1 + cX_2 + d\sin(X_1 X_2)$, problem można rozwiązać tę samą metodą, poprzez sztuczne stworzenie dodatkowych kolumn. Tutaj należy sztucznie wprowadzić $X_3 = \sin(X_1 X_2)$ i rozwiązać problem dla $g(X_1, X_2, X_3) = a + bX_1 + cX_2 + dX_3$. Możliwe problemy wynikające z braku niezależności tak skonstruowanych kolumn zostaną wykryte na etapie oceny przydatności regresji (współczynnik R^2 , rozkład błędów).