

**Poprzednim razem**

- Metoda największej wiarygodności

Funkcja wiarygodności

$$\tilde{L}(\theta_1, \theta_2, \dots, \theta_p) = \prod_{i=1}^n f(x_i; \theta_1, \theta_2, \dots, \theta_p)$$

$$L(\theta_1, \theta_2, \dots, \theta_p) = \ln \tilde{L}(\theta_1, \theta_2, \dots, \theta_p) = \ln \left( \prod_{i=1}^n f(x_i; \theta_1, \theta_2, \dots, \theta_p) \right)$$

Szukamy maksimum tej funkcji w przestrzeni parametrów

- Metoda najmniejszych kwadratów

$$Q^2 = \sum_{i=1}^n w_i (y_i - f(x_i, \vec{\theta}))^2$$

Regresja I rodzaju E(Y|X)

Regresja II rodzaju

$$Q^2 = (\vec{Y} - A\vec{\theta})^T B(\vec{Y} - A\vec{\theta})$$

Rozwiązanie (minimum Q<sup>2</sup>):  $\vec{\theta} = (A^T B A)^{-1} A^T B \vec{Y}$

**Przykład: regresja liniowa zwyczajna**

- regresja liniowa ( $f_1(x_i)=x_i, f_2(x_i)=1$ ); zakładamy  $\sigma_1=\sigma_2=\dots=\sigma_n\equiv\sigma_y$

$$f_1(x_i) = x_i, \quad f_2(x_i) = 1$$

$$y_i = \theta_1 x_i + \theta_2$$

$$\vec{Y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \quad \vec{\theta} = \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix}$$

$$\vec{\theta} = (A^T B A)^{-1} A^T B \vec{Y}$$

$$A_y = \begin{pmatrix} f_1(x_1) & f_2(x_1) & \dots & f_n(x_1) \\ f_1(x_2) & f_2(x_2) & \dots & f_n(x_2) \\ \vdots & \vdots & \ddots & \vdots \\ f_1(x_n) & f_2(x_n) & \dots & f_n(x_n) \end{pmatrix} \rightarrow A_y = \begin{pmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{pmatrix}$$

$$B_{ij} = \delta_{ij} \frac{1}{(\sigma_y)^2} = \begin{pmatrix} (\sigma_y)^{-2} & 0 & \dots & 0 \\ 0 & (\sigma_y)^{-2} & \dots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \dots & (\sigma_y)^{-2} \end{pmatrix} \rightarrow B_{ij} = \frac{1}{\sigma_y^2} \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & \vdots \\ 0 & 0 & 0 & \ddots \\ 0 & \dots & 0 & 1 \end{pmatrix}$$

RPIS 2023/2024 2

**Przykład: regresja liniowa zwyczajna**

$$\vec{\theta} = (A^T B A)^{-1} A^T B \vec{Y}$$

$$\begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} = \begin{pmatrix} x_1 & x_2 & \dots & 1 \\ 1 & 1 & \dots & 1 \end{pmatrix}^{-1} \begin{pmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{pmatrix} \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

W macierzy do odwracania: macierz jednostkową pomijamy

$$\begin{pmatrix} x_1 & x_2 & \dots & 1 \\ 1 & 1 & \dots & 1 \end{pmatrix} \begin{pmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{pmatrix} = \begin{pmatrix} \sum x_i^2 & \sum x_i \\ -\sum x_i & n \end{pmatrix} \quad \begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \frac{1}{ad-bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

$$(A^T B A)^{-1} = \frac{\sigma_y^2}{\det A^T B A} \begin{pmatrix} n & -\sum x_i \\ -\sum x_i & \sum x_i^2 \end{pmatrix} = \frac{\sigma_y^2}{n \sum x_i^2 - (\sum x_i)^2} \begin{pmatrix} n & -\sum x_i \\ -\sum x_i & \sum x_i^2 \end{pmatrix}$$

RPIS 2023/2024 3

**Przykład: regresja liniowa zwyczajna**

$$\vec{\theta} = (A^T B A)^{-1} A^T B \vec{Y}$$

Zatem  $(A^T B A)^{-1} A^T B = \frac{\sigma_y^2}{n \sum x_i^2 - (\sum x_i)^2} \begin{pmatrix} n & -\sum x_i \\ -\sum x_i & \sum x_i^2 \end{pmatrix} \begin{pmatrix} x_1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ x_n & \dots & 1 \end{pmatrix} =$

$$= \frac{1}{n \sum x_i^2 - (\sum x_i)^2} \begin{pmatrix} nx_1 - \sum x_i & \dots & nx_1 - \sum x_i \\ -x_1 \sum x_i + \sum x_i^2 & \dots & -\sum x_i + \sum x_i^2 \end{pmatrix}$$

i ostatecznie

$$\begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} = \frac{1}{n \sum x_i^2 - (\sum x_i)^2} \begin{pmatrix} nx_1 - \sum x_i & \dots & nx_1 - \sum x_i \\ -x_1 \sum x_i + \sum x_i^2 & \dots & -\sum x_i + \sum x_i^2 \end{pmatrix} \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

$$= \frac{1}{n \sum x_i^2 - (\sum x_i)^2} \begin{pmatrix} nx_1 y_1 - y_1 \sum x_i + nx_1 y_2 - y_2 \sum x_i + \dots & \dots & -y_n \sum x_i \\ -x_1 y_1 \sum x_i + y_1 \sum x_i^2 - x_1 y_2 \sum x_i + y_2 \sum x_i^2 - \dots & \dots & \sum x_i + y_n \sum x_i^2 \end{pmatrix}$$

$$= \frac{1}{n \sum x_i^2 - (\sum x_i)^2} \begin{pmatrix} \sum x_i y_i & \dots & \sum y_i \\ -\sum x_i y_i & \dots & \sum y_i \end{pmatrix}$$

RPIS 2023/2024 4

**Regresja krzywoliniowa – przykład liniowy**

$$\vec{\theta} = (A^T B A)^{-1} A^T B \vec{Y}$$

- Wniosek: liniowy związek Y i  $\theta$ .
- Współczynniki  $\theta$  łatwe do wyliczenia (odwracanie macierzy)
- Przykład: regresja liniowa (m=2,  $f_1(x)=x, f_2(x)=1$ ), dodatkowo zakładamy  $\sigma_1=\sigma_2=\dots=\sigma_n\equiv\sigma_y$

$$Y = \theta_1 f_1(X) + \theta_2 f_2(X) = \theta_1 X + \theta_2$$

$$\rightarrow \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} = \frac{1}{n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2} \begin{pmatrix} \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_i^2 \sum_{i=1}^n y_i - \sum_{i=1}^n x_i \sum_{i=1}^n x_i y_i \end{pmatrix}$$

- Jest to regresja zwyczajna.
- Istnieją też inne odmiany regresji: **klasyczna** (nic nie wiemy o  $\sigma_1$ ), **ważona** (różne  $\sigma_i$ ), **efektywna** (uwzględnić także niepewność  $x_i$ ).

RPIS 2023/2024 5

**Prawo przenoszenia błędów**

- Dla liniowych związków pomiędzy Y i  $\theta$  znajdujemy macierz kowariancji C( $\theta$ ) (ćwiczenia)

$$\vec{\theta} = (A^T B A)^{-1} A^T B \vec{Y} \equiv T \vec{Y}$$

$$\theta_i = T_{ij} Y_j$$

$$E(\theta_i) = T_{ij} E(Y_j)$$

$$C_{ij}(\vec{\theta}) \equiv Cov(\theta_i, \theta_j) = \sum_{k=1}^n \sum_{l=1}^n T_{ik} T_{jl} Cov(Y_k, Y_l)$$

$$C(\vec{\theta}) = T C(\vec{Y}) T^T$$

RPIS 2023/2024 6

## Prawo przenoszenia błędów

- W naszym przykładzie

$$C(\vec{Y}) = B^{-1} \rightarrow$$

$$C(\vec{\theta}) = (A^T B A)^{-1} A^T B B^{-1} [(A^T B A)^{-1} A^T B]^T =$$

$$= (A^T B A)^{-1} A^T [A^{-1} B^{-1} (A^T)^{-1} A^T B]^T =$$

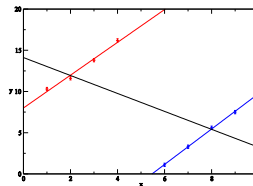
$$= (A^T B A)^{-1} A^T (A^{-1})^T = (A^T B A)^{-1} A^T (A^T)^{-1} = (A^T B A)^{-1}$$

$$C_y(\vec{\theta}) = \begin{pmatrix} \text{var}(\theta_1) & \text{cov}(\theta_1, \theta_2) \\ \text{cov}(\theta_1, \theta_2) & \text{var}(\theta_2) \end{pmatrix} =$$

$$= \sigma_y^2 [A^T A]^{-1} = \frac{\sigma_y^2}{n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2} \begin{pmatrix} n & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{pmatrix}$$

RPIS 2023/2024 7

## Paradoks Simpsona



$$y = 1.99x + 8.0$$

$$y = 2.15x - 11.75$$

$$y = -1.0883x + 14.117$$

- Paradoks Simpsona – połączenie różnych grup może zmienić interpretację wyników. Prowadzi to do praktycznych problemów np. związanych z wyborem kuracji osoby chorej. Rozwiązanie problemu prowadzi do sieci Bayesowskich i teorii grafów i polega na prześledzeniu zależności pomiędzy zmiennymi X, Y i innymi zmiennymi, które mogą mieć wpływ na zmienne X i Y.

RPIS 2023/2024 8

## Miara jakości dopasowania

### – „test Chi-kwadrat”

Zał.  $\sigma_i$  mają rozkład normalny

$$M = \sum_{i=1}^n \frac{(y_i - \sum_{j=1}^p \theta_j f_j(x_i))^2}{\sigma_i^2}$$

- $E(M) = n - p$  – dla rozkładu Chi-kwadrat o  $n - p$  stopniach swobody.
- Gdy  $M > \chi_{\gamma}^{2, n-p}$  ( $n$  - liczba pomiarów,  $p$  liczba parametrów) to na poziomie ufności  $\gamma$  odrzucamy wynik dopasowania.
- Dla regresji liniowej mamy ( $p=2$ )

$$M = \sum_{i=1}^n \frac{(y_i - ax_i - b)^2}{\sigma_i^2}$$

- W praktyce sprawdzamy czy  $M/(n-2) \approx 1$  – prosta dobrze przybliża.
- $M/(n-2) < 1$  prosta, ale za duże  $\sigma_i$ ;  $M/(n-2) > 1$  inna krzywa, za małe  $\sigma_i$  lub niektóre pomiary odstają od prostej.

RPIS 2023/2024 9

## Hipotezy statystyczne

- Testowanie hipotez statystycznych pozwala na sprawdzenie na podstawie wyników próby, przy zadanym poziomie ufności, czy jakieś twierdzenie (hipotezę) dotyczące populacji generalnej jest prawdziwe. Taką procedurę nazywamy **testem statystycznym**.

Przykłady:

- Czy próbka pochodzi z rozkładu normalnego?
- Czy na podstawie próby można powiedzieć, że wartość oczekiwana jest większa niż pewna wybrana liczba?
- Czy wariancje dwóch rozkładów, z których mamy dwie różne próbki, są sobie równe czy nie?

- Jak widać mamy zawsze do czynienia z pewną **hipotezą zerową**  $H_0$ , którą sprawdzamy i **hipotezą alternatywną**  $H_1$ , którą przyjmujemy gdy w wyniku testu odrzucamy  $H_0$ .

- Hipotezy** można podzielić też na **proste** (jednoznacznie określają rozkład, funkcję gęstości prawdopodobieństwa lub dystrybuantę zmiennej losowej) i **złożone** (pozostałe).

- Inny podział: **parametryczna** i **nieparametryczna**.

RPIS 2023/2024 10

## Testy statystyczne

- Schemat postępowania:

- Sformułowanie hipotezy zerowej  $H_0$ .
- Określenie rozmiaru próbki  $n$ .
- Ustalamy **poziom istotności  $\alpha$** , zwany również **poziomem błędzie testu**. Jest to prawdopodobieństwo popełnienia **błędnie I-go rodzaju**, a więc odrzucenia, w wyniku testu, hipotezy  $H_0$ , podczas gdy jest ona w rzeczywistości prawdziwa. Zazwyczaj  $\alpha < 0.1$ .
- Sformułowanie hipotezy alternatywnej  $H_1$ .
- Zakładamy chwilowo, że  $H_0$  jest prawdziwa, i wybieramy statystykę testową  $\delta$ , zależną od próby, o znanym rozkładzie prawdopodobieństwa.
- Wyznaczymy **obszar krytyczny W** dla wartości  $\delta$ , zależny od  $\alpha$ ,  $H_1$  i  $n$ , w taki sposób, że prawdopodobieństwo, że  $\delta \in W$  wynosi  $\alpha$ .
- Wyliczamy, otrzymaną na podstawie pobranej próby, wartość statystyki  $\delta$ . **Gdy należy ona do obszaru krytycznego W to odrzucamy  $H_0$ , w przeciwnym razie twierdzimy, że nie ma podstaw do odrzucenia  $H_0$ .**

RPIS 2023/2024 11

## Testy statystyczne - uwagi

- Określenie obszaru krytycznego na podstawie  $\alpha$  jest niejednoznaczne. Można to poprawić wprowadzając prawdopodobieństwo **błędnie II-go rodzaju**: przyjęciu fałszywej  $H_0$ , gdy w rzeczywistości prawdziwa jest hipoteza alternatywna  $H_1$ . Prawdopodobieństwo to oznaczamy  $\beta$ , a  $1 - \beta$  nazywamy **mocą testu**.
- Jeżeli znamy  $f(\delta|H_0)$  ( $\alpha$ ) i  $f(\delta|H_1)$  ( $\beta$ ) to obszar krytyczny wyznaczamy jako zakres wartości  $\delta$ , który minimalizuje  $\beta$  przy zadanym  $\alpha$ . Zazwyczaj znamy tylko  $f(\delta|H_0)$  i nie znamy mocy testu co oznacza, że nie możemy twierdzić, że  $H_0$  jest prawdziwa, a tylko, że nie odrzucamy  $H_0$ . Test statystyczny, przy nieznajomości  $\beta$ , nazywamy **testem istotności**.
- W praktyce opracowane są zestawy statystyk testowych i zakresy obszaru krytycznego dla różnych rodzajów hipotez  $H_0$  i  $H_1$ . W większości testów zakładamy, że próbka pochodzi z rozkładu normalnego.

RPIS 2023/2024 12