

# 1 Wstęp

## 1.1 Czym są efekty losowe?

### Jednokierunkowa ANOVA

Na poprzednich zajęciach mówiliśmy o modelach liniowych, o jedno- i dwuczynnikowej analizie wariancji. W tych modelach estymowaliśmy nieznanne wartości stałych parametrów, w ANOVIE rolę parametrów pełniły efekty różnych poziomów czynnika. W tym przypadku obserwacje zmiennej odpowiedzi możemy zapisać jako:

$$Y_{ij} = \mu + \alpha_i + \varepsilon_{ij}$$

$$i = 1, 2, \dots, k,$$

$$j = 1, 2, \dots, n$$

$Y_{ij}$  -  $j$ -ta obserwacja na  $i$ -tym poziomie czynnika

$\mu$  - ogólna wartość średnia

$\alpha_i$  - efekt  $i$ -tego poziomu czynnika

$\varepsilon_{ij}$  - niezależne zmienne losowe  $\sim N(0, \sigma^2)$

Obserwacje  $Y_{ij}$ ,  $j = 1, 2, \dots, n$  dla każdej ustalonej wartości wskaźnika  $i$  tworzą **grupę**.

Efekty  $\alpha_i$  są stałymi o nieznanach wartościach, zakładaliśmy przy tym, że  $\sum_{i=1}^k \alpha_i = 0$ . Bez takiego założenia wielkości  $\mu$  i  $\alpha_i$  nie byłyby określone jednoznacznie.

### Dane zgrupowane:

- każda obserwacja należy do jednej grupy i jest tylko jeden czynnik grupujący
- hierarchiczne, zagnieżdżone zbiory danych - bardziej skomplikowane przypadki
- nie możemy zakładać niezależności obserwacji, obserwacje są skorelowane w ramach grup
- efekty losowe - wygodne do modelowania tego typu danych  
Użycie efektów losowych jest popularnym sposobem używanym do modelowania takich danych.

#### 1. **Efekt stały** - nieznaną liczbą, którą estymujemy z danych

Są często używane w modelach liniowych i uogólnionych modelach liniowych

#### 2. **Efekt losowy** - zmienna losowa, estymujemy parametry opisujące rozkład efektu losowego

## 1.2 Przykłady

### 1. Z astronomii:

1861 r. - astronom Airy sformułował model  $Y_{ij} = \mu + \alpha_i + \varepsilon_{ij}$  z losowymi efektami  $\alpha_i$  do opisu obserwowanych przez teleskop obiektów astronomicznych.

Obserwacje  $Y_{ij}$  opisywane są mianowicie jako suma parametru stałego  $\mu$ , określającego średnią wartość obserwacji, losowego efektu  $i$ -tej nocy  $\alpha_i$ , spowodowanego zmiennymi warunkami atmosferycznymi lub innymi czynnikami losowymi charakterystycznymi (i stałymi) dla danej nocy, oraz błędu losowego  $\varepsilon_{ij}$  obserwacji  $Y_{ij}$  spowodowanego czynnikami losowymi różnymi od czynników stałych dla danej nocy, składających się na efekt  $\alpha_i$ .

Jasne, że potraktowanie efektu czynnika nocy jako stałego, a nie losowego, byłoby bardzo grubym błędem.

### 2. Z medycyny:

Zazwyczaj interesuje nas leczenie konkretnymi lekami i traktujemy efekty wywołane przez te leki jako stałe. Sensownie jest jednak uznać efekt pacjentów za losowy. Traktujemy pacjentów poddanych leczeniu, jako wybranych losowo z większego zbioru pacjentów, których cechy chcemy estymować. Zwykle nie interesuje nas wyłącznie ta wąska grupa pacjentów, ale cała populacja pacjentów.

**Kiedy używać efektów losowych?** Czasem w miarę jasne jest, że np. lepiej byłoby użyć efektów losowych, ale w niektórych przypadkach, może być kwestią dyskusji, czy efekty losowe czy stałe będą bardziej odpowiednie. Użycie efektów losowych jest ambitniejsze w takim sensie, że próbujemy powiedzieć coś o szerszej populacji a nie tylko naszej konkretnej próbie.

- gdy efektom danego czynnika nie można przypisać charakteru stałego, należy uznać je za losowe
- gdy chcemy uzyskać informacje o całej populacji, nasze obserwacje traktujemy jako próbkę z całej populacji
- niektórzy statystycy polecają zawsze używanie efektów losowych

## 1.3 Model mieszany

Z efektami, które rozsądnie jest przedstawić jako losowe, spotykamy sięw praktyce nierzadko, zwłaszcza wtedy, gdy mamy do czynienia ze zjawiskami bardziej złożonymi, niż zjawiska dające się opisać za pomocą modelu jednokierunkowej klasyfikacji. Jak wiadomo, modele analizy wariancji możemy przedstawić jako modele liniowej analizy regresji z odpowiednio określoną macierzą planu doświadczenia  $X$ .

- w modelu mieszanym pojawiają się efekty stałe i efekty losowe
- najprostszy przykład - **model dwukierunkowej klasyfikacji**:

$$y_{ijk} = \mu + \tau_i + \nu_j + \varepsilon_{ijk}$$

gdzie

$\mu, \tau_i$  - efekty stałe

efekty losowe, i.i.d.:  $\nu_j \sim N(0, \sigma_\nu^2)$

$\varepsilon_{ijk} \sim N(0, \sigma^2)$

### Przykład - badania biologiczne:

Dla przykładu, w badaniach biologicznych często bada się jakąś cechę potomków zadanej populacji i matek. Powiedzmy, że wartość tej cechy rozsądnie jest uzależnić od płci potomka. Można wówczas zbadać trafność przyjęcia modelu dwukierunkowej klasyfikacji z efektami czynnika płci, jako efektami stałymi oraz efektem j-tej matki,  $j = 1, 2, \dots, l$ , jako efektem Z efektami, które rozsądnie jest przedstawić jako losowe, spotykamy się w praktyce nierzadko, zwłaszcza wtedy, gdy mamy do czynienia ze zjawiskami bardziej złożonymi, niż zjawiska dające się opisać za pomocą modelu jednokierunkowej klasyfikacji.

$Y_{ijk}$  oznacza m-tą obserwację na i-tym poziomie czynnika płci ( $m = 2$ ) i j-tym, losowym poziomie czynnika matka,  $\mu$  jest ogólną wartością średnią,  $\tau_i$  jest (stałym) efektem t-tego poziomu czynnika płci,  $\nu_j$  jest efektem związanym z j-tą matką,  $n_{ij}$  jest liczbą potomków i-tej płci j-tej matki oraz  $\varepsilon_{ijk}$  są niezależnymi zmiennymi losowymi o rozkładzie normalnym  $N(0, \sigma_\varepsilon^2)$ . O efektach stałych zakładamy jak w ANOVIE, że sumują się do zera. Dodatkowo zakładamy zwykle, że efekty  $\nu_j$  są nieskorelowane między sobą i z błędami  $\varepsilon_{ijk}$ , posiadają zerową wartość oczekiwaną i wspólną wariancję  $\sigma_\nu^2 \geq 0$ , niezależną od wskaźnika j.

### Estymacja efektów, przykład

**Hipoteza zerowa** przy estymacji efektów:

- stałych - w modelu mieszanym chcemy estymować  $\tau_i$  i testować hipotezę:  
 $H_0 : \tau_i = 0 \forall i$
- losowych -  $H_0 : \sigma_\nu^2 = 0$

Różnica - w pierwszym przypadku musimy przetestować wiele parametrów efektu stałego, podczas gdy w drugim przypadku estymujemy i testujemy jedynie jeden parametr efektu losowego.

Teraz powiemy więcej o estymacji i testowaniu modeli losowych i mieszanych.

## 2 Estymacja

Estymacja efektów w modelu losowym jest możliwa przy użyciu:

- klasyfikacji jednokierunkowej (ANOVA)
- metody największej wiarygodności (ML)
- metody największej wiarygodności z restrykcjami (REML)

## 2.1 Jednokierunkowa klasyfikacja (ANOVA)

Najprostszy model losowy - **model jednokierunkowej klasyfikacji (ANOVA)**:

$$Y_{ij} = \mu + \alpha_i + \varepsilon_{ij}$$

$$i = 1, \dots, a$$

$$j = 1, \dots, n$$

$$\alpha_i \sim N(0, \sigma_\alpha^2)$$

$$\varepsilon_{ij} \sim N(0, \sigma_\varepsilon^2)$$

- **Komponenty wariancyjne:**  $\sigma_\alpha^2, \sigma_\varepsilon^2$
- **Wewnątrzgrupowy współczynnik korelacji:**  $\rho = \frac{\sigma_\alpha^2}{\sigma_\alpha^2 + \sigma_\varepsilon^2}$

**Estymatory**

- **Dekompozycja wariancji:**

$$\sum_{i=1}^a \sum_{j=1}^n (y_{ij} - \bar{y}_{..})^2 = \sum_{i=1}^a \sum_{j=1}^n (y_{ij} - \bar{y}_{i.})^2 + \sum_{i=1}^a \sum_{j=1}^n (\bar{y}_{i.} - \bar{y}_{..})^2$$

$$SST = SSE + SSA$$

- $E(SSE) = a(n-1)\sigma_\varepsilon^2$   
 $E(SSA) = (a-1)(n\sigma_\alpha^2 + \sigma_\varepsilon^2)$
- **Estymatory:**  
 $\hat{\sigma}_\varepsilon^2 = SSE / (a(n-1)) = MSE$   
 $\hat{\sigma}_\alpha^2 = \frac{SSA / (a-1) - \hat{\sigma}_\varepsilon^2}{n} = \frac{MSA - MSE}{n}$

**Wady korzystania z ANOVY:**

- estymator wariancji może przyjmować ujemne wartości
- skomplikowane obliczenia w przypadku gdy dane są niezbalansowane

## 2.2 Metoda największej wiarygodności (ML)

Jak widzimy ANOVA nie jest doskonała. Szukamy metody, która będzie działała dobrze dla danych niezablansowanych, którą będzie można łatwo zastosować i nie będzie problemów z ujemną wariancją. Te warunki spełnia ML.

- nie ma takich wad jak ANOVA
- musimy założyć, jaki rozkład mają błędy i efekty losowe (zwykle rozkład normalny)  
 ML zadziała też dla innych rozkładów, ale zwykle nie rozważa się innych rozkładów niż rozkład normalny.
- dla modelu wyłącznie z efektami stałymi mamy model:

$$y = X\beta + \varepsilon, y \sim N(X\beta, \sigma^2 I)$$

gdzie

$X$  - macierz planu doświadczenia  $n \times p$

$\beta$  - wektor  $p$  stałych parametrów

$\varepsilon$  - wektor błędów losowych

**Model mieszany** Ww. model możemy uogólnić na model mieszany. Mając dane wartości efektów losowych możemy modelować odpowiedź  $y$ .

- dla modelu mieszanego:

$$y = X\beta + Z\gamma + \varepsilon, y|\gamma \sim N(X\beta + Z\gamma, \sigma^2 I)$$

gdzie

$Z$  - macierz  $n \times q$  opisująca wpływ efektów losowych na obserwacje  $y$

$\gamma$  - wektor  $q$  efektów losowych

$\varepsilon$  - wektor błędów losowych

Przyjmuje się, iż wartości oczekiwane obydwu wektorów losowych są wektorami zerowymi. O macierzy kowariancji wektora  $\varepsilon$  zakłada się zwykle, że jest macierzą diagonalną o postaci  $\sigma^2 I$ , gdzie  $I$  jest macierzą jednostkową.

- dalej zakładamy, że  $\gamma \sim N(0, \sigma^2 D)$ , Na macierz kowariancji wektora efektów losowych  $\gamma$  nie nakłada się w ogólności żadnych ograniczeń (poza oczywistym wymaganiem jej symetryczności i nieujemnej określoności); macierz tę zapisujemy się zwykle jako równą  $\sigma^2 D$  wówczas możemy zapisać rozkład  $y$  (bezwarunkowy)

$$\begin{aligned} \text{var} y &= \text{var} Z\gamma + \text{var} \varepsilon = \sigma^2 Z D Z^T + \sigma^2 I \\ y &\sim N(X\beta, \sigma^2 (I + Z D Z^T)) \end{aligned}$$

- przyjmuje się, że wektory losowe  $\gamma$  i  $\varepsilon$  są niezależne

Przy takich założeniach macierz kowariancji wektora obserwacji  $Y$  ma postać  $\sigma^2 Z D Z^T + \sigma^2 I$ . Najczęściej przyjmuje się ponadto, że zarówno wektor  $\varepsilon$  jak i wektor  $\gamma$  mają rozkłady normalne (o podanych parametrach)

### Estymacja przy użyciu metody największej wiarygodności (ML)

Gdybyśmy znali  $D$ , wówczas moglibyśmy estymować  $\beta$  przy użyciu metody najmniejszych kwadratów. Jednak estymacja komponentów wariacyjnej  $D$  jest jednym z celów naszej analizy. W tej sytuacji jedną z metod estymacji, którą możemy tutaj użyć jest metoda największej wiarygodności.

- jeśli przyjmiemy, że  $V = I + Z D Z^T$ , wówczas łączna gęstość  $y$  wynosi:

$$\frac{1}{(2\pi)^n |\sigma^2 V|^{1/2}} e^{-\frac{1}{2\sigma^2} (y - X\beta)^T V^{-1} (y - X\beta)}$$

- logarytm wiarygodności:

$$l(\beta, \sigma, D|y) = -\frac{n}{2}\log 2\pi - \frac{1}{2}\log|\sigma^2 V| - \frac{1}{2\sigma^2}(y - X\beta)^T V^{-1}(y - X\beta)$$

- estymacji podlega nieznaną wektor stałych parametrów modelu  $\beta$ , wariancja  $\sigma^2$  i macierz komponentów wariancyjnych  $D$

Wszystko wydaje się proste, jednak w praktyce może to stwarzać trudności.

### Wady korzystania z metody największej wiarygodności (ML):

- trudności w przypadku bardziej złożonych modeli, wymagających estymacji większej liczby parametrów efektów losowych pojawiają się trudności.
- problem z estymatorem wariancji, gdy maksimum wiarygodności jest przyjmowane dla ujemnych wartości. Często trzeba w takich sytuacjach przyjąć po prostu, że estymator wariancji wynosi zero (gdy pochodna wiarygodności nie zeruje się np. dla wartości nieujemnych). To komplikuje obliczenia numeryczne.
- estymatory są często dość silnie obciążone (?)

### 2.3 Metoda największej wiarygodności z restrykcjami (REML)

- nie ma wielu wad, które dotyczyły ANOVY i ML próbują rozwiązać problem obciążenia
- dla danych zbalansowanych, estymatory obliczone przy użyciu REML są zwykle takie same jak obliczone przy użyciu ANOVY
- dostajemy estymatory nieobciążone lub o zredukowanym obciążeniu
- idea: bierzemy liniową kombinację  $X$  taką że  $K^T X = 0$ .  
Wówczas:  $K^T y \sim N(0, \sigma^2 K^T V K)$ ,  
następnie możemy korzystać z metody ML (nie mamy już efektów stałych)  
Kiedy już zastymujemy efekty losowe, możemy łatwo estymować parametry efektów stałych

### 2.4 R

- **REML**

Teraz zaprezentujemy estymatory uzyskiwane metodą największej wiarygodności. Używamy pakietu lme4 (dopasowuje modele mieszane; wcześniej nlme).

Model jak widzimy ma efekty stałe i losowe. Efekty losowe tutaj to przecięcie prezentowane przez pierwszą jedynekę w formule modelu. Efekty losowe są reprezentowane przez (1 operator) - zaznaczone jest, że dane są zgrupowane po operatorze, 1 oznacza, że efekt losowy jest stały w każdej grupie.

REML jest domyślną metodą dopasowania (tutaj jest użyty). Widzimy, że daje identyczne estymatory jak ANOVA,  $\hat{\sigma}^2 = 0.106$ ,  $\hat{\sigma}_\alpha^2 = 0.068$  i  $\mu^* = 60.4$ . Dla danych niezbalansowanych estymatory REMLA i ANOVY niekoniecznie są identyczn.

- **ML**

Możemy też skorzystać z tradycyjnej metody największej wiarygodności. Jak widzimy wariancja między grupami wynosi 0.0482 i jest mniejsza niż w REMLu, natomiast wariancja wewnątrz grupy wynosi 0.1118 i jest większa niż poprzednio.

### 3 Testowanie

#### 3.1 Test ilorazu wiarygodności

Używając standardowej teorii wiarygodności, możemy wyprowadzić test, który służy do porównywania dwóch zagnieżdżonych hipotez  $H_0$  i  $H_1$ , obliczając test ilorazu wiarygodności.

Do porównywania hipotez  $H_0$  i  $H_1$  używamy testu **ilorazu wiarygodności**:

$$2(l(\hat{\beta}_1, \hat{\sigma}_1, \hat{D}_1|y) - l(\hat{\beta}_0, \hat{\sigma}_0, \hat{D}_0|y))$$

gdzie

$\hat{\beta}_1, \hat{\sigma}_1, \hat{D}_1$  - estymatory największej wiarygodności (MLE) parametrów przy hipotezie zerowej

$\hat{\beta}_0, \hat{\sigma}_0, \hat{D}_0$  - MLE przy hipotezie alternatywnej

Rozkład zerowy tej statystyki testowej to w przybliżeniu rozkład  $\chi^2$  z liczbą stopni swobody równą różnicy stopni swobody w modelu zerowym i pierwszym (w wymiarach dwóch przestrzeni parametrów (gdy modele są identyfikowalne)).

#### Wady testu ilorazu wiarygodności:

- test jest przybliżony i często przybliżenie to jest nienajlepsze
- wymaga wielu dodatkowych założeń, np. parametry przy hipotezie zerowej nie mogą być na brzegu przestrzeni parametrów  
Jako, że w naszym przypadku  $H_0 : \hat{\sigma}^2 = 0$ , to stanowi to prawdziwy problem. Nawet jeśli te warunki są spełnione, to przybliżenie rozkładu  $\chi^2$  jest bardzo kiepskie.

#### 3.2 Testowanie efektów stałych

- jeśli chcemy użyć testu ilorazu wiarygodności do porównania dwóch zagnieżdżonych modeli, które różnią się tylko efektami stałymi, to nie mo-

zemy użyć metody REML do ich porównania. Przyczyną jest to, że REML estymuje efekty losowe przez rozważanie liniowych kombinacji danych, które likwidują efekty stałe. Jeśli te efekty stałe zostaną zmienione, wówczas wiarygodności dwóch modeli nie będą bezpośrednio porównywalne.

- Jeśli chcemy korzystać z testu ilorazu wiarygodności do testowania efektów stałych używamy ML.
- p-wartości generowane przez test ilorazu wiarygodności dla efektów stałych są przybliżone i często zbyt małe, przez co wyolbrzymiają wagę niektórych efektów.
- do znalezienia dokładniejszej p-wartości testu iloraz wiarygodności możemy użyć metody **parametrycznego bootstrapu**. Generujemy dane (przy założeniu modelu zerowego używając dopasowanych estymatorów przy hipotezie zerowej) a następnie obliczamy test ilorazu wiarygodności dla tak wygenerowanych danych. Powtarzamy to wielokrotnie i używamy otrzymanych wyników do oceny istotności obserwowanej statystyki testowej. Metoda ta zostanie zaprezentowana później.
- alternatywnie: możemy użyć standardowych testów  $F$  lub  $t$  (parametry efektu losowego warunkujemy po estymowanych wartościach)  
Tutaj zakładamy, że kowariancja losowej części modelu,  $D$ , jest równa estymowanej wartości i postępujemy tak jak w metodzie najmniejszych kwadratów.

### 3.3 Testowanie efektów losowych

- hipoteza zerowa zwykle ma postać  $H_0 : \sigma^2 = 0$   
problem - to nie jest wewnątrz przestrzeni parametrów, a standardowe obliczenie asymptotycznego rozkładu  $\chi^2$  dla ilorazu wiarygodności opiera się o hipotezę zerową leżącą we wnętrzu przestrzeni parametrów. To założenie jest złamane, kiedy testujemy, czy wariancja jest równa zero.
- rozkład zerowy jest wtedy w ogólności nieznan, musimy zastosować metody numeryczne (bootstrap), jeśli chcemy dokładniej testować.
- jeśli jednak przyjmiemy, że otrzymujemy rozkład  $\chi^2$  ze standardową liczbą stopni swobody, wówczas otrzymana  $p$ -wartość będzie zwykle większa niż powinna. To oznacza, że jeśli obserwujemy istotny efekt używając przybliżenia  $\chi^2$ , to możemy być na 100 % pewni, że ten efekt jest rzeczywiście istotny.  
Jeśli, otrzymujemy wątpliwe wyniki, powinniśmy skorzystać z dokładniejszej, ale czasochłonnej metody bootstrapu.