

JEDNOCZYNNIKOWA ANALIZA WARIANCJI, ANOVA 1

Obserwowana (badana) cecha — Y

Czynnik wpływający na Y (badany) — A

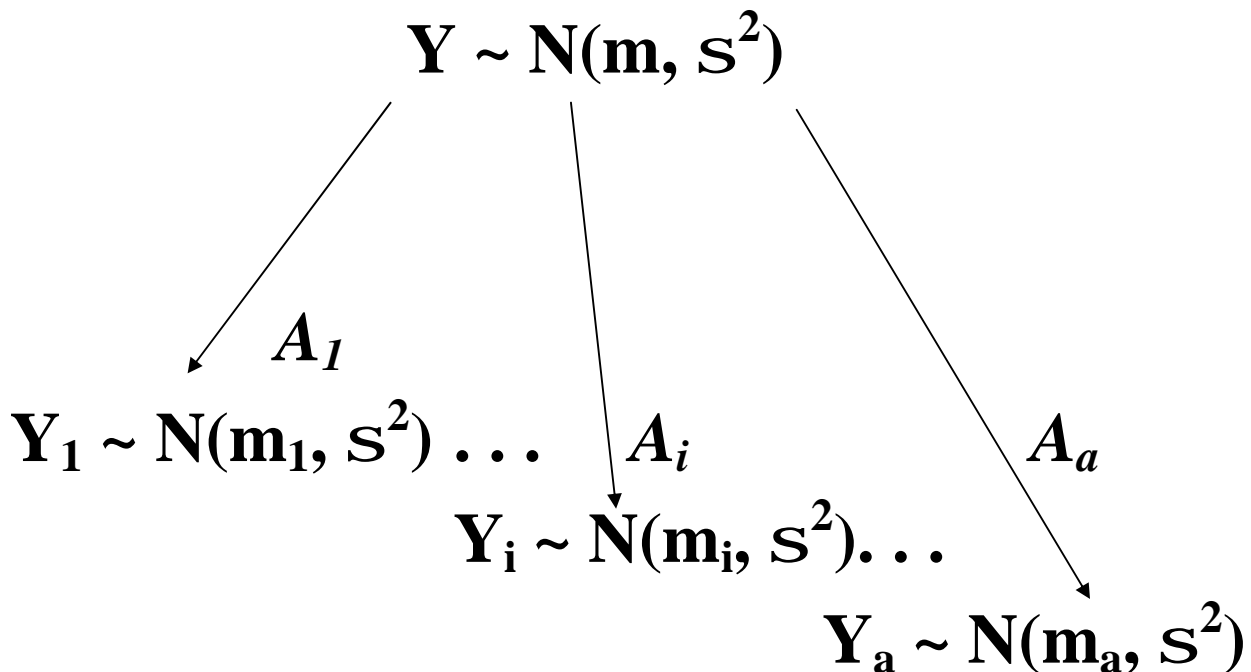
A_i — i -ty poziom czynnika A

a — liczba poziomów ($j=1..a$),

n_i — **liczba powtórzeń** w i -tej populacji

Cecha $Y \sim N(m_i, \sigma^2)$ w **i -tej populacji**
(tej części populacji, która ma czynnik A na
 i -tym poziomie)

Obserwacje są niezależne w ramach powtórzeń
dla poziomu oraz między poziomami



Model liniowy

$$y_{ij} = m + a_i + e_{ij}$$

$$m_i = m + a_i$$

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

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

gdzie

y_{ij} – obserwacja przeprowadzona dla i -tego poziomu czynnika A (A_i), w j – tym powtórzeniu,

m – średnia ogólna,

a_i – efekt i -tego poziomu czynnika A ,

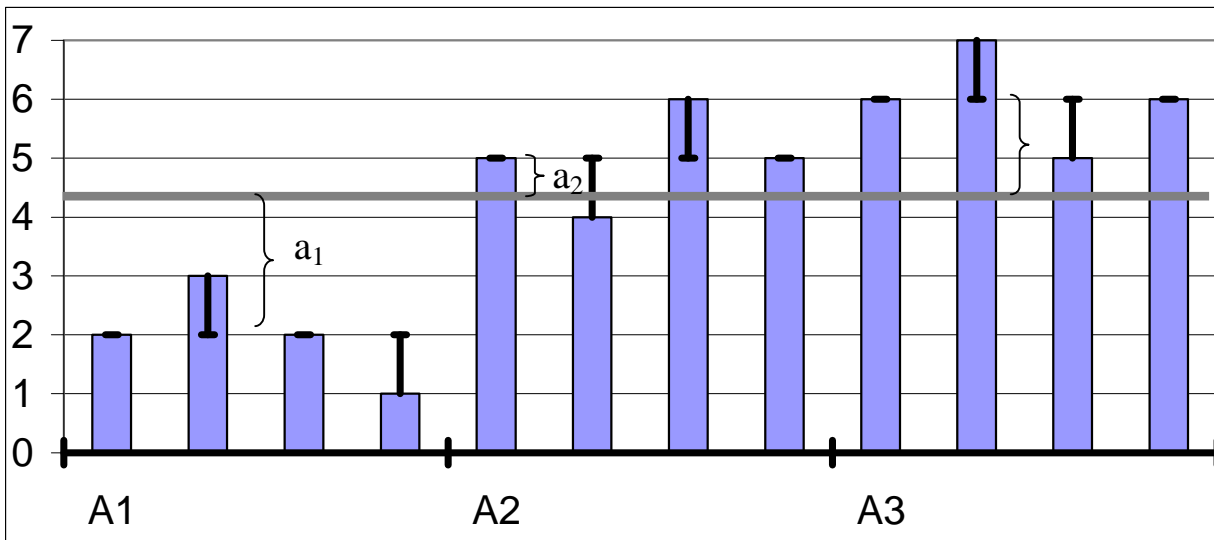
e_{ij} – efekt losowy (błąd losowy)

$\sum_{i=1}^a a_i = 0$ zakładamy dla jednoznaczności

wyniku

Rozkład efektów losowych powinien być rozkładem $N(0, \sigma_{err})$ (rozkłady niezależne dla różnych poziomów czynnika).

rysunek poglądowy dla 3 poziomów czynnika



Średnie populacji (podpopulacji) oceniamy jako 2, 5 i 6 ($\hat{m}_1=2$, $\hat{m}_2=5$, $\hat{m}_3=6$)

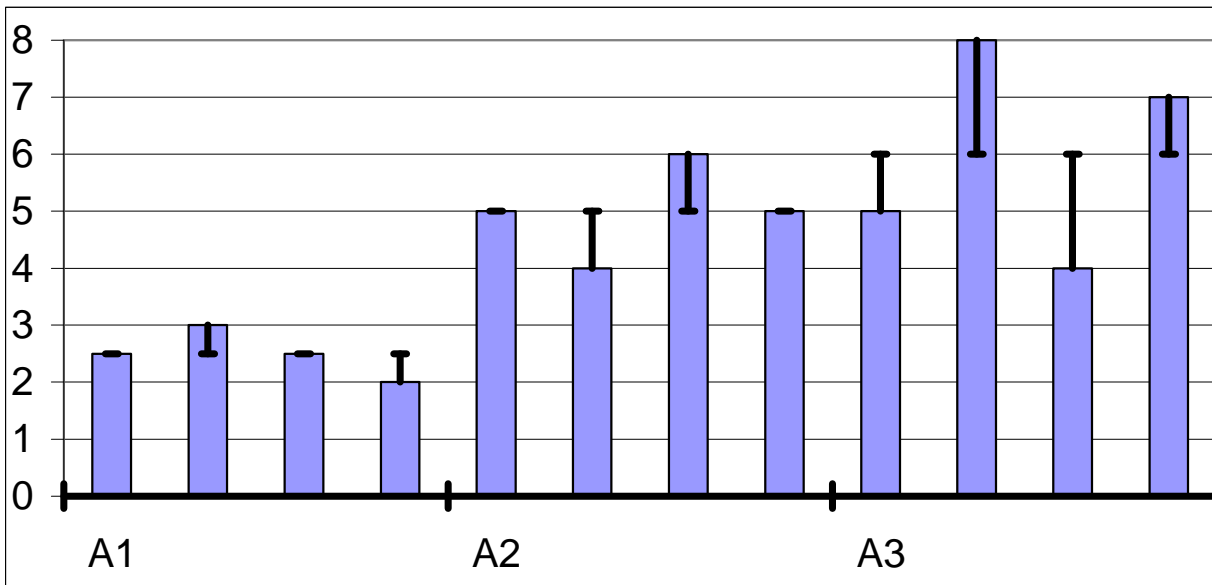
Średnią ogólną oceniamy jako 4,33

Efekty główne czynnika A oceniamy jako: -2,33, 0,66 i 1,66 (estymatory efektów głównych czynnika w całej populacji)

Efekty losowe oznaczono czarną pionową linią np. $e_{11}=0$, $e_{12}=-1$, ..., $e_{22}=1$,

Można przyjąć, że wariancja błędu (σ_{err}) jest równa dla wszystkich podpopulacji. Jej estymator wynosi ok. 0,54.

rysunek niespełnienia założeń o równości wariancji



W tym przypadku wariancje błędów są większe dla trzeciego poziomu czynnika A

Hipoteza: wpływ czynnika A na cechę Y jest nieistotny (czynnik A nie wpływa na wartości cechy Y)

$$\mathbf{H_0 : m_1 = m_2 = \dots = m_a = m}$$

$$(\mathbf{H_0 : \sum_i a_i = 0, \text{ bo } m_i = m + a_i})$$

Hipoteza alternatywna:

$$\mathbf{H_1 : \exists i, m_i \neq m}$$

$$(\mathbf{H_1 : \exists i, a_i \neq 0})$$

Funkcja testowa dla badanej hipotezy:

$$F_{emp} = \frac{S_A^2}{S_{err}^2}$$

Funkcja testowa opiera się na hipotezie o równości wariancji, tzn: zmienność (wariancja) spowodowana czynnikiem A jest taka sama jak zmienność losowa.

Tabela analizy wariancji

Źródło zmienności <i>source</i>	Sumy kwadratów <i>SS</i>	stopnie swobody <i>df</i>	średni kwadrat <i>MS</i>	F_{emp}
Czynnik-A	SS_A	$V_A = a - 1$	$s_A^2 = \frac{SS_A}{a - 1}$	$\frac{s_A^2}{s_e^2}$
Błąd losowy-E	SS_e	$V_e = N - a$	$s_e^2 = \frac{SS_e}{N - a}$	
Zmienność całkowita-T	SS_y	$V_T = N - 1$		

Sumy kwadratów odchyłeń od średnich grupowych zachowują się addytywnie, tzn.

$$SS_y = SS_A + SS_e$$

Jako poziom istotności α wybiera się najczęściej wartości: 0,05 i 0,01 (oznaczając * lub ** odpowiednio).

Symbole występujące w tabeli analizy wariancji liczymy według formuł:

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

$$SS_a = \sum_{i=1}^a n_i (\bar{y}_i - \bar{y})^2 = \frac{1}{n} \sum_{i=1}^a Y_i^2 - \bar{y} Y_{..}$$

$$SS_e = SS_y - SS_A$$

$$\bar{y} = \frac{1}{N} \sum_{i,j} y_{ij}, \quad \bar{y}_i = \frac{1}{n} \sum_{j=1}^n y_{ij}, \quad Y_{..} = \sum_{i,j} y_{ij}$$

$$N = n \cdot a, \quad \text{lub} \quad N = n_1 + n_2 + \dots + n_a$$

Jeżeli zajdzie nierówność $F_{\text{emp}} \geq F_{\alpha, a-1, N-a}$ to na poziomie istotności α hipotezę H_0 należy odrzucić na korzyść hipotezy alternatywnej H_1 . Gdy zajdzie nierówność przeciwna (tzn. $F_{\text{emp}} < F_{\alpha, a-1, N-a}$) to nie ma podstaw do odrzucenia hipotezy H_0 .

P-Value

Ponieważ wnioskowanie często odbywa się na podstawie wartości statystyki określającej **prawdopodobieństwo** (ryzyko) **popęlnienia błędu I rodzaju** przy aktualnie obserwowanych danych (*Pvalue*) to do tabeli analizy wariancji dodaje się kolumnę z wyliczoną wartością tej statystyki.

W przypadku odrzucenia hipotezy zerowej testujemy **szczegółowo** **średnie** podpopulacji, czyli przeprowadzamy podział **średnich** na grupy **jednorodne**.

Szczegółowe porównanie średnich

Grupy jednorodne – podzbiory **średnich**, które można uznać za takie same.

Procedury porównań wielokrotnych – postępowanie statystyczne zmierzające do podzielenia zbioru **średnich** na grupy **jednorodne**.

Ogólna idea procedur porównań wielokrotnych

Są one rozwinięciem metody weryfikacji hipotezy: $H_0: m_1 - m_2 = 0$ przy hipotezie alternatywnej $H_1: m_1 - m_2 \neq 0$, opartej na teście t- Studenta.

NIR – Najmniejsza Istotna Różnica (ang. Least Significant Difference – LSD)

Jeżeli $|\bar{Y}_i - \bar{Y}_{i'}| < NIR$, to uznajemy, że $m_i = m_{i'}$, czyli nie różnią się istotnie.

Jeżeli $|\bar{Y}_i - \bar{Y}_j| < NIR$, $|\bar{Y}_i - \bar{Y}_l| < NIR$,
 $|\bar{Y}_j - \bar{Y}_l| < NIR$, to uznajemy, że $m_i = m_j = m_l$,
czyli średnie te stanowią grupę jednorodną.

Badając w ten sposób wszystkie pary średnich uzyskamy podział zbioru średnich na grupy jednorodne.

Uporządkowanie (malejące lub rosnące) badanych średnich znacznie przyspiesza podział na grupy jednorodne.

Najczęściej stosowane procedury porównań wielokrotnych to procedury: Studenta, **Tukeya** (uwzględnia wielokrotność porównań, gwarantuje jednakowy poziom istotności dla wszystkich porównywanych par), **Dunnetta** (porównanie innych średnich z normą), **Scheffego** (zapewnia łączny poziom istotności przy porównaniu wszystkich par), Bonferroniego (uwzględnia wielokrotność porównań – podzielenie poziomu istotności przez ich liczbę), Duncana i Newmana – Kuelsa.

Procedura t – Studenta

Porównujemy dwie średnie populacyjne (obiektowe): m_i i $m_{i'}$, w oparciu o średnie próbkowe \bar{y}_i i $\bar{y}_{i'}$, uzyskane dla n_i i $n_{i'}$ obserwacji.

Kryterium NIR, oparte na teście t-Studenta, liczymy według wzoru:

$$NIR = t_{a, v_e} \cdot s_r = t_{a, v_e} \sqrt{S_e^2 \left(\frac{1}{n_i} + \frac{1}{n_{i'}} \right)}$$

Procedura Tukey'a

Założenie: $n_1 = \dots = n_a = n$

$$NIR_T = q_{a, a, v_e} \cdot \sqrt{\frac{S_e^2}{n}}$$

Wartość q_{a, a, v_e} pochodzi z tablic wartości krytycznych studentyzowanego rozstępu.

Porównanie z normą test Dunnetta

Zastosowanie: gdy wśród poziomów czynnika A jest wyróżniony poziom kontrolny. Pozostałe poziomy czynnika A porównywane są tylko do kontroli (a nie do siebie nawzajem).

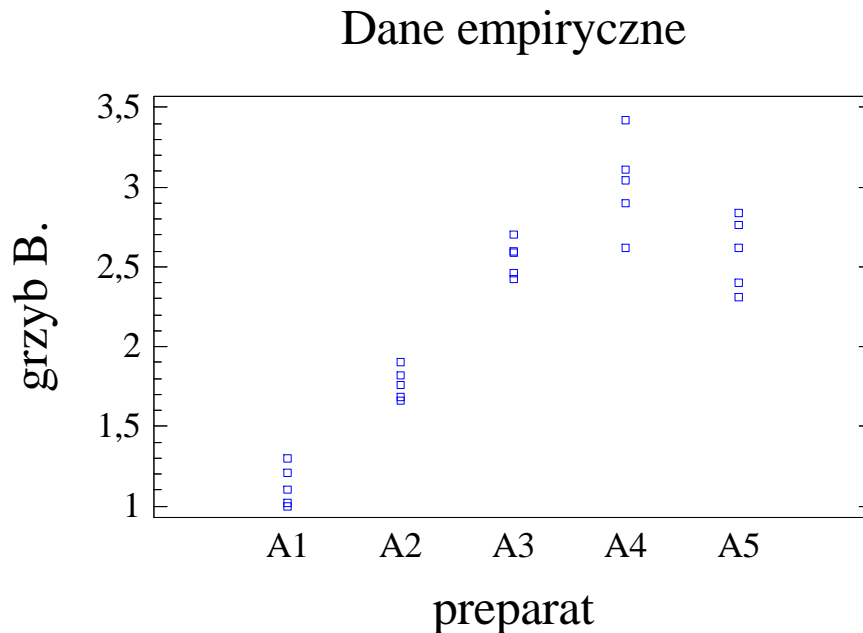
Grupa kontrolna powinna zawierać więcej wyników niż pozostałe grupy (do $\sqrt{k-1}$ razy liczniejsza niż pozostałe grupy).

Przykład:

Badano wpływ pięciu preparatów chemicznych (oznaczonych symbolami A_1 , A_2 , A_3 , A_4 i A_5) na rozwój grzyba *Botrytis cinerea* L. powodującego szarą pleśń. Dla każdego preparatu chemicznego przygotowano pięć jednakowych szalek Petriego z pożywką agarową. Po pewnym czasie zmierzono wielkość kolonii grzyba *Botrytis cinerea* L. (średnicę w cm). Wyniki zawiera tabela:

A_1	1,00	1,21	1,02	1,10	1,30
A_2	1,66	1,90	1,68	1,82	1,76
A_3	2,59	2,46	2,70	2,60	2,42
A_4	2,62	2,90	3,04	3,42	3,11
A_5	2,62	2,31	2,76	2,84	2,40

Graficzna prezentacja danych:



Hipoteza badawcza zakłada, że badane preparaty nie wpływają różnicująco na wielkość kolonii grzyba. Zgodnie z teorią analizy wariancji, odpowiadająca jej hipoteza statystyczna ma postać:

$$H_0 : m_1 = m_2 = \dots = m_5 = m$$

$$(H_1 : \exists i, m_i \neq m)$$

$$F_{emp} = \frac{S_a^2}{S_e^2}$$

Weryfikację hipotezy, czyli konstrukcję funkcji testowej F_{emp} zawiera *tabela analizy wariancji*:

źródło zmienności	sumy kwadratów	stopnie swobody	średnie kwadraty	F_{emp}
preparat-A	$SS_a=11,433$	$v_a = 4$	2,858	80,97
błąd losowy-E	$SS_e=0,706$	$v_e = 20$	0,035	
ogółem-T	$SS_T=12,139$	$v_T = 24$		

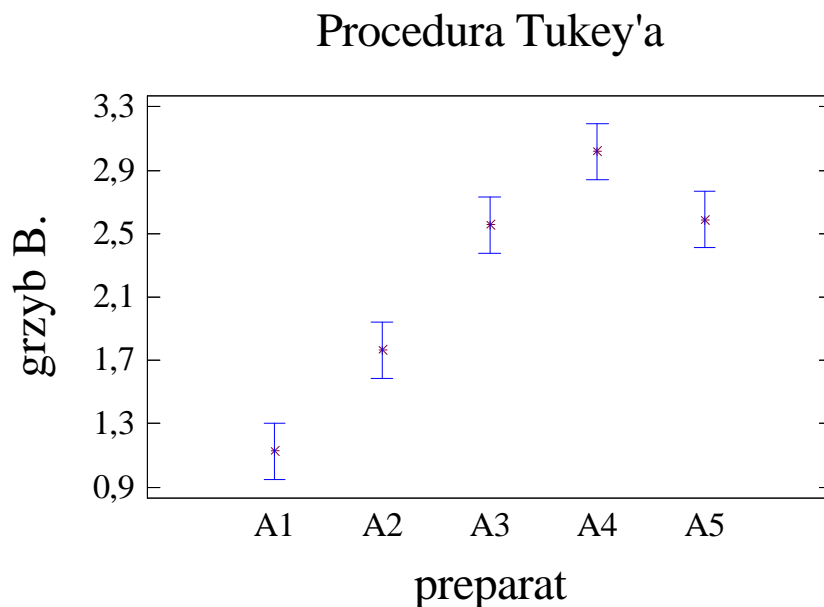
Ponieważ $F_{0.05,4,20} = 2,87$ hipotezę o **równości wszystkich** pięciu **średnich** wielkości kolonii grzyba *Botrytis cinerea* L. odrzucamy. Stwierdzamy istotny wpływ preparatów (zastosowanego preparatu) na wielkość kolonii badanego gatunku grzyba.

Przystępujemy do porównań wielokrotnych średnich obiektowych (populacyjnych), przy pomocy jednej ze znanych procedur podziału na grupy jednorodne. Posłużymy się procedurą Tukey'a:

m_i	\bar{y}_i	
preparat	średnie	gr. jednorodne
A1	1,126	X
A2	1,764	X
A3	2,554	X
A5	2,586	X
A4	3,018	X

$$\text{NIR}_T = 0,355625$$

Graficzna interpretacja testu Tukey'a ma postać:



Niestety test Tukey'a nie gwarantuje rozłącznych podgrup (jak w przypadku powyżej — $\{A1\}$, $\{A2\}$, $\{A3, A5\}$, $\{A4\}$).

preparat	grupy
A1	X
A2	X X
A3	X
A4	X

W tym przypadku najniższą średnią ma $\{A4\}$. Pozostałe 3 grupy nie tworzą jednej podgrupy jednorodnej: A1 i A3 mają istotnie różne średnie. Można powiedzieć, że mamy podgrupę z najwyższą średnią $\{A1, A2\}$ oraz podgrupę środkową $\{A3\}$. Można też stwierdzić, że mamy podgrupę $\{A2, A3\}$ z nieco większą średnią od podgrupy $\{A4\}$, oraz podgrupę $\{A1\}$.

Przykłady **niepoprawnie** założonych doświadczeń:

Badano ilość drobnoustrojów w mleku pewnej firmy. Czy ich ilość jest taka sama w mleku chudym / średniotłustym / tłustym? (3 poziomy czynnika)

Pobrano jedną próbkę z kartonu mleka chudego, próbkę z kartonu mleka średniotłustego i próbkę z kartonu mleka tłustego.

Wada: pojedyncza próbka z każdego poziomu czynnika nie dostarczy informacji o średnim błędzie doświadczalnym.

A jeśli z każdego z 3 kartonów pobrano by po 5 próbek? Wtedy komputer policzy.

Wada: Powtórzenia są z tego samego kartonu, nie z tej samej populacji (kartonów mleka). Komputer policzył, czy te konkretne 3 kartony się różnią, a nie czy różni się mleko chude od średniotłustego od tłustego.

Badano skuteczność różnych środków konserwujących pewien produkt spożywczy. Czy rodzaj środka ma wpływ na stan produktu po 10 dniach?

Wybrano losowo z populacji tego produktu $3 \cdot 10$ próbek, i z 10 losowo wybranych paczek środka konserwującego pobrano odpowiednią ilość środka konserwującego. Każdą partię (10 próbek) konserwowano innym środkiem (A, B lub C).

Ustawiono na parapecie 10 próbek konserwowanych środkiem A, 10 konserwowanych środkiem B i 10 konserwowanych środkiem C.

Wada: Na parapecie warunki przechowywania nie są jednorodne (np. po lewej mogą częściej padać promienie słoneczne). Próbki należało opisać i rozstawić losowo.

Badano zgodność 3 metod pomiaru tłuszczu w mleku. Czy średni wynik metody A jest taki jak metody B i jak metody C?

Z 10 losowo wybranych kartonów mleka pobrano 3 próbki i badano je 3-ma metodami.

Wada: Z kartonu pobrano 3 próbki, więc nie można twierdzić, że te 3 wyniki pomiaru zawartości tłuszczu są niezależne (bo zawartość tłuszczu jest praktycznie identyczna). Analizowanie takich danych za pomocą jednoczynnikowej analizy wariancji nie jest poprawne (inna metoda statystyczna)!