

Cuprins

1	ANOVA - Analiza de varianță	1
1.1	ANOVA uni-factorială. Concepte de bază	1
1.2	Funcționarea și logica analizei de varianță	3
1.3	Calcularea varianței dintre grupuri	6
1.4	Calcularea varianței din interiorul grupurilor	7
1.5	Statistica F. Tabelul ANOVA.	9
1.6	Distribuția de eșantionare F	11
1.7	Asumpțiile analizei de varianță	13
1.8	Glosar de termeni	15

Capitolul 1

ANOVA - Analiza de varianță

1.1 ANOVA uni-factorială. Concepte de bază.

Analiza de varianță continuă seria testelor de semnificație și întărește înțelegerea analizei cauzale, făcând trecerea spre analiza de regresie.

După cum ne reamintim, t este un test de semnificație care testează diferența dintre două medii ale unor variabile metrice. ANOVA este o generalizare a testului t pentru mai mult de două medii, pentru că deseori în cercetarea socială studiem mai mult de două grupuri în același timp. De subliniat este faptul că ANOVA poate să fie folosită cu rezultate foarte bune și la comparația dintre două medii, însă testul își arată adevărata valoare la trei sau mai multe medii.

Exemple de acest tip pot fi multiple:

- putem testa dacă mai multe campanii pro-nataliste diferă unele de altele sub aspectul efectelor acestora (sau în sens contrar, să vedem dacă mai multe campanii diferite de prevenire a sarcinilor nedorite diferă semnificativ unele de altele); unele campanii se pot orienta cu preponderență către spoturi TV, altele către distribuirea de materiale informative tipărite, altele pe consiliere directă ș.a.m.d.
- putem deasemenea să testăm dacă mai multe acțiuni de creștere a participării civice diferă sau nu în ce privește efectele
- putem testa dacă mai multe strategii locale de combatere a sărăciei diferă sau nu între ele, etc.

ANOVA combină și extinde testele t și χ^2 , prin testarea egalității dintre trei sau mai multe medii (pentru trei sau mai multe grupuri). Se testează așadar legătura dintre o variabilă metrică (pentru care se calculează media) și o variabilă calitativă (a cărei valori sau categorii sunt considerate grupuri independente). De asemenea, ANOVA face o introducere clară în analiza cauzală: variabila cauză (independentă) este cea calitativă iar variabila efect (dependentă) este cea metrică.

În exemplul pe care îl vom expune, avem următoarele două variabile:

VÂRSTĂ	vârsta la angajare a unei persoane (variabilă raport, metrică, continuă)
STRATIN	strategia de atragere a tinerilor (variabilă calitativă, nominală, cu mai multe categorii - pot exista mai multe strategii posibile)

În această carte, noi luăm în considerare doar varianta de *ANOVA uni-factorială*, cu un singur factor (în engl. *one-way ANOVA*), deoarece folosim o singură variabilă categorială (denumită în limba engleză “*factor*”) pentru a testa diferențele între mediile grupurilor definite de categoriile acesteia (denumite în engleză “*levels*” - niveluri sau “*treatments*” - tratamente).

În cazul nostru, ne raportăm *doar* la strategia de atragere a tinerilor; este important însă de știut că vârsta la angajare a tinerilor poate fi influențată de diverși factori și că există variante ale analizei de varianță care iau în calcul mai mulți asemenea factori. De pildă, varianta care ia în calcul doi factori se numește *ANOVA bi-factorială* (în engl. *two-way ANOVA*) iar varianta care ia în calcul mai mulți factori se numește *ANOVA multi-factorială* (în engl. *multi-way ANOVA*).

Pentru această secțiune, să presupunem că studiem oportunitățile de acces a tinerilor pe piața forței de muncă și analizăm diferite strategii folosite pentru a atrage tinerii să se angajeze. Ipoteza pe care dorim să o testăm este următoarea: *media de vârstă a persoanelor nou angajate este influențată de strategia de atragere utilizată.*

Setul de ipoteze generale pentru ANOVA este:

$$\begin{cases} H_0 : \mu_1 = \mu_2 = \dots = \mu_k \\ H_A : \text{cel puțin două medii sunt diferite} \end{cases}$$

După cum se poate observa, ipoteza de nul se referă la mediile din populație (notate cu μ). Pentru simplificare, vom lua în considerare doar trei localități, care beneficiază de strategii diferite. Din fiecare localitate vom extrage câte un eșantion (să spunem de 10 persoane), iar ipotezele devin:

$$\begin{cases} H_0 : \mu_1 = \mu_2 = \mu_3 \\ H_A : \text{cel puțin două localități au medii de vârstă sunt diferite} \end{cases}$$

În cuvinte, ipoteza de nul susține că nu este nici o diferență între rezultatele diferitelor strategii (fie sunt toate strategiile foarte bune și atrag mulți tineri, fie sunt toate foarte slabe și nu atrag tineri) iar ipoteza alternativă susține că cel puțin o strategie dă rezultate mai bune decât cel puțin una dintre celelalte (este posibil ca o strategie să aibă un rezultat de mijloc, care nu este semnificativ diferit nici față de strategia de succes maxim, nici față de strategia care atrage cei mai puțini tineri; diferență semnificativă în acest caz există doar între prima și ultima strategie).

După cum se poate vedea în **Tabelul 1.1**, primul grup (cel de control, din localitatea unde nu s-a aplicat nici o strategie) are o medie a vârstei la angajare de 27,8 ani cu o abatere standard de 3,65 ani; al doilea grup (din localitatea unde s-a aplicat prima strategie) are o medie de 23,6 ani cu o abatere standard

Tabelul 1.1: Vârstele la angajare a persoanelor în cadrul a 3 eșantioane independente.

Nr.crt.	Localitate 1	Localitate 2	Localitate 3
1	22	28	20
2	27	22	28
3	32	24	31
4	30	18	26
5	29	21	26
6	27	26	30
7	33	25	21
8	24	20	25
9	24	24	29
10	30	28	27
\bar{x}	27,8	23,6	26,3
s	3,65	3,34	3,59

de 3,34 ani iar al treilea grup o medie de 26,3 ani cu o abatere standard de 3,59 ani.

La o primă vedere, toate cele trei grupuri conțin tineri: există vreo diferență semnificativă între cele trei medii? Cum testăm, mai exact, acest lucru?

În fine, dacă obiectivul principal al acestei analize este de a testa diferențele dintre medii, de ce se numește “*Analiză de varianță*”?

1.2 Funcționarea și logica analizei de varianță

Privind **Tabelul 1.1**, putem extrage câteva informații interesante, care ne vor ajuta în cele ce vor urma.

Avem un eșantion total format din 30 de persoane, deci $n = 30$. Acest eșantion este format din trei grupuri independente de câte 10 persoane fiecare (subeșantioane din trei localități diferite); avem deci: $n_1 = 10$, $n_2 = 10$ și $n_3 = 10$. Pentru fiecare dintre cele trei localități/grupuri putem calcula câte o medie și câte o abatere standard; mai avem așadar: $\bar{x}_1 = 27,8$ cu $s_1 = 3,65$, $\bar{x}_2 = 23,6$ cu $s_2 = 3,34$ și $\bar{x}_3 = 26,3$ cu $s_3 = 3,59$.

În același timp, putem calcula o medie generală pentru eșantionul total (pentru toate cele 30 de observații) $\bar{\bar{x}} = 25,9$ precum și o abatere standard totală $s = 3,714$.

Sintetizând:

$$\begin{array}{lll}
 n = 30 & \bar{\bar{x}} = 25,9 & s = 3,714 \\
 n_1 = 10 & \bar{x}_1 = 27,8 & s_1 = 3,65 \\
 n_2 = 10 & \bar{x}_2 = 23,6 & s_2 = 3,34 \\
 n_3 = 10 & \bar{x}_3 = 26,3 & s_3 = 3,59
 \end{array}$$

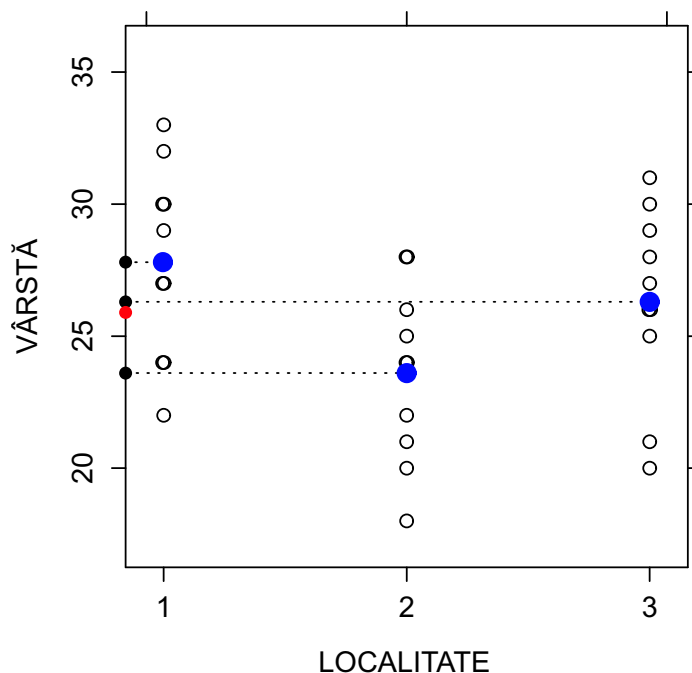
Din faptul că putem calcula abaterile standard, avem un prim indiciu că există o *variație internă* în cadrul fiecărui grup. În **Figura 1.1** se pot observa cele trei grupuri, ale căror observații variază în jurul mediilor reprezentate de punctele colorate în albastru. Această variație poate fi calculată ușor, *pentru fiecare grup în parte*, cu binecunoscuta formulă:

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1} \quad (1.1)$$

Avem așadar trei varianțe în interiorul grupurilor (câte o varianță pentru fiecare din cele trei grupuri).

Apoi, pentru că există trei medii diferite pentru fiecare eșantion, plus o medie generală pentru toate eșantioanele, se poate constata o variație a celor trei medii de grupuri în jurul mediei generale.

Figura 1.1: Scatterplot al vârstelor la angajare pentru cele trei eșantioane.



Celor trei medii le corespund punctele de culoare albastră de pe axa Oy , iar media generală ($\bar{\bar{x}} = 25,9$) este punctul roșu de pe aceeași axă. Poate fi ușor observată variația punctelor albastre în jurul punctului roșu, cu alte cuvinte variația mediilor de grupuri în jurul mediei generale. În cazul nostru, deoarece avem doar trei grupuri, $k = 3$. Adaptând formula 1.1, obținem:

$$s_{\bar{x}}^2 = \frac{\sum_{j=1}^k (\bar{x}_j - \bar{\bar{x}})^2}{k - 1} \quad (1.2)$$

După cum ne aducem aminte de la distribuția de eșantionare a mediei, deviația standard a mediilor în jurul mediei generale este denumită Eroare Standard. **Ecuția 1.2** indică așadar o estimare a Erorii Standard din populație, de unde putem extrage foarte simplu varianța din populație, deoarece:

$$ES = \frac{\sigma}{\sqrt{n}}$$

Ne confruntăm, deci, cu două tipuri de variații: o variație *între grupuri* (variația mediilor de grup în jurul mediei generale) și una *în interiorul grupurilor* (variațiile observațiilor în jurul fiecărei medii de grup). Ambele tipuri de variații sunt folosite ca estimări ale *variației generale* în populație.

IMPORTANT!

Analiza de varianță se bazează pe comparația dintre două estimări ale varianței σ^2 **pentru întreaga populație.**

Logica analizei este următoarea: dacă cele două estimări ale varianței din populație σ^2 sunt aproximativ egale, atunci ipoteza de nul este adevărată (în populație, toate mediile sunt egale). Dacă ipoteza de nul nu este adevărată, atunci cele două estimări ale varianței vor fi semnificativ diferite.

Analiza de varianță se efectuează în trei pași:

- Calcularea primei estimări a varianței în populație: varianța *dintre mediile grupurilor*
- Calcularea celei de a doua estimări a varianței în populație: varianța *din interiorul grupurilor*
- Se compară cele două estimări cu ajutorul statisticii test F. Dacă sunt aproximativ egale (raportul dintre cele două este aproape de valoarea 1), atunci *nu respingem* ipoteza de nul.

Un lucru important de care trebuie să ne aducem aminte (de la măsurile tendinței centrale, capitolul de descriere a variabilelor) este caracterizarea varianței; să mai examinăm încă odată formula:

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}$$

Avem în partea de sus o *sumă de pătrate*, împărțită la un număr de *grade de libertate* ($gl =$ numărul de observații n minus 1, detalii la capitolul de care vorbeam). Ecuția de mai sus mai poate fi scrisă ca:

$$s^2 = \frac{SP}{gl}$$

După cum știm, orice sumă împărțită la numărul de observații se numește *medie*, de unde reiese că varianța nu este nimic altceva decât o medie a unei sume de pătrate.

De aici și denumirile pe care le poartă, în literatura de specialitate, cele două estimări ale varianței din populație σ^2 :

- **MPD** (**M**edia sumei **P**ătratelor **D**intre grupuri) - varianța dintre mediile grupurilor, unde:

$$\text{MPD} = \frac{\text{SPD}}{gl_D}$$

- **MPI** (**M**edia sumei **P**ătratelor din **I**nteriorul grupurilor) - varianța din interiorul grupurilor, unde:

$$\text{MPI} = \frac{\text{SPI}}{gl_I}$$

IMPORTANT!

Atât MPD cât și MPI reprezintă estimări ale varianței în populație, deci pot fi notate amândouă cu $\hat{\sigma}^2$.

1.3 Calcularea varianței dintre grupuri

Aplicăm mai întâi formula din **Ecuția 1.2**:

$$s_{\bar{x}}^2 = \frac{\sum_{j=1}^k (\bar{x}_j - \bar{\bar{x}})^2}{k - 1}$$

Folosind estimarea Erorii Standard, obținem estimarea varianței din populație:

$$\widehat{\text{ES}} = \frac{\hat{\sigma}}{\sqrt{n}} \quad \text{de unde reiese că} \quad \hat{\sigma}^2 = n \cdot \widehat{\text{ES}}^2$$

În această formulă, n este mărimea totală a eșantionului. Cum $s_{\bar{x}}$ este chiar estimarea Erorii Standard, putem înlocui pe $\widehat{\text{ES}}^2$ cu $s_{\bar{x}}^2$ și obținem:

$$\begin{aligned} \hat{\sigma}^2 &= n \cdot s_{\bar{x}}^2 \\ \Rightarrow \hat{\sigma}^2 &= \text{MPD} = \frac{\sum_{j=1}^k n(\bar{x}_j - \bar{\bar{x}})^2}{k - 1} \end{aligned}$$

Numărul de observații n , fiind o constantă, poate sta oriunde: înaintea fracției, înaintea sumei sau chiar în interiorul sumei. Există un motiv special pentru

care n stă în interiorul sumei, pentru că trebuie să luăm în calcul și mărimea grupurilor (în cazul nostru cele trei grupuri au mărime egală, însă de obicei nu este așa); distanța dintre media de grup \bar{x}_j și media generală \bar{x} va fi ponderată cu mărimea grupului respectiv n_j (rezultatul final fiind exact același):

$$\text{MPD} = \frac{\text{SPD}}{gl_D} = \frac{\sum_{j=1}^k n_j (\bar{x}_j - \bar{x})^2}{k - 1} \quad (1.3)$$

Cu alte cuvinte, grupurile mai mari vor avea o pondere mai mare (vor *cântări* mai mult în calcul) decât grupurile mai mici, ceea ce este absolut normal. Se poate observa că *numărul de grade de libertate* pentru MPD este $gl = k - 1$.

Efectuând calculele pentru exemplul nostru:

$$\text{MPD} = \frac{10 \cdot (27,8 - 25,9)^2 + 10 \cdot (23,6 - 25,9)^2 + 10 \cdot (26,3 - 25,9)^2}{3 - 1}$$

$$\Rightarrow \text{MPD} = \frac{36,1 + 52,9 + 1,6}{2} = \frac{90,6}{2}$$

$$\Rightarrow \text{MPD} = 45,3$$

Spunem că varianța *dintre* grupuri este egală cu 45,3.

1.4 Calcularea varianței din interiorul grupurilor

O întrebare pertinentă în acest moment este: *ce înțelegem prin varianța din interiorul grupurilor?* Avem trei grupuri, deci trei varianțe (câte una în interiorul fiecăruia); pe care dintre cele trei o folosim?

Există două răspunsuri posibile:

1. putem folosi oricare dintre cele trei varianțe, dacă ele sunt egale în populație (ceea ce ne duce spre una dintre asumpțiunile acestei analize, prezentată în **Secțiunea 1.7**)
2. putem folosi o medie ponderată a tuturor celor trei varianțe, folosind o procedură derivată de asemenea din formula clasică a varianței.

Având trei grupuri, adunăm trei sume de pătrate:

$$\text{MPI} = \frac{\text{SPI}}{gl} = \frac{\sum_{i=1}^{n_1} (x_{1i} - \bar{x}_1)^2}{n_1 - 1} + \frac{\sum_{i=1}^{n_2} (x_{2i} - \bar{x}_2)^2}{n_2 - 1} + \frac{\sum_{i=1}^{n_3} (x_{3i} - \bar{x}_3)^2}{n_3 - 1}$$

Partea de jos a ecuației este egală cu $n_1 + n_2 + n_3 - 3 = n - 3$ (pierdem trei grade de libertate pentru că avem trei puncte fixe: mediile corespunzătoare celor trei grupuri). La modul general (cu k grupuri) aceasta va fi egală cu $n - k$, iar ecuația poate fi scrisă sub forma unei sume duble:

$$\text{MPI} = \frac{\text{SPI}}{gl_I} = \frac{\sum_{j=1}^k \sum_{i=1}^{n_j} (x_{ji} - \bar{x}_j)^2}{n - k} \quad (1.4)$$

Numărul de grade de libertate pentru MPI este așadar: $gl = n - k$.

Ecuația 1.4 poate fi simplificată și mai mult, folosind formula varianței:

$$s_j^2 = \frac{\sum_{i=1}^{n_j} (x_{ji} - \bar{x}_j)^2}{n_j - 1} \Rightarrow \sum_{i=1}^{n_j} (x_{ji} - \bar{x}_j)^2 = (n_j - 1)s_j^2$$

pentru fiecare dintre cele $j = 1 \dots k$ grupuri, de unde:

$$\text{MPI} = \frac{\text{SPI}}{gl} = \frac{\sum_{j=1}^k (n_j - 1)s_j^2}{n - k} \quad (1.5)$$

La fel ca la MPD, avem și aici o medie ponderată a varianțelor celor k grupuri (aici însă ponderarea s-a realizat prin utilizarea gradelor de libertate ale fiecărui grup): grupurile de mărime mai mare vor avea o pondere mai mare în calcul.

Aplicând **Ecuația 1.5** pentru exemplul nostru cu trei eșantioane:

$$\text{MPI} = \frac{(10 - 1)3,65^2 + (10 - 1)3,34^2 + (10 - 1)3,59^2}{30 - 3}$$

$$\Rightarrow \text{MPI} = \frac{119,6 + 100,4 + 116,1}{27} = \frac{336,1}{27}$$

$$\Rightarrow \text{MPI} = 12,448$$

Spunem că varianța *în interiorul* grupurilor este egală cu 12,448.

Ecuațiile 1.3 și **1.5** sunt cele folosite pentru calcularea celor două estimări ale varianței în populație, în cazul general cu k grupuri și mărimi n_k ale grupurilor. Formulele utilizate nu sunt foarte complicate (chiar dacă așa par la prima vedere), bazându-se exclusiv pe formula clasică a varianței. Din fericire pentru persoanele cu abilități matematice mai scăzute, ele nu trebuie calculate de mână; computerul ne va da automat rezultatele, singura noastră grijă fiind aceea de a le interpreta corect.

1.5 Statistica F. Tabelul ANOVA.

Cel de-al treilea pas în efectuarea analizei de varianță este calcularea statisticii test F, ca raport între cele două estimări ale varianței în populație:

$$F = \frac{\text{Varianța dintre grupuri}}{\text{Varianța în interiorul grupurilor}}$$

sau mai simplu:

$$F = \frac{\text{MPD}}{\text{MPI}} \quad (1.6)$$

După cum vom vedea, există un motiv puternic pentru faptul că MPD se află la numărător, în partea de sus a fracției.

Multe din informațiile prezentate în continuare sunt explicate în detaliu la **Capitolul ??** - Regresia liniară simplă; în general, multe informații în statistică se bazează pe altele, formând un tot unitar. Singura metodă de a înțelege corect toate informațiile este de a citi capitolele care fac referiri unele la altele și de a reciti un capitol cu referințele proaspăt citite.

Să vedem însă cum interpretăm raportul F: MPI, varianța în interiorul grupurilor, este un bun estimator al varianței din populație σ^2 , *indiferent* dacă ipoteza de nul este sau nu adevărată. Aceasta deoarece MPI se bazează pe variațiile din interiorul fiecărui grup, care luate împreună oferă o imagine destul de bună (o estimare destul de bună) a lui σ^2 .

Partea care se află sub lupa testului este însă MPD; dacă ipoteza de nul este adevărată (toate mediile sunt egale) atunci și MPD va fi un bun estimator a lui σ^2 . În această situație MPD va avea o valoare apropiată de cea a lui MPI, iar valoarea lui F va fi aproape de 1 (raportul dintre două cantități egale este egal cu 1). Cu cât valoarea lui F se va apropia mai mult de 1, cu atât va crește probabilitatea de a greși respingând ipoteza de nul.

În cealaltă situație, în care ipoteza de nul nu este adevărată (adică cel puțin una dintre medii este semnificativ diferită), atunci valoarea lui F se va mări considerabil; în același timp, probabilitatea de a greși respingând ipoteza de nul se va micșora pe măsură. Acest lucru se întâmplă deoarece diferențele dintre grupuri tind să mărească MPD.

IMPORTANT!

Statistica F este o măsură care ne arată cât de multă variație se datorează diferențelor dintre grupuri, raportată la variația generată de selecția aleatoare a eșantionului.

Pentru a înțelege și mai bine aceste lucruri, vom introduce încă o sumă de pătrate, ignorată până acum: STP (Suma Totală a Pătratelor); este vorba despre distanța dintre toate observațiile din eșantionul general în jurul mediei generale \bar{x} :

$$\text{STP} = \sum_{i=1}^n (x_i - \bar{x})^2$$

Dacă împărțim STP la $n - 1$ vom obține varianța pentru eșantionul general, a tuturor observațiilor. Deoarece varianța este o măsură a variației iar STP face parte din formula acesteia, rezultă că STP este de asemenea o *bună măsură* a variației totale. Se poate arăta că:

$$\text{STP} = \text{SPD} + \text{SPI} \quad (1.7)$$

Cu alte cuvinte, cantitatea totală de variație este egală cu cantitatea de variație explicată de diferențele dintre grupuri plus cantitatea de variație rămasă neexplicată (erorile aleatoare în jurul mediei); pe scurt, variația totală este egală cu variația explicată plus variația neexplicată, iar valoarea lui F se poate defini ca raport între cele două:

$$F = \frac{\text{VE}}{\text{VN}} \quad (1.8)$$

Cu cât variația explicată va fi mai mare, cu atât va scădea variația neexplicată, iar valoarea lui F va crește spectaculos; invers, cu cât variația explicată va fi mai mică (ceea ce înseamnă că diferențele dintre grupuri sunt foarte mici), cu atât va crește variația neexplicată (datorată erorilor aleatoare) iar valoarea lui F va tinde spre zero (la limită, când grupurile seamănă perfect și nu există absolut nici o diferență între ele, variația explicată va fi egală cu zero).

Diverse programe de analiză statistică pot să difere foarte puțin în modul de prezentare a rezultatelor, însă toate se vor referi la exact același lucru; în general, orice tabel de rezultate va conține următoarele lucruri:

	gl	Suma Pătratelor	Varianța	Valoare F	p
între grupuri	$k - 1$	SPD	MPD	$F = \frac{\text{MPD}}{\text{MPI}}$	$\text{Pr}(> F)$
în interiorul grupurilor	$n - k$	SPI	MPI		
Total	$n - 1$	STP			

Evident că aceasta este o variantă în limba română a tabelului; pentru că cele mai bune programe statistice sunt în limba engleză, rezultatele afișate de calculator vor avea denumirile ca în următorul tabel:

	df	Sum of Squares	Mean Square	F value	p
Between	$k - 1$	SSB	MSB	$F = \frac{\text{MSB}}{\text{MSW}}$	$\text{Pr}(> F)$
Within	$n - k$	SSW	MSW		
Total	$n - 1$	SST			

Uneori pot să apară prescurtări de genul “Sum Sq” sau “Mean Sq”, iar în loc de p se poate găsi frecvent “Sig.” (vom vedea semnificația acestora în următoarea secțiune). Mai departe:

- “df” se referă la gl - gradele de libertate
- “Sum of Squares” înseamnă “Suma Pătratelor”
- SSB (în engl. **Sum of Squares Between**) este echivalent cu SPD
- SSW (în engl. **Sum of Squares Within**) este echivalent cu SPI

- SST (în engl. **Sum of Squares Total**) este echivalent cu STP
- MSB (în engl. **Mean Square Between**) este echivalent cu MPD
- MSW (în engl. **Mean Square Within**) este echivalent cu MPI

Efectuând calculele pentru exemplul nostru, am obținut următorul tabel:

	<i>gl</i>	Suma Pătratelor	Varianța	Valoare F	<i>p</i>
între grupuri	2	90,6	45,3	3,639	0,040
în interiorul grupurilor	27	336,1	12,448		
Total	29	426,7			

Se poate vedea foarte clar modul cum a fost calculat F:

$$F = \frac{45,3}{12,45} = 3,639$$

Valoarea de 3,639 (mult mai mare decât 1) ne sugerează că ipoteza de nul este pe cale de a fi respinsă, pentru că variația explicată de diferențele dintre grupuri este mult mai mare decât variația datorată erorilor aleatoare; existând diferențe majore între grupuri, vor exista cu siguranță și diferențe între mediile acestora.

La fel ca la testele t și χ^2 , decizia se ia după compararea acestei valori cu una critică. Modalitatea alternativă este de a compara valoarea lui p cu pragul de semnificație ales; cum p este mai mic decât $\alpha = 5\%$ (un prag generic, pentru un nivel de încredere de 95%), vom respinge ipoteza de nul: cel puțin una dintre strategii a dat rezultate.

1.6 Distribuția de eșantionare F

În mod similar cu testele t și χ^2 , analiza de varianță folosește o distribuție de eșantionare numită pe scurt “Distribuția F”. Modul de testare a ipotezei de nul este de asemenea similar, prin calcularea unei valori F critice cu care se compară valoarea F obținută în test, sau prin compararea valorii p obținute cu pragul de semnificație α .

Ca la orice distribuție de eșantionare, putem distinge o primă caracteristică a distribuției F: este *continuă* (existând o infinitate de eșantioane posibile), cu un interval de valori care poate varia între 0 și $+\infty$.

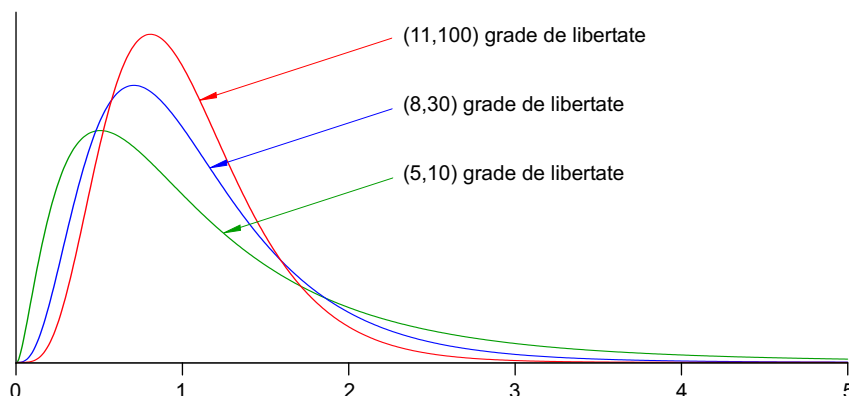
După cum am mai arătat, când eșantioanele sunt perfect similare (media lor este exact aceeași) atunci MPD va fi egal cu 0 iar F va fi egal de asemenea egal cu 0; la celălalt pol, când eșantioanele sunt total diferite (diferența dintre medii este maximă) atunci MPI va fi egal cu 0 iar F va fi egal cu $+\infty$. Astfel, o altă caracteristică a distribuției este aceea că valorile lui F sunt *non-negative* (mai mari sau egale cu zero).

Tot similar cu distribuțiile t și χ^2 , nu există o singură distribuție F, ci o întreagă familie de distribuții; dacă însă până acum forma distribuțiilor depindea de un

singur număr de grade de libertate, la analiza de varianță forma distribuțiilor depinde de o *pereche* de grade de libertate.

Figura 1.2 arată trei asemenea distribuții, unde primul număr reprezintă numărul de grade de libertate de la numărător (din MPD), iar cel de al doilea număr reprezintă numărul de grade de libertate de la numitor (din MPI).

Figura 1.2: Trei distribuții F, cu:



După cum se poate vedea toate curbele sunt mai mult sau mai puțin alungite la dreapta, fiecare având un singur mod. Acestea sunt alte două caracteristici ale distribuției F: este *unimodală* și *alungită la dreapta*. La număr mic al gradelor de libertate pentru numărător (cu alte cuvinte, pentru număr mic de grupuri), curba se apropie din ce în ce mai mult de axa verticală; spre exemplu, la o pereche (1,100) curba va fi chiar lipită de axa Oy . Odată cu creșterea numărului de grade de libertate, cozile distribuțiilor se vor apropia din ce în ce mai mult de axa orizontală Ox , însă nu o vor atinge decât la infinit.

Forma distribuției se modifică deci odată cu creșterea numărului de grade de libertate (atât la numărător cât și la numitor), fiind din ce în ce mai puțin alungită la dreapta. În principiu, creșterea volumului total al eșantionului modifică distribuția până la o formă relativ apropiată de distribuția normală.

Rezumând, proprietățile distribuției F sunt următoarele:

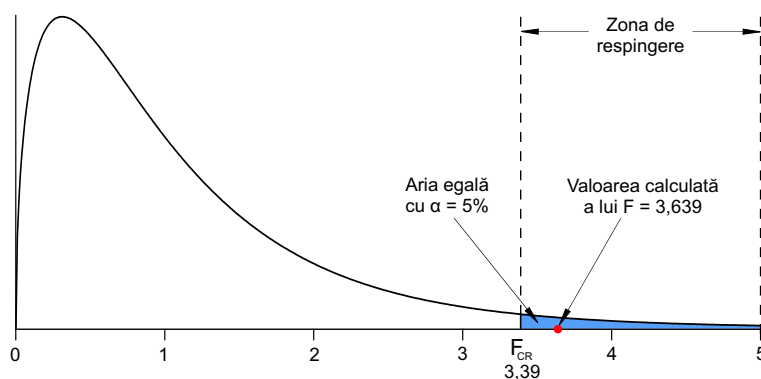
1. este continuă
2. este non-negativă
3. este uni-modală
4. este alungită la dreapta
5. aria de sub curbă este egală cu 1

Testul F este uni-direcțional, *doar* pe coada din dreapta; asta înseamnă că probabilitatea de eroare calculată va fi reprezentată ca o arie sub curbă numai în partea dreaptă. Există o singură valoare critică a lui F, în dreapta căreia se află aria de sub curbă corespunzătoare nivelului de semnificație ales.

Valoarea critică a lui F poate fi găsită cu ajutorul tabelelor de valori care pot fi găsite în anexele oricărui manual de statistică. Există mai multe tabele de valori, câte unul pentru fiecare nivel de semnificație clasic: 10%, 5%, 2,5%, 1% și uneori chiar 0,1%.

Modul de citire a tabelelor este foarte simplu: a) se alege tabelul corespunzător nivelului de semnificație ales; b) se localizează numărul de grade de libertate de la numărător ($k - 1$, de la MPD) pe orizontală, în partea de sus a tabelului; c) se localizează numărul de grade de libertate de la numitor ($n - k$, de la MPI) pe verticală, în partea din stânga a tabelului; d) la intersecția dintre coloana și linia identificate se află valoarea critică a lui F .

Figura 1.3: Probabilitatea de eroare de tipul I, valoarea critică și valoarea calculată a lui F , pe o distribuție cu o pereche (2, 27) grade de libertate



Pentru exemplul nostru, la un nivel de semnificație $\alpha = 5\%$, valoarea critică a lui F este aproximativ egală cu 3,35. Valoarea calculată a lui F este egală cu 3,639 și este mai mare decât valoarea critică, intrând în zona de respingere a ipotezei de nul (colorată cu albastru).

1.7 Asumpțiile analizei de varianță

Pentru a putea utiliza această analiză, trebuie să ne asigurăm că sunt îndeplinite următoarele condiții/asumpții:

1. Fiecare eșantion este extras dintr-o populație cu o distribuție normală.
2. Populațiile din care au fost extrase eșantioanele au aceeași varianță (cu alte cuvinte, toate variază în aceeași măsură).
3. Eșantioanele sunt extrase în mod aleator și independent.

Ca și la testele t și χ^2 , prima asumptie este legată de normalitatea distribuțiilor: populațiile din care au fost extrase eșantioanele trebuie să aibă o distribuție normală. Așa cum se întâmplă mai întotdeauna în practică, însă, asumptia de normalitate este rareori satisfăcută (cazurile în care populația are o distribuție perfect normală sunt foarte rare, dacă nu inexistente).

Există totuși o soluție: violarea acestei asumptii poate fi tolerată, dacă eșantionul este *suficient de mare* (deoarece în acest caz distribuția de eșantionare nu mai depinde de forma distribuției în populație).

Pentru a verifica dacă distribuțiile sunt normale, se construiește câte o histogramă a variabilei metrice pentru fiecare grup (atenție însă: histogramele sunt relevante doar pentru eșantioane mari).

Asumptia de bază a analizei de varianță este cea a *omogenității varianțelor*: această a doua asumptie este cea mai dezbătută de către specialiști. Unii dintre ei afirmă că, dacă varianțele în populațiile din care provin eșantioanele nu sunt egale, atunci ANOVA nu poate fi aplicată. Alții afirmă că acest test este irelevant, deoarece rezultatele lui sunt foarte puternic influențate de forma distribuției în populație (testarea egalității dintre varianțe nu poate fi realizată decât dacă distribuțiile în populație sunt normale).

ANOVA este o analiză destul de robustă, chiar și în cazul în care varianțele nu sunt egale; totul este ca diferența dintre varianțe să nu fie foarte mare (adică o varianță să nu fie de câteva ori mai mare decât alta). Mai mult decât atât, analiza este și mai robustă la încălcarea acestei asumptii dacă eșantioanele sunt de mărime egală ($n_1 = n_2 = n_3$). În concluzie, este bine să evităm aplicarea acestei analize dacă eșantioanele sunt mici, au distribuții puternic deplasate de la normalitate și au varianțe în populație inegale; dacă eșantioanele au însă mărime egală, cu distribuții moderat deplasate și varianțe în populație moderat inegale, atunci putem aplica analiza cu încredere.

În ceea ce ne privește, vom proceda în mod similar cu testul t , unde există o variantă de formulă pentru cazul în care varianțele sunt egale (este vorba de cea clasică, predefinită în orice program de analiză statistică) și o altă variantă de formulă pentru cazul în care varianțele nu sunt egale (testul Welch, care mai este denumit și testul “robust” al egalității mediilor); decizia folosirii uneia sau alteia din variante se ia pe baza valorii lui p din testul Levene de testare a omogenității varianțelor.

Setul de ipoteze din acest test (pentru exemplul nostru particular cu trei eșantioane) este:

$$\begin{cases} H_0 : \sigma_1 = \sigma_2 = \sigma_3 \\ H_A : \text{cel puțin două varianțe sunt diferite} \end{cases}$$

În urma efectuării testului cu datele noastre, a fost obținută o valoare a statisticii test $F = 0,108$ și un $p = 0,898$. După cum se poate judeca din valoarea lui p , dovezile sunt zdrobitoare că varianțele sunt omogene (sunt aproape 90% șanse de a greși afirmând contrariul), drept pentru care vom utiliza testul clasic.

În fine, a treia asumptie arată că toate elementele eșantioanelor trebuie extrase în mod independent, utilizând o tehnică aleatoare. Un rol major îl are metodologia utilizată în cercetare, claritatea cu care a fost făcut instructajul dinaintea cercetării, corectitudinea cu care operatorii de teren aplică instrucțiunile primite etc. Cu cât controlăm mai bine toate aceste detalii, cu atât putem fi mai siguri pe rezultatele noastre. A extrage elemente în mod independent unele de altele înseamnă că între orice pereche de elemente din eșantion nu trebuie să fie nici o legătură (spre exemplu, doi respondenți să nu fie rude).

1.8 Glosar de termeni

Analiza de varianță - ANOVA (în engl. **AN**alysis **O**f **V**ariance). O tehnică statistică utilizată pentru a testa egalitatea dintre trei sau mai multe medii.

Distribuția F (în engl. **F** Distribution). O familie de distribuții de eșantionare folosite pentru a testa diferențele dintre medii sau varianțe, a căror formă depinde de doi parametri (gradele de libertate de la numărător și de la numitor).

Grade de libertate (în engl. **D**egrees of **f**reedom). Număr de observații care pot fi alese în mod liber.

MPD - Varianța dintre grupuri (în engl. **MSB** - **M**ean **S**quare **B**etween sau **Between Group Variance**). Medie a sumei pătratelor dintre grupuri, este o estimare a varianței din populație care calculează variația mediilor de grupuri în jurul mediei generale, împărțind SPD la un număr de grade de libertate.

MPI - Varianța în interiorul grupurilor (în engl. **MSW** - **M**ean **S**quare **W**ithin sau **Within Group Variance**). Medie a sumei pătratelor din interiorul grupurilor, este o estimare a varianței din populație care calculează variația din interiorul tuturor grupurilor (unde grupurile mai mari vor avea o pondere mai mare), împărțind SPI la un număr de grade de libertate.

SPD - Suma Pătratelor Dintre grupuri (în engl. **SSB** - **S**um of **S**quares **B**etween). O măsură a variației dintre grupuri, calculată prin însumarea pătratelor distanțelor de la fiecare medie de grup la media generală.

SPI - Suma Pătratelor în Interiorul grupurilor (în engl. **SSW** - **S**um of **S**quares **W**ithin). O măsură a variației din interiorul tuturor grupurilor, calculată prin însumarea pătratelor distanțelor de la fiecare observație la media grupului de care aparține.

STP - Suma Totală a Pătratelor (în engl. **SST** - **S**um of **S**quares **T**otal). Măsură a variației totale, calculată ca sumă dintre SPD și SPI.

Statistica F sau **Raportul F** (în engl. **F** statistic sau **F** ratio). Este un raport între cele două estimări ale varianței din populație (MPD - variația explicată și MPI - variația neexplicată). Dacă cele două estimări sunt aproximativ egale atunci raportul va fi egal cu 1 iar ipoteza de nul nu poate fi respinsă.

Variație explicată (în engl. **E**xplained **v**ariation). Parte a variației totale, explicată de diferențele dintre grupuri.

Variație neexplicată sau **Eroare** (în engl. **U**nexplained **v**ariation sau **E**rror). Parte a variației totale care nu poate fi explicată și care se datorează erorilor aleatoare ale observațiilor în jurul mediilor de grup.

Variație totală (în engl. **T**otal **v**ariation). Variația generală a tuturor observațiilor din eșantion, egală cu variația explicată plus variația neexplicată.