

Documentație Proiect 1

Probabilități și Statistică

January, 2021

Contents

1	Echipă	2
2	Introducere	2
3	Cerințe	2
3.1	Cerința 1	2
3.2	Cerința 2	2
3.3	Cerința 3	4
3.4	Cerința 4	4
3.5	Cerința 5	4
3.6	Cerința 6	5
3.7	Cerința 7	5
3.8	Cerința 8	8
3.9	Cerința 10	8
3.10	Cerința 11	9
3.11	Cerința 12	11
4	Concluzie	13

1 Echipă

Membrii: Ștefan Radu, Comorașu Ana-Maria, Buhai Darius și Cîiașu Nicoleta.

Grupă: 234

2 Introducere

Construirea unui pachet R pentru lucru cu variabile aleatoare continue.

Folosind documentul suport și orice alte surse de documentare considerați potrivite construiți un pachet R care să permită lucru cu variabile aleatoare continue. Pentru a primi punctaj maxim, pachetul trebuie să implementeze cel puțin 8 din următoarele cerințe (oricare 8):

3 Cerințe

3.1 Cerința 1

Enunț

Fiind dată o funcție f , introdusă de utilizator, să se determine constanta de normalizare. Pentru cazul în care o asemenea constantă nu există, afișarea unui mesaj corespunzător către utilizator.

Rezolvare

Pentru calcularea constantei de normalizare, vom folosi regula: $k \int_{-\infty}^{\infty} f(x)dx = 1$. Astfel, funcția *compute_normalizing_constant* primește 3 parametri, printre care ultimii doi pot avea valori implicite. Se va încerca să se calculeze integrala funcției pe domeniul său - sau pe intervalul $(-\infty, \infty)$. Dacă aceasta a fost calculată cu succes, se va calcula constanta de normalizare drept i^{-1} , pentru i integrala rezultată. Dacă încercarea calculării constantei de normalizare eșuează, se va prinde eroarea și se va afișa mesajul corespunzător, iar valoarea returnată va fi 0.

Cod

```
1 awaiting code here
```

Teste

```
1 f <- function(x){
2   return (exp((-x^2)/2))
3 }
4 compute_normalizing_constant(f, 0, 1)
5 # [1] 1.168737
6 compute_normalizing_constant(f)
7 # [1] 0.3989423
8
9 g <- function(x){
10   x^2
11 }
12 # Caught an error!
13 # <simpleError in integrate(f, lower = lowerBound, upper = upperBound): the integral is
14 # [1] 0 probably divergent>
```

3.2 Cerința 2

Enunț

Verificarea dacă o funcție introdusă de utilizator este densitate de probabilitate.

Rezolvare

Pentru a verifica dacă funcția f este densitate de probabilitate, am testat următoarele 2 proprietăți:

1. $f(x) \geq 0$ (f este nenegativă)
2. $\int_{-\infty}^{\infty} f(x) dx = 1$

În primă parte, pentru a verifica dacă funcția dată este pozitivă, i-am calculat toate valorile din intervalul $[-10000000, 10000000]$ și am verificat dacă sunt pozitive.

Pentru a evita erorile generate în urma integrării funcțiilor divergente, am creat o funcție `safe_integrate` ce returnează `FALSE` în cazul în care funcția dată nu poate fi integrată.

Cod

```
1
2 # Integrate a function and return False if any errors are generated
3 safe_integrate <- function(f, d = c(-Inf, Inf)){
4   tryCatch(integrate(Vectorize(f), d[1], d[2]),
5     error = function(e){
6       print(e)
7       FALSE
8     })
9 }
10
11 # Check if a function is positive in a given interval
12 is_positive <- function(f, d = c(-10000000, 10000000), step=10000000){
13   # Set maximum lower and upper bounds
14   if(d[2] == Inf)
15     d[2] = 10000000;
16   if(d[1] == -Inf)
17     d[1] = -10000000;
18   # Calculate values in that intervals
19   vals <- seq(d[1], d[2], len=step)
20   # Check if all values are positive
21   all(f(vals)>=0)
22 }
23
24 check_pdf <- function(f, d = c(-Inf, Inf)){
25   # Check if the function is positive
26   if(!is_positive(f, d)){
27     print("Function is negative")
28     return(FALSE)
29   }
30   # Safe integrate, handle errors
31   i = safe_integrate(f, d)
32   # Cannot integrate
33   if(typeof(i)=='logical' && i==FALSE)
34     return(FALSE)
35   v = i $ value
36   round(v) == 1
37 }
```

Teste

```
1 f1 <- function(x){
2   if (x > 0 && x < 2){
3     3/8 * (4*x-2*x^2)
4   }else{
5     0
6   }
7 }
8
9 check_pdf(f1)
10 # TRUE
11
```

```

12 f2 <- function(x){
13   x * (4*x-2*x^2)
14 }
15
16 check_pdf(f2)
17 # FALSE

```

Probleme întâlnite

Verificarea funcțiilor pozitive nu este 100% precisă atunci când domeniul nu este definit, deoarece pe un interval infinit, valorile funcției sunt limitate la valori între -10000000 și 10000000.

3.3 Cerința 3

Enunț

Rezolvare

3.4 Cerința 4

Enunț

Rezolvare

3.5 Cerința 5

Enunț

Calculul mediei, dispersiei și a momentelor inițiale și centrate până la ordinul 4 (dacă există). Atunci când unul dintre momente nu există, se va afișa un mesaj corespunzător către utilizator.

Rezolvare

În rezolvarea acestui exercițiu am folosit următoarele formule, preluate de pe [Wikipedia](#):

Momentul Central: $\mu_n = \int_{-\infty}^{\infty} (x - \mu)^n * f(x) dx$

Momentul Inițial: $\mu_n = \int_{-\infty}^{\infty} x^n * f(x) dx$

Iar pentru a determina dacă un moment dat nu există, am folosit funcția `safe_integrate` ce returnează `FALSE` în cazul în care pdf-ul dat nu poate fi integrat, afișând pe ecran un mesaj corespunzător.

Cod

```

1
2 # Calculating central moment of order o
3 central_moment <- function(f, o){
4   m = medium(f)
5   f_new <- function(x){
6     ((x - m) ^ o) * f(x)
7   }
8   res = safe_integrate(f)
9   if(typeof(res)=="logical" && res==FALSE){
10     print("Momentul nu exist !")
11     return(0)
12   }
13   res $ value
14 }
15
16 # Calculating initial moment of order o
17 initial_moment <- function(f, o){
18   f_new <- function(x){
19     x ^ o * f(x)
20   }
21   res = safe_integrate(f)

```

```

22     if(typeof(res)=="logical" && res==FALSE){
23         print("Momentul nu exist !")
24         return(0)
25     }
26     res $ value
27 }

```

Teste

```

1  f = function (x) {
2      if(x<0 || x>1)
3          return(0)
4      x ^ 3
5  }
6
7  initial_moment(f, 2)
8  central_moment(f, 2)

```

3.6 Cerința 6

Enunț

Calculul mediei și dispersiei unei variabile aleatoare $g(X)$, unde X are o repartiție continuă cunoscută iar g este o funcție continuă precizată de utilizator.

Cod

```

1  # am nevoie de functia g, pdf-ul lui X si domeniul de definitie
2  median_and_dispersion <- function(g, fx, domeniu_valori) {
3      # Y = g(X) poate fi vazut ca o v.a.c. noua cu pdf-ul g(f(x)), ambele functii fiind
4      # continue le compun
5
6      # folosesc formula pt media functiilor de x
7      # folosesc integral(din pachetul pracma) deoarece integer imi da eroare cand incerc
8      # sa integrez la +-inf
9      e_y <- integral(Vectorize(function(x){g(x) * fx(x)}),domeniu_valori[1],domeniu_
10     valori[2])
11
12     # pt dispersie, mai folosesc o data formula mediei functiilor de x
13     # pentru a afla X^2
14     e_y2 <- integral(Vectorize(function(x){x^2 * g(x) * fx(x)}),domeniu_valori[1],
15     domeniu_valori[2])
16     dispersie <- e_y2 - e_y^2
17
18     print(paste("Media: ",e_y))
19     print(paste("Dispersia:", dispersie))
20 }

```

Teste

```

1  # exemplu
2  f1 <- function(x)(x^2)
3  f2 <- function(x) (1 * exp(1)^(-1 * x))
4
5  median_and_dispersion(f1,f2, c(0,Inf)) # 2, 20

```

3.7 Cerința 7

Enunț

Crearea unei funcții P care permite calculul diferitelor tipuri de probabilități asociate unei variabile aleatoare continue(similar funcției P din pachetul discreteRV)

Rezolvare

```
1 #####
2 # Cerinta: Crearea unei func ii P care permite calculul diferitelor tipuri de
3 # probabilit i asociate unei variabile aleatoare continue(similar func iei P din
4 # pachetul discreteRV)
5 #
6 # Header functie: myP(f, p)
7 # - unde f este o functie densitate de probabilitate (pdf)
8 # - iar p este un string ce reprezinta probabilitatea (conditionata sau
9 # independenta).
10 #####
11
12 myP <- function(f, p) {
13   operatii_posibile=c("<=", ">=", "=", "<", ">")
14
15   ##### Functii ajutatoare #####
16   # transforma string-ul dat in ceva ce pot utiliza
17   parseaza_expresie <- function(expresie) {
18
19     # scot whitespace
20     expresie <- gsub(" ", "", expresie)
21
22     # iau fiecare operatie posibila pe rand si incerc sa dau split dupa ea, daca pot,
23     # mi-am gasit operatorul si returnez parametri
24     for(op in operatii_posibile) {
25
26       # am dat split corect => in stanga am variabila, in dreapta am bound-ul
27       split <- unlist(strsplit(expresie, op, fixed = TRUE))
28       splitSize <- length(split)
29
30       if (splitSize == 2) {
31         # returnez (v.a.c, operatie, bound)
32         return (c(split[1], op, split[2]))
33       }
34     }
35     #daca am ajuns aici => nu am operator, deci eroare
36     return(c(-1))
37
38   }
39
40   # calculez cdf, adica integrala -inf, bound, adica P(X <= bound)
41   cdf <- function(bound) { return (integrate(f, -Inf, bound) $ value)}
42
43   # transform din string in double
44   compute_bound <-function(bound) {
45     rez <- switch(bound,
46                   "-Inf" = -Inf,
47                   "+Inf" = +Inf,
48                   as.double(bound))
49     return (rez)
50   }
51
52   ## Calculeaza probabilitatea ##
53   evalueaza <- function(operator, bound) {
54
55     # parsez bound-ul, transform in double sau in +-inf
56     bound = compute_bound(bound)
57     integrala <- cdf(bound)
58
59     #pentru >,>= din toata aria(1) scad cdf-ul si obtin restul
60     ans <- switch(
61       operator,
62       "=" = 0,
63       "<=" = integrala,
64       "<" = integrala,
65       ">=" = 1 - integrala,
```

```

66     ">" = 1 - integrala)
67
68     return(ans)
69
70 }
71
72 # daca am o singura expresie e prob independenta
73 prob_independenta <- function(expresie) {
74
75     # scot parametri din expresie
76     parametri <- parseaza_expresie(expresie)
77
78     if(length(parametri) != 3)
79         return("Eroare la parsarea probabilitatii")
80
81     # aici presupun ca expresiile mele sunt mereu de forma x operator bound
82     # ar trebui o verificare, poate, a ordinii
83
84     operator <- parametri[2]
85     bound <- parametri[3]
86
87     print(evalueaza(operator, bound))
88
89 }
90
91 #daca am 2 expresii e prob conditionata
92 prob_conditionata <- function(expresie1, expresie2){
93
94     # scot parametri in variabile ca sa fie readable ce fac mai jos
95     parametri1 <- parseaza_expresie(expresie1)
96     parametri2 <- parseaza_expresie(expresie2)
97     op1 <- parametri1[2]
98     op2 <- parametri2[2]
99     bound1 <- parametri1[3]
100    bound2 <- parametri2[3]
101
102    # am formula  $P(a \text{ depinde de } b) = P(a \text{ intersectat } b)/P(b)$ 
103    # calculez cdf pentru fiecare functie si iau pe cazuri
104    ans1 <- evalueaza(op1, bound1)
105    ans2 <- evalueaza(op2, bound2)
106
107    # daca vreuna e 0, intersectia va da 0
108    if(ans1 == 0)
109        return(0);
110    if(ans2 == 0)
111        return ("Cannot divide by zero")
112
113    ## cazuri
114
115    ## caz in care conditionarea face probabilitatea sa fie imposibila
116    #  $p(x < 3 \mid x > 5) = 0$ 
117    if (op1 %in% c("<=", "<") && op2 %in% c(">=", ">") && bound1 >= bound2)
118        return (0);
119
120    #  $p(x > 5 \mid x < 3) = 0$ 
121    if (op1 %in% c(">=", ">") && op2 %in% c("<=", "<") && bound1 >= bound2)
122        return (0);
123
124    ## caz in care am acelasi fel de operator, facand intersectia defapt doar aleg
125    intervalul cel mai restrans
126    #  $p(x > 3 \mid x > 7) \Rightarrow$  iau  $(-\infty, 3)$ 
127    if(op1 %in% c(">=", ">") && op2 %in% c(">=", ">"))
128        if(bound1 > bound2)
129            return (ans1/ans2)
130        else return (1);
131    #  $p(x < 3 \mid x < 7) \Rightarrow$  iau  $(-\infty, 3)$ 
132    if(op1 %in% c("<=", "<") && op2 %in% c("<=", "<"))
133        if(bound1 < bound2)
134            return (ans1/ans2)
135        else return (1)

```



```

135
136     ## daca nu e niciunul de mai sus, e intersectie de forma x > 5 | x < 7 si fac
diferenta
137     ## din cdf mai mare scad cdf mai mic si mi da bucata de dintre => intersectia
138     return ((cdf(compute_bound(bound2))-cdf(compute_bound(bound1)))/ans2)
139
140 }
141
142 ##### Body functie principala #####
143
144 # parsez parametri
145 parti = unlist(strsplit(p, "|", fixed = TRUE))
146 len = paste(length(parti))
147 switch(len,
148        "0" = return("Eroare"),
149        "1" = return(prob_independenta(p)),
150        "2" = return(prob_conditionata(parti[1],parti[2])),
151        )
152 return ("eroare");
153
154 }

```

Teste

```

1  g <- function (x) {
2    fun <- 0.1*(3*(x^2) + 1)
3    fun[x<0] = 0
4    fun[x>2]=0
5    return ( fun )
6  }
7
8  myP(g, "x>1|x<1.5") #

```

3.8 Cerința 8

Enunț

Rezolvare

3.9 Cerința 10

Enunț

Calculul covarianței și coeficientului de corelație pentru două variabile aleatoare continue.

Rezolvare

În rezolvarea problemei am folosit următoarele formule de calcul ale covarianței și a coeficientului de corelație preluate din Cursurile 8 și 6.

De asemenea, pentru a determina funcțiile pdf ale variabilelor x și y din densitatea comună, am integrat pdf.comun în funcție de x și y pe domeniile aferente.

$$\mu_x = E(x) = \int_a^b x * f(x) dx$$

$$cov(x, y) = (\int_c^d \int_a^b x * y * f(x, y) dx dy) - \mu_x * \mu_y$$

$$cor(x, y) = \frac{cov(x, y)}{\sigma_x * \sigma_y}$$

Cod

```

1
2  # Calculates medium
3  medium <- function(f, d = c(-Inf, Inf)){
4    integrate(function(x){ x * f(x)}, d[1], d[2]) $ value
5  }

```

```

6
7 # Extracts x out of common density of x and y
8 extract_x_marginal <- function(f, dx){
9   Vectorize(function(y){
10     integrate(function(x){f(x, y)}, dx[1], dx[2]) $ value
11   })
12 }
13
14 # Extracts y out of common density of x and y
15 extract_y_marginal <- function(f, dy){
16   Vectorize(function(x){
17     integrate(function(y){f(x, y)}, dy[1], dy[2]) $ value
18   })
19 }
20
21 # Calculates Covariance and Correlation coefficient
22 covariance_and_correlation <- function(pdf, dx, dy){
23
24   fx = extract_x_marginal(pdf, dx)
25   fy = extract_y_marginal(pdf, dy)
26
27   mx = medium(fx, dx)
28   my = medium(fy, dy)
29
30   cov = double_integrate(function(x,y){x*y*pdf(x,y)}, dx, dy) - (mx*my)
31
32   cor = cov / (mx*my)
33
34   list(cov = cov, cor = cor)
35 }

```

Teste

```

1 f <- function (x, y) {
2   return (3/2 * (x^2+y^2))
3 }
4 covariance_and_correlation(f, c(0,1), c(0,3))
5
6 # $cov
7 # [1] -78.04688
8 # $cor
9 # [1] -0.8222222

```

Probleme întâlnite

În rezolvarea problemei am găsit dificil determinarea pdf_x și pdf_y din pdf_comun, lucru pe care l-am rezolvat prin funcțiile extract_x_marginal și extract_y_marginal.

3.10 Cerința 11

Enunț

Pornind de la densitatea comună a două variabile aleatoare continue, construirea densităților marginale și a densităților condiționate.

Rezolvare

Cod

```

1 # conditia 1 pentru o functie de densitate corecta
2 # f trebuie sa aiba valori pozitive pentru
3 # toate valorile din domeniu
4 # cum in domeniu sunt o infinitate de puncte
5 # am ales sa verific 'pe sarite' o parte dintre ele
6 check_positive_2d <- function(f, ab, cd, step) {

```

```

7   for (x in seq(ab[1], ab[2], step)) {
8     for (y in seq(cd[1], cd[2], step)) {
9       if (f(x, y) < 0) return(FALSE)
10    }
11  }
12
13  return(TRUE)
14 }
15
16 # conditia 2 pentru o functie de densitate corecta
17 # f trebuie sa aiba probabilitatea totala egala cu 1
18 # pentru a calcula integrala dubla am folosit functia
19 # 'integral2' din pachetul 'pracma'
20 check_integral_is_1 <- function(f, ab, cd) {
21   integral <- integral2(f, ab[1], ab[2], cd[1], cd[2])$Q
22   return(abs(integral - 1) < 0.0000000001)
23 }
24
25 # intoarce o noua functie g(y) definita
26 # ca f(x, y), unde x este o valoare fixata
27 f_set_x <- function(x, f) function(y) f(x, y)
28
29 # intoarce o noua functie g(x) definita
30 # ca f(x, y), unde y este o valoare fixata
31 f_set_y <- function(y, f) function(x) f(x, y)
32
33 # calculeaza densitatea marginala in X
34 # il setam pe x si integram peste y
35 x_marginal_pdf <- function(f, x, c, d) {
36   return (integrate(f_set_x(x, f), c, d)$value)
37 }
38
39 # calculeaza densitatea marginala in Y
40 # il setam pe y si integram peste x
41 y_marginal_pdf <- function(f, y, a, b) {
42   return (integrate(f_set_y(y, f), a, b)$value)
43 }
44
45 # calculeaza densitatea in X conditionata de Y = y
46 h_x_conditioned_by_y <- function(y, x, f, c, d) {
47   f(x, y) / x_marginal_pdf(f, x, c, d)
48 }
49
50 # calculeaza densitatea in Y conditionata de X = x
51 h_y_conditioned_by_x <- function(x, y, f, a, b) {
52   f(x, y) / y_marginal_pdf(f, y, a, b)
53 }
54
55 marginal_conditional_densities <- function(f, x, y, ab, cd) {
56
57   tryCatch({
58     if (!check_positive_2d(f, ab, cd, 0.01)) {
59       print("functia nu este pozitiva pentru valorile din domeniu")
60       return ()
61     }
62
63     if (!check_integral_is_1(f, ab, cd)) {
64       print("probabilitatea totala este diferita de 1")
65       return ()
66     }
67
68     print(paste("testez cu x=", x, " y=", y, " pe [", ab[1], ",",
69               ab[2], "]x[", cd[1], ",", cd[2], "]", sep=""))
70
71     # densitatea marginala pentru X cu x fixat
72     print(paste("pdf marginal cu x fixat:", x_marginal_pdf(f, x, cd[1], cd[2])))
73
74     # densitatea marginala pentru Y cu y fixat
75     print(paste("pdf marginal cu y fixat:", y_marginal_pdf(f, y, ab[1], ab[2])))
76

```

```

77     # densitatea conditionata in x cand Y = y
78     print(paste("pdf conditionat de y=", y, ": ", h_x_conditioned_by_y(y, x, f, cd
[1], cd[2]), sep=""))
79
80     # densitatea conditionata in y cand X = x
81     print(paste("pdf conditionat de x=", x, ": ", h_y_conditioned_by_x(x, y, f, ab
[1], ab[2]), sep=""))
82     }, error = function(e) {
83         print(paste("A aparut o eroare", e, sep = " "))
84     })
85 }

```

Teste

```

1  test_1 <- function() {
2      f <- function(x, y) (8 / 3) * x ^ 3 * y
3      marginal_conditional_densities(f, 0.5, 1.5, c(0, 1), c(1, 2))
4  }
5
6  test_2 <- function() {
7      f <- function(x, y) (3 / 2) * (x ^ 2 + y ^ 2)
8      marginal_conditional_densities(f, 0, 0.5, c(0, 1), c(0, 1))
9  }
10
11  test_1()
12  test_2()

```

3.11 Cerința 12

Enunț

Construirea sumei și diferenței a două variabile aleatoare continue independente (folosiți formula de convoluție)

Documentare

Întâi a trebuit să mă documentez despre formula de convoluție, așa că m-am orientat spre următoarele surse:

1. Pagina [Convolutions](#) de pe StatLect, unde am găsit formula pentru convoluția a două pdf-uri ale unor variabile aleatoare continue independente.
2. O parte dintr-un curs de pe MIT OpenCourseWare despre [suma variabilelor aleatoare continue independente](#) unde a fost prezentată și o demonstrație pe care voi încerca să o explic mai jos.
3. Acest [video de pe YouTube](#) care vorbește despre formula pentru diferență.

Pentru rezolvarea cerinței, am folosit formulele următoare care nu sunt prezente în curs:

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(z-y)f_Y(y) dy \text{ pentru } Z = X + Y$$

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(y-z)f_Y(y) dy \text{ pentru } Z = X - Y.$$

Cod

```

1  # suma a doua variabile aleatoare continue independente folosind formula de
   convolutie
2
3  convolution_sum <- function(fx,fy) {
4      res <-function(z) {
5          integrate(function(y) {
6              fx(z-y) * fy(y)
7          },-Inf,Inf)
8      }
9  }

```

```

10
11 # diferenta a doua variabile aleatoare continue independente folosind formula de
    convolutie
12 convolution_diff <- function(fx,fy) {
13   res <-function(z) {
14     integrate(function(y) {
15       g(y-z)*f(y)
16     },-Inf,Inf) $ value
17   }
18 }
19

```

Teste

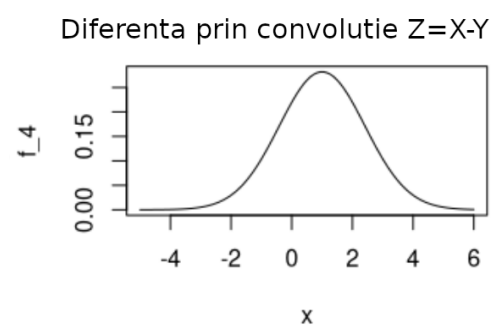
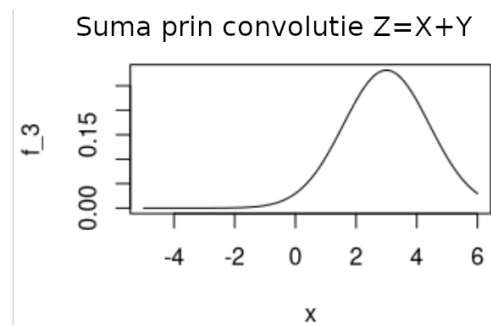
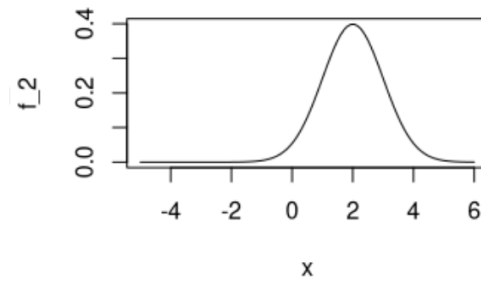
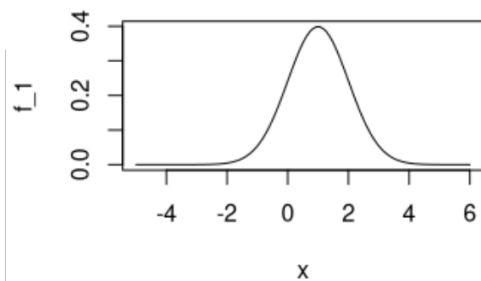
```

1 # doua distributii normale, suma si diferenta facuta cu formula de convolutie
2 f_1 <- function(x)(dnorm(x,mean=1))
3 f_2 <- function(x) (dnorm(x,mean=2))
4 f_3 <- Vectorize(convolution_sum(f_1, f_2))
5 f_4 <- Vectorize(convolution_diff(f_1,f_2))
6
7 plot(f_1,from=-10,to=10,type="l")
8 plot(f_2,from=-10,to=10,type="l")
9 plot(f_3(t),from=-10,to=10,type="l")
10 plot(f_4(t),from=-10,to=10,type="l")

```

Output

Pentru două v.a.c cu distributii normale:



Probleme întâlnite

1. Nu am reușit să integrez spre $+\infty$ cu funcția standard *integrate* fără să obțin erori, așa că am folosit funcția *integral* din pachetul **pracma**.
2. A fost nevoie să vectorizez pdf-urile rezultate prin convoluție ca să pot să le plotez.

4 Concluzie

Prin urmare putem spune că prin rezolvarea acestui proiect, am învățat să aplicăm diferite formule de variabile aleatoare continue în limbajul R, cât și să construim un pachet complet R, cu documentație și manual.