

PROIECT STATISTICĂ

Studenti: Olaru Alexandru - Vergiliu (311)
Matei Elena - Elisabeta (312)
Trache Andrei - Daniel (311)
Negru Liviu - Bogdan (312)

Facultatea de Matematică și
Informatică - Universitatea din București
Specializarea Matematică - Informatică
Februarie 2022

PROBLEMA I

1/ Pentru a determina numărul de elemente aproximativ din intervalul (a, b) a căror sumă depășește $c \gg b$, s-a creat o funcție *average* în care se generează uniform numere din (a, b) . Acestea au fost conținute cu ajutorul variabilei *counter* și însumate în variabila *localSum* până când se îndeplinește criteriul de ieșire din structura repetitivă.

Funcția *average* a fost apelată de 10^9 ori cu ajutorul funcției *replicate*, ce a returnat un vector cu care a fost făcută media aproximațiilor. În plus, pentru cazul

$iterations \leftarrow 10^9$;

$a \leftarrow 0$;

$b \leftarrow 1$;

$c \leftarrow 1$;

S-a observat că aproximația $aproxK = e$, unde e este numărul lui Euler.

```
## 1.
```

```
a <- 1
```

```
b <- 6
```

```
c <- 20
```

```
iterations <- 10^7
```

```
average <- function() { #calculam numarul de elemente aleatorii necesare din (a,b) pentru a depasi c
```

```
  if(a+b >= 0) {  
    localSum <- 0  
    counter <- 0
```

```
    repeat {  
      x <- runif(1, a, b)  
      localSum <- localSum + x  
      counter <- counter + 1  
      if(localSum > c) break
```

```
    }  
    return(counter)
```

```
  } else  
    return(Inf)
```

```
}
```

```
aproxK <- sum(replicate(iterations, average()))/iterations #calculam media
```

2/ $X \sim \text{Unif}(a, b)$, $a \leq b \leq c$, $b \geq 0$.

Am văzut la 1) că K = numărul de variabile aleatoare care însumate dau suma c .

Considerăm distribuția $f_x: (a, b) \rightarrow \mathbb{R}$, $f(x) = \frac{1}{b-a}$.

Media $E[X] = \int_a^b x f_x(x) dx = \int_a^b \frac{x}{b-a} dx = \int_a^b \frac{x}{b-a} dx = \frac{x^2}{2(b-a)} \Big|_a^b = \frac{b^2 - a^2}{2(b-a)}$.

Obținem $E[X] = \frac{a+b}{2}$.

Avem următoarele 3 cazuri:

Cazul 1 $a+b > 0 \Rightarrow E[X] > 0 \Rightarrow K = \left\lceil \frac{c}{E[X]} + 1 \right\rceil = \left\lceil \frac{c}{\frac{a+b}{2}} + 1 \right\rceil = \left\lceil \frac{2c}{a+b} + 1 \right\rceil$.

Deci, pentru $a+b > 0$, avem $K = \left\lceil \frac{2c}{a+b} + 1 \right\rceil$.

Cazul 2 $a+b < 0 \Rightarrow E[X] < 0 \Rightarrow$ în acest caz, me așteptăm ca, în medie, $c \geq 0$ } suma variabilelor să nu atingă c niciodată.

Cazul 3 $a+b = 0 \Rightarrow E[X] = 0 \Rightarrow$ aceeași concluzie ca în cazul 2.
Cum $x \in (a, b)$ și $c \geq b$

După ce a fost calculată valoarea exactă, având deja valoarea aproximativă am creat o variabilă `errorK`, unde am calculat în procente eroarea dintre cele 2 valori.

2.

```
exactK <- ((2*c)/(a+b)) + 1
errorK <- (abs(exactK - aproxK)/exactK)*100 #calculam eroarea in procente
```

Values	
a	1
aproxK	6.2995505
b	6
c	20
errorK	6.17690744680852
exactK	6.71428571428571
iterations	1e+07
Functions	
average	function ()

PROBLEMA II

1/ Pentru a reduce semnificativ din costul computațional, au fost generate o singură dată mulțimile α , reprezentând vectorul probabilităților, și T_i , reprezentând vectorul duratelor de timp ale fiecărei etape, fiind măsurate în algoritmul curent. În plus, $iterations$ reprezintă numărul persoanelor, iar n reprezintă numărul de etape dintr-o secvență.

Pentru a determina timpul în care persoana termină secvența, a fost creată funcția $runStages$ în care se ia o variabilă $totalTime$ pentru a calcula timpul total și o uniformă $nextStage$ pentru a decide dacă se trece la etapa următoare. Aceasta din urmă este comparată cu probabilitatea de trecere $\alpha[i]$. În final, funcția returnează timpul total pe care persoana îl petrece în secvența de etape până la ieșire. Aceasta funcție se execută pentru cele 10^6 persoane, de unde rezultă un vector cu care se calculează timpul mediu petrecut pentru rezolvarea celor n etape.

```
## 1.

iterations <- 10^6 #numarul de persoane
n <- 20 #numarul de etape

alfa <- runif(n-1, 0, 1) #probabilitatea ca o persoana sa termine fiecare activitate, ca vector
Ti <- rexp(n, runif(1, 0, 1)) #timpul necesar pentru fiecare activitate, ca vector

runStages <- function() { #persoana isi parcurge etapele pana la esec
  totalTime <- Ti[1]
  nextStage <- runif(1, 0, 1)

  if(nextStage < alfa[1]) {
    totalTime <- totalTime + Ti[2]

    for (i in 2:(n-1)) {
      nextStage <- runif(1)

      if(nextStage < alfa[i]) {
        totalTime <- totalTime + Ti[i+1]
      }
    }
  }

  return(totalTime)
}

totalTimes <- replicate(iterations, runStages()) #timpul T pentru fiecare persoana, ca vector
averageTime <- sum(totalTimes)/iterations # media timpului necesar fiecarei persoane
```

2/ $i = \overline{1, m}$, $T_i \sim \text{Exp}(\lambda_i)$, $\alpha_i \sim \text{Unif}(0, 1)$, $\alpha_1, \dots, \alpha_{m-1}$ independente
 T = timpul total pentru finalizarea activității
 $E[T] = ?$

Avem pentru $\forall i$, $E[T_i] = \frac{1}{\lambda_i}$, $E[\alpha_i] = \frac{1}{2}$.

De asemenea, avem că T , timpul total petrecut de persoana A, poate lua valorile T_1 , $T_1 + T_2$, $T_1 + T_2 + T_3$, \dots , $T_1 + T_2 + \dots + T_m$.

Așadar, $P(T = T_1) = 1 - \alpha_1$ (A se oprește după etapa 1);

$P(T = T_1 + T_2) = \alpha_1 (1 - \alpha_2)$ (A termină etapa 1, trece mai departe la etapa 2 și se oprește după aceasta).

Analog, $P(T = T_1 + T_2 + \dots + T_k) = \alpha_1 \cdot \alpha_2 \cdot \dots \cdot \alpha_{k-1} (1 - \alpha_k)$, $k = \overline{1, m-1}$.

$$P(T = \sum_{i=1}^m T_i) = \alpha_1 \cdot \alpha_2 \cdot \dots \cdot \alpha_{m-1}$$

$$E[P(T = T_1)] = E[1 - \alpha_1] = 1 - \frac{1}{2} = \frac{1}{2}$$

$$E[P(T = T_1 + T_2)] = E[\alpha_1 (1 - \alpha_2)] = \frac{1}{4}$$

$$\text{Analog, } E[P(T = T_1 + T_2 + \dots + T_k)] = \frac{1}{2^k}, \quad k = \overline{1, m-1}$$

$$E[P(T = T_1 + T_2 + \dots + T_m)] = \frac{1}{2^{m-1}}$$

Notăm $T^{(k)} = T_1 + T_2 + \dots + T_k$ și $P^{(k)} = E[P(T = T^{(k)})]$.

Avem deci $E[T] = \sum_{k=1}^m E[T^{(k)}] \cdot P^{(k)}$

$$= \frac{1}{2} \cdot \frac{1}{\lambda_1} + \frac{1}{2^2} \left(\frac{1}{\lambda_1} + \frac{1}{\lambda_2} \right) + \dots + \frac{1}{2^{m-1}} \left(\frac{1}{\lambda_1} + \dots + \frac{1}{\lambda_{m-1}} \right) + \frac{1}{2^{m-1}} \left(\frac{1}{\lambda_1} + \dots + \frac{1}{\lambda_m} \right)$$

3/ Pentru a afla probabilitatea în care o persoană termină cele m etape ale unei secvențe, s-a făcut produsul $\alpha_i[i]$ -urilor.

3.

```
finishProb <- prod(alfa) #se inmultesc probabilitatile din alfa
```


4/ Pentru a calcula probabilitatea ca o persoană ca o persoană să iasă din secvență înainte de timpul σ , care este luat alator folosind T_i , se adună timpuri până se atinge σ , contorizându-se poziția de ieșire în variabila $counter$, apoi se înmulțesc probabilitățile din α până la $counter-1$. La final, acest rezultat se înmulțește cu $1-\alpha[counter]$, acesta din urmă fiind probabilitatea de ieșire.

4.

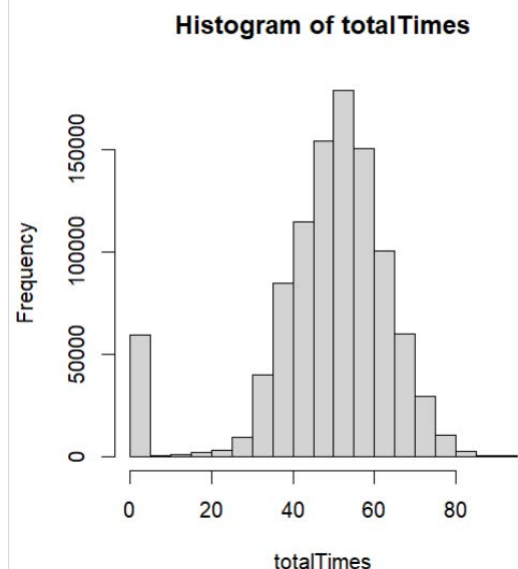
```
sigma <- runif(1, min(cumsum(Ti)), max(cumsum(Ti))) #generam un timp aleator
counter <- 0
timeSum <- 0

if(sigma >= Ti[1]) { #daca sigma e prea mic atunci nu se termin prima etapa
  for(i in 1:n) {
    if((timeSum + Ti[i]) <= sigma) {
      timeSum <- timeSum + Ti[i]
      counter <- counter + 1
    } else {
      break
    }
  }
  sigmaProb <- prod(alfa[1:(counter-1)])*(1-alfa[counter])
} else {
  sigmaProb <- 0
}
```

5/ Pentru a determina timpul minim și timpul maxim în care o persoană termină activitatea, s-au apelat $\min(times)$, respectiv $\max(times)$, unde $times$ reprezintă vectorul celor 10^6 persoane.

5.

```
minTime <- min(totalTimes) #timpul minim
maxTime <- max(totalTimes) #timpul maxim
hist(totalTimes)
```



6/ Pentru a determina probabilitatea de ieșire după o anumită etapă, se alege k etapa și se înmulțesc probabilitățile $\alpha[i]$ până la $k-1$, apoi se înmulțește cu probabilitatea de ieșire $1-\alpha[i]$. Valorile anterioare lui k sunt, de asemenea, adăugate într-un vector *probabilities* cu care se creează un plot. În acest plot se poate observa faptul că este o probabilitate mai mare ca o persoană să iasă înainte să se termine cele k etape.

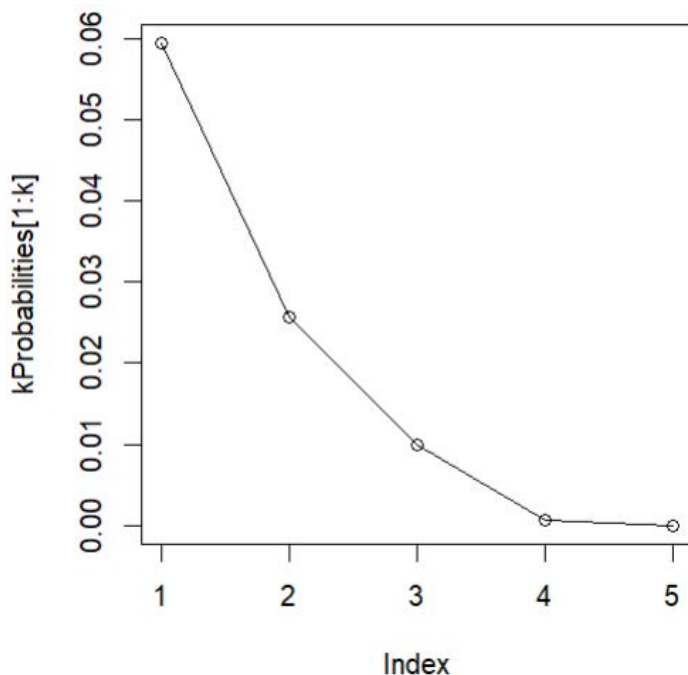
```
## 6.

k <- 5
probabilities <- c(1 - alfa[1])

for (i in 2:(n-1)) {
  tempValue <- prod(alfa[1:(i-1)])*(1-alfa[i])
  probabilities <- c(probabilities, tempValue)
}

kProbabilities <- cumprod(probabilities) #vectorul de probabilitati de a ajunge la k
kProb <- kProbabilities[k] #probabilitatea de ajunge la etapa k

plot(kProbabilities[1:k])
lines(kProbabilities[1:k])
```



Values	
alfa	num [1:19] 0.941 0.54 0.242 0.431...
averageTime	48.7916529679921
counter	14
finishProb	2.51505655620341e-10
i	19L
iterations	1e+06
k	5
kProb	3.23477046646175e-05
kProbabiliti...	num [1:19] 5.94e-02 2.57e-02 9.89...
maxTime	93.9794317960589
minTime	3.17448657228229
n	20
probabilities	num [1:19] 0.0594 0.4326 0.385 0....
sigma	70.8369098514613
sigmaProb	1.56235500269097e-05
tempValue	4.71328684641105e-10
Ti	num [1:20] 3.174 2.636 3.343 0.48...
timeSum	68.8537747138407
totalTimes	Large numeric (1000000 elements, ...)
Functions	
runStages	function ()

PROBLEMA III

1/ Pentru a verifica dacă șirul de variabile aleatoare X_1, X_2, \dots, X_m converge în lege (distributie) la X , aplicăm funcția **check.convergence** pentru ambele distribuții și comparăm graficele rezultate. Ținem cont că toate X_i -urile aparțin unei distribuții **Beta($\frac{1}{m}, \frac{1}{n}$)** iar X este distribuit pe **Binomial($1, \frac{1}{2}$)**. Pentru a genera probabilități aleatoare cu distribuția Beta, ne folosim de funcția **valuesGenBeta** care ne întoarce o listă cu m probabilități. Facem același lucru și pentru variabila repartizată binomial, adică funcția **valuesGenBin** care ne întoarce o probabilitate la care șirul X_i ar trebui să convergă.

Apelăm funcția **check.convergence** pentru fiecare distribuție, cu parametrii respectivi: **nmax** = numărul de puncte n , **M** = numărul de seturi de m valori care să fie generate, **genXm** = funcția care ne generează șirul de probabilități, **mode** = modul de convergență (în cazul nostru, avem "L"). Astfel, comparând cele 2 forme rezultate, observăm că șirul X_i converge în lege la X , când $m \rightarrow \infty$.

Pentru cazul în care X_i sunt în distribuția **Beta($\frac{a}{m}, \frac{b}{m}$)**, cu $a, b > 0$, observăm că nu se aplică convergența în lege deoarece formele rezultate sunt diferite. Convergența se aplică în schimb, doar în cazul în care $a = b$.

```
# III.) #####
require("ConvergenceConcepts")

## 1.

n <- 1000

valuesGenBeta <- function(n) {rbeta(n, 1/n, 1/n)}
valuesGenBin <- function(n) {rbinom(n, 1, 1/2)}

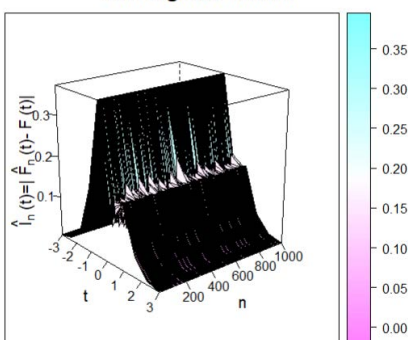
dataBetaL <- check.convergence(nmax = n, M = 5000, genXn = valuesGenBeta, mode = "L")
dataBinL <- check.convergence(nmax = 2, M = 5000, genXn = valuesGenBin, mode = "L")

a <- runif(1, 0, 100)
b <- runif(1, 0, 100)

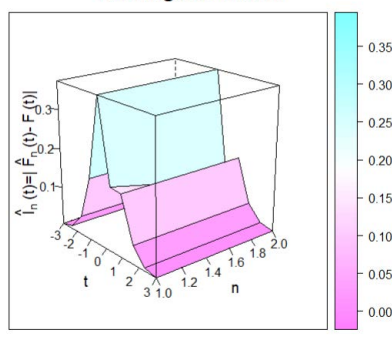
valuesGenBetaC <- function(n) {rbeta(n, a/n, b/n)}

dataBetaLC <- check.convergence(nmax = n, M = 5000, genXn = valuesGenBetaC, mode = "L")
dataBinL <- check.convergence(nmax = 2, M = 5000, genXn = valuesGenBin, mode = "L")
```

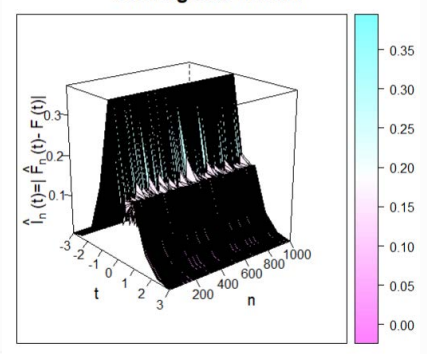
Convergence in law?



Convergence in law?



Convergence in law?



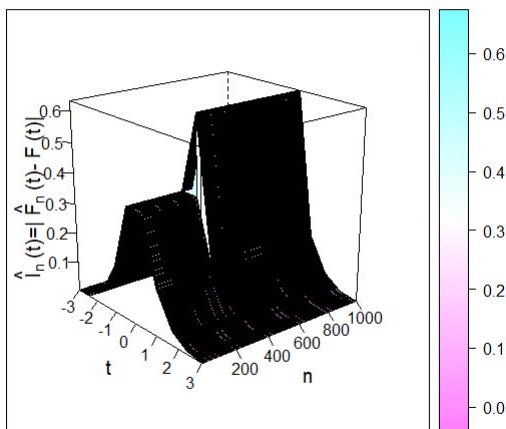
2/ Urmăm un proces asemănător ca la punctul anterior, numai că X_i sunt distribuite uniform pe mulțimea $\{1/m, 2/m, \dots, 1\}$ iar X este uniform distribuit în $(0,1)$. Generăm lista de probabilități cu funcția `valuesGenUnifC`, o variabilă uniformă cu funcția `valuesGenUnif` și comparăm rezultatele. Pentru convergența în distribuție observăm că cele 2 forme rezultate sunt diferite, deci $X_i \not\rightarrow X$, iar pentru convergența în probabilitate, observăm că probabilitatea la care șirul X_i converge pe 0, care este diferită de probabilitatea repartizată uniform X .

2.

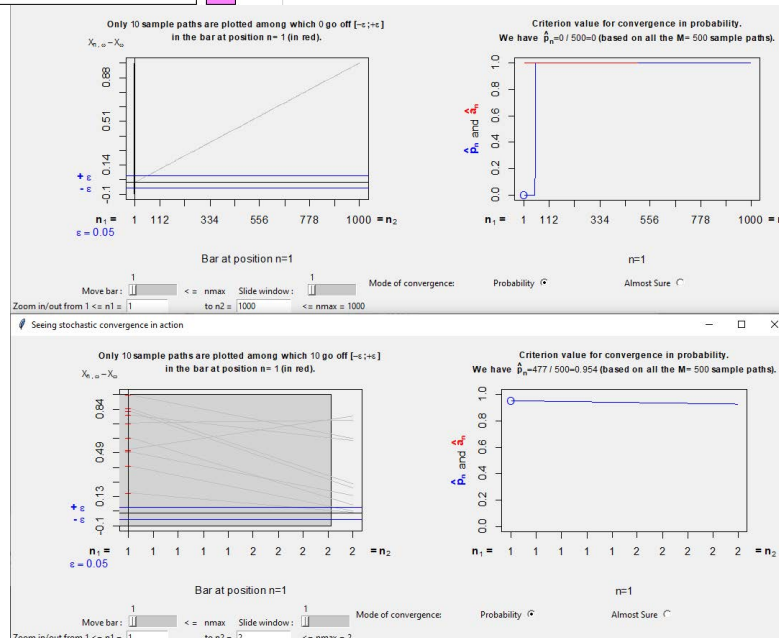
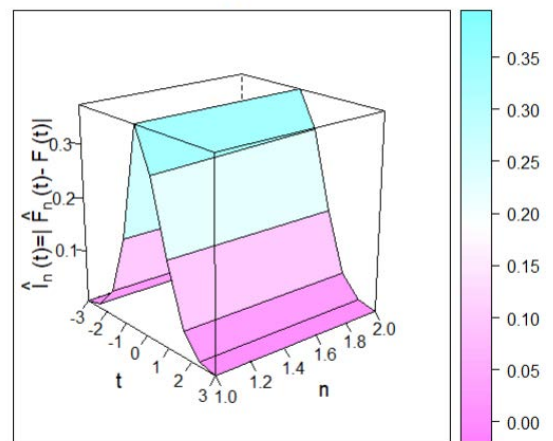
```
valuesGenUnif <- function(n) {runif(n, 0, 1)}
valuesGenUnifC <- function(n) {
  X <- c()
  for(i in 1:n) {
    X <- c(X, (runif(1, i/n, i/n)))
  }
  return(X)
}
```

```
dataUnifLC <- check.convergence(nmax=n, M = 5000, genXn=valuesGenUnifC, mode="L")
dataUnifL <- check.convergence(nmax=2, M = 5000, genXn=valuesGenUnif, mode="L")
dataUnifPC <- check.convergence(nmax=n, M = 5000, genXn=valuesGenUnifC, mode="p")
dataUnifPL <- check.convergence(nmax=2, M = 5000, genXn=valuesGenUnif, mode="p")
```

Convergence in law?



Convergence in law?



PROBLEMA IV

a) Binom (3, p) $\Rightarrow P_K = P(X=K) = C_3^K p^K (1-p)^{3-K}$, $K \in \{0, 1, 2, 3\}$
 $= C_3^K \exp\left(\log\left(\frac{p}{1-p}\right) \cdot K + 3 \log(1-p)\right)$
 $= h(K) \exp(\eta(p) \cdot T(K) - A(p))$

Avem: $h(K) = C_3^K$

$\eta(p) = \log\left(\frac{p}{1-p}\right)$

$T(K) = K$

$A(p) = -3 \log(1-p)$

(considerăm $\theta = p$)

\Rightarrow repartiția binomială face parte din familia exponentială.

b) Binom (m, p) $\Rightarrow P_K = P(X=K) = C_m^K p^K (1-p)^{m-K}$
 $= C_m^K \exp\left(\log\left(\frac{p}{1-p}\right) \cdot K + m \log(1-p)\right)$
 $= h(K) \exp(\eta(p) \cdot T(K) - A(p))$

Avem: $h(K) = C_m^K$

$\eta(p) = \log\left(\frac{p}{1-p}\right) \Rightarrow$ analog a)

$T(K) = K$

$A(p) = -m \log(1-p)$

(considerăm $\theta = p$)

c) Geom(p)

Fie $f: \mathbb{N}^* \rightarrow \mathbb{R}$, $f(x) = (1-p)^{x-1} \cdot p$ și considerăm funcția $I(x) = \begin{cases} 1 & x \in \mathbb{N}^* \\ 0 & \text{altfel} \end{cases}$

Avem că $f(x) = I(x) \cdot p \cdot (1-p)^{x-1}$.

$= I(x) \cdot \left(\frac{p}{1-p}\right) \cdot (1-p)^x$

$= I(x) \cdot \frac{p}{1-p} \exp(x \log(1-p))$

Avem: $h(x) = I(x) \cdot \frac{p}{1-p}$

$\theta = p$

$T(x) = x$

$\eta(\theta) = \log(1-\theta)$

Deci, $\text{Geom}(p) \in$ familiei exponentiale

$$d) \text{Pois}(\lambda) \Rightarrow P_{\theta} = P(X=m) = e^{-\lambda} \cdot \frac{\lambda^m}{m!}$$

$$= \frac{1}{m!} \exp(\log(\lambda) \cdot m - \lambda)$$

$$= h(m) \exp(\eta(\lambda) \cdot T(m) - A(\lambda))$$

$$\text{Averm: } h(m) = \frac{1}{m!}$$

(considerăm $x=m$)

$$\eta(\lambda) = \log(\lambda)$$

(considerăm $\lambda=\theta$)

$$T(m) = m$$

$$A(\lambda) = \lambda$$

\Rightarrow repartiția Poisson face parte din familia exponențială

$$e) \text{Gamma}(\alpha, \beta)$$

$$p(x|\alpha, \beta) = \frac{\alpha^{\beta} x^{\beta-1} e^{-\alpha x}}{\Gamma(\beta)}, \quad 0 < x < \infty$$

$$\text{iar } \Gamma(\beta) = \int_0^{\infty} \alpha^{\beta} x^{\beta-1} e^{-\alpha x} dx$$

$$p(x|\alpha, \beta) = \left[\frac{x^{\beta-1}}{\Gamma(\beta)} \right] \exp(-\alpha x + \beta \log(\alpha))$$

$$= h(x) \exp(\eta T(x) - A(\eta))$$

$$\text{Averm: } h(x) = \left[\frac{x^{\beta-1}}{\Gamma(\beta)} \right]$$

$$\eta = -\alpha$$

$$T(x) = x$$

$$A(\eta) = -\beta \log(\alpha) = -\beta \log(-\eta)$$

\Rightarrow repartiția Gamma face parte din familia exponențială

$$f) \text{Beta}(\alpha, \beta) \Rightarrow \text{fie } f(x) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}$$

$$f(x) = \exp\left(\log\left(\frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}\right)\right) \cdot \frac{1}{x(1-x)} \exp(\alpha \log(x) + \beta \log(1-x))$$

$$= \frac{1}{x(1-x)} \exp\left(\log(\Gamma(\alpha+\beta)) - \log(\Gamma(\alpha)) - \log(\Gamma(\beta)) + \alpha \log(x) + \beta \log(1-x)\right)$$

$$\text{Averm: } h(x) = \frac{1}{x(1-x)}; \quad \theta = (\alpha, \beta)$$

$$\alpha \log(x) + \beta \log(1-x) = \underbrace{(\alpha, \beta)}_{\eta(\theta)} \underbrace{\begin{bmatrix} \log(x) \\ \log(1-x) \end{bmatrix}}_{T(x)}$$

\Rightarrow repartiția Beta face parte din familia exponențială.

$$A(\theta) = \left(\log(\Gamma(\alpha+\beta)) - \log(\Gamma(\alpha)) - \log(\Gamma(\beta)) \right)$$

$$g) \chi^2(v) \text{ îi corespunde lui } \text{Gamma}\left(\frac{v}{2}, \frac{1}{2}\right) \Rightarrow \text{face parte din familia exponențială}$$