

Seminar 11

Partea 3: Mixuri de probabilitate

11.1 Mixturi de probabilitate

Fie $p_1, p_2, \dots, p_n \in (0, 1)$ astfel încât $p_1 + p_2 + \dots + p_n = 1$. Dacă F_1, F_2, \dots, F_n sunt funcțiile de repartiție ale variabilelor aleatoare X_1, X_2, \dots, X_n , atunci funcția

$$F = p_1 F_1 + p_2 F_2 + \dots + p_n F_n$$

este o funcție de repartiție, numită repartiție compusă.

Analog, dacă f_1, f_2, \dots, f_n sunt densitățile de probabilitate ale variabilelor aleatoare X_1, X_2, \dots, X_n , atunci

$$f = p_1 f_1 + p_2 f_2 + \dots + p_n f_n$$

este o densitate de probabilitate, numită densitate compusă.

Definiția 11.1.1 O variabilă aleatoare X ce are densitatea de probabilitate compusă f sau funcția de repartiție compusă F , se numește *mixtură de distribuții de probabilitate* sau, mai simplu, *mixtură de probabilitate*.

Dacă variabila aleatoare X are densitatea $f = p_1 f_1 + p_2 f_2 + \dots + p_n f_n$, aceasta înseamnă că X are densitatea f_1 cu probab. p_1 , \dots , X are densitatea f_n cu probab. p_n . Reprezentând fiecare densitate f_k prin indicele său k , asociem unei densități compuse o variabilă aleatoare discretă:

$$H = \begin{pmatrix} 1 & 2 & \dots & k & \dots & n \\ p_1 & p_2 & \dots & p_k & \dots & p_n \end{pmatrix}. \quad (11.1)$$

Selectând la întâmplare o valoare a lui H este echivalent cu a selecta la întâmplare, cu aceeași probabilitate, o densitate de probabilitate din cele n .

Exemplul 1. O variabilă aleatoare X a cărei densitate de probabilitate

$$f = p_1 f_1 + p_2 f_2 + \dots + p_n f_n, \text{ unde } f_i(x) = \begin{cases} \frac{1}{\theta_i} e^{-x/\theta_i}, & \text{dacă } x \geq 0, \\ 0, & \text{dacă } x < 0, \end{cases} \quad (11.2)$$

este compusa a n densități ale distribuției exponențiale de parametrii $\theta_1, \theta_2, \dots, \theta_n$ se numește variabilă aleatoare hiperexponențială. Variabilele aleat. hiperexponențiale modelează durata serviciului procesorului. Acestea se folosesc în simularea rețelelor de cozi.

Dăm în continuare un exemplu de aplicare a mixturii în analiza algoritmilor:

Exemplul 2. Considerăm instrucțiunea `if-else`:

`if(B) then I1 else I2;`

Fie X_1, X_2 variabilele aleatoare exponențial distribuite care dau timpul de execuție al grupului de instrucțiuni I_1 , respectiv I_2 . Fie $p \in (0, 1)$ probabilitatea ca expresia booleană B să fie adevărată. Densitatea de probabilitate a variabilei ce dă timpul total de execuție al blocului `if-else` este $f = pf_1 + (1 - p)f_2$, unde f_1, f_2 sunt densități exponențiale.

Dacă $p = 0.75$ și $M(X_1) = 20$ milisecunde, $M(X_2) = 40$ milisecunde, să determinăm media timpului de execuție pentru `if-else`.

Se știe că dacă $X \sim \text{Exp}(\theta)$, atunci $M(X) = \theta$. Prin urmare, $\theta_1 = 20$, respectiv $\theta_2 = 40$ și, deci, densitatea de probabilitate a timpului de execuție a lui `if-else` este

$$f(x) = \begin{cases} 0.75 \frac{1}{20} e^{-x/20} + 0.25 \frac{1}{40} e^{-x/40}, & \text{dacă } x \geq 0, \\ 0, & \text{dacă } x < 0. \end{cases}$$

Astfel, media timpului de execuție este

$$\begin{aligned} M(X) &= \int_{-\infty}^{\infty} xf(x)dx \\ &= \int_0^{\infty} x \left(0.75 \frac{1}{20} e^{-x/20} + 0.25 \frac{1}{40} e^{-x/40} \right) dx \\ &= 0.75M(X_1) + 0.25M(X_2). \end{aligned}$$

Mixturile de densități normale, numite și *mixturi Gaussiene*, se folosesc, de exemplu, în analiza și procesarea imaginilor, recunoașterea formelor, în clusterizare etc.

Pe lângă mixturile de distribuții de probabilitate continue se folosesc și mixturi de distribuții discrete.

Exemplul 3. Facebook monitorizează atitudinea unui user față de postările pe wall-uri și îi asociază un număr de reacții ce sunt modelate de o mixtură Poisson. Pentru a înțelege mai ușor, discutăm cazul cel mai simplu când userul are reacția R_1 cu probabilitatea p_1 și reacția R_2 cu probabilitatea p_2 , $p_1 + p_2 = 1$. Reacția R_1 constă din linkuri la articole din *Times New Roman* cu rata λ_1 /oră și reacția R_2 , ce constă din like-uri la pozele amicilor, cu rata λ_2 /oră. Astfel, numărul de reacții/manifestări ale userului pe oră este o variabilă aleatoare X ce are ca distribuție de probabilitate mixtura Poisson:

$$P_X(k) := P(X = k) = p_1 e^{-\lambda_1} \frac{\lambda_1^k}{k!} + p_2 e^{-\lambda_2} \frac{\lambda_2^k}{k!}.$$

ATENȚIE!!! Mixtura de distribuții de probabilitate **nu** înseamnă că variabila X este de forma $X = p_1 X_1 + p_2 X_2$, cu $X_i \sim \text{Poiss}(\lambda_i)$, $i = 1, 2$, ci că X are distribuția Poisson de rată λ_1 cu probab. p_1 , respectiv X are distribuția Poisson de rată λ_2 cu probab. p_2 .

O mixtura modelează foarte bine comportamentul uman, care nu este constant, ci în funcție de diverse circumstanțe "umanul" are o reacție sau alta.

Exercitiu: Dacă $p_1 = 0.6$ și $p_2 = 0.4$, iar $\lambda_1 = 8$, $\lambda_2 = 10$, să se calculeze probabilitatea ca $(X < 3)$ (probabilitatea evenimentului ca într-o oră userul să înregistreze mai puțin de 3 manifestări din cele monitorizate) și numărul mediu de manifestări pe oră.

Mixturile de distribuții m_i -Erlang, $i = \overline{1, n}$, de parametrii $\theta_1, \theta_2, \dots, \theta_n$,

$$f(x) = p_1 \frac{x^{m_1-1} e^{-x/\theta_1}}{\theta^{m_1} (m_1 - 1)!} + p_2 \frac{x^{m_2-1} e^{-x/\theta_2}}{\theta^{m_2} (m_2 - 1)!} + \dots + p_n \frac{x^{m_n-1} e^{-x/\theta_n}}{\theta^{m_n} (m_n - 1)!}$$

se folosesc ca modele pentru aplicații în rețele wireless și sisteme de calcul mobil. O variabilă aleatoare X ce are o astfel de densitate de probabilitate se zice că are distribuție hyper-Erlangen. Exemplu de astfel de model: rețelele wireless de generația a treia oferă servicii integrate de telefonie, date, multimedia etc. Dacă o rețea wireless cu structură celulară oferă n tipuri de servicii și rata medie a sosirii apelurilor de tip i într-o celulă, în unitatea de timp, este λ_i , iar $f_i(x)$ este densitatea de probabilitate a duratei de servire a cererii de tip i într-o celulă, atunci

$$f(x) = \underbrace{\frac{\lambda_1}{\sum_{j=1}^n \lambda_j}}_{p_1} f_1(x) + \underbrace{\frac{\lambda_2}{\sum_{j=1}^n \lambda_j}}_{p_2} f_2(x) + \dots + \underbrace{\frac{\lambda_n}{\sum_{j=1}^n \lambda_j}}_{p_n} f_n(x)$$

este densitatea de probabilitate a duratei de servire a celulei respective.

Exploatând interpretarea dată mai sus distribuției compuse, putem da următoarea modalitate de generare a N valori de observație asupra variabilei aleatoare X ce are distribuție de probabilitate compusă:

```
for i = 1 : N {
  k ← simulatorul variabilei aleatoare discrete H;
  xi ← simulatorul distribuției fk;
}
```

În șirul generat (x_i) , $i = \overline{1, N}$, o proporție de aproximativ $100 p_1\%$ valori vor fi din legea f_1 , $100 p_2\%$ valori vor fi din legea f_2 etc.