

Curs 1/Partea 2: Problema zilei de naștere

Prezentăm în continuare o problemă celebră în teoria probabilităților, numită *problema zilei de naștere*, care a constituit punctul de plecare pentru dezvoltarea unor metode de analiză a algoritmilor, precum și în criptografie.

Exemplul 1. (Problema zilei de naștere) n persoane participă la un reuniune. Care este probabilitatea ca cel puțin două dintre ele să aibă aceeași zi de naștere, în ani diferiți sau nu?

Presupunem că anul are 365 de zile (ignorăm anii bisecți). Notăm cu A mulțimea participanților la reuniune, $A = \{1, 2, \dots, n\}$ și cu $B = \{1, 2, 3, \dots, 365\}$ mulțimea codurilor pentru zilele anului.

- Mulțimea Ω a tuturor posibilităților pentru zilele de naștere ale celor n persoane coincide cu mulțimea n -listelor cu elemente din B sau echivalent cu mulțimea aplicațiilor de la A la B , adică $\Omega = B^A$. De exemplu, dacă $n = 15$ o lista posibilă de zile de naștere este:

$$(32, 24, 125, 51, 73, 84, 279, 330, 23, 127, 71, 95, 199, 211, 5),$$

unde 32 înseamnă 1 februarie, 24 este 24 ianuarie, 125 calculați voi, etc.

Cardinalul lui Ω este $|\Omega| = 365^n$.

- Lucrăm în ipoteza că zilele de naștere sunt aleator (la întâmplare) distribuite în cele 365 zile și n -listele de zile de naștere sunt echiprobabile.

- Fie E_n evenimentul "cel puțin două persoane din cele n ce participă la reuniune au aceeași zi de naștere". Probabilitatea acestui eveniment este $P(E_n) = |E_n|/|\Omega|$. Este mai simplu însă să calculăm probabilitatea complementarei lui E_n .

- Complementul $\mathbb{C}_\Omega E_n$ este evenimentul ca printre cele n persoane să nu existe două cu aceeași zi de naștere, adică:

$$\mathbb{C}_\Omega E_n = \{f \in B^A \mid f \text{ este injectie}\}$$

este mulțimea aplicațiilor injective de la A la B (o injectie asociază la oricare două per-

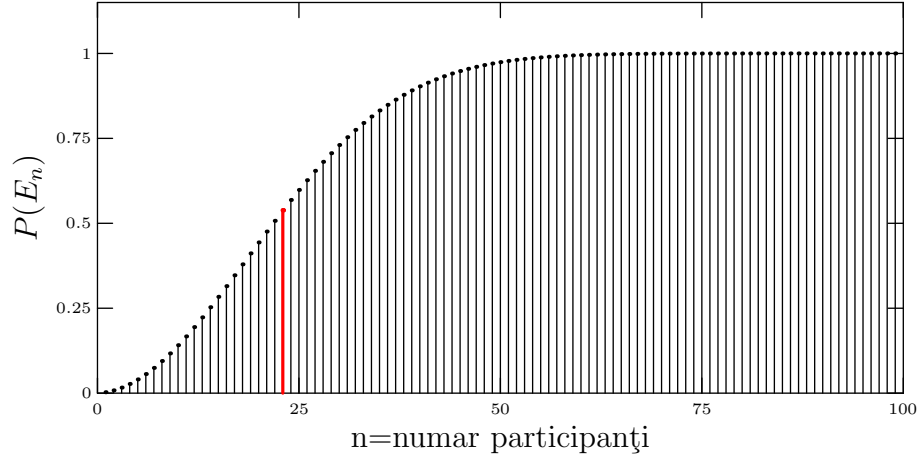


Fig. 1: Ilustrarea prin impulsuri a poziției punctelor $(n, P(E_n))$. Pentru $n \geq 23$, $P(E_n) > 0.5$, $n = 12, \dots, 100$.

soane diferite, zile de naștere diferite). Astfel:

$$\begin{aligned}
 P(\mathbb{C}_\Omega E_n) &= \frac{\text{numărul injectiilor de la A la B}}{\text{numărul aplicațiilor de la A la B}} = \frac{A_{365}^n}{365^n} = \frac{365(365-1) \cdots (365-n+1)}{365^n} \\
 &= \frac{365}{365} \frac{365-1}{365} \cdots \frac{365-(n-1)}{365} = \left(1 - \frac{1}{365}\right) \left(1 - \frac{2}{365}\right) \cdots \left(1 - \frac{n-1}{365}\right) \\
 &= \prod_{k=1}^{n-1} \left(1 - \frac{k}{365}\right)
 \end{aligned}$$

Deci,

$$P(E_n) = 1 - P(\mathbb{C}_\Omega E_n) = 1 - \prod_{k=1}^{n-1} \left(1 - \frac{k}{365}\right)$$

Pentru a extrage informație din această formulă de calcul procedăm astfel: pentru $n = 2, 3, 4, \dots, 100$, calculăm probabilitatea $p_n = P(E_n)$, adică probabilitatea ca într-un grup de n persoane să fie cel puțin două cu aceeași zi de naștere. Într-un sistem de axe ortogonale vizualizăm punctele de coordonate (n, p_n) . Fig.1 ilustrează aceste puncte.

S-a constatat că probabilitatea de a găsi într-un grup de 23 de persoane cel puțin două cu aceeași zi de naștere este mai mare decât $1/2$: $p_{23} = 0.507297$.

Acest rezultat este contrar intuiției umane. Pare nefiresc ca într-un grup de doar 23 persoane să existe cel puțin două cu aceeași zi de naștere, cu o probabilitate mai mare de $1/2$. De aceea, acest fenomen este cunoscut în CS ca *the birthday paradox*.

Generalizare. Problema probabilității coincidenței zilei de naștere apare într-o formă generalizată în criptografie și analiza algoritmilor. Și anume: avem n obiecte numerotate $1, 2, \dots, n$ și N containere (recipienti), $N \geq n$. Pe rând obiectele sunt atribuite

la întâmplare containerelor (sunt aruncate la întâmplare în containere). La întâmplare înseamnă că fiecare obiect poate ateriza în oricare container cu aceeași probabilitate. Dacă într-un container au căzut cel puțin două obiecte spunem că s-a produs o *coliziune*. Probabilitatea să se producă cel puțin o coliziune după aruncarea celor n obiecte în cele N containere se calculează similar cu probabilitatea să participe la reuniune cel puțin două persoane cu aceeași zi de naștere. În cazul coliziunii 365 se înlocuiește cu N . Astfel avem:

Propoziția 1.0.1 Pentru $N \geq n \geq 2$ probabilitatea a cel puțin unei coliziuni $C(N, n)$ este

$$P(C(N, n)) = 1 - \prod_{k=1}^{n-1} \left(1 - \frac{k}{N}\right)$$

și

$$P(C(N, n)) \geq 1 - e^{-n(n-1)/2N}$$

Demonstrație: OPȚIONAL

Probabilitatea de nu se produce nici o coliziune după aruncarea celor n obiecte în N containere este:

$$1 - P(C(N, n)) = \prod_{k=1}^{n-1} \left(1 - \frac{k}{N}\right)$$

Folosind inegalitatea $1 - x \leq e^{-x}$, valabilă pentru orice $x \in [0, 1]$, avem că:

$$\left(1 - \frac{k}{N}\right) \leq e^{-k/N} \text{ și } \prod_{k=1}^{n-1} \left(1 - \frac{k}{N}\right) \leq \prod_{k=1}^{n-1} e^{-k/N} = e^{-\frac{1}{N} - \frac{2}{N} - \dots - \frac{n-1}{N}} = e^{-n(n-1)/2N}$$

Deci probabilitatea producerii a cel puțin unei coliziuni este cel puțin $1 - e^{-n(n-1)/2N}$, adică:

$$P(C(N, n)) \geq 1 - e^{-n(n-1)/2N} \geq 1 - e^{-(n-1)^2/2N}$$

Notând această probabilitate cu $P(n)$ din $P(n) \geq 1 - e^{-(n-1)^2/2N}$ obținem prin inversare

$$n(P) \leq 1 + \sqrt{2N \ln \frac{1}{1-P}} \quad (1)$$

adică numărul de obiecte ce trebuie aruncate în N containere pentru a obține o coliziune cu probabilitatea P este cel mult $1 + \sqrt{2N \ln \frac{1}{1-P}}$. \square

Bazat pe acest rezultat s-a concluzionat că într-o căutare a unui element, obiect, etc, este mult mai ușor să găsim obiecte identice, decât un obiect particular. Astfel s-au dezvoltat algoritmi de coliziune. Ideea de bază a unui astfel de algoritm este să se constituie două liste de elemente și apoi să se caute un element ce apare în ambele, adică să se identifice o coliziune.

Să ilustrăm ce informație obținem din această dependență a numărului n de probabilitatea P , relativ la șansa pe care o are un adversar să ”spargă o valoare hash”.

Atacul birthday asupra unei funcții hash

O funcție hash $h : \{0, 1\}^{\leq L} \rightarrow \{0, 1\}^n$ este o funcție ce asociază unui string de biți de lungime variabilă, mai mică cel mult egală cu L fixat, un string de lungime fixă, n . Hashing înseamnă *chopping and mixing*, adică a ciopârți și a mixa. Principiul de bază al construirii unei funcții hash criptografice este următorul: stringul inițial de biți este divizat în substringuri și acestea sunt mixate prin operații pe biți și rezultatul este un string de lungime n .

Un exemplu de funcție hash este semnătura digitală. Pentru autentificarea expeditorului unui mesaj sau document, de către destinatar, se adaugă mesajului un bloc de date ce constituie așa numita semnătură digitală a expeditorului. Ea se constituie asociind unui string s de lungime variabilă $\leq L$ (ce conține informații despre expeditor), un string de lungime fixă, de n biți, ca rezultat al hashing-ului stringului s .

O funcție hash se construiește astfel încât să poată fi rapid evaluată, dar dacă un adversar interceptează o valoare hash $v \in \{0, 1\}^n$ (semnătura digitală, de exemplu) să fie teoretic imposibil ca într-un timp limitat, să determine unui string s cu proprietatea că $h(s) = v$, chiar dacă acesta dispune de o putere mare de calcul.

Paradoxul zilei de naștere a sugerat și necesitatea asigurării unei alte proprietăți de securitate a funcției hash, numită rezistența la coliziune. Și anume, într-un atac, numit *birthday attack*, adversarul încearcă să găsească două stringuri s, s' ce au aceeași valoare hash: $h(s) = h(s')$.

Să analizăm în cazul funcției cu valori hash pe 64 de biți, câte stringuri trebuie să genereze adversarul, la întâmplare, pentru a obține cu probabilitate mai mare de $1/2$, cel puțin o coliziune, adică două stringuri s, s' ce au aceeași valoare hash. Cu alte cuvinte evaluăm $n(1/2)$ din relația (1) în cazul $N = 2^{64}$, unde N este numărul valorilor hash posibile pentru funcția h :

$$n(1/2) \leq 1 + \sqrt{2^{65} \ln \frac{1}{1 - 1/2}} \approx 1 + 5.056937540686587e + 009$$

Dacă programul adversarului generează într-o secundă un milion de stringuri și valorile lor hash, atunci pentru a evalua $5.056937540686587e + 009$ (neglijăm 1) valori hash are nevoie de

$5.056937540686588e + 003$ secunde ≈ 84 minute = 1 oră și 24 minute (deci nu prea mult!)

În concluzie pentru a preveni succesul unui atac birthday, se construiesc funcții hash cu valori exprimate printr-un număr mare de biți. Funcția hash criptografică, **SHA1**, folosită pentru semnături digitale generează valori hash în $\{0, 1\}^{160}$, în timp ce **MD5** în $\{0, 1\}^{128}$.