

Integracija metodom Monte Carlo

Debak, Matea

Master's thesis / Diplomski rad

2021

Degree Grantor / Ustanova koja je dodijelila akademski / stručni stupanj: **University of Zagreb, Faculty of Science / Sveučilište u Zagrebu, Prirodoslovno-matematički fakultet**

Permanent link / Trajna poveznica: <https://um.nsk.hr/um:nbn:hr:217:635219>

Rights / Prava: [In copyright](#)/[Zaštićeno autorskim pravom.](#)

Download date / Datum preuzimanja: **2024-11-18**



Repository / Repozitorij:

[Repository of the Faculty of Science - University of Zagreb](#)



SVEUČILIŠTE U ZAGREBU
PRIRODOSLOVNO–MATEMATIČKI FAKULTET
MATEMATIČKI ODSJEK

Matea Debak

INTEGRACIJA METODOM MONTE CARLO

Diplomski rad

Voditelj rada:
doc. dr. sc. Snježana Lubura Strunjak

Zagreb, 2021.

Ovaj diplomski rad obranjen je dana _____ pred ispitnim povjerenstvom u sastavu:

1. _____, predsjednik
2. _____, član
3. _____, član

Povjerenstvo je rad ocijenilo ocjenom _____.

Potpisi članova povjerenstva:

1. _____
2. _____
3. _____

Hvala mojoj obitelji, tati, mami, bratu i sestri, na neizmjenoj podršci i ljubavi svih ovih godina! Hvala mojim prijateljima i svim divnim ljudima koji su mi bili poput obitelji tijekom studiranja, a najviše hvala dragom Bogu na tolikim blagoslovima i silnoj ljubavi. Od srca hvala mentorici Snježani na pomoći, savjetima i razumijevanju.

Sadržaj

Sadržaj	iv
Uvod	1
1 Osnovni pojmovi vjerojatnosti i statistike	2
2 Opis problema	8
2.1 Motivacija	8
2.2 Opseg i struktura problema	15
3 Generiranje slučajnih varijabli	18
3.1 Metoda generaliziranog inverza funkcije distribucije	19
3.2 Metoda prihvatanja i odbacivanja	22
4 Integracija metodom Monte Carlo	26
4.1 Monte Carlo metoda	26
4.2 Uzorkovanje po važnosti	33
A Rješenja primjera - R kodovi	43
Bibliografija	51

Uvod

Rješavanje različitih matematičkih problema bez sofisticiranih i praktičnih računalnih metoda u današnje vrijeme gotovo je nezamislivo. Uz iznimne napretke u računalnoj znanosti takav način rješavanja problema postao je izvjestan i nimalo neuobičajen. Međutim, u ne tako davnoj prošlosti, znanstvenici su nailazili na brojne poteškoće prilikom modeliranja problema i određivanja rješenja. Često je dvojba bila sljedeća - modelirali egzaktno problem kojeg nije moguće riješiti kao takvog zbog nedovoljno razvijenih i brzih računala, ili modelirati novi problem unutar nekog već poznatog modela, funkcionalnog i upotrebljivog, koji će dovoljno dobro opisivati inicijalni problem. Iako druga opcija ne omogućava eksplicitno rješenje, znanstvenici su se često odlučivali upravo za tu opciju koja pruža aproksimativno rješenje. Međutim, premda znanost svakim trenutkom sve više napreduje, i dalje se javljaju određeni problemi čija rješenja nerijetko nije moguće odrediti egzaktno, analitičkim putem, jer za njih još uvijek ne postoje efikasni algoritmi. Ponekad, čak i ako postoji algoritam za rješavanje nekog matematičkog problema, zbog njegove nedovoljne brzine ili pak nepraktičnosti izvođenja, isplativije je problem riješiti nekom drugom metodom. Također, aproksimativne metode mogu se primjenjivati u svrhu provjere analitički dobivenih rezultata. Upravo to osnovni su razlozi sve učestalije primjene Monte Carlo statističkih metoda u različitim područjima znanosti. Njihova primjena danas je sasvim uobičajena i poželjna, kako u statistici, tako i u brojnim drugim područjima znanosti, što je posljedica intenzivnijih istraživanja i novih znanstvenih saznanja po tom pitanju.

Cilj ovog rada prvenstveno je približiti samu ideju integracije metodom Monte Carlo. Osim teoretskog dijela opisa i razrade problema, postoji i više motivacijskih primjera koji za cilj imaju na praktičan i opipljiv način pobliže razjasniti i vizualizirati uvedene pojmove te tako postupno pružiti cjelovitu sliku. Ovaj rad zamišljen je kao priručnik za opis i razumijevanje problema integracije metodom Monte Carlo, stoga je naglasak na uvođenju i razradi te ideje, dok su određena predznanja iz područja vjerojatnosti i statistike pretpostavljena i nužna za razumijevanje te navedena u 1. poglavlju. Nakon što u 2. poglavlju opišemo osnovnu ideju ove metode, u 3. poglavlju navest ćemo neke od poznatih metoda generiranja slučajnih varijabli. Konačno, u 4. poglavlju teoretski ćemo razraditi integraciju metodom Monte Carlo.

Poglavlje 1

Osnovni pojmovi vjerojatnosti i statistike

U nastavku su dani neki osnovni pojmovi iz područja vjerojatnosti i statistike čije je razumijevanje neophodno u daljnjem uvođenju i razradi problema. Također, ovdje su navedeni i iskazi određenih matematičkih teorema i propozicija iz područja vjerojatnosti i statistike čiji su rezultati korišteni kasnije u radu.

Sljedeće definicije i teoremi preuzeti su iz [3] i [6].

Definicija 1.0.1. *Neka je Ω neprazan skup. Familija podskupova \mathcal{F} od Ω zove se σ -algebra (ili σ -algebra događaja) ako vrijede sljedeća tri svojstva:*

- (i) $\Omega \in \mathcal{F}$;
- (ii) Ako je $A \in \mathcal{F}$, onda je i $A^c \in \mathcal{F}$ (zatvorenost na komplement);
- (iii) Ako su $A_j \in \mathcal{F}$, $j \in \mathbb{N}$, onda je i $\bigcup_{j=1}^{\infty} A_j \in \mathcal{F}$ (zatvorenost na prebrojive unije).

Uređen par (Ω, \mathcal{F}) zove se izmjeriv prostor.

Definicija 1.0.2. *Vjerojatnost na izmjerivom prostoru (Ω, \mathcal{F}) je funkcija $\mathbb{P} : \mathcal{F} \rightarrow [0, 1]$ koja zadovoljava sljedeća tri aksioma:*

- (A1) (nenegativnost) Za sve $A \in \mathcal{F}$, $\mathbb{P}(A) \geq 0$;
- (A2) (normiranost) $\mathbb{P}(\Omega) = 1$;
- (A3) (σ -aditivnost) Za svaki niz $(A_j)_{j \in \mathbb{N}}$ po parovima disjunktnih događaja $A_j \in \mathcal{F}$ (za $i \neq j$, $A_i \cap A_j = \emptyset$) vrijedi $\mathbb{P}\left(\bigcup_{j=1}^{\infty} A_j\right) = \sum_{j=1}^{\infty} \mathbb{P}(A_j)$.

Uređena trojka $(\Omega, \mathcal{F}, \mathbb{P})$ zove se vjerojatnosni prostor.

Definicija 1.0.3. Neka je $(\Omega, \mathcal{F}, \mathbb{P})$ vjerojatnosni prostor. Slučajna varijabla na $(\Omega, \mathcal{F}, \mathbb{P})$ je svaka funkcija $X : \Omega \rightarrow \mathbb{R}$ takva da vrijedi

$$\{X \leq x\} = \{\omega \in \Omega : X(\omega) \leq x\} \in \mathcal{F}, \quad x \in \mathbb{R}.$$

Definicija 1.0.4. Funkcija distribucije slučajne varijable X je funkcija $F : \mathbb{R} \rightarrow [0, 1]$ definirana formulom

$$F(x) = \mathbb{P}(X \leq x), \quad x \in \mathbb{R}.$$

Teorem 1.0.5. Neka je X slučajna varijabla definirana na vjerojatnosnom prostoru $(\Omega, \mathcal{F}, \mathbb{P})$ te neka je F pripadna funkcija distribucije. Tada vrijedi:

- (i) F je neopadajuća;
- (ii) F je neprekidna zdesna u svakoj točki $x \in \mathbb{R}$;
- (iii) F ima limes s lijeva u svakoj točki $x \in \mathbb{R}$;
- (iv) $F(-\infty) := \lim_{x \rightarrow -\infty} F(x) = 0$ i $F(+\infty) := \lim_{x \rightarrow +\infty} F(x) = 1$.

Definicija 1.0.6. Slučajna varijabla $X : \Omega \rightarrow \mathbb{R}$ je apsolutno neprekidna ako postoji funkcija $f : \mathbb{R} \rightarrow [0, \infty)$ takva da za sve $x \in \mathbb{R}$ vrijedi

$$F(x) = \mathbb{P}(X \leq x) = \int_{-\infty}^x f(t) dt.$$

Funkcija f zove se funkcija gustoće od X .

Definicija 1.0.7. Neka je X neprekidna slučajna varijabla s funkcijom gustoće

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \quad x \in \mathbb{R},$$

gdje su $\mu \in \mathbb{R}$ i $\sigma^2 > 0$ parametri. Slučajna varijabla X zove se normalna slučajna varijabla s parametrima μ i σ . Oznaka je $X \sim N(\mu, \sigma^2)$.

Ako vrijedi $\mu = 0$ i $\sigma^2 = 1$, kažemo da je X standardna normalna slučajna varijabla i pišemo $X \sim N(0, 1)$.

Definicija 1.0.8. Neka je X neprekidna slučajna varijabla s funkcijom gustoće

$$f(x) = \begin{cases} \lambda e^{-\lambda x}, & x \geq 0 \\ 0, & x < 0, \end{cases}$$

gdje je $\lambda > 0$ parametar. Slučajna varijabla X zove se eksponencijalna slučajna varijabla s parametrom λ . Oznaka je $X \sim \text{Exp}(\lambda)$.

Definicija 1.0.9. Neka je X neprekidna slučajna varijabla s funkcijom gustoće

$$f(x) = \begin{cases} \frac{1}{b-a}, & a < x < b \\ 0, & \text{inače,} \end{cases}$$

gdje su $a, b \in \mathbb{R}$, $a < b$, parametri. Slučajna varijabla X zove se uniformna slučajna varijabla na intervalu $[a, b]$. Oznaka je $X \sim \text{Unif}[a, b]$.

Definicija 1.0.10. Neka je X neprekidna slučajna varijabla s funkcijom gustoće

$$f(x) = \begin{cases} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}, & 0 < x < 1 \\ 0, & \text{inače,} \end{cases}$$

gdje su $\alpha > 0$ i $\beta > 0$ parametri, a $\Gamma(\alpha) = \int_0^\infty t^{\alpha-1} e^{-t} dt$ gama funkcija. Slučajna varijabla X ima beta distribuciju s parametrima α i β . Oznaka je $X \sim \text{Beta}(\alpha, \beta)$.

Definicija 1.0.11. Neka je X neprekidna slučajna varijabla s funkcijom gustoće

$$f(x) = \frac{1}{\pi(1+x^2)}, \quad x \in \mathbb{R}.$$

Slučajna varijabla X zove se Cauchyjeva slučajna varijabla.

Definicija 1.0.12. Neka je X apsolutno neprekidna slučajna varijabla s funkcijom gustoće f . Ako je $\int_{-\infty}^{\infty} |x|f(x) dx < \infty$, onda postoji matematičko očekivanje od X koje definiramo s

$$\mathbb{E}(X) := \int_{-\infty}^{\infty} xf(x) dx.$$

Propozicija 1.0.13. Neka je X slučajna varijabla s funkcijom gustoće f i očekivanjem $\mathbb{E}(X)$ te neka su $a, b \in \mathbb{R}$. Tada vrijedi $\mathbb{E}(aX + b) = a\mathbb{E}(X) + b$.

Definicija 1.0.14. Neka je X slučajna varijabla s funkcijom gustoće f i očekivanjem $\mathbb{E}(X)$. Varijanca od X definira se kao

$$\text{Var}(X) := \mathbb{E}[(X - \mathbb{E}(X))^2].$$

Standardna devijacija od X definirana je kao $\sigma(X) := \sqrt{\text{Var}(X)}$.

Uočimo, $0 \leq \text{Var}(X) \leq \infty$ i $\text{Var}(X) = \mathbb{E}(X^2) - (\mathbb{E}(X))^2$.

Propozicija 1.0.15. Ako je $\mathbb{E}(X^2) < \infty$, onda za sve $a, b \in \mathbb{R}$ vrijedi $\text{Var}(aX+b) = a^2\text{Var}(X)$.

Vrijedi sljedeće:

Ako je X apsolutno neprekidna slučajna varijabla s funkcijom gustoće f te $g : \mathbb{R} \rightarrow \mathbb{R}$ neka funkcija, tada je (uz neke dodatne uvjete na funkciju g , npr. g ima najviše prebrojivo mnogo prekida) $Y := g \circ X = g(X)$ također slučajna varijabla.

Teorem 1.0.16. Neka je X apsolutno neprekidna slučajna varijabla s funkcijom gustoće f_X i neka je $g : \mathbb{R} \rightarrow \mathbb{R}$ neka funkcija (koja zadovoljava određena svojstva, npr. g ima najviše prebrojivo mnogo prekida). Ako je $\int_{-\infty}^{\infty} |g(x)|f_X(x) dx < \infty$, onda $Y := g \circ X = g(X)$ ima matematičko očekivanje i vrijedi

$$\mathbb{E}(Y) = \mathbb{E}(g(X)) = \int_{-\infty}^{\infty} g(x)f_X(x) dx.$$

Definicija 1.0.17. Neka su $X_k : \Omega \rightarrow \mathbb{R}$, $k = 1, \dots, n$, neprekidne slučajne varijable na vjerojatnosnom prostoru $(\Omega, \mathcal{F}, \mathbb{P})$. Uređena n -torka $X := (X_1, X_2, \dots, X_n)$ zove se n -dimenzionalan neprekidan slučajni vektor. $X : \Omega \rightarrow \mathbb{R}^n$ je preslikavanje s prostora elementarnih događaja Ω u \mathbb{R}^n .

Definicija 1.0.18. Neka je $X = (X_1, X_2, \dots, X_n)$ neprekidan slučajni vektor na vjerojatnosnom prostoru $(\Omega, \mathcal{F}, \mathbb{P})$. Funkcija $F_X : \mathbb{R}^n \rightarrow [0, 1]$ definirana s

$$F_X(x) = \mathbb{P}(X \leq x) = \mathbb{P}\left(\bigcap_{i=1}^n \{X_i \leq x_i\}\right), \quad \text{za } x = (x_1, \dots, x_n) \in \mathbb{R}^n,$$

je funkcija distribucije slučajnog vektora X .

Definicija 1.0.19. Slučajni vektor $X = (X_1, X_2, \dots, X_n)$ je neprekidan ako postoji funkcija $f : \mathbb{R}^n \rightarrow [0, \infty)$ takva da za sve $x = (x_1, \dots, x_n) \in \mathbb{R}^n$ vrijedi

$$F_X(x) = \int_{-\infty}^{x_n} \int_{-\infty}^{x_{n-1}} \dots \int_{-\infty}^{x_1} f(t_1, t_2, \dots, t_n) dt_1 dt_2 \dots dt_n.$$

Funkcija f zove se funkcija gustoće neprekidnog slučajnog vektora X .

Definicija 1.0.20. Neka su X i Y neprekidne slučajne varijable definirane na vjerojatnosnom prostoru $(\Omega, \mathcal{F}, \mathbb{P})$, s funkcijama gustoće f_X i f_Y . Nadalje, neka je f funkcija gustoće od (X, Y) . Kažemo da su X i Y nezavisne ako vrijedi $f(x, y) = f_X(x)f_Y(y)$ za sve $x, y \in \mathbb{R}$.

Analogno gornjoj definiciji, kažemo da su slučajne varijable X_1, X_2, \dots, X_n nezavisne ako za svaki izbor broja k ($2 \leq k \leq n$) i svaki izbor k -člane kombinacije $X_{i_1}, X_{i_2}, \dots, X_{i_k}$ tog niza varijabli vrijedi da je

$$f_{X_{i_1}, \dots, X_{i_k}}(x_1, \dots, x_k) = f_{X_{i_1}}(x_1) \cdots f_{X_{i_k}}(x_k) \quad \text{za sve } x_1, \dots, x_k \in \mathbb{R}.$$

Ako je (X, Y) neprekidan slučajni vektor s funkcijom gustoće $f_{X,Y}$ i funkcija $g : \mathbb{R}^2 \rightarrow \mathbb{R}$ takva da je $g(X, Y) = g \circ (X, Y)$ neprekidna slučajna varijabla, tada je

$$\mathbb{E}[g(X, Y)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y) f_{X,Y}(x, y) dx dy.$$

Definicija 1.0.21. Neka je X slučajna varijabla s funkcijom distribucije F definirana na vjerojatnosnom prostoru $(\Omega, \mathcal{F}, \mathbb{P})$. Označimo s C_F skup svih točaka u \mathbb{R} u kojima je F neprekidna. Niz funkcija distribucije $(F_n)_{n \in \mathbb{N}}$ konvergira funkciji distribucije F ako vrijedi

$$F(x) = \lim_{n \rightarrow \infty} F_n(x), \quad x \in C_F.$$

Niz slučajnih varijabli $(X_n)_{n \in \mathbb{N}}$ konvergira po distribuciji slučajnoj varijabli X ako niz pripadnih funkcija distribucije konvergira funkciji distribucije od X .

Teorem 1.0.22 (Centralni granični teorem). Neka je $(X_n)_{n \in \mathbb{N}}$ niz nezavisnih jednako distribuiranih slučajnih varijabli sa zajedničkim očekivanjem μ i zajedničkom varijancom σ^2 . Za $n \in \mathbb{N}$, definirajmo $S_n := X_1 + \dots + X_n$. Tada za sve $x \in \mathbb{R}$ vrijedi

$$\lim_{n \rightarrow \infty} \mathbb{P}\left(\frac{S_n - n\mu}{\sigma \sqrt{n}} \leq x\right) = \Phi(x),$$

gdje je Φ funkcija distribucije standardne normalne slučajne varijable. Drugim riječima, niz $((S_n - n\mu)/(\sigma \sqrt{n}))_{n \in \mathbb{N}}$ konvergira po distribuciji ka $N(0, 1)$.

Definicija 1.0.23. Za niz slučajnih varijabli $(X_n)_{n \in \mathbb{N}}$ kažemo da konvergira po vjerojatnosti slučajnoj varijabli X ako za svaki $\varepsilon > 0$ vrijedi

$$\lim_{n \rightarrow \infty} \mathbb{P}(|X_n - X| > \varepsilon) = 0.$$

Pišemo $X_n \xrightarrow{\mathbb{P}} X$ (ili $(\mathbb{P}) \lim_{n \rightarrow \infty} X_n = X$). Niz slučajnih varijabli $(X_n)_{n \in \mathbb{N}}$ konvergira gotovo sigurno slučajnoj varijabli X ako

$$\mathbb{P}(\lim_{n \rightarrow \infty} X_n = X) = 1.$$

Pišemo $X_n \xrightarrow{\text{g.s.}} X$.

Teorem 1.0.24. *Konvergencija gotovo sigurno povlači konvergenciju po vjerojatnosti.*

Teorem 1.0.25 (Jaki zakon velikih brojeva). *Neka je $(X_n)_{n \in \mathbb{N}}$ niz nezavisnih jednako distribuiranih slučajnih varijabli takvih da je $\mathbb{E}(X_n) = \mu$. Tada, kada $n \rightarrow \infty$ vrijedi*

$$\frac{X_1 + \dots + X_n}{n} \rightarrow \mu, \quad \text{g.s.}$$

Definicije i teorem u nastavku preuzeti su iz [4].

Definicija 1.0.26. *Za procjenitelj T_n kažemo da je nepristran procjenitelj za parametar τ ako vrijedi*

$$\mathbb{E}[T_n] = \tau.$$

Definicija 1.0.27. *Za procjenitelj T_n parametra τ za koji je*

$$(\mathbb{P}) \lim_{n \rightarrow \infty} T_n = \tau$$

kažemo da je (slabo) konzistentan. T_n je jako konzistentan za τ ako je

$$\mathbb{P}(\lim_{n \rightarrow \infty} T_n = \tau) = 1.$$

Teorem 1.0.28. *Neka je X_1, X_2, \dots, X_n slučajni uzorak s konačnom varijancom i neka je μ parametar očekivanja. Tada je aritmetička sredina \bar{X}_n :*

- (1) *nepristran procjenitelj za μ ;*
- (2) *konzistentan procjenitelj za μ .*

Sada smo uveli osnovne pojmove i tvrdnje koje se koriste u nastavku i neophodne su za razumijevanje poglavlja koja slijede. Ovdje su navedeni samo iskazi određenih teorema i propozicija, a u [3], [4] i [6] mogu se pronaći dokazi te dodatne tvrdnje i definicije koje mogu biti od pomoći. U dijelovima koji slijede nećemo se pozivati na iskazane tvrdnje, već se podrazumijeva da su one usvojene.

Poglavlje 2

Opis problema

2.1 Motivacija

U statistici, unatoč brojnim sofisticiranim računalnim metodama, nerijetko i danas dolazi do poteškoća prilikom modeliranja i rješavanja problema. Već pri odabiru vjerojatnosnog modela i varijabli koje opisuju problem može se dogoditi da, u nastojanju da ga što bolje i detaljnije opišemo, dobijemo previše složen model kojeg je vrlo teško parametarski opisati. Također, složenost modela može otežati ili onemogućiti statističko zaključivanje, primjerice testiranje hipoteza i procjenu parametara. Upravo zbog te pretjerane složenosti i zamršenosti, ponekad je do zaključka i pronalaska rješenja moguće doći samo simulacijskim metodama, odnosno uz pomoć računala dovoljno velikim brojem izračuna i iteracija pokušati predvidjeti ponašanje složenijih sustava. Naš primarni cilj je pokazati na koji se način simulacijske metode danas mogu koristiti u procesu pronalaska određenih rješenja, konkretno u metodi integracije. U nastavku je stoga dan jedan uvodni primjer koji za cilj ima intuitivno približiti ideju integracije uz pomoć Monte Carlo metode.

Primjer 2.1.1 (Motivacijski primjer). ¹ Neka je zadana funkcija

$$h(x) = [\cos(50x) + \sin(20x)]^2.$$

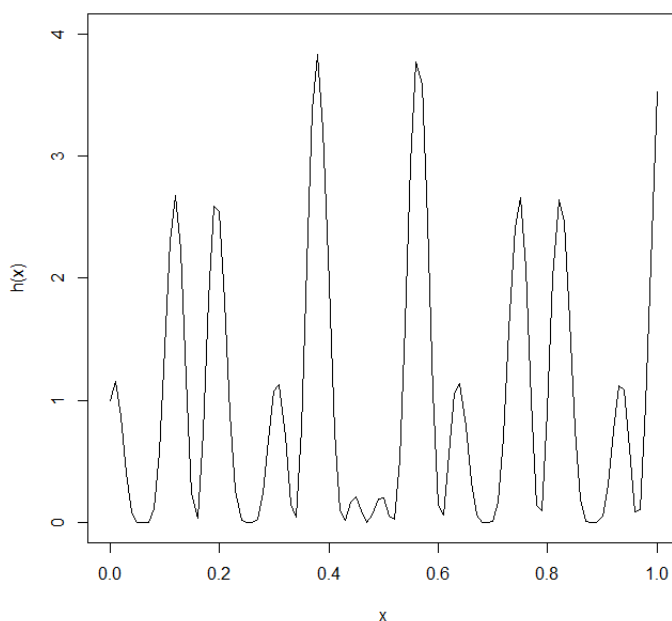
Želimo izračunati integral

$$I = \int_a^b h(x)dx.$$

Integral I moguće je izračunati standardnim metodama supstitucije i parcijalne integracije, kao i numeričkim metodama integracije, koristeći primjerice trapeznu integracijsku

¹Ovaj primjer preuzet je iz [2].

formulu. Međutim, u praksi je najčešće potrebno izračunati integrale složenijih funkcija, čime se oba spomenuta pristupa rješavanju mogu bitno zakomplicirati, pogotovo u višedimenzionalnim problemima. U nastavku ovog primjera pokazat ćemo još jedan mogući pristup rješavanju ovog integrala uz pomoć simulacije slučajnih varijabli.



Slika 2.1: Funkcija h

Umjesto integracije ove funkcije analitičkim putem, uočimo sljedeće:

Za slučajnu varijablu X s uniformnom distribucijom na segmentu $[a, b]$, $a, b \in \mathbb{R}$, $a < b$, $X \sim \text{Unif}[a, b]$, vrijedi:

$$\mathbb{E}[h(X)] = \int_a^b h(x)f_X(x)dx,$$

gdje je f_X pripadna funkcija gustoće uniformne slučajne varijable $X \sim \text{Unif}[a, b]$ dana s:

$$f_X(x) = \begin{cases} \frac{1}{b-a}, & a < x < b \\ 0, & \text{inače.} \end{cases}$$

Dakle,

$$\mathbb{E}[h(X)] = \int_a^b h(x)\frac{1}{b-a}dx = \frac{1}{b-a} I,$$

iz čega slijedi

$$I = (b - a) \mathbb{E}[h(X)] = \mathbb{E}[(b - a)h(X)]. \quad (2.1)$$

Sada je ideja očita, računanje integrala I zapravo svodimo na računanje matematičkog očekivanja $\mathbb{E}[h(X)]$ za $X \sim \text{Unif}[a, b]$. Zapravo, osnovna ideja je simulirati niz od n nezavisnih i jednako distribuiranih (*n.j.d.*) slučajnih varijabli te iskoristiti Jaki zakon velikih brojeva (JZVB).

Za niz nezavisnih i jednako distribuiranih slučajnih varijabli X_1, X_2, \dots, X_n , $X_i \sim \text{Unif}[a, b]$, vrijedi

$$\mathbb{E}[(b - a)h(X_i)] < \infty, \quad i = 1, \dots, n,$$

tj. postoji matematičko očekivanje $\mathbb{E}[(b - a)h(X_i)]$ jer je funkcija h ograničena na segmentu $[a, b]$. Neka je $\mathbb{E}[(b - a)h(X_i)] = \mu$, $\mu \in \mathbb{R}$. Tada, po Jakom zakonu velikih brojeva kada $n \rightarrow \infty$ vrijedi:

$$\frac{\sum_{i=1}^n (b - a)h(X_i)}{n} \rightarrow \mathbb{E}[(b - a)h(X_1)] = \mu, \quad \text{g.s.}$$

Zaključujemo, za procjenitelj

$$\tilde{I}_n := \frac{1}{n} \sum_{i=1}^n (b - a)h(X_i)$$

po Jakom zakonu velikih brojeva vrijedi, kada $n \rightarrow \infty$,

$$\tilde{I}_n \rightarrow \mathbb{E}[(b - a)h(X_1)] = I \quad \text{g.s.}$$

Za primjer odaberimo $a = 0$ i $b = 1$. Želimo izračunati vrijednost integrala

$$I = \int_0^1 h(x)dx.$$

Sada vrijedi sljedeće:

$$\tilde{I}_n = \frac{1}{n} \sum_{i=1}^n h(X_i) \rightarrow \mathbb{E}[h(X_1)] = I \quad \text{g.s.},$$

odnosno, kada $n \rightarrow \infty$,

$$\tilde{I}_n \approx I.$$

Dodatno, vrijedi i $\sigma^2 := \text{Var}[h(X_1)] < \infty$. Iz pokazanog slijedi da je \tilde{I}_n nepristan i konzistentan procjenitelj za I .

\tilde{I}_n , odnosno procjenu intervala I , moguće je izračunati u programskom paketu R , standardnom alatu za statističke izračune i analizu podataka, koji pruža mogućnost jednostavne simulacije slučajnih varijabli iz već poznatih distribucija.

Primjerice, ako simuliramo $n = 10\,000$ uniformno distribuiranih slučajnih varijabli $X_i \sim \text{Unif}[0, 1]$ i izračunamo aritmetičku sredinu

$$\tilde{I}_n := \frac{1}{n} \sum_{i=1}^n h(X_i), \quad (2.2)$$

zapravo smo izračunali procjenu \tilde{I}_n za I . Jedna takva procjena dobivena u programu R jednaka je $\tilde{I}_n = 0.9638259^2$, dok egzaktna vrijednost integrala iznosi $I \approx 0.96520$. Kako se radi o procjeni stvarne vrijednosti, također je potrebno ocijeniti grešku. To je moguće koristeći Centralni granični teorem (CGT), budući da nezavisne i jednako distribuirane slučajne varijable $h(X_i)$ imaju matematičko očekivanje μ i konačnu varijancu σ^2 , po kojem vrijedi:

$$\frac{\sum_{i=1}^n h(X_i) - n\mu}{\sigma \sqrt{n}} \xrightarrow{D} Z \sim N(0, 1), \quad n \rightarrow \infty,$$

odnosno, uz oznaku uvedenu u (2.2) vrijedi

$$\sqrt{n} \cdot \frac{\tilde{I}_n - \mu}{\sigma} \xrightarrow{D} Z \sim N(0, 1), \quad n \rightarrow \infty,$$

odakle konačno slijedi

$$\sqrt{n} \cdot \frac{\tilde{I}_n - I}{\sigma} \xrightarrow{D} Z \sim N(0, 1), \quad n \rightarrow \infty.$$

Intuitivno, vrijedi

$$\tilde{I}_n \stackrel{D}{\approx} N\left(I, \frac{\sigma^2}{n}\right), \quad \text{kada } n \rightarrow \infty,$$

odnosno, možemo reći - procjenitelj \tilde{I}_n je asimptotski normalan s asimptotskom standardnom devijacijom $\frac{\sigma}{\sqrt{n}}$.

Za $\alpha \in (0, 1)$ sada znamo odrediti i aproksimativni $(1 - \alpha)\%$ pouzdani interval. Uzmimo za primjer $\alpha = 0.05$, i odredimo 95% pouzdani interval za I .

Tražimo $z_{\alpha/2}$ takav da je

$$\mathbb{P}(-z_{\alpha/2} \leq Z \leq z_{\alpha/2}) = 0.95,$$

²Vrijednost procjene \tilde{I}_n razlikuje se za svaki različiti niz od n simuliranih slučajnih varijabli. Iz tog razloga se i granice pouzdanog intervala za I razlikuju ovisno o simuliranom uzorku.

gdje je $z_{\alpha/2}$ $\left(1 - \frac{\alpha}{2}\right)$ -kvantil standardne normalne distribucije.³
Stoga je

$$\mathbb{P}\left(-z_{\alpha/2} \leq \sqrt{n} \cdot \frac{\tilde{I}_n - I}{\sigma} \leq z_{\alpha/2}\right) \approx 0.95,$$

odnosno,

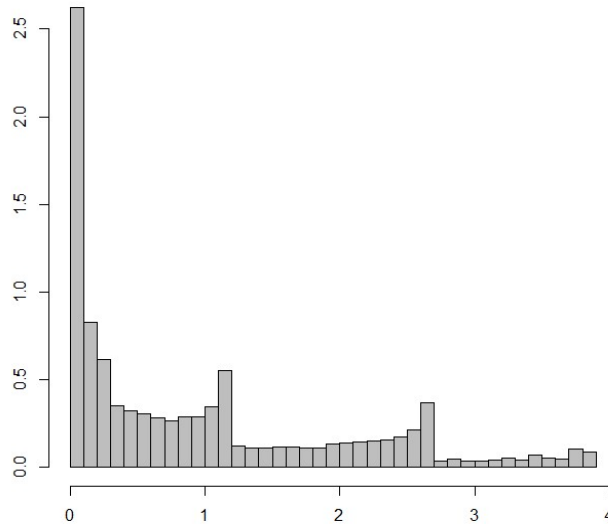
$$\mathbb{P}\left(\tilde{I}_n - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \leq I \leq \tilde{I}_n + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}\right) \approx 0.05.$$

Dakle, 95% pouzdani interval za I je

$$\left[\tilde{I}_n - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}, \tilde{I}_n + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}\right].$$

Očito, pouzdani interval je manji što je n veći, tj. greška procjene teži u 0 kada $n \rightarrow \infty$.
U našem primjeru, uvrštavanjem izračunatih vrijednosti dobijemo procjenu 95% pouzdanog intervala za I koja iznosi $[0.9434476, 0.9842042]$ ⁴.

Izračunate vrijednosti u nastavku su prikazane grafički.



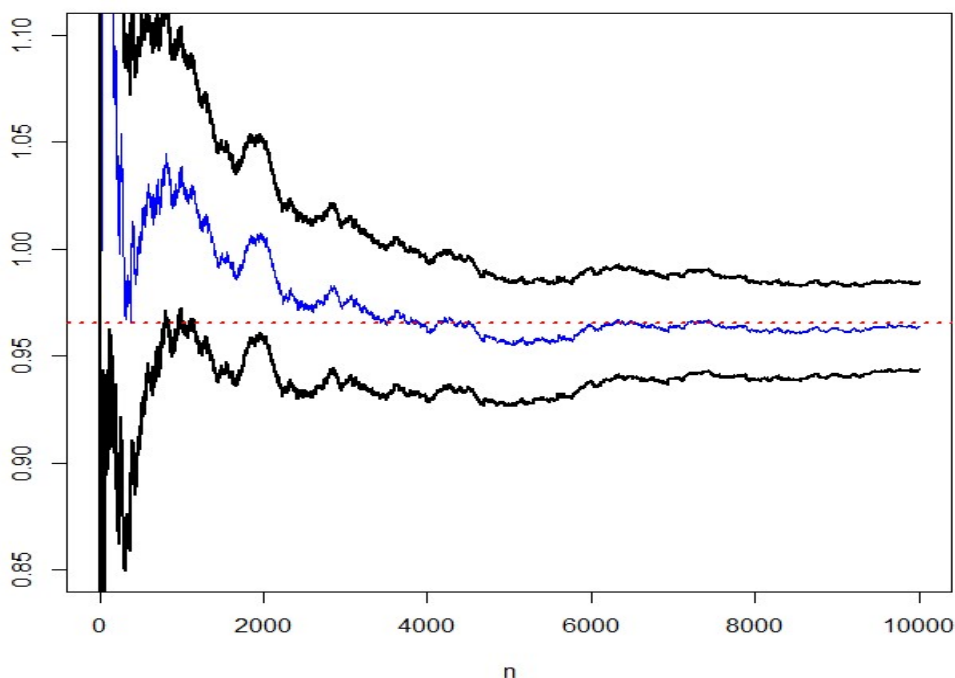
Slika 2.2: Histogram simuliranog uzorka podataka $h(X_i)$

³Odnosno, mora vrijediti $\mathbb{P}(Z \leq z_{\alpha/2}) = 1 - \alpha/2$, za $\alpha \in (0, 1)$.

⁴Varijancu σ^2 možemo procijeniti uzoračkom varijancom S_n^2 . Vrijedi $S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (h(X_i) - \tilde{I}_n)^2 = \frac{1}{n-1} (\sum_{i=1}^n (h(X_i))^2 - n\tilde{I}_n^2)$, iz čega direktno dobijemo procjenu za standardnu devijaciju σ .

Na slici 2.2 prikazan je histogram podataka $h(X_i)$, za $n = 10^4$ simuliranih slučajnih varijabli $X_i \sim \text{Unif}[0, 1]$ i $h(x) = [\cos(50x) + \sin(20x)]^2$.

Već smo zaključili, što je n veći, širina pouzdanog intervala je manja, odnosno greška je manja. Grafički to izgleda ovako:

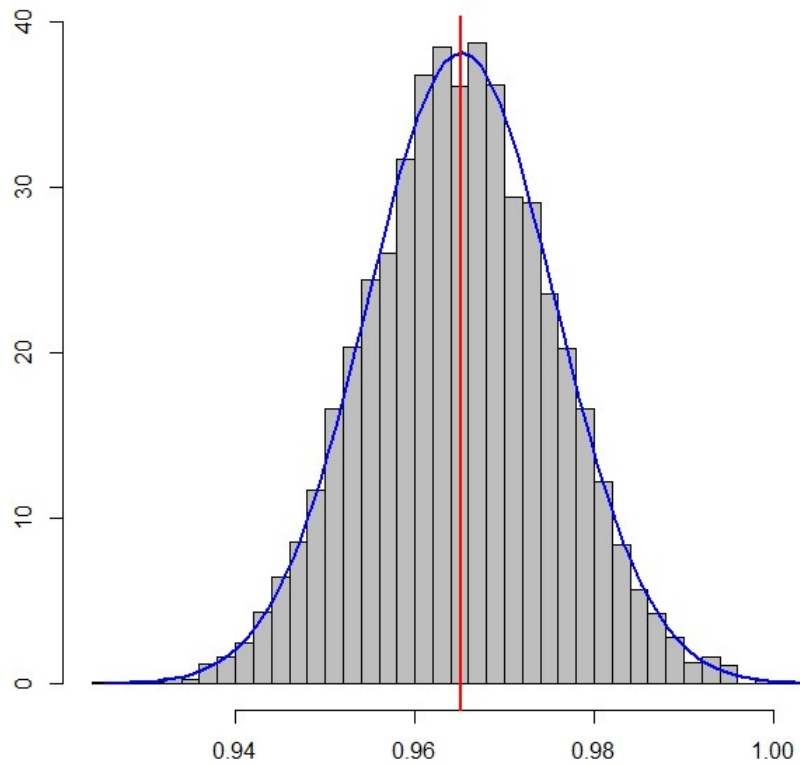


Slika 2.3: Ocjena greške u ovisnosti o n

Vrijednost procjenitelja \tilde{I}_j , za $j = 1, \dots, n$, na slici 2.3 prikazana je plavom linijom, dok crna linija predstavlja pripadni 95% pouzdani interval za I . Stvarna vrijednost integrala I prikazana je crvenom isprekidanom linijom. Dakle, širina intervala pouzdanosti teži u 0 kada $n \rightarrow \infty$, odnosno, veći n daje bolju procjenu vrijednosti integrala I ilustriranom Monte Carlo metodom.

Prikažimo još grafički asimptotsku normalnost procjenitelja \tilde{I}_n , odnosno da vrijedi $\tilde{I}_n \stackrel{D}{\approx} N(I, \sigma^2/n)$, kada $n \rightarrow \infty$, te tako pokažimo da i u ovom primjeru vrijedi Centralni granični teorem. Ideja je zapravo izračunati k realizacija procjenitelja \tilde{I}_{n_j} , $j = 1, \dots, k$, te pokazati da su normalno distribuirani, s matematičkim očekivanjem I i varijancom σ^2/n .

Odaberimo primjerice $k = 10^4$. Na slici 2.4 prikazan je histogram izračunatih k procjenitelja, zajedno s funkcijom gustoće normalne razdiobe $N \sim (I, \sigma^2/n)$ te stvarnom vrijednosti integrala I (crvena linija).



Slika 2.4: Histogram 10^4 realizacija procjenitelja \tilde{I}_n

Pritom, primijetimo sljedeće: stvarnu vrijednost standardne devijacije σ od $h(X)$, za $X \sim \text{Unif}[0, 1]$, možemo izračunati iz

$$\begin{aligned} \text{Var}[h(X)] &= \mathbb{E}[h(X)^2] - (\mathbb{E}[h(X)])^2 \\ &= \mathbb{E}[h(X)^2] - I^2, \end{aligned}$$

gdje za $\mathbb{E}[h(X)^2]$ vrijedi:

$$\mathbb{E}[h(X)^2] = \int_0^1 h(x)^2 \underbrace{f(x)}_{=1} dx = \int_0^1 h(x)^2 dx.$$

□

Ovaj uvodni primjer od iznimne je važnosti za razumijevanje cijelog procesa integracije metodom Monte Carlo. Naime, na jednom konkretnom primjeru vidjeli smo osnovnu ideju ovakve integracije - izračun vrijednosti integrala svodimo na računanje matematičkog očekivanja. Pokazali smo na koji način određujemo procjenitelj \tilde{I}_n za I koristeći Jaki zakon velikih brojeva, te kako možemo ocijeniti grešku procjene koristeći Centralni granični teorem.

Točnije, vidjeli smo da za neku funkciju $h : \mathbb{R} \rightarrow \mathbb{R}$ (neprekidnu) i za neprekidnu slučajnu varijablu X s pripadnom funkcijom gustoće f_X vrijedi

$$\mathbb{E}[h(X)] = \int_{\mathbb{R}} h(x)f_X(x)dx. \quad (2.3)$$

Uočimo da, ako na prije opisani način znamo izračunati $\mathbb{E}[h(X)]$, istu ideju možemo iskoristiti i za izračun $\mathbb{E}[X]$ - u tom slučaju funkcija h je funkcija identiteta, tj. $h(x) = x$. Također, u posebnom slučaju moguće je primjerice procijeniti vjerojatnost $\mathbb{P}(X \in [a, b])$, za $a, b \in \mathbb{R}$, $a < b$. Naime, znamo da vrijedi $\mathbb{P}(X \in [a, b]) = \int_a^b f_X(x)dx$, pa možemo iskoristiti jednakost (2.3) za indikatorsku funkciju $h(x) = \mathbb{1}_{(a,b)}(x)$.

Kasnije ćemo pokazati kako se opisana metoda lagano može poopćiti i primijeniti u višedimenzionalnim problemima, odnosno za računanje integrala funkcija više varijabli. Dok se kod numeričke integracije situacija bitno zakomplicira za dimenzije $d > 1$ (u slučaju $d = 1$ numerička integracija vrlo je efikasna), integracija metodom Monte Carlo u slučaju većih dimenzija i dalje "radi". Štoviše, greška procjene, tj. širina pouzdanog intervala ovisi o faktoru $\frac{1}{\sqrt{n}}$, koji ne ovisi o dimenziji d .

2.2 Opseg i struktura problema

Osnovna tema ovog rada je integracija uz pomoć Monte Carlo metode. Međutim, u tom procesu integracije posebno se ističu dvije značajnije podteme:

- generiranje slučajnih varijabli, odnosno simulacija
- tzv. uzorkovanje po važnosti (*Importance Sampling*).

U uvodnom primjeru trebali smo generirati n slučajnih varijabli iz uniformne distribucije $\text{Unif}[a, b]$, što je u programskom jeziku R poprilično jednostavno, koristeći funkciju `runif`. Općenito, potrebno je znati simulirati niz nezavisnih slučajnih varijabli iz različitih distribucija. Osim za uniformnu distribuciju, u R -u za mnoge druge distribucije također

već postoje pripadne funkcije za generiranje slučajnih varijabli. No, nameće se pitanje što u slučaju kada je potrebno generirati slučajni uzorak iz distribucije za koju ne postoji takva odgovarajuća funkcija implementirana u R -u. Postoji nekoliko metoda za generiranje slučajnih varijabli u tom slučaju, kao što su primjerice metoda inverza te metoda prihvatanja i odbacivanja (*Accept-Reject metoda*), koje ćemo kasnije spomenuti.

Osim tim pitanjem, dodatno ćemo se fokusirati i na pitanje izbora distribucije slučajne varijable X i funkcije h . Naime, taj izbor nije jedinstven, a o njemu ovisi varijanca Monte Carlo procjenitelja. Naravno, želimo da varijanca procjenitelja bude što manja pa nam je cilj odabrati takvu distribuciju od X i funkciju h koje će maksimalno smanjiti varijancu. Taj problem nejedinstvenosti izbora h i X najbolje ćemo opisati jednostavnim primjerom:

Primjer 2.2.1. Integral $I = \int_{-\infty}^{+\infty} e^{-x^2} x^2 dx$ prikažimo kao $\mathbb{E}[h(X)]$, za neku neprekidnu slučajnu varijablu X i funkciju $h: \mathbb{R} \rightarrow \mathbb{R}$.

Vrijedi sljedeće:

$$I = \int_{-\infty}^{+\infty} e^{-x^2} x^2 dx = \int_{-\infty}^{+\infty} \underbrace{\sqrt{2\pi} e^{-x^2/2} x^2}_{h_1(x)} \underbrace{\frac{1}{\sqrt{2\pi}}}_{f_{X_1}(x)} e^{-x^2/2} dx,$$

odnosno

$$I = \int_{-\infty}^{+\infty} h_1(x) f_{X_1}(x) dx = \mathbb{E}[h_1(X_1)],$$

za $h_1(x) = \sqrt{2\pi} e^{-x^2/2} x^2$ te $X_1 \sim N(0, 1)$, gdje smo iskoristili činjenicu da je funkcija gustoće f_X slučajne varijable $X \sim N(\mu, \sigma^2)$ oblika

$$f_X(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \quad x \in \mathbb{R}.$$

No, primijetimo da koristeći zamjenu varijabli možemo dobiti i sljedeće jednakosti:

$$\begin{aligned} I &= \int_{-\infty}^{+\infty} e^{-x^2} x^2 dx = \left[x = \frac{z}{\sqrt{2}} \quad dx = \frac{dz}{\sqrt{2}} \right] = \int_{-\infty}^{+\infty} e^{-z^2/2} \frac{z^2}{2} \frac{1}{\sqrt{2}} dz \\ &= \int_{-\infty}^{+\infty} \underbrace{\sqrt{\pi} \frac{z^2}{2}}_{h_2(z)} \underbrace{\frac{1}{\sqrt{2\pi}}}_{f_Z(z)} e^{-z^2/2} dz. \end{aligned}$$

Iz ovog očito slijedi da je

$$I = \int_{-\infty}^{+\infty} h_2(z) f_Z(z) dz = \mathbb{E}[h_2(Z)],$$

za $h_2(x) = \sqrt{\pi} \frac{z^2}{2}$ i $Z \sim N(0, 1)$, čime smo pokazali da odabir funkcije h nije jedinstven.

Pogledajmo još jedan raspis integrala I . Ideja je taj integral "na silu" zapisati u obliku

$$\int_{-\infty}^{+\infty} \frac{1}{\sigma \sqrt{2\pi}} e^{-x^2/(2\sigma^2)} h_3(x) dx.$$

Uočimo da vrijedi $e^{-x^2/(2\sigma^2)} = e^{-x^2}$ za $\sigma^2 = \frac{1}{2}$, tj. za $\sigma = \frac{1}{\sqrt{2}}$, iz čega slijedi sljedeći raspis:

$$I = \int_{-\infty}^{+\infty} \frac{1}{\sigma \sqrt{2\pi}} e^{-x^2/(2\sigma^2)} \underbrace{\sigma \sqrt{2\pi} x^2}_{h_3(x)} dx, \quad \text{gdje je } \sigma = \frac{1}{\sqrt{2}}.$$

Konačno slijedi

$$I = \int_{-\infty}^{+\infty} h_3(x) f_{X_3}(x) dx = \mathbb{E}[h_3(X_3)],$$

za $h_3(x) = \sqrt{\pi} x^2$ i $X_3 \sim N\left(0, \frac{1}{2}\right)$. Dakle, ni odabir distribucije slučajne varijable X nije jedinstven. Upravo taj izbor distribucije slučajne varijable X i funkcije h , koji za cilj ima smanjivanje varijance procjenitelja, opisan je u nastavku metodom uzorkovanja po važnosti.

□

Poglavlje 3

Generiranje slučajnih varijabli

U prethodnom poglavlju vidjeli smo kako je sama ideja integracije metodom Monte Carlo utemeljena na simuliranju slučajnih varijabli. Ova tema poprilično je opširna i složena jer, osim već postojećih algoritama za simuliranje slučajnih varijabli, u novije vrijeme razvijene su i brojne nadograđene verzije poznatih algoritama koje za cilj imaju postizanje još optimalnijih performansi. Upravo zato, u nastavku ćemo navesti neke od najpoznatijih metoda generiranja slučajnih varijabli te ih ukratko opisati, dok se u [1] može pronaći iscrpna analiza metoda za generiranje slučajnih varijabli.

Osnovna ideja ovog poglavlja opisati je na koji način možemo generirati slučajne varijable iz distribucija za koje nemamo odgovarajuće funkcije za generiranje implementirane u standardnim statističkim programskim paketima. Vidjet ćemo da se takve simulacije iz nestandardnih distribucija zapravo svode na generiranje uniformno distribuiranih slučajnih varijabli, pri čemu se ističu uniformne slučajne varijable na intervalu $[0, 1]$ koje pružaju osnovni vjerojatnosni prikaz slučajnosti. Možemo se zapitati je li uopće moguće računalnim metodama producirati deterministički niz vrijednosti koji oponaša niz nezavisnih i jednako distribuiranih uniformnih slučajnih varijabli na intervalu $[0, 1]$, no to je već pitanje za neku filozofsku raspravu. Mi ćemo se fokusirati na metode koje koriste deterministički proces za generiranje niza slučajnih varijabli, gdje se slučajnost generiranog niza očituje u pojedinačnom uzorku X_1, X_2, \dots, X_n kada $n \rightarrow \infty$.

Kako bismo izbjegli spomenuta pitanja filozofske prirode, uvedimo sljedeću definiciju:

Definicija 3.0.1. ¹ *Uniformni generator pseudo-slučajnih brojeva je algoritam koji, počevši od inicijalne vrijednosti u_0 i funkcije transformacije D , producira niz $(u_i) = (D^i(u_0))$ vrijednosti iz intervala $[0, 1]$. Za sve n , vrijednosti (u_1, u_2, \dots, u_n) oponašaju ponašanje uzorka (V_1, V_2, \dots, V_n) nezavisnih i jednako distribuiranih (uniformnih) slučajnih varijabli.*

¹Definicija je preuzeta iz [1].

Provjera valjanosti algoritma svodi se na testiranje statističke hipoteze

$$H_0 : U_1, U_2, \dots, U_n \text{ su n.j.d. slučajne varijable uniformne na } [0, 1],$$

za niz U_1, U_2, \dots, U_n , pa je gornja definicija više funkcionalne prirode.

3.1 Metoda generaliziranog inverza funkcije distribucije

Za početak, neka je $([0, 1], \mathcal{B}([0, 1]), \mathbb{P}_X)$ inducirani vjerojatnosni prostor slučajnom varijablom $X \sim \text{Unif}[0, 1]$, gdje je $\mathcal{B}([0, 1])$ σ -algebra Borelovih skupova na intervalu $[0, 1]$. Slučajne varijable X na tom vjerojatnosnom prostoru su funkcije $X : [0, 1] \rightarrow \mathbb{R}$.

Definicija i Lema navedene ispod preuzete su iz [1].

Definicija 3.1.1. Za funkciju distribucije F na \mathbb{R} , generalizirani inverz F^- od F je funkcija $F^- : [0, 1] \rightarrow \mathbb{R}$ definirana ² s

$$F^-(u) = \inf\{x : F(x) \geq u\}, \quad u \in [0, 1]. \quad (3.1)$$

Ovdje napomenimo da u slučaju kad je funkcija F bijekcija možemo definirati njen inverz $F^{-1} : [0, 1] \rightarrow \mathbb{R}$, te tada vrijedi $F^- = F^{-1}$.

Sljedeća lema pokazuje kako je bilo koju slučajnu varijablu moguće prikazati kao inverz uniformne slučajne varijable.

Lema 3.1.2. Ako je $U \sim \text{Unif}[0, 1]$, tada slučajna varijabla $F^-(U)$ ima distribuciju F (pišemo $F^-(U) \sim F$).

Dokaz. Za sve $u \in [0, 1]$ i za sve $x \in F^-([0, 1])$, generalizirani inverz zadovoljava jednakosti

$$F(F^-(u)) \geq u \quad \text{i} \quad F^-(F(x)) \leq x.$$

Stoga vrijedi

$$\{(u, x) : F^-(u) \leq x\} = \{(u, x) : F(x) \geq u\},$$

odnosno

$$\mathbb{P}(F^-(U) \leq x) = \mathbb{P}(U \leq F(x)) = F(x).$$

□

²Primijetimo da generalizirani inverz F^- uvijek možemo definirati, neovisno o tom je li funkcija F bijekcija.

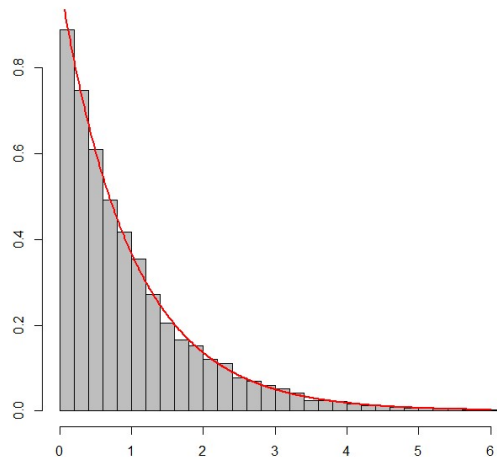
Dakle, ako želimo generirati slučajnu varijablu $X \sim F$, dovoljno je generirati uniformnu slučajnu varijablu $U \sim \text{Unif}[0, 1]$ te izvršiti transformaciju $x = F^{-}(u)$. Drugim riječima, ako znamo generirati iz uniformne distribucije na intervalu $[0, 1]$ (što znamo!), jednostavnom metodom inverza znamo generirati i iz distribucije F . Ilustrirajmo to na sljedećem jednostavnom primjeru.

Primjer 3.1.3 (Generiranje eksponencijalnih slučajnih varijabli). ³ Neka je X eksponencijalna slučajna varijabla s parametrom $\lambda = 1$, tj. $X \sim \text{Exp}(1)$. Tada je pripadna funkcija distribucije F dana s

$$F(x) = 1 - e^{-x}, \quad x \geq 0.$$

Ako jednadžbu $u = 1 - e^{-x}$ riješimo po x , dobijemo $x = -\log(1 - u)$, odnosno izraz za generalizirani inverz $F^{-}(u)$ (ovdje je \log oznaka za prirodni logaritam \ln). Zaključujemo, ako je $U \sim \text{Unif}[0, 1]$, tada slučajna varijabla $X = -\log(U)$ ima eksponencijalnu razdiobu s parametrom $\lambda = 1$ (jer su slučajne varijable U i $1 - U$ jednako distribuirane, tj. $1 - U \sim \text{Unif}[0, 1]$).

Sada u programskom paketu R jednostavno možemo generirati niz od n nezavisnih i jednako distribuiranih slučajnih varijabli $X_i \sim \text{Exp}(1)$ koristeći opisanu metodu inverza. Na slici 3.1 prikazan je histogram generiranog uzorka iz $\text{Exp}(1)$ razdiobe, zajedno s pripadnom funkcijom gustoće.



Slika 3.1: Histogram generiranog uzorka duljine $n = 10^4$ iz $\text{Exp}(1)$ razdiobe, zajedno s pripadnom funkcijom gustoće

□

³Ovaj primjer preuzet je iz [1].

Ovakav način generiranja slučajnih varijabli koristan je za razumijevanje ostalih metoda generiranja, no u praksi se metoda inverza temeljena na Lemi 3.1.2 rijetko koristi. Naime, za implementaciju ove metode potrebno je znati eksplicitan oblik funkcije distribucije F , što u većini problema nije slučaj. Neke druge metode, kao primjerice metoda prihvaćanja i odbacivanja, primjenjive su u općenitijim slučajevima (npr. kod generiranja u dimenzijama većim od 1).

Dodatno, napomenimo da kada postoji relativno jednostavan način na koji do distribucije F možemo doći koristeći neku drugu distribuciju, često možemo razviti algoritam za generiranje iz F na temelju veze tih dviju distribucija. Stoga, postoji više različitih metoda za generiranje slučajnih varijabli iz distribucija koje nisu uniformne, a koje su alternativa spomenutoj metodi inverza. Međutim, takve metode dosta su specifične jer ovise o svojstvima dane distribucije, te se teško mogu generalizirati. Više primjera takvih općih metoda transformacije navedeno je i pobliže opisano u [1]. Ovdje ćemo istaknuti poznatu Box-Muller metodu za generiranje iz standardne normalne razdiobe $N(0, 1)$, koja je opisana na sljedeći način:

Ako su U_1 i U_2 nezavisne slučajne varijable uniformno distribuirane na $[0, 1]$, tada su slučajne varijable X_1 i X_2 definirane s

$$X_1 = \sqrt{-2 \log(U_1)} \cos(2\pi U_2), \quad X_2 = \sqrt{-2 \log(U_1)} \sin(2\pi U_2)$$

nezavisne slučajne varijable iz standardne normalne razdiobe $N(0, 1)$.

Box-Muller algoritam

1. generiraj nezavisne $U_1, U_2 \sim \text{Unif}[0, 1]$;
2. definiraj $x_1 = \sqrt{-2 \log(u_1)} \cos(2\pi u_2)$, $x_2 = \sqrt{-2 \log(u_1)} \sin(2\pi u_2)$;
3. uzmi x_1 i x_2 kao dvije nezavisne simulacije iz $N(0, 1)$.

3.2 Metoda prihvaćanja i odbacivanja

Pretpostavimo da želimo generirati slučajnu varijablu iz distribucije F na \mathbb{R} koju je teško invertirati. Postoji puno takvih distribucija, a u nekim slučajevima funkciju distribucije uopće nije moguće prikazati u eksplicitnom obliku pa ne možemo primijeniti spomenutu metodu inverza. U takvim okolnostima koristimo drugačiju metodu generiranja - metodu prihvaćanja i odbacivanja (*Accept-Reject metodu*), za čiju je primjenu potrebno poznavati samo oblik funkcije gustoće f . Osnovna ideja ove metode zapravo je generiranje iz neke druge gustoće g , iz koje je jednostavnije generirati nego iz f , a odgovarajući Accept-Reject algoritam zasnovan je na jednostavnoj uniformnoj distribuciji. Ovaj odjeljak za cilj ima dati kratak pregled Accept-Reject algoritma⁴.

Sljedeća lema osnovni je temelj za Accept-Reject algoritam.

Lema 3.2.1.⁵ *Pretpostavimo da za funkcije gustoće f i g na \mathbb{R} te $M \geq 1$ vrijedi*

$$f(x) \leq Mg(x), \quad x \in \mathbb{R}.$$

Neka su $X \sim g$ i $Y \sim f$, te $U \sim \text{Unif}(0, 1)$ nezavisna od X . Tada vrijedi

$$\mathbb{P}\left(U \leq \frac{f(X)}{Mg(X)}\right) = \frac{1}{M}, \quad (3.2)$$

te

$$\mathbb{P}\left(X \leq t \mid U \leq \frac{f(X)}{Mg(X)}\right) = \mathbb{P}(Y \leq t), \quad t \in \mathbb{R}. \quad (3.3)$$

Odnosno, uvjetno na $U \leq \frac{f(X)}{Mg(X)}$, X ima istu distribuciju kao Y .

Dokaz. Neka je $t \in \mathbb{R}$ proizvoljan i $E = \{x \in \mathbb{R} : g(x) = 0\}$. Po pretpostavci vrijedi da su X i U nezavisne pa slijedi

$$\begin{aligned} \mathbb{P}\left(X \leq t \mid U \leq \frac{f(X)}{Mg(X)}\right) &= \int_{-\infty}^t g(x) \mathbb{1}_E(x) \int_0^{\frac{f(x)}{Mg(x)}} du dx = \int_{-\infty}^t \frac{f(x)}{Mg(x)} g(x) \mathbb{1}_E(x) dx \\ &= \frac{1}{M} \int_{-\infty}^t f(x) \mathbb{1}_E(x) dx = \frac{1}{M} \int_{-\infty}^t f(x) dx = \frac{1}{M} \mathbb{P}(Y \leq t). \end{aligned}$$

⁴Razrada i analiza same Accept-Reject metode posebna je tema, koja je detaljnije opisana u [1].

⁵Lema 3.2.1. preuzeta je iz [7].

Zbog $\mathbb{P}(X \in E^c) = 0$ te $\frac{f(x)}{Mg(x)} \leq 1$ za $x \in E$ vrijedi prva jednakost. Četvrta jednakost slijedi iz $f(x) = 0$ za sve $X \in E^c$. Sada se lako pokaže da vrijedi (3.2) puštanjem $t \rightarrow \infty$, iz čega slijedi i (3.3). \square

Accept-Reject algoritam opisan je sljedećim pseudokodom:

Accept-Reject algoritam

1. generiraj $X \sim g$ i $U \sim \text{Unif}[0, 1]$;
2. ako je $U \leq \frac{f(X)}{Mg(X)}$ idi na 3., inače idi na 1.;
3. vrati $Y := X$.

Primijetimo da je $\mathbb{P}\left(U \leq \frac{f(X)}{Mg(X)}\right)$ zapravo vjerojatnost prihvatanja generirane slučajne varijable $X \sim g$ koja po Lemi 3.2.1 iznosi $\frac{1}{M}$, odnosno očekivani broj pokušaja do prvog prihvatanja u algoritmu je M .

Dakle, za danu funkciju gustoće g želimo odabrati M najmanji mogući, ali takav da zadovoljava jednakost $f(x) \leq Mg(x)$, $x \in \mathbb{R}$, jer je u tom slučaju algoritam najefikasniji.

Također, uočimo da je uvjet $f(x) \leq Mg(x)$, $x \in \mathbb{R}$ ekvivalentan sa sljedeća dva uvjeta:

- (1) $g(x) = 0 \Rightarrow f(x) = 0$, za sve $x \in \mathbb{R}$
- (2) $f(x)/g(x) \leq M$ za sve $x \in \mathbb{R}$ takve da $g(x) > 0$.

Napomena 3.2.2. Bitno je primijetiti da ovim algoritmom možemo generirati iz bilo koje gustoće f koju znamo samo do na multiplikativnu konstantu, tj. ako je $f = cf^*$, za nepoznatu konstantu $c > 0$ i poznatu funkciju gustoće f^* . Naime, ako za $M^* > 0$ vrijedi $f^* \leq M^*g$, tada vrijedi $f \leq Mg$ za $M := cM^*$ te

$$\frac{f(x)}{Mg(x)} = \frac{cf^*(x)}{cM^*g(x)} = \frac{f^*(x)}{M^*g(x)}, \text{ za } x \in \mathbb{R}.$$

Uočimo da u Accept-Reject algoritmu neke realizacije simulirane iz g odbacujemo ("reject") pa je jedna od kritika ovog algoritma da generira "beskorisne" slučajne varijable

pri odbacivanju. Konačno, napomenimo da se Accept-Reject algoritam može primijeniti i kada su f i g gustoće na \mathbb{R}^d na analogan način.

Sljedećim primjerom, preuzetim iz [2], pokažimo primjenu ove metode.

Primjer 3.2.3. *Iskoristimo Accept-Reject algoritam za generiranje slučajnih varijabli iz Beta(2.7, 6.3) razdiobe. Za $\alpha, \beta > 0$ gustoća beta razdiobe Beta(α, β) je*

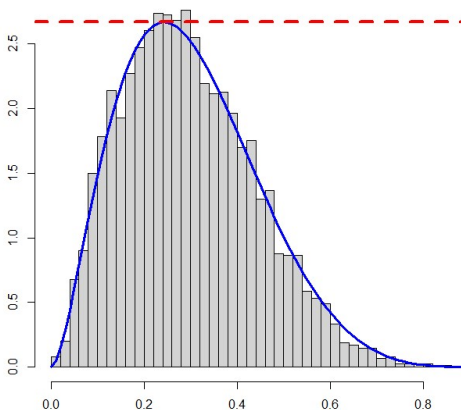
$$f(x) = \begin{cases} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}, & 0 < x < 1 \\ 0, & \text{inače.} \end{cases}$$

Uočimo da za $\alpha > 1$ i $\beta > 1$ vrijedi $f(x) \leq \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}$ za sve $x \in (0, 1)$, tj. funkcija f je ograničena. U tom slučaju, za izbor gustoće g takve da zadovoljava uvjet $f(x) \leq Mg(x)$ prirodno se nameće funkcija $g(x) = \mathbb{1}_{(0,1)}(x)$, odnosno funkcija gustoća uniformne slučajne varijable $X \sim \text{Unif}[0, 1]$. Za gornju ogradu M očito možemo uzeti $M = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}$. Međutim, kako efikasnost Accept-Reject algoritma ovisi o vjerojatnosti prihvatanja $1/M$ (što je ta vjerojatnost veća, manje je "beskorisnih" generiranja iz g), želimo M što manji.

Odredimo takav najmanji M za $\alpha = 2.7$ i $\beta = 6.3$.

Tada je f oblika $f(x) = c \cdot x^{1.7} (1-x)^{5.3}$, gdje je $c = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} = \frac{\Gamma(9)}{\Gamma(2.7)\Gamma(6.3)}$ konstanta. Funkciju f koja je neprekidna na intervalu $(0, 1)$ možemo proširiti na segment $[0, 1]$ s $f(0) = f(1) = 0$, pa znamo da f poprima minimum i maksimum na $[0, 1]$. Uz pomoć funkcije `optimize` u `R`-u (ili deriviranjem funkcije f) lako dobijemo da se maksimum postiže u 0.2428608 i iznosi 2.669744 (približne vrijednosti).

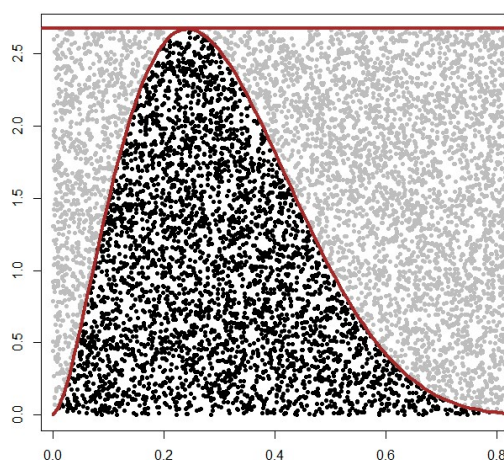
Dakle, pronašli smo najmanji M za koji vrijedi $\frac{f(x)}{g(x)} = f(x) \leq M$ i on iznosi $M \approx 2.67$.



Slika 3.2: Histogram generiranog uzorka duljine $n = 10^4$ iz Beta(2.7, 6.3) razdiobe, zajedno s pripadnom funkcijom gustoće f i funkcijom Mg

Na slici 3.2 prikazan je histogram uzorka duljine $n = 10^4$ iz beta razdiobe $\text{Beta}(2.7, 6.3)$ koji je generiran u R-u uz pomoć *Accept-Reject* algoritma, odnosno histogram prihvaćenih generiranih varijabli $X_i \sim g$, zajedno s funkcijom gustoće f $\text{Beta}(2.7, 6.3)$ razdiobe i funkcijom Mg .

Primijetimo da, uvjetno na $X = x$, slučajna varijabla $U^* := UMg(x)$ ima uniformnu distribuciju na $[0, Mg(x)]$ pa umjesto $U \sim \text{Unif}[0, 1]$, u 1. koraku algoritma možemo generirati $U^* \sim \text{Unif}[0, Mg(x)]$, te tada generirani $X = x$ prihvaćamo ako je $U^* \leq f(x)$.



Slika 3.3: Grafički prikaz prihvaćenih (crnih) i odbačenih (sivih) točaka generiranih *Accept-Reject* algoritmom

Na slici 3.3 prikazane su točke $(X_i, U_iMg(X_i))$, $i = 1, \dots, n$, za $X_i \sim g$ i $U_i \sim \text{Unif}[0, 1]$. Točke za koje vrijedi $U_iMg(X_i) \leq f(X_i)$ su točke prihvaćene u *Accept-Reject* algoritmu (one koje se nalaze ispod grafa funkcije gustoće f), a točke iznad grafa od f su odbačene točke koje ne zadovoljavaju taj uvjet.

U ovom konkretnom primjeru vjerojatnost prihvaćanja generirane slučajne varijable $X_i \sim g$ je $1/M \approx 0.3745677$, a očekivani broj pokušaja do prvog prihvaćanja u algoritmu je $M \approx 2.67$. □

Poglavlje 4

Integracija metodom Monte Carlo

Cilj ovog poglavlja teoretski je opisati i sažeti ideju integracije metodom Monte Carlo. Za razumijevanje osnovnog koncepta Monte Carlo metode motivacijski primjer iz Poglavlja 2 od iznimne je koristi i pomoći. Opisat ćemo način na koji računalno generirane slučajne varijable možemo iskoristiti za izračun aproksimativne vrijednosti jednostrukih i višestrukih integrala. U prethodnom poglavlju rekli smo nešto više o samom generiranju slučajnih varijabli, na kojem se zasniva cijeli proces integracije ovom metodom.

4.1 Monte Carlo metoda

Primarni cilj ovog odjeljka je pružiti detaljniju teoretsku analizu problema i opisati opći pristup rješavanju istog, da bismo zatim mogli pokazati primjenu ove metode na nekim praktičnim primjerima. Dva teorema iz teorije vjerojatnosti, Jaki zakon velikih brojeva i Centralni granični teorem, odnosno njihovi rezultati, nužni su u analizi koja slijedi.

Neka je $h : \mathbb{R}^d \rightarrow \mathbb{R}$ izmjeriva funkcija i X d -dimenzionalan neprekidan slučajni vektor s pripadnom funkcijom gustoće $f_X : \mathbb{R}^d \rightarrow [0, 1]$. Pretpostavimo da slučajna varijabla $h(X)$ ima matematičko očekivanje, odnosno da vrijedi $\mathbb{E}[|h(X)|] < \infty$. Tada znamo da za $\mathbb{E}[h(X)]$ vrijedi

$$\mathbb{E}[h(X)] = \int_{\mathbb{R}^d} h(x)f_X(x) dx =: I, \quad I \in \mathbb{R}. \quad (4.1)$$

Kao što smo vidjeli u uvodnom primjeru, za izračun aproksimativne vrijednosti integrala u (4.1) prirodno se nameće korištenje simuliranog slučajnog uzorka X_1, X_2, \dots, X_n od n nezavisnih i jednako distribuiranih X_i (generiranih iz gustoće f_X , tj. $X_i \sim X$).

Definirajmo prvo

$$\tilde{I}_n = \frac{1}{n} \sum_{i=1}^n h(X_i). \quad (4.2)$$

Tada, po Jakom zakonu velikih brojeva vrijedi

$$\frac{h(X_1) + h(X_2) + \dots + h(X_n)}{n} \xrightarrow{g.s.} \mathbb{E}[h(X_i)],$$

odnosno \tilde{I}_n gotovo sigurno konvergira prema $\mathbb{E}[h(X)]$, tj. prema integralu I . Kažemo da je \tilde{I}_n Monte Carlo procjenitelj za I .

Nadalje, kako slučajna varijabla $|h(X)|$ ima konačno očekivanje, moguće je procijeniti brzinu konvergencije nepristanog i konzistentnog procjenitelja \tilde{I}_n , budući da varijancu $\text{Var}[\tilde{I}_n]$ za koju vrijedi

$$\begin{aligned} \text{Var}[\tilde{I}_n] &= \text{Var}\left[\frac{1}{n} \sum_{i=1}^n h(X_i)\right] = \text{Var}\left(\frac{1}{n}h(X_1) + \dots + \frac{1}{n}h(X_n)\right) = \frac{1}{n^2} \text{Var}[h(X_1)] + \dots + \frac{1}{n^2} \text{Var}[h(X_n)] \\ &= \frac{1}{n} \text{Var}[h(X_1)] = \frac{1}{n} \mathbb{E}\left[(h(X_1) - \mathbb{E}[h(X_1)])^2\right] = \frac{1}{n} \int_{\mathbb{R}^d} (h(x) - \mathbb{E}[h(X_1)])^2 f_X(x) dx \end{aligned}$$

možemo procijeniti s

$$v_n^2 = \frac{1}{n^2} \left(\sum_{i=1}^n (h(X_i))^2 - n\tilde{I}_n^2 \right) \quad (4.3)$$

uz pomoć uzoračke varijance S_n^2 . Naime, prvo uočimo da varijancu $\sigma^2 := \text{Var}[h(X_i)] = \text{Var}[h(X)]$ možemo jednostavno procijeniti uzoračkom varijancom S_n^2 :

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (h(X_i) - \tilde{I}_n)^2 = \frac{1}{n-1} \left(\sum_{i=1}^n (h(X_i))^2 - n\tilde{I}_n^2 \right). \quad (4.4)$$

Želimo ocijeniti grešku procjene integrala I koristeći Centralni granični teorem, po kojem slijedi

$$\sqrt{n} \cdot \frac{\tilde{I}_n - I}{\sigma} \xrightarrow{D} Z \sim N(0, 1), \quad n \rightarrow \infty.$$

Drugim riječima, niz $(\sqrt{n}(\tilde{I}_n - I)/\sigma)_{n \in \mathbb{N}}$ konvergira po distribuciji ka $N(0, 1)$. Kažemo da je \tilde{I}_n asimptotski normalan procjenitelj za I s asimptotskom standardnom devijacijom $\frac{\sigma}{\sqrt{n}}$,

tj. vrijedi $\tilde{I}_n \stackrel{D}{\approx} N\left(I, \frac{\sigma^2}{n}\right)$, kada $n \rightarrow \infty$.

Sada vidimo da je $\text{Var}[\tilde{I}_n]$ jednaka $\frac{\sigma^2}{n}$, iz čega direktno slijedi procjena varijance $\text{Var}[\tilde{I}_n]$ u (4.3), uz činjenicu da σ^2 možemo procijeniti uzoračkom varijancom kao u (4.4).

Također, ranije smo pokazali da je aproksimativni $(1 - \alpha)\%$ pouzdani interval, za $\alpha \in (0, 1)$, dan s

$$\left[\tilde{I}_n - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}, \tilde{I}_n + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \right],$$

pri čemu za procjenu standardne devijacije σ opet možemo iskoristiti (4.4).

Dakle, pokazali smo ono što je zapravo intuitivno jasno i očekivano - što je duljina niza n generiranih slučajnih varijabli veća, opisana Monte Carlo metoda daje bolju procjenu \tilde{I}_n vrijednosti integrala I , odnosno širina pouzdanog intervala je manja.

Dodatno, uočimo sljedeće:

Ako trebamo odrediti vrijednost nekog nepoznatog parametra $\theta \in \mathbb{R}$ kojeg je moguće prikazati u obliku

$$\theta = \mathbb{E}[h(X)],$$

za $h : \mathbb{R}^d \rightarrow \mathbb{R}$ izmjerivu funkciju i X d -dimenzionalan neprekidan slučajni vektor, tada po Jakom zakonu velikih brojeva na iznad opisan način možemo izračunati procjenu vrijednosti parametra θ koristeći Monte Carlo procjenitelj $\tilde{\theta}_n$ definiran s

$$\tilde{\theta}_n := \frac{1}{n} \sum_{i=1}^n h(X_i),$$

gdje je $h(X_1), h(X_2), \dots, h(X_n)$ niz nezavisnih i jednako distribuiranih slučajnih varijabli ($X_i \sim X$) i n velik, te ocijeniti grešku procjene koristeći Centralni granični teorem.

Primjer 4.1.1 (Funkcija distribucije Φ standardne normalne razdiobe). ¹ *Promotrimo funkciju distribucije slučajne varijable $X \sim N(0, 1)$ iz standardne normalne razdiobe*

$$\Phi(t) = \mathbb{P}(X \leq t) = \int_{-\infty}^t \frac{1}{\sqrt{2\pi}} e^{-y^2/2} dy.$$

¹Primjer je preuzet iz [1].

Uočimo da je tablicu vrijednosti funkcije distribucije standardne normalne razdiobe moguće izračunati koristeći metodu Monte Carlo, jer se gornji integral može prikazati kao matematičko očekivanje $\mathbb{E}[h(X)]$ za $h(x) = \mathbb{1}_{x \leq t}(x)$ i $X \sim N(0, 1)$, to jest

$$\Phi(t) = \int_{-\infty}^t \frac{1}{\sqrt{2\pi}} e^{-y^2/2} dy = \int_{-\infty}^t \underbrace{1}_{h(y)} \cdot \underbrace{\frac{1}{\sqrt{2\pi}} e^{-y^2/2}}_{f_{N(0,1)}(y)} dy.$$

Prema tome, definiramo Monte Carlo procjenitelj

$$\tilde{\Phi}(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{X_i \leq t}.$$

Vidjeli smo već kako se u R -u vrlo jednostavno mogu simulirati slučajne varijabli iz osnovnih distribucija, pa tako i iz $N(0, 1)$. Izračunajmo stoga neke realizacije definiranog Monte Carlo procjenitelja $\tilde{\Phi}(t)$.

$n \setminus t$	0.0	0.67	0.84	1.28	1.65	2.32	2.58	3.09	3.72
10^2	0.485	0.74	0.77	0.9	0.945	0.985	0.995	1	1
10^3	0.4925	0.7455	0.801	0.902	0.9425	0.9885	0.9955	0.9985	1
10^4	0.4962	0.7425	0.7941	0.9	0.9498	0.9896	0.995	0.999	0.9999
10^5	0.4995	0.7489	0.7993	0.9003	0.9498	0.9898	0.995	0.9989	0.9999
10^6	0.5001	0.7497	0.8	0.9002	0.9502	0.99	0.995	0.999	0.9999
10^7	0.5002	0.7499	0.8	0.9001	0.9501	0.99	0.995	0.999	0.9999
10^8	0.5	0.75	0.8	0.9	0.95	0.99	0.995	0.999	0.9999

Tablica 4.1: Vrijednosti nekih procjenitelja $\tilde{\Phi}(t)$, u ovisnosti o broju simuliranih slučajnih varijabli n i vrijednosti t

Tablica 4.1 prikazuje razvoj vrijednosti procjenitelja $\tilde{\Phi}(t)$ izračunatog metodom Monte Carlo u R -u, za odabrane t , u ovisnosti o n . Uočimo da simulacije za $n = 10^8$ u posljednjem retku daju najbolju procjenu vrijednosti funkcije distribucije $\Phi(t)$ za pripadni parametar t , tj. širina pouzdanog intervala se smanjuje kako n raste. \square

Sljedeći primjer poznata je ilustracija primjene Monte Carlo metode (spomenuto u [9]). Pokazuje na koji način je moguće procijeniti površinu kruga D radijusa $r = 1$, za koju znamo da vrijedi $P(D) = \pi$, $D \subseteq \mathbb{R}^2$.

Primjer 4.1.2. Za krug $D \subseteq \mathbb{R}^2$ radijusa $r = 1$ procijenimo površinu $P(D)$ koristeći Monte Carlo simulacije. Ideja je definirati vjerojatnost p na sljedeći način:

$$p := \mathbb{P}(X \in D) = \frac{P(D)}{P([-1, 1]^2)},$$

gdje je $X \sim \text{Unif}[-1, 1]^2$ slučajan vektor. Iz tog očito slijedi da je $p = \frac{\pi}{4}$.

Također, vrijedi:

$$p = \mathbb{P}(X \in D) = \mathbb{E}[\mathbb{1}_{\{X \in D\}}],$$

odnosno $p = \mathbb{E}[h(X)]$, za $h(x) = \mathbb{1}_D(x)$, $h : \mathbb{R}^2 \rightarrow \mathbb{R}$.

Sada je očito da vjerojatnost p možemo procijeniti koristeći Monte Carlo procjenitelj

$$\tilde{p}_n = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{X_i \in D\}}.$$

Dakle, generiramo niz od n nezavisnih i jednako distribuiranih slučajnih vektora X_1, X_2, \dots, X_n , $X_i \sim \text{Unif}[-1, 1]^2$. Primijetimo da ovako definiran procjenitelj \tilde{p}_n zapravo predstavlja postotak generiranih točaka, odnosno slučajnih vektora X_i koji "upadaju" u D .

Pritom, uočimo da je $\mathbb{1}_{\{X \in D\}}$ Bernoullijeva slučajna varijabla s vjerojatnošću uspjeha p , tj. vrijedi

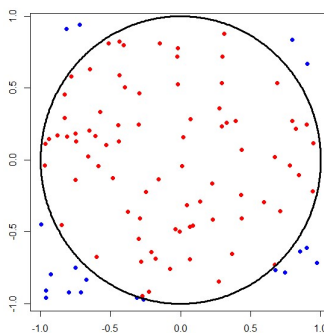
$$\mathbb{1}_{\{X \in D\}} \sim \begin{pmatrix} 0 & 1 \\ 1-p & p \end{pmatrix},$$

odnosno

$$\text{Var}[\mathbb{1}_{\{X \in D\}}] = p(1-p),$$

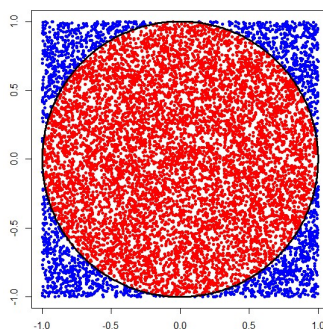
pa varijancu $\text{Var}[h(X)] = \text{Var}[\mathbb{1}_{\{X \in D\}}]$ možemo procijeniti s $\tilde{p}_n(1 - \tilde{p}_n)$.

Konačno, površinu π procijenimo kao $\pi = 4\tilde{p}_n$.



Slika 4.1: Grafički prikaz simuliranog uzorka duljine $n = 10^2$

Procjenu površine π sada lako možemo izračunati u R -u. Na slici 4.1 grafički je prikazan položaj simuliranih točaka na kvadratu $[-1, 1]^2$ za $n = 10^2$, dok je na slici 4.2 isto prikazano za $n = 10^4$. Ovaj ilustrativan prikaz pokazuje kako je za velike n Monte Carlo metodom moguće procijeniti površinu jediničnog kruga u \mathbb{R}^2 .



Slika 4.2: Grafički prikaz simuliranog uzorka duljine $n = 10^4$

□

Pokažimo još kratko na primjeru ispod kako se ova ideja integracije može poopćiti za rješavanje trostrukog integrala.

Primjer 4.1.3. Koristeći Monte Carlo procjenitelj procijenimo vrijednosti integrala I_1 i I_2 :

(a) Neka je zadan integral

$$I_1 = \int_0^1 \int_0^2 \int_0^3 xyz \, dx \, dy \, dz.$$

Uočimo sljedeće:

Neka su $X \sim \text{Unif}[0, 3]$, $Y \sim \text{Unif}[0, 2]$ i $Z \sim \text{Unif}[0, 1]$ nezavisne slučajne varijable. Tada za slučajan vektor (X, Y, Z) i funkciju $h(x, y, z) = 6xyz$ vrijedi

$$\begin{aligned} \mathbb{E}[h(X, Y, Z)] &= \int_0^1 \int_0^2 \int_0^3 h(x, y, z) f_{(X,Y,Z)}(x, y, z) \, dx \, dy \, dz \\ &\stackrel{\text{nez.}}{=} \int_0^1 \int_0^2 \int_0^3 h(x, y, z) f_X(x) f_Y(y) f_Z(z) \, dx \, dy \, dz, \\ &= \int_0^1 \int_0^2 \int_0^3 6 \cdot xyz \cdot \frac{1}{3} \cdot \frac{1}{2} \cdot \frac{1}{1} \, dx \, dy \, dz = I_1, \end{aligned}$$

gdje je $f_{(X,Y,Z)}$ funkcija gustoće slučajnog vektora (X, Y, Z) , a f_X, f_Y i f_Z funkcije gustoće uniformno distribuiranih komponenti X, Y i Z slučajnog vektora (X, Y, Z) . Druga jednakost slijedi iz nezavisnosti komponenti slučajnog vektora, jer u tom slučaju vrijedi

$$f_{(X,Y,Z)}(x, y, z) = f_X(x)f_Y(y)f_Z(z), \quad (x, y, z) \in \mathbb{R}^3.$$

Dakle, za $h(x, y, z) = 6xyz$ te $X \sim \text{Unif}[0, 3]$, $Y \sim \text{Unif}[0, 2]$ i $Z \sim \text{Unif}[0, 1]$ nezavisne slučajne varijable vrijedi

$$I_1 = \mathbb{E}[h(X, Y, Z)],$$

stoga vrijednost integrala I_1 možemo procijeniti koristeći Monte Carlo procjenitelj

$$\tilde{I}_{n_1} = \frac{1}{n} \sum_{i=1}^n h(X_i, Y_i, Z_i),$$

tako da za svaku komponentu slučajnog vektora (X, Y, Z) simuliramo uzorak od n nezavisnih i jednako distribuiranih slučajnih varijabli iz odgovarajuće distribucije.

U tablici 4.2 prikazane su neke realizacije procjenitelja \tilde{I}_{n_1} izračunate u R-u u ovisnosti o duljini uzoraka n , zajedno s pripadnim 95% pouzdanim intervalima. Stvarna vrijednost ovog integrala je $I_1 = \frac{9}{2}$.

n	\tilde{I}_{n_1}	95% pouzdani interval za I_1
10^4	4.429479	[4.327818, 4.531140]
10^6	4.508728	[4.498377, 4.519079]
10^8	4.500044	[4.499012, 4.501076]

Tablica 4.2: Vrijednosti nekih procjenitelja \tilde{I}_{n_1} za I_1 i pripadni 95% pouzdani intervali, u ovisnosti o n

(b) Isti postupak ponovimo i za procjenu integrala

$$I_2 = \int_{-1}^1 \int_0^1 \int_{-3}^2 (\sin(x)y^2 - z) dx dy dz.$$

U ovom slučaju vrijedi

$$I_2 = \mathbb{E}[h(X, Y, Z)],$$

za $h(x, y, z) = 10(\sin(x)y^2 - z)$ te $X \sim \text{Unif}[-3, 2]$, $Y \sim \text{Unif}[0, 1]$ i $Z \sim \text{Unif}[-1, 1]$. Stvarna vrijednost integrala I_2 iznosi $I_2 \approx -0.38256$, a u tablici 4.3 prikazane su neke realizacije procjenitelja \tilde{I}_{n_2} , zajedno s pripadnim 95% pouzdanim intervalima.

n	\tilde{I}_{n_2}	95% pouzdani interval za I_2	širina intervala pouzdanosti
10^4	-0.2240785	$[-0.35492612, -0.09323092]$	0.2616952
10^6	-0.378237	$[-0.3912608, -0.3652131]$	0.02604771
10^8	-0.3819667	$[-0.3832701, -0.3806632]$	0.002606904

Tablica 4.3: Vrijednosti nekih procjenitelja \tilde{I}_{n_2} za I_2 i pripadni 95% pouzdani intervali, u ovisnosti o n

Primijetimo da je u (b) tek za $n = 10^8$ dobivena procjena s točne dvije značajne znamenke jer procjenjujemo malu vrijednost, pa je u takvim slučajevima često potreban jako velik n . Naravno, što je n veći, vrijeme potrebno za generaciju slučajnih varijabli i izračun se produljuje.

□

Napomena 4.1.4. Uočimo da, koristeći zaključke iz prethodna dva primjera, ovom metodom između ostalog znamo odrediti i volumen tijela u \mathbb{R}^3 (npr. volumen jedinične kugle).

4.2 Uzorkovanje po važnosti

Nakon nekoliko primjera integracije metodom Monte Carlo koje smo dosad vidjeli, mogli bismo zaključiti da takvom metodom možemo lagano doći do vrijednosti bilo kojeg integrala, odnosno njegove procjene. U pravilu to zaista vrijedi te ovom jednostavnom metodom možemo dobiti dobre procjene vrijednosti integrala i ocijeniti njihovu grešku. Međutim, kao što smo već vidjeli u Primjeru 2.2.1, bitno je uočiti da izbor funkcije h i distribucije od X nije jedinstven, a upravo o tom izboru ovisi varijanca procjenitelja σ , odnosno kvaliteta same procjene. Nije teško zaključiti da je cilj odabrati h i X takve da varijanca σ bude što manja. Upravom tim problemom bavi se metoda uzorkovanja po važnosti (*Importance Sampling*), najraširenija metoda za smanjivanje varijance procjenitelja.

Osvrnimo se još jednom na prikaz integrala preko matematičkog očekivanja:

$$\mathbb{E}[h(X)] = \int_{\mathbb{R}^d} h(x)f_X(x) dx, \quad (4.5)$$

te primijetimo da je ovaj integral moguće na beskonačno mnogo načina prikazati kao očekivanje od $h(X)$, za različite h i X . Stoga optimalne h i X , odnosno one za koje je varijanca pripadnog procjenitelja minimalna, ne možemo tražiti izračunom svih mogućih

procjenitelja te usporedbom njihove varijance.

Osnovna ideja ove metode je očekivanje $\mathbb{E}[h(X)]$ prikazati na sljedeći alternativan način:

$$\mathbb{E}[h(X)] = \int_S h(x)f(x) dx = \int_S h(x) \frac{f(x)}{g(x)} g(x) dx, \quad (4.6)$$

gdje su f i g funkcije gustoće na nekom skupu $S \subseteq \mathbb{R}^d$ za koje vrijedi $f(x) \neq 0 \Rightarrow g(x) \neq 0$, odnosno $g(x) = 0 \Rightarrow f(x) = 0$, $x \in \mathbb{R}^d$. Prikaz očekivanja od $h(X)$ kao u (4.6) zovemo *temeljni identitet uzorkovanja po važnosti*.

Prema tome, za X s funkcijom gustoće f i \tilde{X} s funkcijom gustoće g (pišemo $X \sim f$, $\tilde{X} \sim g$) vrijedi

$$\mathbb{E}[h(X)] = \int_S h(x)f(x) dx = \int_S h(x) \frac{f(x)}{g(x)} g(x) dx = \mathbb{E} \left[\frac{hf}{g}(\tilde{X}) \right]. \quad (4.7)$$

Tada, ako vrijedi

$$\text{Var} \left[\frac{hf}{g}(\tilde{X}) \right] < \text{Var}[h(X)],$$

i ako znamo generirati slučajne varijable iz gustoće g , jednakost (4.7) nudi alternativno aproksimativno rješenje za $\mathbb{E}[h(X)]$. Dodatno, za gustoće f i g dovoljno je da vrijedi $g(x) = 0 \Rightarrow h(x)f(x) = 0$, $x \in \mathbb{R}^d$.

Dakle, umjesto procjenitelja \tilde{I}_n definiranog u (4.2), definiramo novi procjenitelj

$$\tilde{J}_n := \frac{1}{n} \sum_{i=1}^n h(\tilde{X}_i) \frac{f(\tilde{X}_i)}{g(\tilde{X}_i)}, \quad (4.8)$$

za $\tilde{X}_1, \tilde{X}_2, \dots, \tilde{X}_n$ slučajan uzorak generiran iz gustoće g .

Tako definiran procjenitelj \tilde{J}_n konvergira u $\int_S h(x)f(x) dx$ po Jakom zakonu velikih brojeva, kao i \tilde{I}_n , i to neovisno o izboru distribucije od \tilde{X} , odnosno neovisno o izboru g (uz uvjet $g(x) = 0 \Rightarrow f(x) = 0$, $x \in \mathbb{R}^d$).

Definicija 4.2.1. ² *Metoda uzorkovanja po važnosti je procjena vrijednosti integrala*

$\mathbb{E}[h(X)] = \int_S h(x)f(x) dx$ pomoću procjenitelja

$$\frac{1}{n} \sum_{i=1}^n h(\tilde{X}_i) \frac{f(\tilde{X}_i)}{g(\tilde{X}_i)},$$

gdje je f funkcija gustoće od X , g funkcija gustoće od \tilde{X} , a $\tilde{X}_1, \tilde{X}_2, \dots, \tilde{X}_n$ slučajan uzorak generiran iz gustoće g (tj. $\tilde{X}_i \sim \tilde{X}$).

²Definicija je preuzeta iz [1].

Kako metoda uzorkovanja po važnosti zahtijeva vrlo malo ograničenja na funkciju g , jasno je da je to značajna i često korištena metoda za smanjivanje varijance procjenitelja (g možemo odabrati iz neke distribucije iz koje znamo lako simulirati). Također, isti uzorak simuliran iz g može se koristiti za različite funkcije h i f , što je prednost u brojnim analizama.

Iako je moguće izabrati više različitih gustoća g , postavlja se pitanje kako izabrati što bolju g . Naime, očito svi izbori nisu jednako dobri jer se varijanca pripadnih procjenitelja razlikuje ovisno o g . Zbog te činjenice želimo usporediti različite odabire gustoće g .

Primijetimo sljedeće:

Iako procjenitelj \tilde{J}_n gotovo sigurno konvergira u $\int_S h(x)f(x) dx$, da bi njegova varijanca bila konačna mora vrijediti

$$\mathbb{E}_g \left[h^2(\tilde{X}) \frac{f^2(\tilde{X})}{g^2(\tilde{X})} \right] = \mathbb{E}_f \left[h^2(\tilde{X}) \frac{f(\tilde{X})}{g(\tilde{X})} \right] = \int_S h^2(x) \frac{f^2(x)}{g(x)} dx < \infty.$$

Dakle, zaključujemo da gustoću g treba odabrati tako da f/g bude ograničena, jer u suprotnom varijanca procjenitelja \tilde{J}_n u većini slučajeva neće bit konačna.

Sljedeći teorem³ pokazuje da je, od onih gustoća g koje daju konačnu varijancu procjenitelja, moguće odabrati optimalnu g za zadanu funkciju h i fiksnu gustoću f .

Teorem 4.2.2. *Gustoća g koja minimizira varijancu procjenitelja dana je s*

$$g^*(x) = \frac{|h(x)| f(x)}{\int_S |h(z)| f(z) dz}.$$

Dokaz. Prvo primijetimo da vrijedi

$$\text{Var} \left[\frac{h(\tilde{X}) f(\tilde{X})}{g(\tilde{X})} \right] = \mathbb{E}_g \left[\frac{h^2(\tilde{X}) f^2(\tilde{X})}{g^2(\tilde{X})} \right] - \left(\mathbb{E}_g \left[\frac{h(\tilde{X}) f(\tilde{X})}{g(\tilde{X})} \right] \right)^2,$$

te da izraz $\left(\mathbb{E}_g \left[\frac{h(\tilde{X}) f(\tilde{X})}{g(\tilde{X})} \right] \right)^2$ ne ovisi o g . Dakle, problem minimizacije varijance svodi se na minimizaciju izraza $\mathbb{E}_g \left[\frac{h^2(\tilde{X}) f^2(\tilde{X})}{g^2(\tilde{X})} \right]$.

³Teorem 4.2.2. preuzet je iz [1].

Po Jensenovoj nejednakosti⁴ slijedi

$$\mathbb{E}_g \left[\frac{h^2(\tilde{X}) f^2(\tilde{X})}{g^2(\tilde{X})} \right] \geq \left(\mathbb{E}_g \left[\frac{|h(\tilde{X})| f(\tilde{X})}{g(\tilde{X})} \right] \right)^2 = \left(\int_S |h(x)| f(x) dx \right)^2,$$

što daje donju granicu neovisnu o izboru g . Sada se lako provjeri da se ta donja granica postiže upravo za $g = g^*$. \square

Napomena 4.2.3. Uočimo, za odabir gustoće g po prethodnom teoremu potrebno je znati vrijednost integrala $\int_S |h(z)| f(z) dz$, što je zapravo inicijalni integral kojem želimo izračunati aproksimativnu vrijednost, pa je ovaj teorem više informativan.

U praksi, Teorem 4.2.2 možemo shvatiti na sljedeći način: trebamo tražiti gustoću g za koju je $\frac{|h|f}{g}$ gotovo konstantna s konačnom varijancom. Zahtjev ograničenosti varijance nije nužan za konvergenciju procjenitelja, no u suprotnom metoda uzorkovanja po važnosti daje poprilično loše rezultate te se u tim slučajevima ne preporučuje koristiti.

Ilustrirajmo ovu metodu na primjeru spomenutom u [1]:

Primjer 4.2.4. Neka je $Z \sim N(0, 1)$. Želimo odrediti $\mathbb{P}(Z > 4.5)$ (jako mala vrijednost!). Koristeći Monte Carlo metodu možemo generirati n slučajnih varijabli $Z_i \sim N(0, 1)$ i izračunati procjenitelj

$$\mathbb{P}(Z > 4.5) = \mathbb{E}[\mathbb{1}_{\{Z > 4.5\}}] \approx \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{Z_i > 4.5\}}.$$

No, u tom slučaju čak i za npr. $n = 10^6$ dobijemo tek jednu značajnu znamenku u rezultatu procjene, dok su sve ostale znamenke nula, jer procjenjujemo vjerojatnost jako rijetkog događaja. Tada, za prihvatljivo rješenje dobiveno "naivnom" Monte Carlo metodom potreban je jako velik broj n . Pokažimo kako metoda uzorkovanja po važnosti može poboljšati točnost procjene.

Umjesto procjene vjerojatnosti $\mathbb{P}(Z > 4.5)$ koristeći jednakost

$$\mathbb{P}(Z > 4.5) = \mathbb{E}[h(Z)], \quad \text{za } h(x) = \mathbb{1}_{\{x > 4.5\}}(x),$$

pogledajmo drugačiji raspis.

⁴Vrijedi: $\varphi(\mathbb{E}[X]) \leq \mathbb{E}[\varphi(X)]$, za φ konveksnu funkciju i X takvu da je $\mathbb{E}[|X|] < \infty$.

Za $\tilde{Z} \sim N(\mu, 1)$ vrijedi:

$$\begin{aligned} \mathbb{P}(Z > 4.5) &= \mathbb{E}[h(Z)] = \mathbb{E}[\mathbb{1}_{\{Z > 4.5\}}] = \mathbb{E}\left[\mathbb{1}_{\{\tilde{Z} > 4.5\}} \cdot \frac{\frac{1}{\sqrt{2\pi}} e^{-\tilde{Z}^2/2}}{\frac{1}{\sqrt{2\pi}} e^{-(\tilde{Z}^2 - \mu)/2}}\right] \\ &= \mathbb{E}\left[\mathbb{1}_{\{\tilde{Z} > 4.5\}} \cdot e^{-\mu(\mu - 2\tilde{Z})/2}\right] = \mathbb{E}\left[\frac{hf}{g}(\tilde{Z})\right], \end{aligned}$$

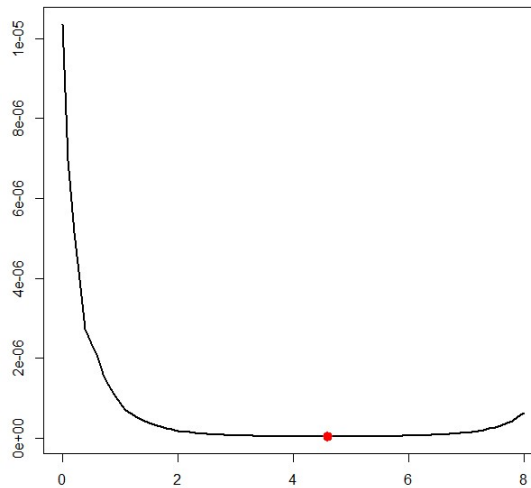
gdje je g funkcija gustoće od $\tilde{Z} \sim N(\mu, 1)$. Sada želimo izračunati procjenu za $\mathbb{P}(Z > 4.5)$ koristeći dobivenu jednakost

$$\mathbb{P}(Z > 4.5) = \mathbb{E}\left[\frac{hf}{g}(\tilde{Z})\right],$$

i to u ovisnosti o parametru μ . Na taj način možemo odrediti za koji μ dobijemo najbolju procjenu, odnosno za koji μ je pripadni 95% pouzdani interval najuži.

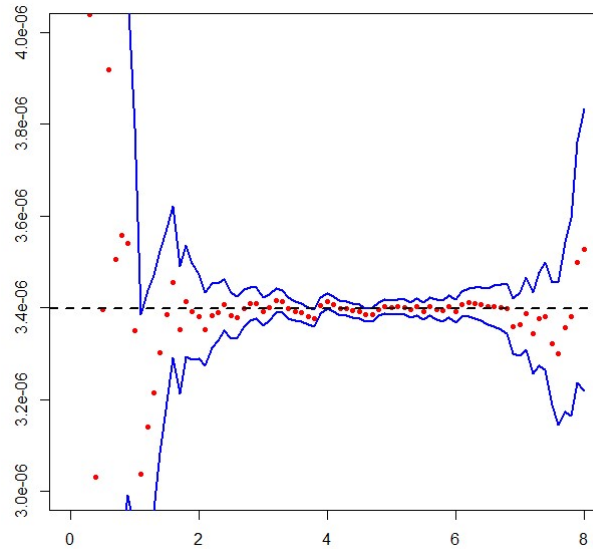
U R-u za $n = 10^6$ izračunajmo realizacije procjenitelja u ovisnosti o parametru $\mu \in [0, 8]$. Za jedan simulirani slučajni uzorak iz $N(\mu, 1)$, za razne μ , interval pouzdanosti je najuži za $\mu = 4.5$. Za takav optimalan μ procjena iznosi $3.384107 \cdot 10^{-6} = 0.000003384107$ (minimalna širina intervala u ovom slučaju iznosi $2.987477 \cdot 10^{-8}$).

Na slici 4.3 prikazana je širina intervala pouzdanosti u ovisnosti o $\mu \in [0, 8]$, te točka minimuma širine intervala koji se postiže u $\mu = 4.5$.



Slika 4.3: Širina intervala pouzdanosti u ovisnosti o $\mu \in [0, 8]$

Na slici 4.4 dodatno su prikazati svi izračunati procjenitelji za $\mu \in [0, 8]$, pripadni 95% pouzdani intervali, kao i stvarna vrijednost vjerojatnosti $\mathbb{P}(Z > 4.5)$.



Slika 4.4: Svi izračunati procjenitelji za $\mu \in [0, 8]$ (crvene točke), pripadni 95% pouzdani intervali (plave linije) i stvarna vrijednost vjerojatnosti $\mathbb{P}(Z > 4.5)$ (crna isprekidana linija)

Primijetimo da $\mathbb{P}(Z > 4.5)$ možemo raspisati i na sljedeći način:
Neka je Y slučajna varijabla s funkcijom gustoće

$$f_Y(y) = e^{-(y-4.5)} / \int_{4.5}^{\infty} e^{-x} dx.$$

Kažemo da je Y slučajna varijabla iz eksponencijalne distribucije (lijevo) skraćena na 4.5 s parametrom $\lambda = 1$. Ako generiramo iz f_Y i koristimo metodu uzorkovanja po važnosti, dobijemo

$$\mathbb{P}(Z > 4.5) \approx \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{Y_i > 4.5\}} \frac{\frac{1}{\sqrt{2\pi}} e^{-Y_i^2/2}}{f_Y(Y_i)},$$

što je još jedan procjenitelj za $\mathbb{P}(Z > 4.5)$. Jedna realizacija tog procjenitelja iznosi 0.000003377. \square

Možemo zaključiti da metoda uzorkovanja po važnosti daje bolju procjenu vrijednosti od "naivne" metode Monte Carlo. Po Teoremu 4.2.2 znamo da postoji optimalniji izbor procjenitelja od Monte Carlo procjenitelja definiranog u (4.2), odnosno "isplati" se slučajni uzorak generirati iz neke gustoće g različite od f . Pritom je bitno odabrati gustoću g za koju je varijanca procjenitelja konačna, odnosno takvu da vrijedi $\mathbb{E}_f \left[\left| \frac{f(X)}{g(X)} \right| \right] < \infty$.

Na kraju, pokažimo na jednom primjeru kako na vrijednost procjene integrala može utjecati odabir metode generiranja slučajnih varijabli te odabir Monte Carlo procjenitelja.

Primjer 4.2.5. Želimo izračunati vrijednost procjene integrala

$$I = \int_4^{\infty} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx.$$

Stvarna vrijednost tog integrala iznosi $I \approx 3.167124 \cdot 10^{-5}$. Kako se radi o jako maloj vrijednosti, za duljinu niza generiranih slučajnih varijabli u R -u odabrat ćemo $n = 10^6$ u svim metodama generiranja koje ćemo provesti u nastavku. Prvo prikažimo I kao očekivanje neke slučajne varijable i definirajmo pripadni Monte Carlo procjenitelj.

Vrijedi:

$$I = \int_4^{\infty} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx = \mathbb{E}[h(X)],$$

za funkciju $h(x) = \mathbb{1}_{\{x > 4\}}(x)$ i slučajnu varijablu $X \sim N(0, 1)$, pa procjenu za I možemo izračunati koristeći Monte Carlo procjenitelj

$$\tilde{I}_n = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{X_i > 4\}},$$

gdje su $X_i \sim N(0, 1)$, $i = 1, \dots, n$, nezavisne i jednako distribuirane varijable. Osim Monte Carlo procjenitelja \tilde{I}_n kojeg smo odredili "naivnom" Monte Carlo metodom, iskoristimo metodu uzorkovanja po važnosti i definirajmo još jedan Monte Carlo procjenitelj \tilde{J}_n , koji ima manju varijancu od \tilde{I}_n .

Ako funkciju gustoće slučajne varijable $X \sim N(0, 1)$ označimo s f , vrijedi sljedeće:

$$\mathbb{E}[h(X)] = \mathbb{E}[\mathbb{1}_{\{X > 4\}}] = \mathbb{E} \left[\mathbb{1}_{\{\tilde{X} > 4\}} \cdot \frac{\frac{1}{\sqrt{2\pi}} e^{-\tilde{X}^2/2}}{\frac{1}{\sqrt{2\pi}} e^{-(\tilde{X}^2 - \mu)/2}} \right] = \mathbb{E} \left[\frac{hf}{g}(\tilde{X}) \right],$$

za $\tilde{X} \sim N(\mu, 1)$ i g funkciju gustoće od \tilde{X} .

Dakle, za izračun procjene za I također možemo iskoristiti novi procjenitelj

$$\tilde{J}_n := \frac{1}{n} \sum_{i=1}^n h(\tilde{X}_i) \frac{f(\tilde{X}_i)}{g(\tilde{X}_i)},$$

gdje su $\tilde{X}_i \sim N(\mu, 1)$. Kao i u prethodnom primjeru, možemo odrediti parametar μ za kojeg je pripadni 95%-pouzdana interval najuži, odnosno greška najmanja. Takav optimalan μ iznosi $\mu = 4.1$.

Da bismo odredili vrijednosti procjenitelja \tilde{I}_n i \tilde{J}_n , moramo znati efikasno generirati slučajne varijable iz distribucija $N(0, 1)$ i $N(4.1, 1)$. To ćemo napraviti na nekoliko načina te ovisno o primijenjenoj metodi generiranja odrediti vrijednosti procjenitelja \tilde{I}_n i \tilde{J}_n .

- **Accept-Reject algoritam**

Za primjenu Accept-Reject algoritma prvo je potrebno pronaći funkciju g i konstantu M takve da zadovoljavaju uvjet $f \leq Mg$. Uočimo sljedeće⁵:

Neka je $g(x) = \frac{1}{\pi(1+x^2)}$, $x \in \mathbb{R}$, funkcija gustoće Cauchyjeve slučajne varijable i f funkcija gustoće standardne normalne slučajne varijable. Tada je funkcija f/g ograničena te postoji najmanji M takav da vrijedi $f \leq Mg$ koji iznosi $M = \sqrt{2\pi/e}$. To se lako pokaže deriviranjem funkcije f/g (stacionarne točke su $-1, 0, 1$, a maksimum se postiže u -1 i 1 te poprima vrijednost $\sqrt{2\pi/e}$). Sada za takve M i g lako provedemo Accept-Reject algoritam i znamo da vjerojatnost prihvatanja generiranog $X \sim g$ iznosi $1/M \approx 0.6577$. To znači da je od ukupno $n = 10^6$ generiranih slučajnih varijabli u algoritmu približno njih 667 700 iz $N(0, 1)$ razdiobe. Međutim, zbog usporedbe s ostalim metodama generiranja, želimo ukupno $n = 10^6$ slučajnih varijabli iz $N(0, 1)$ razdiobe generiranih ovom metodom. Zato ćemo, umjesto početnih 10^6 , generirati $m = 1\,520\,000$ slučajnih varijabli iz g ($m \cdot \frac{1}{M} \approx 10^6 = n$). Pritom, za generiranje Cauchyjevih slučajnih varijabli koristit ćemo funkciju `rcauchy` implementiranu u R-u.

- **funkcija `rnorm` u R-u**

Za generiranje slučajnih varijabli iz $N(0, 1)$ razdiobe korištenje funkcije `rnorm` u R-u uvijek je opcija.

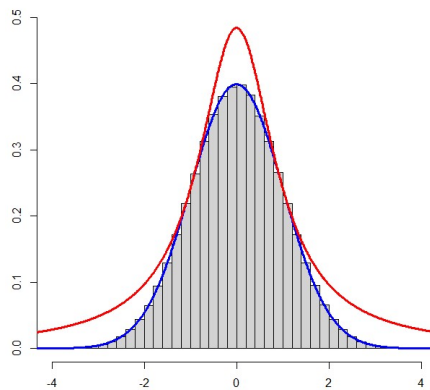
- **Box-Muller algoritam**

U odjeljku 3.1. opisali smo jednostavan Box-Muller algoritam za generiranje iz $N(0, 1)$ razdiobe, a ovdje ćemo ga iskoristiti za još jedan primjer generiranja.

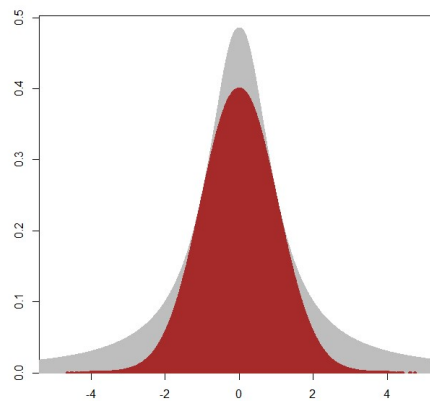
⁵Ovakav način simuliranja iz $N(0, 1)$ razdiobe spomenut je u problemu 2.34 u [1].

Primijetimo da metodu inverza u ovom slučaju ne možemo primijeniti zbog oblika funkcije distribucije F standardne normalne razdiobe.

Na slikama ispod prikazani su histogram uzorka generiranog Accept-Reject algoritmom iz $N(0, 1)$ razdiobe te grafički prikaz prihvaćenih i odbačenih točaka. Primijetimo da je u ovom primjeru funkcija g puno bolje odabrana nego primjerice u Primjeru 3.2.3., što je vidljivo sa slike 4.6, ali i iz veće vjerojatnosti prihvaćanja $1/M$ koja ovdje iznosi približno 0.6577. Očekivani broj pokušaja do prvog prihvaćanja u algoritmu je $M \approx 1.52$.



Slika 4.5: Histogram uzorka generiranog Accept-Reject algoritmom iz $N(0, 1)$ razdiobe, zajedno s funkcijom gustoće f i funkcijom Mg



Slika 4.6: Grafički prikaz prihvaćenih (smeđih) i odbačenih (sivih) točaka generiranih Accept-Reject algoritmom

Konačno, pokažimo u tablicama dobivene rezultate procjena vrijednosti integrala I zajedno s širinama pripadnih 95% pouzdanih intervala u ovisnosti o metodi generiranja slučajnog uzorka i odabiru Monte Carlo procjenitelja. Uočimo, ako je $X \sim N(0, 1)$, tada je $X + a \sim N(a, 1)$, $a \in \mathbb{R}$, pa za generiranu slučajnu varijablu $X \sim N(0, 1)$ lako znamo odrediti $X + 4.1 \sim (4.1, 1)$.

metoda generiranja	procjenitelj \tilde{I}_n za I	širina 95% pouzdanog intervala za I
Accept-Reject algoritam	$2.898788 \cdot 10^{-5}$	$1.711821 \cdot 10^{-5}$
funkcija <code>rnorm</code>	$2.8 \cdot 10^{-5}$	$2.074203 \cdot 10^{-5}$
Box-Muller algoritam	$2.3 \cdot 10^{-5}$	$1.879911 \cdot 10^{-5}$

Tablica 4.4: Vrijednosti procjenitelja \tilde{I}_n za I i širina pripadnog 95% pouzdanog intervala u ovisnosti o metodi generiranja slučajnih varijabli (rezultati dobiveni "naivnom" Monte Carlo metodom)

metoda generiranja	procjenitelj \tilde{J}_n za I	širina 95% pouzdanog intervala za I
Accept-Reject algoritam	$3.161039 \cdot 10^{-5}$	$2.622262 \cdot 10^{-7}$
funkcija <code>rnorm</code>	$3.163275 \cdot 10^{-5}$	$2.625811 \cdot 10^{-7}$
Box-Muller algoritam	$3.171749 \cdot 10^{-5}$	$2.626585 \cdot 10^{-7}$

Tablica 4.5: Vrijednosti procjenitelja \tilde{J}_n za I i širina pripadnog 95% pouzdanog intervala u ovisnosti o metodi generiranja slučajnih varijabli (rezultati dobiveni metodom uzorkovanja po važnosti)

Iako smo generirali slučajne uzorke duljine $n = 10^6$, rezultati dobiveni "naivnom" Monte Carlo metodom nemaju točnu nijednu značajnu znamenku, a najveću grešku, odnosno najširi pouzdani interval, dobili smo koristeći funkciju `rnorm`. Za razliku od "naivne" Monte Carlo metode, metoda uzorkovanja po važnosti za sve tri metode generiranja daje točnije procjene i manje greške, odnosno uže pouzdane intervale (red veličine greške se smanjio), pa zaključujemo da se itetako isplati provesti metodu uzorkovanja po važnosti. \square

Dodatak A

Rješenja primjera - R kodovi

U ovom dodatku nalaze se linije koda iz programskog paketa *R* za rješavanje nekih primjera iz prethodnih poglavlja, koje pokazuju na koji način i uz pomoć kojih naredbi možemo doći do određenih izračuna.

Primjer 2.1.1.

```
n=10^4
x<-runif(n,0,1)
h<-function(x)
{(cos(50*x)+sin(20*x))^2}
In<-mean(h(x))
#ocjena greske:
alfa<-0.05
Sn2<-(1/(n-1))*(sum((h(x))^2)-n*In^2)
Sn<-sqrt(Sn2)
#ili: Sn<-sd(h(x))
#aproksimativni pouzdani interval:
z<-qnorm(1-alfa/2)
l<-In-(Sn/sqrt(n))*z
d<-In+(Sn/sqrt(n))*z
c(l,d)
#histogram uzorka h(Xi):
hist(h(x), probability=T, col="grey", breaks=30, xlab="")
#greška u ovisnosti o n:
I_j<-cumsum(h(x))/(1:n)
Sn_j<-(cumsum((h(x))^2)-(1:n)*I_j^2)/c(1,1:(n-1))
Sn_j<-sqrt(Sn_j)
```



```

plot(1:n,I_j,xlab="n", type="l", ylim=c(0.85,1.1), col="blue")
lines(I_j+ Sn_j*z/sqrt(1:n), lwd=2)
lines(I_j- Sn_j*z/sqrt(1:n), lwd=2)
I<-integrate(h,0,1)$value
abline(h=I, lty=3, lwd=2, col="red")

#k=10000 realizacija procjenitelja:
k<-10000
y<-matrix(runif(k*n,0,1),ncol=k)
In<-apply(h(y),2,mean)

#mozemo koristiti i for petlju:
n<-10^4
k<-10^4
In2<-numeric(k)
for(j in 1:k)
  In2[j]<-mean(h(y[,j]))

#stvarna vrijednost integrala I=E[h(X)]:
I<-integrate(h,0,1)$value
#stvarna vrijednost standardne devijacije od h(X):
#a:=E[h(X)^2]
g<-function(x) {h(x)^2}
a<-integrate(g,0,1)$value
sigma<-sqrt(a-I^2)
hist(In, probability=T, col="grey", main="Histogram 10000 realizacija
      Monte Carlo procjenitelja", breaks=30, xlab="")
t<-seq(min(In),max(In),by=0.001)
lines(t,dnorm(t,I,sigma/sqrt(n)), col="blue", lwd=2)
abline(v=I, col="red",lwd=2)

```

Primjer 3.1.3.

```

n<-10^4
F_inv<-function (u)
{ -log(u) }
u<-runif(n,0,1)
x<-F_inv(u)
hist(x, probability = TRUE, breaks=60, col="grey", xlim=c(0,6),
      ylim=c(0,0.9))

```

```
t<-seq(0,6,by=0.01)
points(t,dexp(t, 1) , type = 'l' , col="red", lwd=2)
```

Primjer 3.2.3.

```
optimize(f=function(x){dbeta(x,2.7,6.3)},interval=c(0,1),maximum=T)
M<-optimize(f=function(x){dbeta(x,2.7,6.3)},interval=c(0,1),
            maximum=T)$objective
```

#Accept-Reject algoritam:

```
n<-10^4
x<-runif(n,0,1)
u<-runif(n,0,1)
uvjet<-u<(dbeta(x,2.7,6.3)/M)
x_accepted<-x[uvjet]
hist(x_accepted, probability = TRUE, breaks = 30)
t<-seq(0,1, by = 0.01)
points(t, dbeta(t,2.7,6.3), type = "l",col="blue",lwd=3)
abline(h=M,col="red",lwd=4,lty=2)
#grafički prikaz prihvacenih i odbacenih točaka:
plot(x[uvjet], u[uvjet]*M, col="black", pch = 20)
points(x[!uvjet], u[!uvjet]*M, col="grey", pch = 20)
points(t, dbeta(t,2.7,6.3), type = "l", col="brown", lwd =4)
abline(h=M,col="brown",lwd=4)
postotak_prihvacenih<-length(x_accepted)/n
vjerojatnost_prihvacanja<-1/M
```

Primjer 4.1.2.

```
n<-10^4
x<-matrix(runif(2*n,-1,1), n, 2)
h<-function(x)
{sum(x^2)<=1}
y<-apply(x,1,h)
pn<-mean(y)
4*pn
#grafička ilustracija:
par(pty="s")
plot(x[!y,1],x[!y,2], col="blue", pch=19)
points(x[y,1],x[y,2], col="red", pch=19)
theta<-seq(0,2*pi,0.01)
points(cos(theta),sin(theta), type="l", col="black", lwd=3)
```

Primjer 4.1.3.

```
#(a)
n<-10^4
x<-runif(n,0,3)
y<-runif(n,0,2)
z<-runif(n,0,1)
h<-function(x,y,z)
{6*x*y*z}
A<-matrix(c(x,y,z),ncol=3)
c<-numeric(n)
for(i in 1:n)
c[i]<-h(A[i,1],A[i,2],A[i,3])
In<-mean(c)
#ocjena greske:
alfa<-0.05
Sn<-sd(c)
#aproksimativni pouzdani interval:
z<-qnorm(1-alfa/2)
l<-In-(Sn/sqrt(n))*z
d<-In+(Sn/sqrt(n))*z
c(l,d)

#(b)
n<-10^8
x<-runif(n,-3,2)
y<-runif(n,0,1)
z<-runif(n,-1,1)
h<-function(x,y,z)
{10*(sin(x)*y^2-z)}
A<-matrix(c(x,y,z),ncol=3)
c<-numeric(n)
for(i in 1:n)
c[i]<-h(A[i,1],A[i,2],A[i,3])
In<-mean(c)
#ocjena greske:
alfa<-0.05
Sn<-sd(c)
#aproksimativni pouzdani interval:
z<-qnorm(1-alfa/2)
```

```
l<-ln-(Sn/sqrt(n))*z
d<-ln+(Sn/sqrt(n))*z
c(l,d)
```

Primjer 4.2.4.

```
#mi+Z uzorak iz N(mi,1)
rezultat_procjene<-function(n=10^4,Z=rnorm(n),mi=0) {
  ImpSampl<-(mi+Z>4.5)*dnorm(mi+Z,0,1)/dnorm(mi+Z,mi,1)
  procjenitelj<-mean(ImpSampl)
  sigma<-sd(ImpSampl)
  alfa<-0.05
  z<-qnorm(1-alfa/2)
  l<-procjenitelj-z*sigma/sqrt(n)
  d<-procjenitelj+z*sigma/sqrt(n)
  return(c(procjenitelj,d-l,l,d))
}
n<-10^6
Z<-rnorm(n,0,1)
mi<-seq(0,8,0.1)
d<-length(mi)
širine_pouz_d_int<-numeric(d)
indeks_min<-0
min_sirina<-1
for(i in 1:d) {
  širine_pouz_d_int[i]<-rezultat_procjene(n,Z,mi[i])[2]
  if(širine_pouz_d_int[i]<min_sirina) {
    min_sirina<-širine_pouz_d_int[i]
    indeks_min<-i
  }
}
opt_mi<-mi[indeks_min]
procj<-rezultat_procjene(n,Z,opt_mi)

plot(mi,širine_pouz_d_int, type="l", lwd=2)
points(mi[indeks_min], širine_pouz_d_int[indeks_min],
       col="red", pch=19, lwd=5)
matrica_rezultata<-matrix(d*4,nrow=d,ncol=4)
for(i in 1:d)
{ matrica_rezultata[i,]<-rezultat_procjene(n,Z,mi[i]) }
```

```

min(matrica_rezultata[,3])
max(matrica_rezultata[,4])
plot(mi, matrica_rezultata[,1], col="red", ylab="", pch=20,
      ylim=c(0.000003,0.000004))
points(mi,matrica_rezultata[,3],type="l",col="blue",lwd=2)
points(mi,matrica_rezultata[,4],type="l",col="blue",lwd=2)
p<-1-pnorm(4.5)
abline(h=p,lty=2,lwd=2)

```

Primjer 4.2.5.

```

f<-function(x)
{ 1/sqrt(2*pi)*exp(-x^2/2)}
I<-integrate(f,4,Inf)$value
#odredjivanje optimalnog parametra mi:
rezultat_procjene<-function(n,Z=rnorm(n),mi=0) {
  ImpSampl<-(mi+Z>4)*dnorm(mi+Z,0,1)/dnorm(mi+Z,mi,1)
  procjenitelj<-mean(ImpSampl)
  sigma<-sd(ImpSampl)
  alfa<-0.05
  z<-qnorm(1-alfa/2)
  l<-procjenitelj-z*sigma/sqrt(n)
  d<-procjenitelj+z*sigma/sqrt(n)
  return(c(procjenitelj,d-l,l,d))
}
n<-10^6
x<-rnorm(n,0,1)
mi<-seq(0,8,0.1)
length(mi)      #81
širine_pouzd_int<-numeric(81)
indeks_min<-0
min_sirina<-1
for(i in 1:81) {
  širine_pouzd_int[i]<-rezultat_procjene(n,x,mi[i])[2]
  if(širine_pouzd_int[i]<min_sirina) {
    min_sirina<-širine_pouzd_int[i]
    indeks_min<-i }
}

```

```

min_sirina
indeks_min
opt_mi<-mi[indeks_min]
plot(mi,širine_pouz_d_int, type="l", lwd=2)
points(mi[indeks_min], širine_pouz_d_int[indeks_min],
       col="red", pch=19, lwd=5)

#generiranje slučajnih varijabli:
n<-10^6
#1)Accept-Reject algoritam
M<-sqrt(2*pi)*exp(-1/2)
m<-1520000
x<-rcauchy(m,0,1)
u<-runif(m,0,1)
uvjet<-u<(dnorm(x,0,1)/(M*dcauchy(x,0,1)))
x_accepted<-x[uvjet]
hist(x_accepted,probability=T,breaks=30,ylim=c(0,0.5),xlim=c(-4,4))
t<-seq(-5,5,0.01)
points(t,dnorm(t,0,1),type="l",lwd=3,col="blue")
points(t,M*dcauchy(t,0,1),type="l",xlab="x",ylab="",col="red",lwd=3)
length(x_accepted)/n
1/M
#graficki prikaz:
y<-u*M*dcauchy(x,0,1)
plot(t,M*dcauchy(t,0,1),type="l",xlab="x",ylab="",lwd=3)
points(t,dnorm(t,0,1),type="l", lwd=3)
points(x[!uvjet],y[!uvjet],col="grey",pch=20)
points(x[uvjet],y[uvjet], col="brown", pch = 20)
length(x_accepted)
#"naivna" MC metoda:
h<-function(x)
{ x>4 }
In_acc<-mean(h(x_accepted))
In_acc
Sn_acc<-sd(h(x_accepted))
alfa<-0.05
z<-qnorm(1-alfa/2)
l<-In_acc-(Sn_acc/sqrt(m))*z
d<-In_acc+(Sn_acc/sqrt(m))*z
c(l,d,d-1)

```

```
#uzorkovanje po vaznosti:
duljina<-length(x_accepted)
procj<-rezultat_procjene(duljina,x_accepted,opt_mi)

#2) rnorm funkcija
#x<-rnorm(n,0,1) je već generiran za odredjivanje parametra mi
#"naivna MC metoda:
In<-mean(h(x))
In
Sn<-sd(h(x))
alfa<-0.05
z<-qnorm(1-alfa/2)
l<-In-(Sn/sqrt(n))*z
d<-In+(Sn/sqrt(n))*z
c(l,d,d-1)
#uzorkovanje po vaznosti:
procj<-rezultat_procjene(n,x,opt_mi)

#3) Box-Muller algoritam
u_i<-runif(n,0,1)
x_i<-numeric(n)
for(i in 1:(n-1)) {
  x_i[i]<-sqrt(-2*log(u_i[i]))*cos(2*pi*u_i[i+1])
  x_i[i+1]<-sqrt(-2*log(u_i[i]))*sin(2*pi*u_i[i+1])
  i<-i+2
}
length(x_i)
#"naivna MC metoda:
In_BM<-mean(h(x_i))
In_BM
Sn_BM<-sd(h(x_i))
alfa<-0.05
z<-qnorm(1-alfa/2)
l<-In_BM-(Sn_BM/sqrt(n))*z
d<-In_BM+(Sn_BM/sqrt(n))*z
c(l,d,d-1)
#uzorkovanje po vaznosti:
procj<-rezultat_procjene(n,x_i,opt_mi)
```

Bibliografija

- [1] Christian P. Robert, George Casella, *Monte Carlo Statistical Methods*, Springer, New York, 2004
- [2] Christian P. Robert, George Casella, *Introducing Monte Carlo Methods with R*, Springer, New York, 2010
- [3] Z. Vondraček, N. Sandrić, *Vjerojatnost*, predavanja, PMF-MO, 2019., https://www.pmf.unizg.hr/images/50023697/vjer_predavanja.pdf
- [4] M. Huzak, *Matematička statistika*, predavanja, PMF-MO, 2020., <https://web.math.pmf.unizg.hr/nastava/ms/index.php?sadrzaj=predavanja.php>
- [5] A. Leko, *Neprekidan slučajan vektor*, završni rad, Osijek, 2013., <http://www.mathos.unios.hr/~mdjunic/uploads/diplomski/LEK05.pdf>
- [6] M. Huzak, *Vjerojatnost i matematička statistika*, predavanja, PMF-MO, 2006., <http://aktuari.math.pmf.unizg.hr/docs/vms.pdf>
- [7] B. Basrak, H. Planinić, *Financijski praktikum*, predavanja, PMF-MO, 2021., https://moodle.srce.hr/2020-2021/pluginfile.php/5038949/mod_resource/content/3/SimSlVar_rejection.pdf
- [8] P. Glasserman, *Monte Carlo Methods in Financial Engineering*, Springer, New York, 2003
- [9] *Monte Carlo Integration*, https://en.wikipedia.org/wiki/Monte_Carlo_integration

Sažetak

Pri rješavanju problema u brojim područjima znanosti nerijetko se pojavljuju integrali čije je vrijednosti potrebno izračunati. Međutim, sam proces integracije složen je te ponekad vrijednost integrala nije moguće odrediti egzaktno. Iako postoje različite numeričke metode integracije koje pružaju efikasne algoritme za aproksimativno rješenje, u određenim slučajevima, primjerice u višedimenzionalnim problemima, zbog nedovoljne brzine ili nepraktičnosti izvođenja one se mogu bitno zakomplicirati. U ovom radu opisali smo jednu drugačiju metodu integracije - Monte Carlo metodu, kojom u pravilu možemo dobiti dobre procjene vrijednosti integrala i ocijeniti njihovu grešku. Ova aproksimativna metoda koristi se u slučajevima kada problem nije moguće (ili nije jednostavno) riješiti analitički, a posebno je korisna u rješavanju višedimenzionalnih problema jer točnost procjene ne ovisi o dimenziji prostora. Stoga je primjena ove metode u novije vrijeme sasvim uobičajena i poželjna. Primarni cilj ovog rada bio je opisati ideju integracije metodom Monte Carlo i teoretski ju razraditi. Zaključili smo kako se cijeli proces integracije ovom metodom temelji na mogućnosti efikasnog generiranja slučajnih varijabli iz određene distribucije te naveli neke od metoda generiranja. Također, pokazali smo kako je u tom procesu integracije ipak potrebno pažljivo odabrati Monte Carlo procjenitelj, odnosno distribuciju iz koje ćemo generirati slučajne varijable. Naime, upravo o tom izboru ovisi varijanca procjenitelja koju želimo maksimalno smanjiti. Taj proces smanjivanja varijance Monte Carlo procjenitelja opisali smo metodom uzorkovanja po važnosti (*Importance Sampling*) i zaključili da ta metoda daje bolju procjenu vrijednosti integrala od "naivne" Monte Carlo metode.

Summary

Solving problems in multiple areas of science often requires calculating the values of integrals. However, the process of integration is so complex that sometimes it is not possible to determine the exact value of an integral. There are various numerical integration methods that offer efficient algorithms for calculating the approximate solution. Still, in certain cases (e.g. multidimensional problems), due to insufficient speed or impracticality of the performance they can get significantly complicated. In this thesis, we described a different method of integration - the Monte Carlo method. This method allows us to get appropriate estimated values of integrals and their approximation error. This approximate method is also used in situations when the problem is too difficult (or impossible) to solve analytically. It is especially useful when solving multidimensional problems because the accuracy of the estimated value does not depend on the dimension of space. Hence, using this method is both common and desirable in modern times. Primary goal of this thesis was to describe Monte Carlo integration method and to elaborate on it theoretically. It was concluded that the whole process of using this method is based on the idea of efficiently generating random variables from a certain distribution. In addition, we named some of the methods used for generating such variables. It is worth pointing out that it is necessary to carefully select the Monte Carlo estimator i.e. the distribution from which we generate random variables. Namely, it is on that selection of the distribution that the variance of the estimator which we are trying to maximally reduce depends on. This process of reduction of the variance of the Monte Carlo estimator was described using the Importance Sampling method which has proven to give a better estimated value of the integral than the “naive” Monte Carlo method.

Životopis

Rođena sam 25. listopada 1995. godine u Splitu. Odrasla sam u Ljubitovici, malom mjestu pored Trogira, gdje sam pohađala Osnovnu školu kralja Zvonimira. Srednjoškolsko obrazovanje nastavljam u Trogiru, gdje upisujem Srednju školu Ivana Lucića, smjer opća gimnazija. Završetkom srednje škole, 2014. godine upisujem preddiplomski studij Matematika na Prirodoslovno-matematičkom fakultetu u Zagrebu, kojeg završavam 2018. godine i stječem titulu sveučilišne prvostupnice matematike te zatim, na istom fakultetu, upisujem diplomski studij Financijska i poslovna matematika. Od studenog 2020. godine radim u OTP Osiguranju kao aktuarski suradnik.