

Sisukord

1	Pseudojuhuslikud arvud. Lineaarsed kongruentsed generaatorid	3
1.0.1	Lineaarne kongruentne generaator (LKG)	3
2	Generaatorite testimine	5
2.1	Testimine ühe sündmuse abil	5
2.2	Mitme sündmuse toimumissageduste üheaegne vaatlemine - χ^2 -test	6
2.2.1	Hii-ruut test	6
2.3	Empiirilise ja teoreetilise jaotusfunktsiooni võrdlemisel põhinevad testid .	6
2.3.1	Kolmogorov-Smirnovi test	8
2.3.2	Oomega-ruut test	9
2.4	Sõltumatuse testimine	9
2.4.1	Sõltumatute paaride test	9
2.4.2	Vahemike test	10
2.4.3	Seeriaste test	11
2.4.4	Pokker-test	11
3	Etteantud jaotusele vastavate pseudojuhuslike arvude genereerimine	12
3.1	Jaotusfunktsiooni pööramise meetod	12
3.1.1	Lõpliku arvu väärtustega diskreetse juhusliku suuruse esitamine ühtlase jaotuse kaudu	14
3.2	Valikumeetodid	14
3.2.1	Lihtne valikumeetod	14
3.2.2	Üldine valikumeetod	16
3.3	Genereerimine segujaotusest	16
3.4	Pideva juhusliku suuruse teisendused	17
3.5	Genereerimine tinglikust jaotusest	18
4	Pseudojuhuslike vektorite genereerimine	19
4.1	Pidevate juhuslike vektorite teisendused	19
4.1.1	Tihedusfunktsiooni teisenemine lineaarsete teisenduste korral	19
4.2	Normaaljaotusega juhuslike suuruste ja juhuslike vektorite genereerimine .	19
4.2.1	Box-Mülleri meetod	20
4.2.2	Jaotusega $N(\mu, \Sigma)$ juhuslike vektorite genereerimine	20
4.3	Simuleerimine tinglike jaotuste abil	21

4.4	Metropolis-Hastings algoritm	21
4.5	Gibbsi valik	22
5	Taasvaliku meetodid: bootstrap ja jackknife	23
5.1	Bootstrap meetod	23
5.1.1	Nihke hindamine Bootstrap meetodil	23
5.1.2	Bootstrap usaldusvahemik	24
5.2	Jackknife meetod	24
6	Integraalide ja keskväärtuste arvutamine MC meetodiga. Dispersiooni vähendamise meetodid	26
6.1	Keskvärtuse hindamine MC meetodil	26
6.1.1	MC meetodite võrdlemine	27
6.2	Dispersiooni vähendamise meetodid	28
6.2.1	Antiteetilised (<i>antithetic</i>) juhuslikud suurused	28
6.2.2	Kontrollmuutujate meetod	29
6.2.3	Olulise valimi meetod	29
6.2.4	Kihtvalimi meetod	30
6.2.5	Kihtvalimi meetodi optimaalsed proportsioonid	31
6.2.6	Kihtvalimi meetodi rakendamine	32
7	Kasutatavad tulemused tõenäosusteooriast	33

Peatükk 1

Pseudojuhuslikud arvud. Lineaarsed kongruentsed generaatorid

Suur osa statistika rakendusi on kas otseselt või kaudselt seotud sõltumatute katsete sooritamisega ja saadud katsetulumuste põhjal mitmesuguste hinnangute leidmisega.

Definitsioon 1 Arve x_1, \dots, x_n nimetatakse juhuslikule suurusele X vastavateks juhuslikeks arvudeks, kui nad on saadud juhusliku suuruse X väärtustena juhusliku katse sõltumatul kordamisel.

Sõltumatuid katseid on raske arvutis piisavalt kiiresti sooritada, seetõttu kasutatakse arvutisimulatsioonides mittejuhuslike arve, mis imiteerivad sõltumatute katsete tulemusi piisavalt hästi.

Definitsioon 2 Arve x_1, \dots, x_n nimetatakse juhuslikule suurusele X vastavateks pseudojuhuslikeks arvudeks, kui nad on saadud mingi algoritmi või eeskirja kohaselt ja neid ei õnnestu eristada juhuslikule suurusele X vastavatest juhuslikest arvudest teatud komplekti statistiliste testide abil.

Eelnev definitsioon ei ole matemaatiliselt päris korrektne, sest kasutatav komplekt statistilisi teste ei ole täpselt määratletud, kuid annab siiski aimu sellest, mida arvutisimulatsioonides pseudojuhuslike arve kasutades peab silmas pidama.

1.0.1 Lineaarne kongruentne generaator (LKG)

Järgnevas tähistame täisarvude x ja $y > 0$ korral x jagamisel y -ga tekkivat jääki kujul $x \bmod y$.

Definitsioon 3 Etteantud täisarvude m , $0 < a < m$, $0 \leq b < m$ ja $0 \leq x_0 < m$ nimetame eeskirja

$$x_i := (ax_{i-1} + b) \bmod m, \quad u_i := \frac{x_i}{m}$$

abil lõiku $[0, 1]$ kuuluvaid arve u_i , $i \geq 1$ tekitavat algoritmi lineaarseks kongruentseks generaatoriks (tähis $LKG(a, b, m)$). Kui eelnevas algoritmis kehtib $b = 0$, siis on tegemist multiplikatiivse generaatoriga (inglise keeles ka Lehmer random number generator).

Lihtne on näha, et LKG väljastab maksimaalselt m erinevat arvu ja seejärel hakkab ennast kordama (alustab uut tsükli). Samas ebaõnnestunud parameetrite valiku korral võib generaatori tsükli pikkus olla tunduvalt lühem.

Teoreem 4 (*Hull-Dobell teoreem*) Lineaarse kongruentse meetodi korral saavutatakse tsükli maksimaalne pikkus m parajasti siis, kui on täidetud järgmised tingimused:

1. suurim ühistegur $\text{SYT}(b, m) = 1$
2. kui m jagub p -ga, siis $(a - 1)$ jagub p -ga
3. kui m jagub 4 -ga, siis $(a - 1)$ jagub 4 -ga

Multiplikatiivse generaatori korral on selge, et kui $x_i = 0$ mingi i korral, siis kõik edasised väärtused on samuti võrdsed nulliga. Seega maksimaalne tsükkel ei saa sisaldada arvu 0 ja selle pikkus ei saa olla seega suurem kui $m - 1$.

Teoreem 5 *Multiplikatiivne generaatori tsükkel on maksimaalse pikkusega $m - 1$ parajasti siis, kui*

1. m on algarv ja
2. a on algjuur, st iga $m - 1$ algteguri p korral $a^{(m-1)/p}$ ei jagu arvuga m

Oma lihtsuse tõttu on olnud laialdaselt kasutusel, aga praegu ei soovitata enam kasutada nt Monte-Carlo simulatsioonide puhul ning krüpteerimise eesmärgil. Alternatiivid: Mersenne Twister (\mathbb{R} vaikegeneraator, vt Wikipedia algoritmi osas), Blum-Blum-Shub generaator (näide võimalikust krüptograafiliselt turvaliseks loetavast generaatorist, vt Wikipedia).

Peatükk 2

Generaatorite testimine

Tahame kontrollida, et generaatori väärtused käituvad nagu soovitud jaotusele vastavad juhuslikud arvud. Vaatleme kahte tüüpi teste - jaotuse kontrollimise testid (kas väärtused paiknevad reaalteljel nii, nagu soovitava jaotusega juhusliku suuruse väärtused peaksid käituma) ja sõltumatuse testid (kas väärtused järgnevad üksteisele nii nagu sõltumatute katsete tulemused peaks käituma).

2.1 Testimine ühe sündmuse abil

Varasematest kursustest teame, et fikseeritud sündmuse A toimumiste arv n sõltumatus katses on binoomjaotusega $B(n, p)$, kus $p = P(A)$ on A toimumise tõenäosus ühel katsel. Kui meil on teada näiteks juhusliku suuruse X jaotusfunktsioon F_X ja A on kujul

$$A = \{a < X \leq b\},$$

siis jaotusfunktsiooni definitsioonist järeldub võrdus

$$P(A) = F_X(b) - F_X(a).$$

Kui me sooritame sõltumatuid katseid n korda ja loeme kokku A toimumiste arvu N_A , siis eelneva põhjal $N_A \sim B(n, p)$. Enamasti on põhjust etteantud jaotusele vastavuses kahelda siis, kui sündmuse A toimumiste arv erineb oluliselt oodatavast keskväärtusest $E(N_A) = np$. Kuna generaatoris kahtlemine toob kaasa lisatööd (tuleb otsida vigu, leida parem generaator vms), ei taha me enamasti liiga kergekäeliselt katsetulemusi kahtlaseks lugeda. Seega valime küllalt väikese α korral sellised arvud x_1 ja x_2 , et vaadeldavale jaotusele vastavate juhuslike arvude korral kehtib

$$P(-x_1 < N_A - np \leq x_2) \geq 1 - \alpha$$

ning loeme kahtlaseks tulemused, mis jäävad väljaspoole neid piire. Sobiv valik x_1 ja x_2 jaoks on selline, et

$$P(N_A \leq np - x_1) = \frac{\alpha}{2}$$

ja

$$P(N_A \leq np + x_2) = 1 - \frac{\alpha}{2},$$

mis puhul eelnev tingimus on samaväärne katsetulemuste põhjal leitud N_A väärtuse jäämisega vaadeldava binoomjaotuse $\frac{\alpha}{2}$ ja $(1 - \frac{\alpha}{2})$ -kvantiilide vahele.

2.2 Mitme sündmuse toimumissageduste üheaegne vaatlemine - χ^2 -test

Eeldame, et fikseeritud on m vastastikku välistavat sündmust A_1, \dots, A_m , mille tõenäosused vaadeldava jaotuse korral on p_1, \dots, p_m ($p_i > 0 \forall i$) ja mis katavad ära kõik võimalused (täissüsteem). Näiteks võime vaadelda sündmuseid $A_i = \{x_{i-1} < X \leq c_i\}$, kus

$$-\infty = c_0 < c_1 < \dots < c_{m-1} < c_m = \infty$$

ja punktid c_i , $i = 1, 2, \dots, m-1$ on valitud nii, et kõik sündmused on vaadeldava jaotuse puhul võimalikud (st positiivse tõenäosusega). Loomulik on vaadelda nende sündmuste esinemiste arvu erinevusi oodatavatest esinemiste arvudest (vastavate juhuslik suuruste keskväertustest) np_i , kusjuures väga suured erinevused on mõistlik lugeda kahtlasteks. Samas on selge, et erinevusi keskmistest ei ole hea vaadata üksikhaaval, sest vähegi suurema sündmuste komplekti korral on ka õigest jaotusest pärinevate katsetulemuste korral küllalt suur tõenäosus, et mõne sündmuse esinemissagedus erineb üsna palju oodatavast. Seega on hea vaadelda mingit summaarset oodatavatest sagedustest erinevuse mõõdikut, mille jaotus (või vähemalt asümptootiline jaotus piisavalt suure katsete arvu korral) on vaadeldavast jaotusest pärinevate juhuslike arvude korral teada. Üks selline mõõdik on χ^2 -statistik.

2.2.1 Hii-ruut test

Fikseerime m sündmusest koosneva vastastikku välistavate sündmuste süsteemi, mille toimumise tõenäosused ühel katsel on p_i , $i = 1, \dots, m$. Olgu katseseeria põhjal (ehk n genereeritud väärtuse põhjal) leitud toimumiste arvud n_i , $i = 1, \dots, m$. **Pearsoni χ^2 -statistiku** väärtus on sel juhul defineeritud valemiga

$$\chi^2 := \sum_{i=1}^m \frac{(n_i - np_i)^2}{np_i}.$$

On teada, et piiril $n \rightarrow \infty$ on see statistik χ^2 jaotusega vabadusastmete arvuga $m-1$. *NB! Enamasti (vt näiteks Wikipedia testile pühendatud artikkel inglisekeelses versioonis) loetakse, et piirjaotuse kasutamine on õigustatud ainult siis, kui n on nii suur, et $np_i \geq 5 \forall i$.*

Testi nullhüpoteesiks H_0 on, et katsetulemused vastavad juhuslikele arvudele vaadeldavast jaotusest. Alternatiivne hüpotees H_1 on see, et katsetulemused ei vasta vaadeldavale jaotusele. Otsustuskriteeriumiks on vaatlusaluste sündmuste teoreetiliste tõenäosusete abil arvutatud χ^2 -statistiku p -väärtus: kui see on väiksem kui meie poolt valitud olulisuse niivo α , siis loeme nullhüpoteesi ümberlükatuks; vastasel korral aga jääme nullhüpoteesi juurde (st generaatori testimise kontekstis vaadeldavad katsetulemused ei anna piisavalt alust generaatori õigsuses kahtlemiseks).

Ühtlasele jaotusele vastavuse testimise puhul on sellel testile oma nimi - **sageduste test**, kui sündmuseid defineerivateks reaaltelje jaotuspunktideks on $c_i = \frac{i}{m}$, $i = 1, 2, \dots, m$.

2.3 Empiirilise ja teoreetilise jaotusfunktsiooni võrdlemisel põhinevad testid

Mitme sündmuse koosvaatlemisel on raske otsustada, milliseid piirkondi valida. Et selle küsimusega mitte oma pead vaevata, on loomulik vaadelda korraga komplekti lõpmatult

paljudest piirkondadest kujul $A_x = \{X \leq x\}$, kus x reaalarv. Sellisesse piirkonda sattumise tõenäosus vaadeldava jaotuse korral on $F_X(x)$ (kus F_X on X jaotusfunktsioon). Valimi x_1, \dots, x_n korral piirkonda A_x kuuluvate vaatluste suhteline sagedus vastab **empiirilisele jaotusfunktsioonile**

$$F_n(x) := \frac{\#\{i \mid x_i \leq x\}}{n},$$

kus $\#B$ tähistab hulga B elementide arvu. Seega konkreetse x korral iseloomustab vahe $F_n(x) - F(x)$ sündmuse A_x empiirilise esinemissageduse erinevust teoreetilisest tõenäosusest ning suur erinevus võiks anda alust kahtlusele, et andmed ei vasta vaadeldava jaotusega juhuslikele arvudele. Sündmuste korruga vaatlemiseks aga on vaja kasutada mingit mõõdikut, mis kõikidele reaalarvudele x vastavad erinevused kokku võtaks ning selleks on mitmeid võimalusi. Vaatleme neist kahte, kuid eelnevalt tõestame ühe tulemuse, mis on järgnevate testide tööpõhimõttest arusaamisel äärmiselt oluline

Teoreem 6 *Olgu F pidev jaotusfunktsioon ning olgu X sellele vastav juhuslik suurus. Siis $Y = F(X)$ on ühtlase jaotusega $U(0, 1)$*

Tõestus. Näitame, et juhusliku suuruse Y jaotusfunktsioon F_Y on ühtlase jaotuse jaotusfunktsioon, st avaldub kujul

$$F_Y(y) = \begin{cases} 0, & \text{kui } y \leq 0, \\ y, & \text{kui } 0 < y < 1 \\ 1 & \text{kui } y \geq 1. \end{cases}$$

Selleks paneme tähele, et

$$F_Y(y) \stackrel{\text{def}}{=} P(\{Y \leq y\}) = P(\{F(X) \leq y\}) = P(\{X \in A_y\}),$$

kus hulk A_y on selliste reaalarvude hulk, mille korral F väärtus ei ole suurem y väärtusest, ehk

$$A_y = \{x \in \mathbb{R} : F(x) \leq y\}.$$

Vaatleme eraldi juhte $0 < y < 1$, $y \leq 0$ ja $y \geq 1$

juht $0 < y < 1$: Jaotusfunktsiooni omadusest 3 (Vaata teoreem 21) ja F pidevusest järeldub, et leidub selline reaalarv a_y , mille korral $F(a_y) = y$, seega A_y ei ole tühihulk. Jaotusfunktsiooni monotoonse kasvamise omadusest (Teoreem 21, omadus 2) järeldub võrratusest $x \leq a_y$ võrratus $F(x) \leq y$, seega

$$(-\infty, a_y] \subset A_y.$$

Tõenäosuse monotoonsuse omadusest (Teoreem 18, omadus 3) järeldub seega võrratus

$$y = F(a_y) = P(\{X \in (-\infty, a_y]\}) \leq P(\{X \in A_y\}) = F_Y(y).$$

Veendume ka vastupidise võrratuse kehtimises. Olgu $\varepsilon > 0$ selline reaalarv, mille korral $y + \varepsilon < 1$. Jällegi saame F pidevuse tõttu leida sellise $a_{y+\varepsilon}$, et $F(a_{y+\varepsilon}) = y + \varepsilon$. Kuna iga $x \geq a_{y+\varepsilon}$ korral kehtib F monotoonse kasvamise tõttu $F(x) \geq y + \varepsilon > y$, siis $A_y \subset (-\infty, a_{y+\varepsilon}]$. Seega jällegi tõenäosuse monotoonsuse põhjal

$$F_Y(y) = P(\{X \in A_y\}) \leq P(\{X \in (-\infty, a_{y+\varepsilon}]\}) = F(a_{y+\varepsilon}) = y + \varepsilon.$$

Piiril $\varepsilon \rightarrow 0$ saame siit $F_Y(y) \leq y$ ning eelnevat vastupidist võrratust arvestades olema näidanud, et kehtib võrdus

$$F_Y(y) = y, \quad 0 < y < 1.$$

juht $y \leq 0$: iga reaalarve $\varepsilon \in (0, 1)$ korral kehtib jaotusfunktsiooni monotoonsuse tõttu vaadeldaval juhul võrratus $F_Y(y) \leq F_Y(\varepsilon) = \varepsilon$. Piiril $\varepsilon \rightarrow 0$ saame siit $F_Y(y) \leq 0$. Jaotusfunktsiooni mittenegatiivsusest järeldub siit, et kehtib võrdus

$$F_Y(y) = 0, \quad y \leq 0.$$

juht $y \geq 1$: Kuna iga jaotusfunktsioon rahuldab tingimust $F(x) \leq 1 \quad \forall x \in \mathbb{R}$, siis A_y on kogu reaaltelg ja seetõttu X väärtuse sattumine hulka A_y on kindel sündmus. Järelikult $y \geq 1$ korral

$$F_Y(y) = 1.$$

Sellega on teoreemi väide tõestatud. \square

2.3.1 Kolmogorov-Smirnovi test

Kolmogorov-Smirnovi test võimaldab otsustada, kas valim pärineb etteantud pidevale jaotusfunktsioonile F vastavast jaotusest. Testi nullhüpoteesiks on see, et valim koosneb jaotusfunktsioonile F vastavatest juhuslikest arvudest (sõltumatute katsete tulemustest). Alternatiivne hüpotees H_1 on see, et valim ei koosne jaotusfunktsioonile F vastavatest juhuslikest arvudest. Otsustamiseks kasutatakse statistiku

$$K_n := \sup_{x \in \mathbb{R}} |F_n(x) - F(x)|$$

valimi põhjal arvutatud väärtuse p -väärtust ehk tõenäosust, et jaotusfunktsioonile F vastavate juhuslike arvude korral saadakse vähemalt sama suur väärtus, kui vaadeldava valimi korral. Kui p -väärtus on väiksem, kui meie valitud olulisuse nivoo α , loeme tõestatuks alternatiivse hüpoteesi.

Veendume, et statistiku jaotus nullhüpoteesi kehtimisel ei sõltu vaadeldavast jaotusfunktsioonist F . Tähistagu $x_{(i)}$ valimi suuruselt i -ndat elementi, kusjuures defineerime lisaks

$$x_{(0)} = -\infty, \quad x_{(n+1)} = \infty.$$

Kuna

$$\sup_{x \in \mathbb{R}} |F_n(x) - F(x)| = \max_{i=0, \dots, n} \sup_{x \in [x_{(i)}, x_{(i+1)})} |F_n(x) - F(x)|,$$

siis leiame järgmiseks supreemumid üle osalõikude.

Empiirilise jaotusfunktsiooni definitsioonist järeldub, et

$$F_n(x) = \frac{i}{n}, \quad i = 0, \dots, n, \quad x \in [x_{(i)}, x_{(i+1)}),$$

seega

$$\sup_{x \in [x_{(i)}, x_{(i+1)})} |F_n(x) - F(x)| = \sup_{x \in [x_{(i)}, x_{(i+1)})} |F(x) - \frac{i}{n}|.$$

Funktsiooni F mittekahanemisest järeldub, et

$$F(x_{(i)}) - \frac{i}{n} \leq F(x) - \frac{i}{n} \leq F(x_{(i+1)}) - \frac{i}{n}, \quad x \in [x_{(i)}, x_{(i+1)}),$$

kusjuures F pidevuse tõttu

$$\sup_{x \in [x_{(i)}, x_{(i+1)})} F(x) - \frac{i}{n} = F(x_{(i+1)}) - \frac{i}{n}.$$

Analüüsidest juhte $F(x_{(i+1)}) < \frac{i}{n}$, $F(x_{(i)}) < \frac{i}{n}$, $F(x_{(i+1)}) \geq \frac{i}{n}$, $F(x_{(i)}) \geq \frac{i}{n}$, saame kõigil juhtudel, et

$$\sup_{x \in [x_{(i)}, x_{(i+1)})} |F_n(x) - F(x)| = \max\left\{\frac{i}{n} - F(x_{(i)}), F(x_{(i+1)}) - \frac{i}{n}\right\}.$$

Seega kehtib võrdus

$$K_n = \max_{i=1, \dots, n} \left\{F(x_{(i)}) - \frac{i-1}{n}, \frac{i}{n} - F(x_{(i)})\right\}.$$

Kuna aga $F(X)$ on ühtlase jaotusega $U(0, 1)$ juhuslik suurus, siis $y_i = F(x_{(i)})$ on nullhüpoteesi korral sõltumatu valim jaotusest $U(0, 1)$, $y_{(i)} = F(x_{(i)})$ on selle valimi suurusel i -s element ning seega statistiku jaotus ei sõltu vaadeldavast jaotusfunktsioonist, vaid on leitav statistiku käitumise uurimisega jaotusest $U(0, 1)$ pärineva valimi korral.

2.3.2 Oomega-ruut test

Erinevust $F_n(x) - F(x)$ saab mõõta ka teisiti. Kui X on pidev juhuslik suurus jaotusfunktsiooniga F ja tihedusfunktsiooniga f , siis Cramer - von Mises test ehk oomega-ruut test põhineb järgmisel F ja F_n erinevust mõõtvast statistikul:

$$\omega_n^2 := n \int_{-\infty}^{\infty} (F(x) - F_n(x))^2 f(x) dx.$$

Saab näidata (vt Kollo õpik, lk 35-36), et kehtib

$$\omega_n^2 = \frac{1}{12n} + \sum_{i=1}^n \left(F(x_{(i)}) - \frac{i-0.5}{n} \right)^2,$$

seega jällegi ei sõltu statistiku jaotus nullhüpoteesi kehtimisel vaadeldavast jaotusfunktsioonist F , vaid on leitav ühtlasest jaotusest pärineva valimi analüüsimise teel. Otsustamise eeskiri (hüpoteesid) on samad, kui KS testi puhul.

2.4 Sõltumatuse testimine

Eelnevad testid arvestavad ainult valimis olevate arvude paiknemist reaalteljel, kuid ei arvesta nende järgnevusega seotud infot. Näiteks saaksime samad testitulemused, kui eelnevalt sorteeriksime valimi kasvamise järjekorras. Samas sõltumatute katsete tulemuste puhul on oluline ka järgnevus, st eelmiste väärtuste teadmine ei tohi anda mingit infot järgmiste väärtuste kohta ja kui generaator väljastab väärtuseid kasvavas järjekorras, siis ei vasta need sõltumatute katsete tulemustele.

Järgnevalt vaatleme teste, mis arvestavad nii soovitud jaotusega seotud tõenäosuseid kui ka sõltumatute katsete puhul kehtivaid järjestikuste katsetulemustega seotud omadusi.

2.4.1 Sõltumatute paaride test

Sõltumatute katsete põhiomaduseks on see, et neile vastavate sündmuste koos toimumise tõenäosus on vastavate sündmuste tõenäosuste korrutis. Sõltumatute paaride testi korral vaadeldakse järjestikuste väärtuste paare ehk liitkatset, kus iga kord tekib arvupaar. Testi sooritamise algoritm on järgmine:

1. Fikseerime m vastastikku välistavat sündmust A_1, \dots, A_m , mille tõenäosused vaadeldava jaotuse korral on p_1, \dots, p_m ja mis katavad ära kõik võimalused (täissüsteem)
2. Igale katsetulemusele seame vastavusse sündmuse numbri, mis sellel katsel toimus. Saame numbritest $1, \dots, m$ koosneva jada k_1, k_2, \dots, k_n .
3. Vaatleme paare $(k_1, k_2), (k_3, k_4), \dots$. Kui tegemist on sõltumatute katsetega, siis paari (i, j) saamise tõenäosus igal liitkatsel on $p_i \cdot p_j$
4. Kontrollime χ^2 -testiga, kas arvupaaride (i, j) , $i, j \in \{1, \dots, m\}$ esinemissagedused on kooskõlas nende saamise tõenäosusega.

Selgituseks: kui meil on $2n$ väärtusest koosnev valim, siis võime nende põhjal tekitatud arvupaare (k_{2i-1}, k_{2i}) , $i = 1, \dots, n$ vaadelda kui valimit n liitkatse tulemustest, mille korral sündmused A_{ij} = "liitkatse puhul saadakse paar (i, j) " moodustavad täissüsteemi. Seega on tegemist just sellise juhuga, mille korral saab hii-ruut testi rakendada.

Kui testi p -väärtus on väiksem kui meie olulisuse nivoo, siis loeme tõestatuks, et katsetulemused ei vasta vaadeldava jaotusega juhuslikele (või pseudojuhuslikele) arvudele. Mittevastavuse põhjuseks võib olla kas etteantud jaotusele mittevastamine (tõenäosused ei klapi) või siis see, et katsetulemused ei ole sõltumatud.

Märkus: χ^2 -testil on olemas versioon, mis kontrollib ainult sõltumatust. Kontroll seisneb andmete põhjal sündmuste A_i , $i = 1, \dots, m$ tõenäosuste hindamises ja kontrollis, kas arvupaaride (i, j) esinemissagedused on kooskõlas hinnatud tõenäosuste korrutisega hii-ruut statistiku abil, kus teoreetilised tõenäosused on asendatud hinnatud tõenäosustega. Kuna siin kursusel me kontrollime teadaolevale jaotusele vastamist, siis seda testi versiooni me ei kasuta.

2.4.2 Vahemike test

Vaatleme järgmist tegevust:

- Fikseerime mingi sündmuse A , mille toimumise tõenäosus ühel katsel on p
- teostame katseid, kuni sündmus A toimub, paneme kirja selleks kuluvate katsete arvu
- teostame jälle katseid, kuni sündmus A toimub, paneme kirja selleks kuluvate katsete arvu
- jne

Vaadeldes valimit kui sellise tegevuse teostamisel kirjepandud katsetulemusi, saame teha kindlaks, mitu eelpool kirjeldatud katseseeriat teostati (viimane võib olla poolik ja seda me ei arvesta) ning samuti ka iga seeria korral saadud tulemuse (vajaläinud katsete arvu). Seega saame esialgsest valimist uue (lühema) valimi liitkatsete tulemustega. Teame, et katsete arv kuni fikseeritud sündmuse toimumiseni on geomeetrilise jaotusega $Geom(p)$, seega uus valim peaks vastama valimile jaotusest $Geom(p)$. Fikseerides nüüd uue valimi korral ℓ sündmust kujul "A tulekuks kulub i katset", kus $i = 1, 2, \dots, \ell - 1$ ja "A tulekuks kulub vähemalt ℓ katset", saame hii-ruut testika kontrollida tulemuste vastavust jaotusele $Geom(p)$. (NB! Leia geomeetrilise jaotuse definitsiooni põhjal viimati defineeritud sündmuse tõenäosus!)

2.4.3 Seeriaste test

Vahemike test nõuab ühe sündmuse fikseerimist. Samas võime vaadelda suurema arvu sündmuste korral seda, kui pikad on sama sündmuse järjestikustel katsel kordumiste jadad. Võime mõelda järgmisest tegevusest:

- Fikseerime m vastastikku välistavat võrdvõimalikku sündmust A_1, \dots, A_m
- Teeme ühe katse, paneme kirja selle korral toimunud sündmuse numbriga
- Teeme katseid, kuni tuleb kirjapandud numbriga sündmusest erinev sündmus. Fikseerime selleks kulunud katsete arvu (ehk kirjapandud numbriga sündmuste tulemise seeria pikkuse) ja paneme kirja viimasel katsel saadud sündmuse numbriga
- Kordame eelmises kirjeldatud tegevusi soovitud arvu kordi.

Kuna iga seeria korral teeme sõltumatuid katseid, kuni toimub kirjapandud numbriga sündmusest erinev sündmus, siis on igal katsel seeria lõppemise tõenäosus $\frac{m-1}{m}$ ja seega on vajaminevate katsete arv jaotusega $Geom(\frac{m-1}{m})$. Kui vaadelda algset valimit kui sellise tegevuse käigus kirjapandud katsetulemusi, saame nende põhjal kindlaks teha teostatud katsevoorude arvu ja nende käigus leitud seeriaste pikkused. Kui esialgne valim koosneb vaadeldavale jaotusele vastavatest juhuslikest arvudest, siis leitud seeriaste pikkused moodustavad valimi jaotusest $Geom(\frac{m-1}{m})$. Saadud seeriaste pikkuste valimi vastavust geomeetrilisele jaotusele saame testida samamoodi nagu vahemike testi puhul

2.4.4 Pokker-test

Pokker-testi algoritm on järgmine:

1. Fikseerime m vastastikku välistavat võrdvõimalikku sündmust A_1, \dots, A_m
2. Igale katsetulemusele seame vastavusse sündmuse numbriga, mis sellel katsel toimus
3. vaatleme viisikuid $(k_1, k_2, k_3, k_4, k_5), \dots$
4. loendame erinevate mustrite esinemisi, mustriteks on: kõik erinevad, üks paar, kaks paari, kolmik, paar ja kolmik, nelik, viisik
5. Kasutame χ^2 -testi mustrite esinemissageduste võrdlemisel teoreetiliste sagedustega

Mustrite teoreetilised esinemissagedused saab leida kombinatorika valemeid kasutades (vt Kollo õpikust)

Peatükk 3

Etteantud jaotusele vastavate pseudojuhuslike arvude genereerimine

Järgnevalt eeldame, et me oskame genereerida ühtlasele jaotusele $U(0,1)$ vastavaid pseudojuhuslike arve või meil on füüsiline generaator, mis väljastab sellele jaotusele vastavaid juhuslike arve. Vaatleme võimalusi, kuidas nende abil saab tekitada suvalisele jaotusele vastavaid (pseudo)juhuslike arve.

3.1 Jaotusfunktsiooni pööramise meetod

Kuna eelneva põhjal teame, et suvalise pideva jaotusfunktsiooniga juhuslikust suuruselt saame sellele jaotusfunktsiooni rakendades ühtlase jaotusega $U(0,1)$ juhusliku suuruse, siis võib loota, et jaotusfunktsiooni pöördfunktsiooni abil saab ühtlasest jaotusest soovitud jaotusega juhusliku suuruse. Osutub, et see on alati õige, kui pöördfunktsiooni mõistet sobivalt üldistada juhtudele, kus tavalist pöördfunktsiooni ei leidu.

Definitsioon 7 Olgu $F : R \rightarrow [0,1]$ mittekahanev funktsioon. Siis selle üldistatud pöördfunktsiooniks nimetatakse funktsiooni $F^- : R \rightarrow \bar{R}$ (kus $\bar{R} = R \cup \{-\infty, \infty\}$), mis on defineeritud võrdusega

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

Eelnevas definitsioonis on kasutatud kokkulepet, et infimum tühjast hulgast on $+\infty$.

Eelneva definitsiooni üle järele mõeldes on küllalt lihtne veenduda, et juhul, kui vaadeldava y korral võrrandil $F(x) = y$ on ühene lahend $x = x_y$, siis $F^-y = x_y$, seega pöördfunktsiooni olemasolul langevad üldistatud pöördfunktsioon ja tavaline pöördfunktsioon kokku. Kui selliseid lahendeid on aga mitu, siis pöördfunktsiooni väärtuseks võetakse lahendite hulga alumine raja ning kui lahendeid pole, siis on pöördfunktsiooni väärtuseks see x väärtus, mille korral funktsiooni graafik "hüppab" läbi taseme y .

Järgnev tulemus lihtsustab meil põhitulemuse tõestamist.

Lemma 8 Olgu $F : R \rightarrow [0,1]$ mittekahanev ja paremalt pidev funktsioon. Siis selle üldistatud pöördfunktsiooni korral on võrratus $F^-(y) \leq x$ samaväärne võrratusega $y \leq F(x)$.

Tõestus. Samaväärsuse näitamiseks on vaja näidata, et esimesest võrratusest järeldub teine ning et teisest järeldub esimene.

- a) Olgu x ja y sellised, et kehtib võrratus $F^-(y) \leq x$. Funktsiooni F^- definitsiooni kohaselt leidub siis selline monotoonselt kahanev jada a_n , et $\lim_{n \rightarrow \infty} a_n \leq x$ ning $F(a_n) \geq y$. Defineerides $b_n = \max(a_n, x)$ saame uue monotoonselt kahaneva jada, mis funktsiooni F mittekahanevuse tõttu rahuldab samuti tingimust $F(b_n) \geq y \forall n$ ning samuti kehtib $\lim_{n \rightarrow \infty} b_n = x$. Kuna F on paremalt pidev, kehtib seega ka

$$F(x) = \lim_{n \rightarrow \infty} F(b_n) \geq y.$$

- b) Olgu x ja y sellised, et kehtib võrratus $y \leq F(x)$. Siis

$$x \in \{a : F(a) \geq y\}$$

ning seetõttu

$$F^-(y) = \inf\{a : F(a) \geq y\} \leq x.$$

Lemma on tõestatud. \square

Üldistatud pöördfunktsiooni abil saab esitada järgmise olulise tulemuse.

Teoreem 9 *Olgu F mingi tõenäosusjaotuse jaotusfunktsioon ning olgu F^- selle üldistatud pöördfunktsioon. Siis juhusliku suuruse $X = F^-(Y)$, kus $Y \sim U[0, 1]$, jaotusfunktsiooniks on F .*

Tõestus. Leiame juhusliku suuruse X jaotusfunktsiooni. Definitsiooni kohaselt

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

Olgu elementaarsündmuste ruumiks Ω . Kuna lemma 8 kohaselt

$$\{\omega \in \Omega : F^-(Y(\omega)) \leq x\} = \{\omega \in \Omega : Y(\omega) \leq F(x)\},$$

siis

$$F_X(x) = P(\{Y \leq F(x)\}) = F_Y(F(x)).$$

Võttes arvesse, et ühtlase jaotusega $U[0, 1]$ juhusliku suuruse Y jaotusfunktsiooni korral $F_Y(y) = y$, $0 \leq y \leq 1$, siis olemegi näidanud, et

$$F_X(x) = F(x) \quad \forall x \in \mathbb{R}. \square$$

Märkus: Tõenäosusteooria seisukohalt ei ole vahet, kas punkt valitakse juhuslikult lõigust $[0, 1]$ või vahemikust $(0, 1)$, sest otspunktide saamise tõenäosus on esimesel juhul 0 ja seetõttu ei muutu tõenäosusarvutuse seisukohalt midagi, kui need võimalused ära jätta. Et aga $F^-(0) = -\infty$ ja $F^-(1)$ võib olla $+\infty$, siis on jaotusfunktsiooni pöördfunktsiooni meetodi rakendamisel mugavam kasutada jaotusega $U(0, 1)$ (otspunktid välistatud) juhuslikke suuruseid

3.1.1 Lõpliku arvu väärtustega diskreetse juhusliku suuruse esitamine ühtlase jaotuse kaudu

Vaatleme juhuslikku suurust, mille võimalikeks väärtusteks on $c_1 < c_2 \dots < c_m$, olgu p_1, p_2, \dots, p_m vastavad tõenäosused. Sellise juhusliku suuruse jaotusfunktsioon on tükiti konstantne, katkevuspunktideks on c_1, \dots, c_m ja hüpete kõrguseks p_1, \dots, p_m . Seega avaldub üldistatud pöördfunktsioon vahemikus $y \in (0, 1)$ kujul

$$F^-(y) = \begin{cases} c_1, & 0 < y \leq p_1, \\ c_2, & p_1 < y \leq p_1 + p_2 \\ \dots & \\ c_m, & \sum_{i=1}^{m-1} p_i < y < 1. \end{cases}$$

Seega saame soovitud jaotusega juhusliku suuruse väärtusi genereerida, genereerides ühtlase jaotusega $U(0, 1)$ juhuslikke suuruseid ning rakendades neile eelneva valemi abil defineeritud üldistatud pöördfunktsiooni.

3.2 Valikumeetodid

Selle asemel, et kohe soovitud jaotusest arve genereerida, võime proovida seda tegevust jaotada lihtsamateks osadeks: genereerime punkte mõnest tuntud jaotusest ja siis "hõrendame" neid nii, et erinevatesse piirkondadesse jäävate punktide osakaalud vastaksid soovitud jaotusele. See ongi valikumeetodite idee.

3.2.1 Lihtne valikumeetod

Vaatleme pidevat juhuslikku suurust tihedusfunktsiooniga f . Eeldame, et f on 0 väljaspool lõiku $[a, b]$ ning samuti eeldame, et f on ülalt tõkestatud konstandiga c , st $f(x) \leq c \forall x$. Valikumeetodi algoritmi intuiitiivne kirjeldus on järgmine

1. Pommitame riskülikut $[a, b] \times [0, c]$ selles juhuslikult valitud punktidega (vastavalt ühtlasele jaotusele)
2. Jätame alles ainult need, mis satuvad $y = f(x)$ alla X väärtusteks võtame saadud punktide x -koordinaadid

Osutub, et nii saame valimi tihedusfunktsioonile f vastavast jaotusest. Veelgi enam, algoritmi võib rakendada suvalise mittenegatiivse tõkestatud funktsiooni f korral ning tulemuseks saame valimi jaotusest, mille tihedusfunktsioon on proportsionaalne funktsiooniga f . Tõestame vastava tulemuse.

Teoreem 10 *Eeldame, et mittenegatiivne funktsioon f on 0 väljaspool lõiku $[a, b]$ ning et f on ülalt tõkestatud konstandiga c , st $f(x) \leq c \forall x$. Olgu $Y_1 \sim U(a, b)$ ja $Y_2 \sim U(0, c)$ sõltumatud juhusliud suurused. Defineerime juhusliku suuruse X järgmiselt:*

$$X = Y_1 \text{ tingimusel, et } Y_2 \leq f(Y_1)$$

Siis X on pidev juhuslik suurus tihedusfunktsiooniga, mis on proportsionaalne funktsiooniga f .

Tõestus. Leiame juhusliku suuruse X jaotusfunktsiooni avaldise. Definitsiooni kohaselt

$$F_X(x) = P(X \leq x) = P(Y_1 \leq x \mid Y_2 \leq f(Y_1)) = \frac{P(Y_1 \leq x, Y_2 \leq f(Y_1))}{P(Y_2 \leq f(Y_1))}.$$

Tähistame

$$D_x = \{(y_1, y_2) \in \mathbb{R}^2 : y_1 \leq x, y_2 \leq f(y_1)\},$$

siis tingimus $\{Y_1 \leq x, Y_2 \leq f(Y_1)\}$ on esitatav kujul $(Y_1, Y_2) \in D_x$. Teame, et

$$f_{Y_1}(y_1) = \begin{cases} \frac{1}{b-a}, & y_1 \in [a, b], \\ 0, & y_1 \notin [a, b], \end{cases} \quad f_{Y_2}(y_2) = \begin{cases} \frac{1}{c}, & y_2 \in [0, c], \\ 0, & y_2 \notin [0, c]. \end{cases}$$

Kuna Y_1 ja Y_2 sõltumatud, siis juhusliku vektori (Y_1, Y_2) tihedusfunktsioon avaldub kujul

$$f_{Y_1, Y_2}(y_1, y_2) = f_{Y_1}(y_1)f_{Y_2}(y_2), \quad y_1, y_2 \in \mathbb{R}$$

ning pideva juhusliku vektori tihedusfunktsiooni omaduse 3 (vt lemma 24) kohaselt

$$P((Y_1, Y_2) \in D_x) = \int_{-\infty}^x \left(\int_{-\infty}^{f(y_1)} f_{Y_2}(y_2) dy_2 \right) f_{Y_1}(y_1) dy_1 = \int_{-\infty}^x \frac{f(y_1)}{c} f_{Y_1}(y_1) dy_1.$$

Arvestades, et $f_{Y_1}(y_1)$ on nullist erinev ainult piirkonnas $y_1 \in [a, b]$, tuleb eelneva integraali arvutamisel vaadelda kolme juhtu: $x < a$, $x \in [a, b]$, $x > b$. Seega saame

$$P((Y_1, Y_2) \in D_x) = \begin{cases} 0, & x < a, \\ \int_a^x \frac{f(y_1)}{c(b-a)} dy_1, & x \in [a, b], \\ \int_a^b \frac{f(y_1)}{c(b-a)} dy_1, & x > b. \end{cases}$$

Tähistades

$$c_f = \int_a^b f(y) dy,$$

saame

$$P(Y_2 \leq f(Y_1)) = P((Y_1, Y_2) \in D_\infty) = \frac{c_f}{c(b-a)}$$

ning seega

$$F_X(x) = \begin{cases} 0, & x < a, \\ \int_a^x \frac{f(y)}{c_f} dy, & x \in [a, b], \\ 1, & x > b \end{cases}$$

Arvestades, et

$$f_X(x) = F'_X(x) = \begin{cases} \frac{f(x)}{c_f}, & x \in [a, b], \\ 0, & x \notin [a, b], \end{cases}$$

on teoreemi väide tõestatud. \square

Märkus. Kasutades teadmist, et $U \sim U(0, 1)$ korral $Y_1 = a + (b-a)U$ puhul $Y_1 \sim U(a, b)$ ja $Y_2 = cU$ puhul $Y_2 \sim U(0, c)$, saab valikumeetodit teostada, lähtudes jaotusele $U(0, 1)$ vastavatest pseudojuhuslikest arvudest.

Valikumeetodi *efektiivsuseks* nimetatakse X defineerimisel kasutatud sündmuse toimumise tõenäosust, st lihtsa valikumeetodi puhul arvu

$$P(Y_2 \leq f(Y_1)) = \frac{c_f}{c(b-a)}$$

3.2.2 Üldine valikumeetod

Lihtne valikumeetod võimaldab genereerida ainult tõkestatud väärtustega juhuslikke suuruseid. Suur osa huvipakkuvatest jaotustest on aga sellised, mille tihedusfunktsioon on nullist erinev kuitahes suurte (või väikeste) argumentide väärtuste korral, mistõttu ka vastava juhusliku suuruse väärtuste piirkond ei ole tõkestatud. Osutub, et valikumeetodit saab üldistada selliste juhtudega tegelemiseks

Teoreem 11 *Olgu f ja g mittenegatiivsed ning lõpliku graafikualust pindala omavad funktsioonid ja c selline reaalarv, et kehtib tingimus $f(x) \leq cg(x) \forall x$. Olgu Y_1 ja Y_2 sõltumatud juhuslikud suurused, kusjuures Y_1 on tihedusfunktsioon on kujul*

$$f_{Y_1}(y_1) = \frac{g(y_1)}{c_g}, \quad c_g = \int_{-\infty}^{\infty} g(x) dx$$

ja Y_2 on ühtlase jaotusega $U(0, 1)$. Defineerime juhusliku suuruse

$$X = Y_1 \mid \{g(Y_1) > 0, Y_2 \leq \frac{f(Y_1)}{c \cdot g(Y_1)}\}.$$

Siis juhuslik suurus X on funktsiooniga f proportsionaalsele tihedusfunktsioonile vastava jaotusega.

Tõestus. Tõestuse on analoogiline lihtsa valikumeetodi puhul toodud tõestusega, mistõttu selle läbikirjutamine jääb harjutuseks lugejale. \square

Lihtne on veenduda, et üldise valikumeetodi efektiivsus avaldub kujul

$$P(Y_2 \leq \frac{f(Y_1)}{cg(Y_1)}) = \frac{c_f}{c \cdot c_g},$$

kus

$$c_f = \int_{-\infty}^{\infty} f(x) dx.$$

3.3 Genereerimine segujaotusest

Sageli võib andmete jaotusi uurides täheldada mitut "küüru" või lihtsalt erinevat tüüpi käitumist erinevates väärtuste piirkondades, mistõttu ei vasta nende käitumine lihtsatele tuntud jaotustele. Tihti on selline käitumine põhjustatud üldkogumi mittehomogeensusest ehk sellest, et ülkogumis on mitu erinevat rühma (näiteks mehed/naised), milles mõõdetav tunnus käitub erinevalt. Olgu meil m rühma, kusjuures igas rühmas käitugu uuritav tunnus vastavalt pidevale juhuslikule suurusele tihedusfunktsiooniga f_i , $i = 1, 2, \dots, m$. Olgu iga rühma osakaal üldkogumis p_i , siis täistõenäosuse valemi (vaata lemma 19) abil on lihtne näidata, et uuritava tunnuse tihedusfunktsioon üldkogumist juhusliku liikme valikul avaldub kujul

$$f(x) = \sum_{i=1}^m p_i f_i(x), \quad -\infty < x < \infty.$$

Tõestame selle valemi.

Teoreem 12 Olgu üldkogumis m rühma, kusjuures igas rühmas käitugu uuritav tunnus vastavalt pidevale juhuslikule suurusele tihedusfunktsiooniga f_i , $i = 1, 2, \dots, m$. Siis uuritava tunnuse tihedusfunktsioon üldkogumist juhusliku liikme valikul avaldub kujul

$$f(x) = \sum_{i=1}^m p_i f_i(x), \quad -\infty < x < \infty.$$

Tõestus Olgu X uuritava tunnuse väärtus üldkogumist juhuslikult valitud liikme korral ning olgu A_i sündmus, et valitud liige kuulub rühma i . Tähistagu F_i tihedusfunktsioonile f_i vastavat jaotusfunktsiooni. Täistõenäosuse valemi kohaselt

$$F_X(x) = P(X \leq x) = \sum_{i=1}^m P(A_i)P(X \leq x | A_i).$$

Ülesande teksti põhjal on juhusliku suuruse X tinglik jaotus tingimusel, et valituks osutub rühma A_i liige, antud jaotusfunktsiooniga F_i , seega

$$F_X(x) = \sum_{i=1}^m p_i F_i(x).$$

Kuna tihedusfunktsioon on jaotusfunktsiooni tuletis, saame nüüd

$$f_X(x) = \sum_{i=1}^m p_i F_i'(x) = \sum_{i=1}^m p_i f_i(x), \quad -\infty < x < \infty.$$

Teoreem on tõestatud. \square

Siit saame algoritmi etteantud tihedusfunktsiooniga juhusliku suuruse väärtuste genereerimiseks segujaotuse meetodil:

1. Leiame etteantud tihedusfunktsiooni f esituse tuntud tihedusfunktsioonide f_i kaudu kujul

$$f(x) = \sum_{i=1}^m p_i f_i(x), \quad x \in \mathbb{R}$$

2. Genereerime i väärtuse jaotusest (j, p_j) , $j = 1, \dots, m$
3. Genereerime X väärtuse tihedusele f_i vastavast jaotusest.

3.4 Pideva juhusliku suuruse teisendused

Mõnikord on võimalik soovitud jaotusega juhuslikku suurust lihtsalt saada nii, et rakendada tuntud jaotusega juhuslikule suurusele mingit funktsiooni. Rangelt monotoonsete funktsioonide rakendamisel on võimalik anda lihtsa seose esialgse ja teisendatud juhusliku suuruse tihedusfunktsioonide vahel.

Lemma 13 Kui X on pidev juhuslik suurus tihedusfunktsiooniga f ja g on rangelt monotoonne funktsioon pöördfunktsiooniga h (st $h(g(x)) = x \quad \forall x$), siis juhusliku suuruse $Y = g(X)$ tihedusfunktsioon avaldub kujul

$$f_Y(y) = \begin{cases} f_x(h(y))|h'(y)|, & y \in g(R) \\ 0, & y \notin g(R) \end{cases}$$

Tõestus. Tõestus on harjutuseks lugejale. Soovitus: kõigepealt saab leida seose jaotusfunktsioonide vahel ja seejärel tuletise abil seose tihedusfunktsioonide vahel. Eraldi tuleb käsitleda rangelt kasvavat funktsiooni g ja rangelt kahanevat funktsiooni g , sest võrratusele rangelt kahaneva funktsiooni rakendamisel saame esialgsuga vastupidise võrratuse. \square

3.5 Genereerimine tinglikust jaotusest

Vaatleme juhtu, kus me tahame genereerida juhusliku suuruse väärtuseid tingimusel, et need väärtused jäävad teatud vahemikku. Leiame juhusliku suuruse $Y = X \mid \{a < X \leq b\}$ jaotusfunktsiooni:

$$F_Y(y) = P(Y \leq y) = P(X \leq y \mid a < X \leq b) = \frac{P(a < X \leq y)}{P(a < X \leq b)}.$$

Kuna

$$P(a < X \leq b) = F_X(b) - F_X(a)$$

ning

$$\{X \leq y, a < X \leq b\} = \begin{cases} \emptyset, & y \leq a, \\ \{a < X \leq y\}, & y \in (a, b], \\ \{a < X \leq b\} & y > b \end{cases},$$

siis

$$F_Y(y) = \begin{cases} 0, & y \leq a, \\ \frac{F_X(y) - F_X(a)}{F_X(b) - F_X(a)}, & a < y \leq b, \\ 1, & y > b. \end{cases}$$

Kui X on pidev juhuslik suurus, siis siit järeldub, et tingliku jaotuse tihedusfunktsioon avaldub kujul

$$f_Y(y) = \begin{cases} 0, & y \leq a, \\ \frac{f_X(y)}{F_X(b) - F_X(a)}, & a < y \leq b, \\ 0, & y > b. \end{cases}$$

Osutub, et tinglikust jaotusest on lihtne väärtuseid genereerida, kasutades jaotusfunktsiooni pöördfunktsiooni.

Teoreem 14 *Olgu X juhuslik suurus jaotusfunktsiooniga F_X . Olgu U juhuslik suurus jaotusega $U(F_X(a), F_X(b))$, kus $a < b$ on mingid reaalarvud. Siis juhuslik suurus $Y = F_X^{-1}(U)$, kus F_X^{-1} on funktsiooni F_X üldistatud pöördfunktsioon, on sama jaotusega kui $X \mid \{a < X \leq b\}$*

Tõestus. Leiame juhusliku suuruse Y jaotusfunktsiooni. Lemmat 8 kasutades saame

$$F_Y(y) = P(F_X^{-1}(U) \leq y) = P(U \leq F_X(y)).$$

Ühtlase jaotuse jaotusfunktsiooni avaldist kasutades saame nüüd

$$F_Y(y) = \begin{cases} 0, & F_X(y) \leq F_X(a), \\ \frac{F_X(y) - F_X(a)}{F_X(b) - F_X(a)}, & F_X(a) < F_X(y) \leq F_X(b), \\ 1, & F_X(y) > F_X(b). \end{cases}$$

Arvestades funktsiooni F_X monotoonset kasvamist, on lihtne veenduda, et eelnev valem annab iga y korral sama väärtuse, mis tingliku jaotuse $X \mid \{a < X \leq b\}$ jaotusfunktsiooni valem, seega on teoreem tõestatud. \square

Peatükk 4

Pseudojuhuslike vektorite genereerimine

4.1 Pidevate juhuslike vektorite teisendused

Teoreem 15 Olgu X pidev juhuslik vektor, mis võtab väärtusi piirkonnas $G \subset \mathbb{R}^k$ ja mille tihedusfunktsioon on f_X . Olgu $g = (g_1, g_2, \dots, g_k) : G \rightarrow H$, $H \subset \mathbb{R}^k$ pööratav teisendus, mille kõik komponendid g_i on diferentseeruvad ja mille pöördteisendus on $g^{-1} = (h_1, h_2, \dots, h_k)$. Sellisel juhul juhusliku vektori $Y = g(X)$ tihedusfunktsioon on

$$f_Y(y) = f_X(g^{-1}(y)) \operatorname{abs} \left(\begin{vmatrix} \frac{\partial h_1(y)}{\partial y_1} & \frac{\partial h_1(y)}{\partial y_2} & \dots & \frac{\partial h_1(y)}{\partial y_k} \\ \frac{\partial h_2(y)}{\partial y_1} & \frac{\partial h_2(y)}{\partial y_2} & \dots & \frac{\partial h_2(y)}{\partial y_k} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial h_k(y)}{\partial y_1} & \frac{\partial h_k(y)}{\partial y_2} & \dots & \frac{\partial h_k(y)}{\partial y_k} \end{vmatrix} \right), \text{ kui } y \in H.$$

Absoluutväärtus maatriksist tähistab eelnevas teoreemis selle maatriksi determinanti.

4.1.1 Tihedusfunktsiooni teisenemine lineaarsete teisenduste korral

Olgu X pidev m -mõõtmeline juhuslik vektor ja f_X selle tihedusfunktsioon. Vaatleme dimensiooniga $m \times m$ pööratava maatriksi A ja vektori $b = (b_1, \dots, b_m)^T$ poolt defineeritud teisenduse

$$g(x) = Ax + b$$

rakendamisel saadavat juhuslikku suurust $Y = AX + b$. Kuna pöördteisendus avaldub sel juhul kujul

$$h(y) = A^{-1}(y - b),$$

saame teoreemi 15 põhjal avaldada Y tihedusfunktsiooni kujul

$$f_Y(y) = f_X(A^{-1}(y - b)) \operatorname{abs}(|A^{-1}|).$$

4.2 Normaaljaotusega juhuslike suuruste ja juhuslike vektorite genereerimine

Osutub, et mõnikord on lihtsam genereerida sõltumatute juhuslike suuruste paare kui üksikuid väärtuseid. Üks näide selle kohta on Box-Mülleri meetod sõltumatute standardsete

normaaljaotusega juhuslike suuruste genereerimiseks.

4.2.1 Box-Mülleri meetod

Meetod põhineb teadmisel, et standardse kahemõõtmelise normaaljaotusega vektori komponendid on sõltumatud normaaljaotusega juhuslikud suurused. Algoritm on järgmine:

1. genereeri kaks pseudujuhuslikku arvu u_1, u_2 jaotusest $U(0, 1)$,
2. arvuta $r = \sqrt{-2 \cdot \ln(u_1)}$, $\theta = 2\pi \cdot u_2$,
3. arvuta $x_1 = r \cdot \cos(\theta)$, $x_2 = r \cdot \sin(\theta)$.

Saadud väärtused x_1 ja x_2 vastavad sõltumatute jaotusega $N(0, 1)$ juhuslike suuruste väärtustele.

Veendume eelnevas väites, kasutades teoreemi 15. Alustame sellest, et $U = (U_1, U_2)$ on ühtlase jaotusega ühikruudul, seega

$$f_U(u_1, u_2) = \begin{cases} 1, & 0 < u_1 < 1, 0 < u_2 < 1, \\ 0 & \text{mujal.} \end{cases}$$

Box-Mülleri meetod vastab sellele vektorile teisenduse

$$g(u_1, u_2) = (\sqrt{-2 \cdot \ln(u_1)} \cdot \cos(2\pi u_2), \sqrt{-2 \cdot \ln(u_1)} \cdot \sin(2\pi u_2))$$

rakendamises. Teisendatud juhusliku suuruse tihedusfunktsiooni avaldise leidmiseks läheb meil vaja pöördteisenduse valemeid. Pöördteisenduse esimene funktsioon on lihtsalt leitav:

$$h_1(x_1, x_2) = e^{-\frac{x_1^2 + x_2^2}{2}}$$

Pöördteisenduse teise funktsiooni puhul tuleb arvestada, et funktsioon \arctan väärtused on piirkonnas $(-\frac{\pi}{2}, \frac{\pi}{2})$, seega tuleb olla hoolikas erinevatele tasandi veeranditele vastavate punktide teisendamisel. Tulemuseks on:

$$h_2(x_1, x_2) = \frac{1}{2\pi} \begin{cases} \arctan\left(\frac{x_2}{x_1}\right), & x_1 > 0, x_2 > 0, \\ \frac{\pi}{2}, & x_1 = 0, x_2 > 0, \\ \pi + \arctan\left(\frac{x_2}{x_1}\right), & x_1 < 0, \\ \frac{3\pi}{2}, & x_1 = 0, x_2 < 0, \\ 2\pi + \arctan\left(\frac{x_2}{x_1}\right), & x_1 > 0, x_2 < 0. \end{cases}$$

Edasine kontroll (jakobiaani leidmine ja veendumine, et teisendatud vektori tihedusfunktsioon on tõepoolest standardse kahemõõtmelise normaaljaotuse tihedusfunktsion) on harjutus lugejale.

4.2.2 Jaotusega $N(\mu, \Sigma)$ juhuslike vektorite genereerimine

Soovime genereerida juhuslikke vektoreid jaotusega $N(\mu, \Sigma)$, kus μ on m -mõõtmeline vektor ja Σ on pööratav $m \times m$ kovariatsioonimaatriks. Vaadeldes lineaarteisenduse $Y = AX + \mu$ abil saadud juhusliku vektori tihedusfunktsiooni, on lihtne veenduda, et standardse normaaljaotuse teisendamisel saame normaaljaotuse, mille keskvärtus on μ ja kovariatsioonimaatriks on $\Sigma = AA^T$. Seega tuleb etteantud Σ korral leida selline A , et $\Sigma = AA^T$. Üks võimalus selleks on defineerida maatriksi V , mille veergudeks on Σ normeeritud omavektorid ja diagonaalmaatriksi Λ , kus diagonaalil on Σ omaväärtused. Siis maatriksi $A = V\Lambda^{\frac{1}{2}}$ korral kehtib $AA^T = \Sigma$ (veendu selles!)

4.3 Simuleerimine tinglike jaotuste abil

Teame tõenäosuse korrutamise reeglit

$$P(A_1 A_2 \cdots A_m) = P(A_1) P(A_2 | A_1) \cdots P(A_m | A_1 \cdots A_{m-1})$$

Osutub, et sarnane valem kehtib ka pideva juhusliku vektori tihedusfunktsiooni jaoks

$$f_X(x_1, \dots, x_m) = f_{X_1}(x_1) f_{X_2|X_1}(x_2 | x_1) \cdots f_{X_m|X_1, X_2, \dots, X_{m-1}}(x_m | x_1, x_2, \dots, x_{m-1}),$$

kus tinglikud tihedusfunktsioonid saab leida samm-sammult seostest

$$\begin{aligned} f_{X_1}(x_1) f_{X_2|X_1}(x_2 | x_1) \cdots f_{X_k|X_1, X_2, \dots, X_{k-1}}(x_k | x_1, x_2, \dots, x_{k-1}) \\ = \underbrace{\int \cdots \int}_{m-k} f(x_1, \dots, x_m) dx_{k+1} \cdots dx_m, \quad 1 \leq k \leq m, \end{aligned}$$

Siin 0-kordne integraal tähistab lihtsalt integraalialust funktsiooni. Seega, kui suudame nendele tinglikele tihedustele vastavaid juhuslikke suuruseid genereerida, saame soovitud jaotusega juhusliku vektori tekitada samm-sammult: kõigepealt genereerime X_1 väärtuse vastavalt tihedusele f_{X_1} , seejärel genereerime X_2 vastavalt tihedusele $f_{X_2|X_1}$ jne.

4.4 Metropolis-Hastings algoritm

Mitmemõõtmelise jaotuse puhul sageli väga raske leida tinglikke jaotuseid, samuti on raske valida head lähtejaotust valikumeetodi jaoks. Samas juhuslike protsesside puhul sageli kehtib tulemus, et pärast küllalt pikka ajavahemikku on protsessi mingile konkreetsele ajamomendile vastav väärtus teatud kindla jaotusega (nn statsionaarse jaotusega) juhuslik suurus sõltumata sellest, milline oli protsessi algväärtus. Siit tuleb idee vaadelda protsessi, kus m -mõõtmeline punkt eksleb juhuslikult nii, et liikumine soodustab sattumist sinna, kuhu soovitud tõenäosusega juhusliku vektori väärtused peaks sattuma suurema tõenäosusega ja harvemini satutakse sinna, kuhu vaadeldava vektori väärtused satuvad väikese tõenäosusega. Osutub, et liikumist saab organiseerida nii, et statsionaarseks jaotuseks on meie soovitud jaotus.

Vaatleme algoritmi erijuhtu, kus ekslemisel uue asukoha määramisel kasutame normaaljaotust (nn Metropolise algoritm). Eeldame, et teame soovitud tihedusfunktsiooniga proportsionaalset funktsiooni f . Algoritm on järgmine:

1. Valime algpunkti x_1 (NB! peab olems selline, et $f(x_1) > 0$)
2. iga $i = 1, 2, \dots, n - 1$ korral
 - (a) leiame uue asukoha kandidaadi $x \sim N(x_i, \Sigma)$
 - (b) genereerime $u \sim U(0, 1)$
 - (c) Kui $u \leq \frac{f(x)}{f(x_i)}$, siis $x_{i+1} = x$, muidu $x_{i+1} = x_i$
3. Väljastame punktid $x_{n_0}, x_{n_0+\ell}, x_{n_0+2\ell}$ jne, kus n_0 ja ℓ on sobivalt valitud naturaalarvud

Algoritmi puhul on oluline mõista, et kuigi pärast piisavalt pikka aega n_0 vastavad väärtused x_i soovitud jaotusega juhuslike suuruste väärtustele, ei vasta järjestikused x_i väärtused sõltumatute katsete tulemustele. Sellest on tingitud täiendava parameetri ℓ valik, mis peab olema piisavalt suur, et sõltuvus eelmisest valitud väärtusest oleks praktiliselt olematu.

4.5 Gibbsi valik

Mõnikord on raske leida tinglikke jaotuseid juhul, kui üle ühe vektori koordinaatidest on fikseeritud, kuid ainult ühe koordinaadi fikseerimise korral on lihtne kindlaks teha, milline on vaadeldavale jaotusele vastav tinglik jaotus. Sel juhul on võimalik kasutada Gibbsi valiku algoritmi.

Eeldame, et suudame genereerida väärtuseid tinglikest jaotustest

$$X_i \mid X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_m, \quad i = 1, \dots, m$$

Algorit on järgmine:

1. Valime algpunkti $(x_{11}, x_{12}, \dots, x_{1m})$
2. iga $i = 2, 3, \dots, n$ korral genereerime
 - (a) x_{i1} jaotusest $X_1 \mid X_2 = x_{i-1,2}, \dots, X_m = x_{i-1,m}$,
 - (b) x_{i2} jaotusest $X_2 \mid X_1 = x_{i,1}, X_3 = x_{i-1,3}, \dots, X_m = x_{i-1,m}$,
 - (c) x_{i3} jaotusest $X_3 \mid X_1 = x_{i,1}, X_2 = x_{i,2}, X_4 = x_{i-1,4}, \dots, X_m = x_{i-1,m}$,
 - (d) ...
 - (e) x_{im} jaotusest $X_m \mid X_1 = x_{i,1}, X_2 = x_{i,2}, \dots, X_{m-1} = x_{i,m-1}$,

Jällegi tuleb arvestada, et järjestikused punktid ei vasta enamasti sõltumatutele juhuslikele vektoritele, seetõttu tuleks kasutada neist ainult piisavalt hõredalt (indeksi i mõttes) valitud punkte!

Peatükk 5

Taasvaliku meetodid: bootstrap ja jackknife

Sageli ei ole selge, millised on õiged eeldused jaotuse kohta, kust andmed pärinevad. Seetõttu on küllalt populaarsed meetodid, mis võimaldavad vähemalt ligikaudselt hinnata huvipakkuvaid suuruseid (hinnangute täpsust, nihet jms) ainult olemasolevate andmete põhjal. Selliseid meetodeid nimetatakse taasvaliku meetoditeks.

5.1 Bootstrap meetod

Sageli on meil antud ainult valim mahuga n ja me ei tea, mis jaotusest see pärineb. Valimi põhjal hindame mingit üldkogumi parameetrit (näiteks dispersiooni, asümmeetriakordajat vms) ning lisaks hinnangule tahame tavaliselt teada hinnangu täpsust (või hinnangufunktsiooni jaotust). Kui me teaks üldkogumi jaotust, siis saaks hinnangu täpsust kindlaks teha näiteks uusi valimeid genereerides. Õiget üldkogumi jaotust me ei tea, kuid me teame valimile vastavat empiirilist jaotust.

Kui juhuslik suurus on jaotusest jatusfunktsiooniga F_X , siis valim on juhuslik vektor jaotusfunktsiooniga

$$F_{X_1, \dots, X_n}(x_1, \dots, x_n) = F_X(x_1)F_X(x_2) \dots F_X(x_n)$$

Kui kasutame simulatsioonimeetodeid statistiku jaotuse leidmisel, genereerime vektoreid sellest jaotusest. Intuitiivselt, kui muuta jaotusfunktsioone piisavalt vähe, muutub ka enamuse statistikute käitumine küllalt vähe.

Teame, et valimi empiiriline jaotusfunktsioon F_n läheneb valimimahu kasvades jaotusfunktsioonile F_x , seega, asendades F_X funktsiooniga F_n , võib loota, et valimi jaotusfunktsioon ei muutu väga palju. Seetõttu genereerides jaotusfunktsioonile F_n vastavaid sõltumatu valimeid ja leides huvipakkuva statistiku käitumise, võib loota, et see on lähedane statistiku käitumisele õige jaotuse korral. Jaotusest F_n sõltumatu valimi genereerimine on samaväärne tagasipanekuga juhusliku valimi moodustamisega esialgsesse valimisse sattunud väärtustest.

5.1.1 Nihke hindamine Bootstrap meetodil

Nihke hindamiseks vaatame, kuidas hinnang käitub eeldusel, et valimi empiiriline jaotus vastab üldkogumi jaotusele. Kui me teame üldkogumi jaotust, siis enamasti on võimalik leida huvipakkuva suuruse (parameetri) üldkogumile vastav täpne väärtus, olgu selleks θ_B .

Seejärel leiame statistiku väärtuse m bootstrap valimi korral (tagasipanekuga juhuslik valik esialgse valimi väärtustest, valimi suurus võrdne esialgse valimiga), saame hinnangud $\hat{\theta}_i$, $i = 1, \dots, m$

Sageli nii saadud hinnangute keskmine ei koonu empiirilise jaotuse põhjal arvatud täpseks väärtuseks, st hinnang on nihkega. Nihke hinnanguks võtame $\frac{1}{m} \sum_{i=1}^m \hat{\theta}_i - \theta_B$

Saadud nihke hinnangu abil võime parandada esialgsel valimil statistikut rakendades saadud hinnangut, lahutades selle nihkest maha.

Et saada paremini aru, mis tegelikult tehakse, vaatleme juhtu, kus hindame standardhälvet tavapärase valemiga ja mei valimis on kaks arvu, 1 ja 3. Siis valimi põhjal arvatud standardhälbe hinnang on $\sqrt{2}$, Bootstrap jaotuse põhjal arvatud täpne standardhälve on 1 ja kui valida bootstrap jaotusest kaheelemendilisi valimeid ja arvutada nende põhjal saadud standardhälbe hinnangute keskmine, siis piisavalt suure m korral saame (kas oskad põhjendada, miks see nii on?)

$$\frac{1}{m} \sum_{i=1}^m \hat{\theta}_i \approx \frac{\sqrt{2}}{2}.$$

Seega bootstrap meetodil leitud nihke suuruseks on $\frac{\sqrt{2}}{2} - 1$ ja parandatud standardhälbe hinnanguks saame $\sqrt{2} + 1 - \frac{\sqrt{2}}{2}$.

5.1.2 Bootstrap usaldusvahemik

Parameetri hinnangu usaldusvahemikke saab bootstrap meetodiga leida mitmel viisil, vaatleme ühte, mis arvestab ka võimalikku nihet. Vaadeldav hinnang põhineb ideel kirjutada otsitav parameeter kujul

$$\theta = \hat{\theta} + (\theta - \hat{\theta}),$$

kus θ on hinnatav parameeter ja $\hat{\theta}$ on valimi põhjal hinnatud väärtus. Leiame vahe $\theta - \hat{\theta}$ jaotuse eeldusel, et lähtejaotus on bootstrap jaotus, st θ asemel kasutame bootstrap jaotuse korral leitud täpset väärtust θ_B . Seega vahe $\frac{\alpha}{2}$ -kvantiil on

$$\theta_B - q_{1-\frac{\alpha}{2}},$$

kus q_α tähistab hinnangu $\hat{\theta}$ bootstrap jaotuse korral leitud α -kvantiil:

$$P(\theta_B - \hat{\theta} \leq \theta_B - q_{1-\frac{\alpha}{2}}) = P(\hat{\theta} \geq q_{1-\frac{\alpha}{2}}) = 1 - P(\hat{\theta} < q_{1-\frac{\alpha}{2}}) = \frac{\alpha}{2},$$

Vahe $(1 - \frac{\alpha}{2})$ -kvantiil on $\theta_B - q_{\frac{\alpha}{2}}$, seega bootstrap usaldusvahemik olulisuse nivool α on

$$(\hat{\theta} + \theta_B - q_{1-\frac{\alpha}{2}}, \hat{\theta} + \theta_B - q_{\frac{\alpha}{2}}).$$

5.2 Jackknife meetod

Meetodi motiveerimiseks vaatleme suuruse $f(EX)$ hindamist valimi põhjal, kasutades statistikut $S = f(\bar{X})$, kus \bar{X} tähistab valimi aritmeetilist keskmist. Kasutades funktsiooni f Taylori rittaarendust kohal $x = EX$, saame

$$f(\bar{X}) = f(EX) + f'(EX)(\bar{X} - EX) + \frac{f''(EX)}{2}(\bar{X} - EX)^2 + \dots$$

Seega nihe avaldub kujul

$$E(f(\bar{X}) - f(EX)) = \frac{f''(EX)}{2n}DX + \dots$$

ja ei ole üldjuhul võrdne nulliga.

Tähistame nüüd kujul $S(i)$ sama statistikut, mis on hinnatud ilma i -nda vaatluseta ja defineerime uue statistiku valemiga

$$S_{(i)}^* = nS - (n-1)S(i).$$

Rittaarendusi kasutades on nüüd lihtne veenduda, et $S_{(i)}^*$ on väiksema nihkega, st rittaarenduses rohkem liikmeid kaovad ära (veendu selles!). Kuna selliseid hinnanguid saab arvutada valimi iga elemendi ärajätmisel, siis kokkuvõttes peaks kõige parema hinnangu saama siis, kui leida nende keskmine. Seega **Jacknife hinnang** on kujul

$$S^* = \frac{1}{n} \sum_{i=1}^n S_{(i)}^* = nS - \frac{n-1}{n} \sum_{i=1}^n S(i).$$

Siit saab leida ka jackknife hinnangu nihkele, milleks on

$$nihe_{jack} = S - S^* = (n-1) \left(\frac{1}{n} \sum_{i=1}^n S(i) - S \right).$$

Lisaks väiksema nihkega hinnangule on enamasti vaja teada ka uue hinnangu usaldusväärsus, st selle standardhälvet ja usaldusintervalle.

Jacknife hinnangu standardhälbe hinnang on

$$s_{jack}^* = \sqrt{\frac{n-1}{n} \sum_{i=1}^n \left(S(i) - \frac{1}{n} \sum_{j=1}^n S(j) \right)^2} = \sqrt{\frac{1}{n(n-1)} \sum_{i=1}^n (S_{(i)}^* - S^*)^2}$$

ning Turkey (1958) valem jackknife hinnangu usaldusintervallile olulisusnivool α on kujul

$$(S^* - t_{n-1, 1-\frac{\alpha}{2}} s_{jack}^*, S^* + t_{n-1, 1-\frac{\alpha}{2}} s_{jack}^*),$$

kus $t_{k, \alpha}$ tähistab vabadusastmete arvuga k t-jaotuse α -kvantiili.

Peatükk 6

Integraalide ja keskväärtuste arvutamine MC meetodiga. Dispersiooni vähendamise meetodid

Sageli pakuvad praktilist huvi teatud juhuslike suuruste keskväärtused, näiteks oodatav tulu investeerimisstrateegialt, finantsoptsoonide hinnad, kindlustusfirma järgmise viie aasta jooksul laostumise tõenäosus jne. Enamasti esituvad huvipakkuvad keskväärtused kujul kujul $g(X)$, kus g on mingi funktsioon ja X on juhuslik suurus või vektor.

Kui X on pidev juhuslik suurus või vektor, vastab keskväärtuse arvutamine (ühe- või mitmekordse) integraali arvutamisele

$$E g(X) = \int_{\mathbb{R}^m} g(x) f_X(x) dx,$$

kus f_X on juhusliku suuruse (või vektori) X tihedusfunktsioon. Samas integraalidel on palju muid rakendusi, näiteks erinevate objektide pindalad/ruumalad, kehale mõjuva jõu poolt tehtud töö keha liikumisel, muutuva tihedusega keha mass jne.

Tuleb välja, et iga integraali saab esitada keskväärtusena. Vaatleme integraali

$$I = \int_D f(x) dx, \quad D \subset \mathbb{R}^m.$$

Valime ühe tihedusfunktsiooni f_X , mille korral kehtib $f_X(x) > 0$, $x \in D$ ning defineerime

$$g(x) = \begin{cases} \frac{f(x)}{f_X(x)}, & x \in D, \\ 0, & x \notin D. \end{cases}$$

Siis

$$I = \int_D \frac{f(x)}{f_X(x)} f_X(x) dx = \int_{\mathbb{R}^m} g(x) f_X(x) dx = E g(X).$$

Jaotuse f_X valikul tuleb aga mõelda sellele, et tekkiv funktsioon g mingite argumentide korral väga halvasti käituma (st mingite x väärtuste korral liiga kiiresti pluss- või miinuslõpmatusse minema) ei hakka, sest halvasti valitud f_X korral võib $g(X)$ dispersioon olla lõpmatu ja siis koondub valimi keskmine keskväärtuseks väga aeglaselt.

6.1 Keskväärtuse hindamine MC meetodil

Vaatleme juhuslikku suuruse $Y = g(X)$ keskväärtuse arvutamist Monte Carlo meetodil. Algoritm on järgmine:

1. Genereerime n -elemendilise valimi X jaotusest, saame x_1, \dots, x_n
2. Rakendame funktsiooni g , saame valimi Y jaotusest: $g(x_1), \dots, g(x_n)$
3. Hindame Y keskväärtust valimi keskmise abil

$$EY \approx \bar{y} = \frac{1}{n} \sum_{i=1}^n g(x_i)$$

4. **Eeldusel, et Y dispersioon on lõplik**, saame valimi keskmise jaoks erinevad hinnangud (vaja osata ka tuletada tsentraalsest piirteoreemist!)

- (a) ligikaudne usaldusintervall (z_α tähistab standardse normaaljaotuse α -kvantiili)

$$\left(\bar{y} + z_{\frac{\alpha}{2}} \frac{\sigma_Y}{\sqrt{n}}, \bar{y} - z_{\frac{\alpha}{2}} \frac{\sigma_Y}{\sqrt{n}}\right)$$

- (b) Tõenäosusega $1 - \alpha$ kehtiv (ligikaudne) veahinnang:

$$-z_{\frac{\alpha}{2}} \frac{\sigma_Y}{\sqrt{n}}$$

- (c) **Tõenäoline viga** on tõenäosusega $\frac{1}{2}$ kehtiv veahinnang

5. Veahinnangute arvutamisel asendame σ_Y valimi põhjal leitud hinnanguga

6.1.1 MC meetodite võrdlemine

Ühte ja sama integraali või keskväärtust saab MC meetodiga arvutada väga erinevaid lähenemisi kasutades. Näiteks integraalide arvutamisel on võimalik valida lõpmatult palju erinevaid tihedusfunktsioone ja iga valik annab erineva juhusliku suuruse, mille keskväärtust MC meetodil arvutada. Kuigi kõigil juhtudel on sama keskväärtus, võivad dispersioonid olla vägagi erinevad (osadel juhtudel dispersioon võib ka puududa).

Kui me tahame erinevaid meetodeid omavahel võrrelda, siis selleks on mitmeid võimalusi:

- Tööaeg etteantud täpsuse saavutamiseks
- Etteantud täpsuse saavutamiseks vajaminev genereerimiste arv
- täpsus fikseeritud genereerimiste arvu korral

Kuigi tööaeg etteantud täpsuse saavutamiseks on kõige olulisem, sõltub see kasutatavast arvutist ja programmikoodi efektiivsusest, mistõttu on seda raske täpselt hinnata. Seetõttu keskendume genereeritavate juhuslike suuruste arvule.

Eeldame järgnevalt, et Y on lõpliku dispersiooniga, siis täpsus etteantud n korral sõltub ainult fikseeritud α väärtusest ja Y standardhälbest. Seega väiksem standardhälve annab suurema täpsuse.

MC meetodi veahinnangu valemist saame, et genereerimiste arv etteantud täpsuse ϵ saavutamiseks on kujul $c_{\alpha,\epsilon}DY$, kus

$$c_{\alpha,\epsilon} = \frac{\epsilon^2}{z_{\frac{\alpha}{2}}^2}$$

Tööaeg etteantud täpsuse saavutamiseks on seega $t c_{\alpha,\epsilon}DY$, kus t on ühe väärtuse genereerimise aeg. Edaspidi keskendume DY suuruse hindamisele erinevate meetodite korral. Oluliselt erinevate meetodite võrdlemisel tuleks võimalusel arvestada ka ühe väärtuse genereerimise aegade erinevust.

6.2 Dispersiooni vähendamise meetodid

6.2.1 Antiteetilised (*antithetic*) juhuslikud suurused

Antiteetlisteks (vastandlikeks) nimetame juhuslikke suuruseid, mis on sama jaotusega, kuid negatiivselt korreleeritud. Lihtne näide on järgmine: kui $X_1 \sim U(0, 1)$, siis $X_2 = 1 - X_1$ puhul on X_1 ja X_2 antiteetilised juhuslikud suurused (veendu selles!).

Üldiselt, kui pideva juhusliku suuruse X_1 tihedusfunktsioon on sümmeetriline μ suhtes (st $f_{X_1}(\mu - a) = f_{X_1}(\mu + a) \forall a \in \mathbb{R}$), siis $X_2 = 2\mu - X_1$ korral on X_1 ja X_2 antiteetilised juhuslikud suurused (veendu selles!).

Lisaks on oluline tähele panna, et kui X_1 ja X_2 on sama jaotusega, siis suvalise funktsiooni g korral on $g(X_1)$ ja $g(X_2)$ sama jaotusega, kuid ka antiteetiliste X_1 ja X_2 korral ei pruugi alati olla antiteetilised juhuslikud suurused.

Kui on soovi leida MC meetodiga EY , kus $Y = g(X)$ ja X jaotus on sümmeetriline $x = \mu$ suhtes, siis võime defineerida $X_2 = 2\mu - X$ ja

$$\tilde{Y} = \frac{g(X) + g(X_2)}{2}.$$

Siis $D\tilde{Y} = DY \cdot \frac{1+\rho}{2}$, kus $\rho = \text{cor}(g(X), g(X_2))$. Suuruse \tilde{Y} dispersioon ei ole kunagi suurem, kui DY . Kuna aga \tilde{Y} ühe väärtuse arvutamine on umbes 2 korda töömahukam kui Y ühe väärtuse arvutamine, võidame tööajas ainult siis, kui $\rho < 0$ ehk kui $g(X)$ ja $g(X_2)$ on antiteetilised.

Eelneva põhjal on meetodi kasulikkuse jaoks väga oluline, et $g(X)$ ja $g(X_2)$ oleks antiteetilised. Osutub, et kui g on monotoonne funktsioon, siis see tingimus on täidetud.

Teoreem 16 *Eeldame, et X pidev juhuslik suurus, mille tihedusfunktsioon on sümmeetriline punkti $x = \mu$ suhtes ning et g on monotoonne (mittekahanev või mittekasvav) funktsioon, mille korral $D(g(X)) < \infty$. Siis*

$$\rho = \text{cor}(g(X), g(X_2)) \leq 0.$$

Tõestus. Teoreemi tõestamiseks piisab näidata, et

$$\text{cov}(g(X), g(X_2)) \leq 0.$$

Paneme tähele, et suvaliste juhuslike suuruste X ja Y ja suvalise reaalarvu a korral kehtib

$$\begin{aligned} \text{cov}(X, Y) &= E[(X - EX)(Y - EY)] = E[(X - EX)(Y - a) + (X - EX)(a - EY)] \\ &= E[(X - EX)(Y - a)]. \end{aligned}$$

Tähistame $\mu_g = E g(X)$. Vaatleme mittekahaneva g juhtu. Kuna juhusliku suuruse keskväertus on tema väärtuste infimumi ja supremumi vahel, saame g mittekahanemise tõttu leida leida sellise arvu x^* , et $g(x) \leq \mu_g$, $x < x^*$ ja

$$g(x) \geq \mu_g, \quad x > x^*.$$

Valime $a = g(2\mu - x^*)$, siis

$$\text{cov}(g(X), g(X_2)) = E[(g(X) - \mu_g)(g(X_2) - a)].$$

Kui $X < x^*$, siis x^* definitsiooni kohaselt $g(X) - \mu_g \leq 0$. Samas kehtib siis $X_2 = 2\mu - X > 2\mu - x^*$ ning g mittekahanemise tõttu saame

$$g(X_2) \geq g(2\mu - x^*) = a,$$

seega vaadeldaval juhul kehtib

$$(g(X) - \mu_g)(g(X_2) - a) \leq 0.$$

Sarnane arutelu annab ka juhul $X > x^*$, et kehtib võrratus $(g(X) - \mu_g)(g(X_2) - a) \leq 0$ ning juhul $X = x^*$ on viimane teguritest a definitsiooni kohaselt võrdne nulliga, mistõttu saime, et juhuslik suurus $(g(X) - \mu_g)(g(X_2) - a)$ on alati mittepositiivne. Keskväertuse monotoonsuse omadusest järeldub nüüd, et vaadeldav kovariatsioon on mittepositiivne ja teoreemi väide kehtib.

Mittekasvava g korral on arutelu analoogiline (mõttele see läbi!). \square

6.2.2 Kontrollmuutujate meetod

Sageli on võimalik koos huvipakkuva suurusega Y genereerida lihtsam juhuslik suurus Z , mille keskväertus EZ on teada või mille keskväertust on lihtne arvutada. Seega saame moodustada uue juhusliku suuruse

$$\tilde{Y} = Y - a(Z - EZ),$$

kus a on meie valitud reaalarv. Lihtne on näha, et $E\tilde{Y} = EY$, nii et Y keskväertuse asemel võime leida \tilde{Y} keskväertuse. Kordaja a tuleks valida nii, et \tilde{Y} dispersioon on minimaalne, mis taandub ruutvõrrandi (a suhtes) miinimumkoha leidmisele. Lahendades selle võrrandi, saame

$$a = \frac{\text{cov}(Y, Z)}{DZ}.$$

Enamasti hindame a väärtuse valimi põhjal, kasutades kovariatsiooni ja dispersiooni hinnanguid.

Asetades optimaalse a väärtuse \tilde{Y} dispersiooni avaldisse, saame

$$D\tilde{Y} = DY \cdot (1 - \rho^2),$$

kus $\rho = \text{cor}(Y, Z)$.

6.2.3 Olulise valimi meetod

Sageli pakuvad praktikas huvi juhuslikud suurused kujul $Y = g(X)$, kus g omandab suuri (olulisi) väärtuseid piirkonnas kuhu X väärtused satuvad küllalt väikese tõenäosusega. Näiteks finantsmatemaatikas huvitavad investoreid sageli optsioonid, mis toovad omanikule raha sisse ainult juhul, kui aktsiahind oluliselt tõuseb (nt 50% praegusega võrreldes 3 kuu jooksul). Sageli sellisel juhul on dispersioon suur, mis annab idee muuta X jaotust nii, et näeme olulisi väärtuseid rohkem. Osutub, et seda saab alati teha. Nimelt saame valida juhusliku suuruse Z nii, et $f_Z(x) > 0$, kui $g(x)f_X(x) \neq 0$ Kuna alati kehtib võrdus (veendu selles!)

$$E g(X) = E[g(Z) \frac{f_X(Z)}{f_Z(Z)}],$$

siis saame valida sellise jaotusega Z , mille korral g olulised väärtused on nähtavad suurema tõenäosusega kui X korral.

Oluline on märkida, et seda ideed saab kasutada ka mitmemõõtmelisel juhul: kui sõltumatute juhuslike suuruste X_1 ja X_2 korral kasutada X_1 asemel suurust Z_1 ja X_2 asemel suurust Z_2 , siis

$$E g(X_1, X_2) = E[g(Z_1, Z_2) \frac{f_{X_1}(Z_1) f_{X_2}(Z_2)}{f_{Z_1}(Z_1) f_{Z_2}(Z_2)}].$$

Kuna normaaljaotus esineb stohhastilistes mudelites väga sageli, leiame suhte $\frac{f_X(z)}{f_Z(z)}$ juhul $X \sim N(\mu_1, \sigma)$, $Z \sim N(\mu_2, \sigma)$. Arvutuste tulemuseks (tee need läbi!) saame

$$\frac{f_X(z)}{f_Z(z)} = e^{\frac{2(\mu_1 - \mu_2)z + \mu_2^2 - \mu_1^2}{2\sigma^2}}.$$

6.2.4 Kihtvalimi meetod

Olgu B_1, \dots, B_m mingi sündmuste täissüsteem (vastastikku välistavad, positiivse tõenäosusega, üks alati toimub igal katsel). Teame, et keskväärtust saab esitada tinglike keskväärtuste kaudu:

$$EY = \sum_{i=1}^m P(B_i) E(Y | B_i).$$

Seega, kui oskame Y väärtuseid genereerida vastavatest tinglikest jaotustest, võime ühe keskväärtuse asemel arvutada MC meetodiga m keskväärtust

Tähistame $p_i = P(B_i)$ ning olgu Y_{ij} sõltumatud juhuslikud suurused tingliku jaotusega $Y|B_i$. Olgu n kasutatavate juhuslike suuruste arv. Siis tavalise MC meetodi korral on EY hinnangufunktsiooniks

$$H_n = \frac{1}{n} \sum_{i=1}^n Y_i,$$

kus Y_i on sõltumatud suurusega Y sama jaotusega juhuslikud suurused ja seega hinnangufunktsiooni dispersioon on $D(H_n) = \frac{DY}{n}$. Vaatame, kas sama arvu juhuslike suuruseid kasutades on võimalik saada väiksema dispersiooniga hinnang

Olgu n_i , $i = 1, \dots, m$ sellised positiivsed täisarvud, et $\sum_{i=1}^m n_i = n$. Defineerime kihtvalimi hinnangufunktsiooni

$$\tilde{H}_n = \sum_{i=1}^m p_i \frac{\sum_{j=1}^{n_i} Y_{ij}}{n_i}.$$

Leiame kihtvalimi hinnangufunktsiooni dispersiooni

$$D\tilde{H}_n = \sum_{i=1}^m \frac{p_i^2}{n_i} D(Y | B_i).$$

Põhiline küsimus on, kuidas valida n_i nii, et kihtvalimi hinnangufunktsiooni dispersioon oleks väiksem kui tavalise MC oma?

Teoreem 17 Olgu B_1, \dots, B_m sündmuste täissüsteem. Tähistame $p_i = P(B_i)$. Kui $n_i = p_i n$, siis

$$DH_n - D\tilde{H}_n = \frac{1}{n} \sum_{i=1}^m p_i (\mu_i - EY)^2,$$

kus $\mu_i = E(Y | B_i)$.

Eelneva teoreemi kohaselt ei ole seega tõenäosustega proportsionaalse valiku korral kihtvalimi hinnangu dispersioon sama arvu juhuslike suuruste kasutamisel kunagi suurem kui tavalise MC hinnangu dispersioon, st tulemus on enamasti täpsem. Valemist on samuti näha, et mida erinevad on juhusliku suuruse tinglikud keskväärtused erinevad üldisest keskmisest, seda suurem on võit täpsuses, seega kasutatavad sündmused võiksid olla hästi seotud vaadeldava juhusliku suuruse väärtustega.

6.2.5 Kihtvalimi meetodi optimaalsed proportsioonid

Nägime, et valik $n_i = p_i n$ on hea (annab enamasti parema meetodi). Aga kas saab paremini?

Kokku peab kasutatavate juhuslike suuruste summa olema n , seega võime otsida proportsioone q_i , $i = 1, \dots, m$, mis rahuldavad

$$q_i > 0, \quad i = 1, \dots, m, \quad \sum_{i=1}^m q_i = 1.$$

Otsime selliseid proportsioone, mille korral valiku $n_i = q_i n$ puhul on $D\tilde{H}_n$ minimaalne. Tähistame

$$\sigma_i^2 = D(Y | B_i).$$

Saame tingliku ekstreemumi leidmise ülesande: minimiseerida muutujate q_1, \dots, q_m suhtes suurust

$$f(q_1, \dots, q_m) = \frac{1}{n} \sum_{i=1}^m \frac{p_i^2 \sigma_i^2}{q_i}$$

tingimustel

$$q_i > 0, \quad i = 1, \dots, m, \quad \sum_{i=1}^m q_i = 1.$$

Lahendame ülesande Lagrange kordajate meetodil.

Lagrange kordajate meetodil (vt. näiteks Kõrgem matemaatika II loengukonspekti) peavad tingimust $\sum_{i=1}^m q_i = 1$ rahuldavates ekstreemumpunktides olema funktsiooni

$$f(q_1, \dots, q_m) + \lambda \cdot \left(\sum_{i=1}^m q_i - 1 \right)$$

osatuletised muutujate q_1, \dots, q_m järgi võrduma nulliga ning lisaks peab olema rahuldatud kitsendus $\sum_{i=1}^m q_i = 1$. Siit saame võrrandid

$$-\frac{p_i^2 \sigma_i^2}{n q_i^2} + \lambda = 0, \quad i = 1, 2, \dots, m,$$

kust järeldub, et positiivsetele proportsioonidele vastavad ekstreemumpunktid esituvad kujul

$$q_i = \frac{p_i \sigma_i}{\sqrt{n \lambda}}$$

Tingimusest $\sum_{i=1}^m q_i = 1$ järeldub nüüd, et

$$\sqrt{n \lambda} = \sum_{i=1}^m p_i \sigma_i,$$

mistõttu

$$q_i = kp_i\sigma_i, \quad i = 1, \dots, m,$$

$$\text{kus } k = \frac{1}{\sum_{i=1}^m p_i\sigma_i}.$$

Seega parima tulemuse saame, kui igas kihis on kasutatud väärtuste arv proportsionaalne korrutisega $p_i\sigma_i$. Hinnangu dispersioon on sel juhul

$$D\tilde{H}_n = \frac{(\sum_{i=1}^m p_i\sigma_i)^2}{n}$$

Kuna iga kihtvalimi meetodi puhul on Suure n korral on hinnang ligikaudu normaaljaotusega (iga kihi keskmine on ligikaudu normaaljaotusega tsentraalse piirteoreemi kohaselt ning sõltumatute normaaljaotusega juhuslike suuruste summa on normaaljaotusega), seega saame ligikaudu leida tõenäosusega $1 - \alpha$ kehtiva vea kujul

$$-z_{\frac{\alpha}{2}} \sqrt{D\tilde{H}_n}$$

6.2.6 Kihtvalimi meetodi rakendamine

=====

Kõigepealt tuleb valida m ja fikseerida sündmused B_1, \dots, B_m . Teoreetiliselt, mida suurem m , seda suuremat võitu on võimalik saada. Praktiliselt peame iga kihi korral hindama tinglikke dispersioone ja normaalse täpsusega hinnangu jaoks on vaja piisavalt suur arv (nt 100) genereeritud väärtust, seetõttu minimaalne n_i võiks olla vähemalt 100. Suure m korral aga võib see tähendada väga suurt n väärtust.

Sündmused B_i peaks olema juhusliku suuruse käitumisega seotud. Kui $Y = g(X)$, siis on hea valik $B_i = \{a_{i-1} < X \leq a_i\}$, kus

$$-\infty = a_0 < a_1 < \dots < a_m = \infty,$$

sest selliste juhuslike suuruste korral oskame me jaotusfunktsiooni pöördfunktsiooni meetodil juhuslike suuruste väärtuseid tinglikust jaotusest genereerida. Sageli on mõistlik defineerida a_i juhusliku suuruse X kvantiilide kaudu, siis saame määrata lihtsalt sündmuste tõenäosused (nt tagada, et $p_i = \frac{1}{m} \forall i$).

Peatükk 7

Kasutatavad tulemused tõenäosusteooriast

Teoreem 18 (Tõenäosuse omadused) Olgu (Ω, \mathcal{F}, P) mingi tõenäosusruum. Siis kehtivad järgnevad omadused:

1. $P(\emptyset) = 0$;

2. kui $A_i \in \mathcal{F}$, $i = 1, 2, \dots, n$ on vastastikku välistavad, st $A_i \cap A_j = \emptyset$, $i \neq j$, siis kehtib võrdus

$$P\left(\bigcup_{i=1}^n A_i\right) = \sum_{i=1}^n P(A_i);$$

3. kui $A, B \in \mathcal{F}$, $A \subset B$, siis $P(A) \leq P(B)$ (monotoonsus).

4. $P(\bar{A}) = 1 - P(A)$;

5. $P(A \setminus B) = P(A) - P(A \cap B) \forall A, B \in \mathcal{F}$;

6. $P(A) \leq 1 \forall A \in \mathcal{F}$;

7. $P(A \cup B) = P(A) + P(B) - P(A \cap B) \forall A, B \in \mathcal{F}$;

$$P\left(\bigcup_{i=1}^n A_i\right) = \sum_{i=1}^n P(A_i) - \sum_{1 \leq i < j \leq n} P(A_i \cap A_j) + \sum_{1 \leq i < j < k \leq n} P(A_i \cap A_j \cap A_k) - \dots + (-1)^{n-1} P(A_1 \cap A_2 \cap \dots \cap A_n), \quad A_i \in \mathcal{F}, \quad i = 1, 2, \dots, n;$$

8. $P\left(\bigcup_{i=1}^{\infty} A_i\right) \leq P(A) + P(B) \forall A, B \in \mathcal{F}$;

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) \leq \sum_{i=1}^{\infty} P(A_i), \quad A_i \in \mathcal{F}, \quad i \in \mathbb{N};$$

9. Tõenäosuse pidevus:

- $A_i \in \mathcal{F}$, $i \in \mathbb{N}$, $A_1 \subset A_2 \subset A_3 \subset \dots \Rightarrow \lim_{n \rightarrow \infty} P(A_n) = P\left(\bigcup_{i=1}^{\infty} A_i\right)$;
- $A_i \in \mathcal{F}$, $i \in \mathbb{N}$, $A_1 \supset A_2 \supset A_3 \supset \dots \Rightarrow \lim_{n \rightarrow \infty} P(A_n) = P\left(\bigcap_{i=1}^{\infty} A_i\right)$;

Lemma 19 (Täistõenäosuse valem). Rahuldagu sündmused $B_i \in \mathcal{F}$, $i = 1, \dots, n$ tingimusi

$$P(B_i) \neq 0, \quad i = 1, 2, \dots, n; \quad B_i \cap B_j = \emptyset, \quad i \neq j; \quad \bigcup_{i=1}^n B_i = \Omega. \quad (7.1)$$

Siis iga sündmuse $A \in \mathcal{F}$ korral kehtib võrdus

$$P(A) = \sum_{i=1}^n P(B_i)P(A|B_i).$$

Definitsioon 20 Juhusliku suuruse X jaotusfunktsiooniks nimetatakse funktsiooni

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

Teoreem 21 Olgu X juhuslik suurus ning F tema jaotusfunktsioon. Siis kehtivad järgnevad omadused.

1. $0 \leq F(x) \leq 1$ iga $x \in \mathbb{R}$ korral.
2. F on monotoonselt kasvav: kui $x_1 < x_2$, siis $F(x_1) \leq F(x_2)$.
3. Kehtivad piirväärtused

$$\lim_{x \rightarrow -\infty} F(x) = 0, \quad \lim_{x \rightarrow \infty} F(x) = 1.$$

4. F on paremalt pidev:

$$\lim_{x > a, x \rightarrow a} F(x) = F(a) \quad \forall a \in \mathbb{R}.$$

5. Kehtib võrdus

$$P(\{X = a\}) = F(a) - \lim_{x < a, x \rightarrow a} F(x).$$

6. Kehtib võrdus $P(\{a < X \leq b\}) = F(b) - F(a)$.

Definitsioon 22 Juhusliku vektori (X, Y) jaotusfunktsiooniks (ehk juhuslike suuruste X ja Y ühisjaotuse jaotusfunktsiooniks) nimetatakse funktsiooni

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

Definitsioon 23 Juhuslikku vektorit (X, Y) nimetatakse pidevaks, kui tema jaotusfunktsioon avaldub kujul

$$F_{X,Y}(x, y) = \int_{-\infty}^x \left(\int_{-\infty}^y f_{X,Y}(u, v) dv \right) du, \quad x, y \in \mathbb{R}$$

mingi funktsiooni $f_{X,Y} : \mathbb{R}^2 \rightarrow \mathbb{R}$ korral. Funktsiooni $f_{X,Y}$ nimetatakse sel juhul juhusliku vektori (X, Y) tihedusfunktsiooniks (ehk juhuslike suuruste X ja Y ühistiheduseks).

Lemma 24 (Tihedusfunktsiooni omadused) Olgu (X, Y) pidev juhuslik vektor jaotusfunktsiooniga $F_{X,Y}$ ja tihedusfunktsiooniga $f_{X,Y}$. Siis kehtivad järgmised omadused:

1. Funktsioon $f_{X,Y}$ on mittenegatiivne, st $f_{X,Y}(x, y) \geq 0 \quad \forall (x, y) \in \mathbb{R}^2$;
2. kehtivad võrdused

$$\begin{aligned} f_X(x) &= \int_{-\infty}^{\infty} f_{X,Y}(x, y) dy, \\ f_Y(y) &= \int_{-\infty}^{\infty} f_{X,Y}(x, y) dx, \\ \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) dx dy &= 1 \end{aligned}$$

3. Kui $D \subset \mathbb{R}^2$ on Boreli σ -algebra suhtes mõõtv hulk (st esitatav loenduva arvu ristkõikude abil kasutades ühendeid, ühisosasid ja täiendeid), siis

$$P(\{(X, Y) \in D\}) = \iint_D f_{X,Y}(x, y) dx dy.$$

4. Kui $g : \mathbb{R}^2 \rightarrow \mathbb{R}$ on "piisavalt heade omadustega" funktsioon (nt pidev või selline, mille valemit me oskame kirja panna) ning

$$\iint_{\mathbb{R}^2} |g(x, y)| f_{X,Y}(x, y) dx dy < \infty,$$

siis

$$E(g(X, Y)) = \iint_{\mathbb{R}^2} g(x, y) f_{X,Y}(x, y) dx dy.$$

5. Kui $F_{X,Y}$ on diferentseeruv punktis (x, y) , siis

$$f_{X,Y}(x, y) = \frac{\partial^2 F_{X,Y}}{\partial x \partial y}(x, y)$$

Lemma 25 Pidevad juhuslikud suurused X ja Y on sõltumatud parajasti siis, kui nende ühisjaotuse tihedusfunktsioon avaldub kujul

$$f_{X,Y}(x, y) = f_X(x) f_Y(y) \quad \forall x, y \in \mathbb{R}.$$