

## 3A Keskihajonta ja korrelaatio

### Tuntitehtävät

**3A1** (Korrelaatio ja riippuvuus) Diskreettien satunnaismuuttujien  $X$  ja  $Y$  yhteisjakauma on esitetty allaolevana taulukkona:

	Y		
X	-1	0	1
-1	0	$\frac{1}{6}$	$\frac{1}{6}$
0	$\frac{1}{3}$	0	0
1	0	$\frac{1}{6}$	$\frac{1}{6}$

- (a) Määritä  $X$ :n jakauma, odotusarvo ja keskihajonta.
- (b) Määritä  $Y$ :n jakauma, odotusarvo ja keskihajonta.
- (c) Laske  $X$ :n ja  $Y$ :n korrelaatio.
- (d) Selvitä, ovatko  $X$  ja  $Y$  riippuvat vai riippumattomat.

### Ratkaisu.

- (a) Yhteisjakauman taulukon rivisummat laskemalla saadaan  $X$ :n jakaumaksi

$k$	-1	0	1
$P(X = k)$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$

Näin ollen  $X$ :n odotusarvoksi saadaan

$$E(X) = \sum_k k P(X = k) = (-1) \times \frac{1}{3} + 0 \times \frac{1}{3} + 1 \times \frac{1}{3} = 0.$$

Keskihajonta saadaan varianssin neliöjuurena, ja varianssi puolestaan on ehkä kätevintä laskea kaavalla  $\text{Var}(X) = E(X^2) - E(X)^2$ . Koska

$$E(X^2) = \sum_k k^2 P(X = k) = (-1)^2 \times \frac{1}{3} + 0^2 \times \frac{1}{3} + 1^2 \times \frac{1}{3} = \frac{2}{3},$$

havaitaan, että  $X$ :n keskihajonta on

$$\text{SD}(X) = \sqrt{\frac{2}{3} - 0^2} = \sqrt{E(X^2) - E(X)^2} = \sqrt{\frac{2}{3}} \approx 0.82.$$

(b) Yhteisjakauman taulukon sarakesummat laskemalla saadaan  $Y$ :n jakaumaksi

$k$	-1	0	1
$P(Y = k)$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$

Tästä nähdään, että  $X$  ja  $Y$  ovat samoin jakautuneita. Näin ollen  $Y$ :n odotusarvo ja keskihajonta ovat samat mitä  $X$ :lläkin, eli  $E(Y) = 0$  ja  $SD(Y) = \sqrt{2/3} \approx 0.82$ .

(c) Korrelaatio saadaan normittamalla kovarianssi, ja kovarianssi yleensä saadaan helpon laskukaavasta  $\text{Cov}(X, Y) = E(XY) - E(X)E(Y)$ . Yhteisjakauman taulukon avulla tulee laskea

$$E(XY) = \sum_i \sum_j ij P(X = i, Y = j).$$

Ylläolevasta summasta voidaan jättää pois kaikki termit, joissa  $i = 0$  tai  $j = 0$  (koska niiden kontribuutio summaan on nolla), joten

$$\begin{aligned} E(XY) &= \sum_{i \neq 0} \sum_{j \neq 0} ij P(X = i, Y = j) \\ &= (-1) \times (-1) \times 0 + (-1) \times 1 \times \frac{1}{6} + 1 \times (-1) \times 0 + 1 \times 1 \times \frac{1}{6} \\ &= 0. \end{aligned}$$

Näin ollen  $X$ :n ja  $Y$ :n korrelaatio on

$$\text{Cor}(X, Y) = \frac{\text{Cov}(X, Y)}{SD(X)SD(Y)} = \frac{E(XY) - E(X)E(Y)}{SD(X)SD(Y)} = 0.$$

(d)  $X$  ja  $Y$  ovat stokastisesti riippuvat, sillä esimerkiksi

$$P(Y = -1 | X = 0) = 1,$$

kun taas

$$P(Y = -1) = 1/3$$

Tämän tehtävän viesti on, että korreloimattomuus *ei* takaa riippumattomuutta. Stokastinen riippumattomuus on vahva ominaisuus (kertolaskukaavan pitäisi päteä joka kohdassa yhteisjakaumaa; tai toisin sanottuna,  $Y$ :n ehdollisten jakaumien pitäisi olla samat jokaisella arvolla  $X = x$ ). Stokastisesta riippumattomuudesta seuraa kyllä korreloitumattomuus.

Korrelaatio puolestaan mittaa vain tiettytyypistä, luonteeltaan lineaarista, stokastista riippuvuutta. Muuttujilla voi olla epälineaarinen stokastinen riippuvuus, kuten tässä tehtävässä, ja niiden korrelaatio voi silti olla nolla.

**3A2** (Nopanheittojen keskiarvo) Tavallista noppaa heitetään monta kertaa peräkkäin. Heittojen tuloksia merkitään  $X_1, X_2, \dots$  ja ne ovat keskenään riippumattomat. Ensimmäisten  $n$ :n tuloksen keskiarvoa merkitään  $A_n = \frac{1}{n}(X_1 + \dots + X_n)$ .

- (a) Laske satunnaismuuttujan  $X_1$  odotusarvo ja keskihajonta.
- (b) Määritä satunnaismuuttujan  $A_2 = \frac{1}{2}(X_1 + X_2)$  jakauma.  
Vihje: Selvitä ensin  $A_2$ :n arvojoukko. Tutki sitten, aluksi pienillä arvoilla  $a$ , milloin eli millä parin  $(X_1, X_2)$  arvoilla toteutuu  $A_2 = a$ . Koeta yleistää päättelysi.
- (c) Laske satunnaismuuttujan  $A_2$  odotusarvo ja keskihajonta. Vertaa muuttujien  $A_2$  ja  $X_1$  keskihajontoja laskemalla niiden suhde.
- (d) Laske odotusarvo ja keskihajonta satunnaismuuttujalle

$$A_{100} = \frac{1}{100}(X_1 + X_2 + \dots + X_{100}).$$

Vertaa muuttujien  $A_{100}$  ja  $X_1$  keskihajontoja laskemalla niiden suhde.

Opastus. C-kohdan voit laskea joko b-kohdassa lasketusta jakaumasta, tai voit hyödyntää odotusarvon ja varianssin yhteenlaskuominaisuuksia. D-kohdassa tarkan jakauman selvittäminen olisi hyvin työlästä, joten on parempi käyttää yhteenlaskuominaisuuksia.

Lauseesta 4.9 seuraa, että  $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2 \cdot \text{Cov}(X, Y)$ . **Riippumattomilla** satunnaismuuttujilla em. kaavasta jää kovarianssitermi nolaksi, ja usealle **riippumattomalle** muuttujalle kaava yleistyy muotoon  $\text{Var}(X_1 + \dots + X_n) = \text{Var}(X_1) + \dots + \text{Var}(X_n)$ . Summien variansseihin tutustutaan tarkemmin luentomonisteen luvussa 5.

Huomaa, että yhteenlaskukaava pätee nimenomaan variansseille, ei keskihajonnoille.

## Ratkaisu.

- (a)  $X_1$  noudattaa joukon  $\{1, 2, \dots, 6\}$  tasajakaumaa, joten sen odotusarvo saadaan kaavasta

$$E(X_1) = \sum_k k P(X_1 = k) = \sum_{k=1}^6 k \times \frac{1}{6} = \frac{7}{2} = 3.5.$$

Keskihajonnan laskemiseksi lasketaan ensin  $X_1$ :n neliön odotusarvo

$$E(X_1^2) = \sum_k k^2 P(X_1 = k) = \sum_{k=1}^6 k^2 \times \frac{1}{6} = \frac{91}{6}.$$

Tämän jälkeen keskihajonta saadaan kaavasta

$$\text{SD}(X_1) = \sqrt{E(X_1^2) - E(X_1)^2} = \sqrt{\frac{91}{6} - \left(\frac{7}{2}\right)^2} = \sqrt{\frac{35}{12}} \approx 1.7078.$$

- (b) Tarkastelemalla mitä arvoja  $X_1$  ja  $X_2$  voivat saada, havaitaan että  $A_2$ :n mahdollisten arvojen joukko on  $\{1.0, 1.5, 2.0, 2.5, \dots, 5.5, 6.0\}$ . Kyseisten arvojen todennäköisyydet voidaan määrittää kohta kohdalta:

- Tapahtuma  $\{A_2 = 1\}$  sattuu täsmälleen silloin, kun  $X_1 = 1$  ja  $X_2 = 1$ . Näin ollen  $P(A_2 = 1) = \left(\frac{1}{6}\right)^2 = \frac{1}{36}$ .

- Tapahtuma  $\{A_2 = 1.5\}$  sattuu, jos  $(X_1, X_2) = (1, 2)$  tai  $(X_1, X_2) = (2, 1)$ . Näin ollen  $P(A_2 = 1.5) = 2 \times (\frac{1}{6})^2 = \frac{2}{36}$ .
- ...

Näin etenemällä saadaan  $A_2$ :n jakauma määritettyä taulukkoon:

$k$	1.0	1.5	2.0	2.5	3.0	3.5	4.0	4.5	5.0	5.5	6.0
$P(A_2 = k)$	$\frac{1}{36}$	$\frac{2}{36}$	$\frac{3}{36}$	$\frac{4}{36}$	$\frac{5}{36}$	$\frac{6}{36}$	$\frac{5}{36}$	$\frac{4}{36}$	$\frac{3}{36}$	$\frac{2}{36}$	$\frac{1}{36}$

- (c) **Tapa 1.** Satunnaismuuttujan  $A_2$  odotusarvo ja keskihajonta voidaan määrittää b)-kohdan taulukosta kaavoilla:

$$E(A_2) = \sum_x x P(A_2 = x) = 1.0 \times \frac{1}{36} + 1.5 \times \frac{2}{36} + \dots = 3.5$$

ja

$$E(A_2^2) = \sum_x x^2 P(A_2 = x) = 1.0^2 \times \frac{1}{36} + 1.5^2 \times \frac{2}{36} + \dots \approx 13.71,$$

josta

$$SD(A_2) = \sqrt{E(A_2^2) - (E(A_2))^2} \approx \sqrt{13.71 - 3.5^2} \approx 1.2.$$

Lukuarvot voi laskea taskulaskimella tai esim. R-komennoilla

```
x <- seq(1.0, 6.0, by=0.5)
p <- c(1:6, 5:1)/36
m1 <- sum(x*p)
m2 <- sum(x^2*p)
mu <- m1
sigma <- sqrt(m2-m1^2)
```

**Tapa 2.** Koska  $X_2$ :llä on sama jakauma kuin  $X_1$ :llä, sillä on myös sama odotusarvo ja varianssi, jotka laskettiin a-kohdassa. Odotusarvon laskusäännöillä

$$E(A_2) = E\left(\frac{1}{2}(X_1 + X_2)\right) = \frac{1}{2}(E(X) + E(Y)) = \frac{1}{2}(3.5 + 3.5) = \mathbf{3.5}.$$

Varianssin laskusäännöillä (vrt. opastus)

$$\text{Var}(X_1 + X_2) = \text{Var}(X_1) + \text{Var}(X_2) + 2 \cdot \text{Cov}(X_1, X_2) = (35/12) + (35/12) + 2 \cdot 0 = 35/6,$$

missä kovarianssi oli nolla, koska  $X_1$  ja  $X_2$  ovat riippumattomat. Lasketaan vielä keskihajonta ja lopuksi skaalataan kertoimella puoli.

$$SD(X_1 + X_2) = \sqrt{35/6}$$

ja

$$\text{SD}(A_2) = \text{SD}\left(\frac{1}{2}(X_1 + X_2)\right) = \frac{1}{2}\text{SD}(X_1 + X_2) = \frac{1}{2}\sqrt{35/6} = \sqrt{35/24} \approx \mathbf{1.2076}.$$

Lasketaan vielä kysytty suhde.

$$\frac{\text{SD}(A_2)}{\text{SD}(X_1)} = \frac{\sqrt{35/24}}{\sqrt{35/12}} = \frac{1}{\sqrt{2}} \approx \mathbf{0.7071}.$$

- (d) Satunnaismuuttujan  $A_{100}$  arvojoukko on  $\{1.00, 1.01, 1.02, \dots, 5.99, 6.00\}$ . Tämä 501 rationaaliluvun arvojoukko voitaisiin periaatteessa listata taulukkoon ja laskea jokaisen lukuarvon todennäköisyys. Ei ole kuitenkaan aivan yksinkertaista laskea esim. tapahtuman  $\{A_{100} = 3.97\}$  todennäköisyyttä käymällä läpi kaikkia 100 nopan kombinaatioita, jotka tuottavat keskiarvoksi 3.97, sillä sadan nopan tulokombinaatioita on tähtitieteellisen suuri määrä:  $6^{100} \approx 6 \cdot 10^{77}$  kappaletta.

Satunnaismuuttujan  $A_{100}$  odotusarvo voidaan helposti laskea lineaarisuutta käyttämällä:

$$E(A_{100}) = E\left(\frac{1}{100} \sum_{i=1}^{100} X_i\right) = \frac{1}{100} \sum_{i=1}^{100} E(X_i) = \frac{1}{100} \sum_{i=1}^{100} 3.5 = \mathbf{3.5}.$$

(Summan termit olivat samat, koska kaikki noppatulokset  $X_1, X_2, \dots$  noudattavat samaa lukujoukon  $\{1, \dots, 6\}$  tasajakaumaa.)

Satunnaismuuttujan  $A_{100}$  varianssi voidaan laskea yleisten varianssin laskusääntöjen avulla tuntematta  $A_{100}$ :n jakaumaa, sillä tulokset  $X_1, X_2, \dots$  ovat toisistaan stokastisesti riippumattomat. Lisäksi tiedetään että jokaisella heittotuloksella  $X_i$  on sama jakauma ja siis sama varianssi kuin tuloksella  $X_1$ . Näin ollen

$$\text{Var}(A_{100}) = \text{Var}\left(\frac{1}{100} \sum_{i=1}^{100} X_i\right) = \left(\frac{1}{100}\right)^2 \sum_{i=1}^{100} \text{Var}(X_i) = \frac{\text{Var}(X_1)}{100}$$

ja  $A_{100}$ :n keskihajonnaksi saadaan (a)-kohdan avulla

$$\text{SD}(A_{100}) = \sqrt{\frac{\text{Var}(X_1)}{100}} = \frac{\text{SD}(X_1)}{10} = \frac{\sqrt{35/12}}{10} \approx \mathbf{0.1708}.$$

Nähdään, että kysytty suhde on

$$\frac{\text{SD}(A_{100})}{\text{SD}(X_1)} = \mathbf{0.10}$$

eli sadan heittotuloksen keskiarvon keskihajonta on **kymmenesosa** yhden heittotuloksen keskihajonnasta.

Tehtävä osoittaa, että riippumattomien ja samoin jakautuneiden satunnaismuuttujien keskiarvoistaminen *säilyttää odotusarvon ennallaan* mutta *pienentää keskihajontaa*. Tämä tärkeä havainto selittää, miten riskejä voi pienentää hajauttamalla. Tähän myös pohjautuu suurten finanssi- ja vakuutusyhtiöiden sekä kasinojen toiminta. Toisaalta tähän perustuu myös tilastotieteellinen estimointi, koska otoksen suurentaminen saa otoksen keskiarvon osumaan lähemmäs jakauman odotusarvoa, eli tieto odotusarvosta tarkentuu.

## Kotitehtävät

**3A3** (Lämpötilamalli) Meteorologi mallintaa tämän ja huomisen päivän lämpötilojen  $T_0$  ja  $T_1$  välistä yhteyttä kaavalla

$$T_1 = T_0 + \Delta T$$

jossa  $\Delta T$  kuvaa lämpötilojen muutosta. Satunnaismuuttujat  $T_0$  ja  $\Delta T$  oletetaan toisistaan riippumattomiksi. Lisäksi tiedetään, että  $E(T_0) = \mu$  ja  $\text{Var}(T_0) = \sigma^2$  sekä  $E(\Delta T) = 0$  ja  $\text{Var}(\Delta T) = \theta^2$ . Mallin parametrit  $\mu$ ,  $\sigma$  ja  $\theta$  oletetaan ennalta tunnetuiksi, ja lisäksi  $\sigma > 0$  ja  $\theta \geq 0$ .

- Määritä  $E(T_1)$ .
- Määritä  $\text{SD}(T_1)$ . *Vrt. tehtävään 3A2. Miten lasketaan summan varianssi?*
- Määritä  $\text{Cov}(T_1, T_0)$ . *Käytä kovarianssin bilineaarisuutta.*
- Määritä  $\text{Cor}(T_1, T_0)$ . Ennen kuin lasket korrelaation tarkan arvon, yritä intuitiivisesti päätellä korrelaation tulkinnan kautta miten se käyttäytyy tapauksissa  $\theta = 0$  ja  $\theta \gg \sigma$  (eli  $\theta$  paljon suurempi kuin  $\sigma$ ).

Tulokset tulee ilmoittaa mahdollisimman yksinkertaisina lausekkeina mallin parametreista  $\mu$ ,  $\sigma$  ja  $\theta$ .

### Ratkaisu.

- Odotusarvon lineaarisuuden perusteella

$$E(T_1) = E(T_0 + \Delta T) = E(T_0) + E(\Delta T) = \mu + 0 = \mu.$$

Siis huomisen lämpötilan odotusarvo on sama kuin tämän päivän lämpötilan odotusarvo.

- Satunnaismuuttujien summan varianssin kaavan perusteella

$$\begin{aligned}\text{Var}(T_1) &= \text{Var}(T_0 + \Delta T) \\ &= \text{Var}(T_0) + 2 \text{Cov}(T_0, \Delta T) + \text{Var}(\Delta T) \\ &= \sigma^2 + 0 + \theta^2 \\ &= \sigma^2 + \theta^2,\end{aligned}$$

missä  $\text{Cov}(T_0, \Delta T) = 0$  koska  $T_0$  ja  $\Delta T$  ovat riippumattomia. Siis  $\text{SD}(T_1) = \sqrt{\sigma^2 + \theta^2}$ .

- Käyttäen kovarianssin laskusääntöjä havaitaan, että

$$\begin{aligned}\text{Cov}(T_1, T_0) &= \text{Cov}(T_0 + \Delta T, T_0) \\ &= \text{Cov}(T_0, T_0) + \text{Cov}(\Delta T, T_0) \\ &= \text{Var}(T_0) + 0 \\ &= \sigma^2.\end{aligned}$$

- (d) Korrelaatio mittaa sitä, kuinka hyvin satunnaismuuttujalla voi (linearisesti) ennustaa toisen satunnaismuuttujan arvoja. Muutoksen keskihajonnan  $\theta$  ollessa lähellä nollaa, pitäisi huomisen lämpötilan olla paremmin ennustettavissa tämän päivän lämpötilasta. Muutoksen keskihajonnan  $\theta$  kasvattaminen suureksi puolestaan tarkoittaa, että huomisen lämpötila on keskimäärin kauempana tämän päivän lämpötilasta, ja täten heikommin ennustettava. Näin ollen on oletettavaa, että korrelaatio lähestyy kohti nollaa muutoksen keskihajonnan  $\theta$  kasvaessa suureksi. Varmistetaan tämä laskemalla, kovarianssin bilineaarisuuden avulla:

$$\text{Cor}(T_1, T_0) = \frac{\text{Cov}(T_1, T_0)}{\text{SD}(T_1) \text{SD}(T_0)} = \frac{\sigma^2}{(\sqrt{\sigma^2 + \theta^2}) \cdot \sigma} = \frac{\sigma}{\sqrt{\sigma^2 + \theta^2}}.$$

Viimeinen lauseke on tosiaan  $\theta$ :n suhteen vähenevä funktio, vahvistaen aiemman intuition. Pisteessä  $\theta = 0$  lauseke antaa korrelaation arvoksi täsmälleen  $+1$ , mikä vastaa tarkkaa lineaarista riippuvuutta. Jos taas  $\theta$  on hyvin suuri, niin lausekkeen arvo on lähellä nollaa.

**3A4** (Kustannusfunktion minimointi) Abel ja Bertta ovat töissä eri sääennusteyrityksissä. He ovat kumpikin simuloineet tietokoneella mahdollisia säätiloja ja laskeneet ennusteen huomisen keskipäivän lämpötilalle  $X$  celsiusasteina. Kumpikin on päätenyt samaan tulokseen, että lämpötila on satunnaismuuttuja, jolla on kolmion muotoinen tiheysfunktio  $f(x) = x/18$  välillä  $[0, 6]$  (ja nolla muualla).

Kummankaan työnantaja ja sääennusteita odottavat asiakkaat eivät kuitenkaan halua kuulla mistään “jakaumista”, vaan he haluavat yksinkertaisesti *piste-ennusteen* eli yhden lämpötilaluvun. Abelin on siis valittava ennusteekseen jokin luku  $a \in [0, 6]$ , samoin Bertan on valittava jokin luku  $b \in [0, 6]$ . He voivat valita eri luvut, jos haluavat.

- (a) Tarkista integroimalla, että  $f$  tosiaan on kelvollinen jatkuvan jakauman tiheysfunktio.
- (b) Laske  $X$ :n *odotusarvo*  $\mu = E(X)$  sekä *mediaani* eli sellainen luku  $m$ , että  $P(X \leq m) = \frac{1}{2}$ .
- (c) Abelin työnantaja kannustaa häntä osumaan lähelle oikeaa siten, että hänen palkastaan vähennetään *neliöllinen sakko*  $(X - a)^2$ , jossain sopivissa rahayksiköissä, kun  $X$  on huomenna havaittu lämpötila ja  $a$  on Abelin ennustama lämpötila. Koska  $X$ :n arvo ei ole vielä tiedossa, Abel pyrkii valitsemaan  $a$ :n siten, että hänen sakkonsa odotusarvo olisi mahdollisimman pieni, ts. hän haluaa minimoida funktion  $q(a) = E((X - a)^2)$ . Sievennä funktio  $q$  sellaiseen muotoon, että se on yksinkertainen  $a$ :n polynomi, jossa ei ole E-merkkejä.

**Vihje:** Voit esim. aloittaa kertomalla binomin neliön auki. Vaihtoehtoisesti voit aloittaa huomamalla, että  $(X - a)^2$  on eräs  $X$ :n muunnos, ja esittää sen odotusarvon integraalina muunnoksen odotusarvon kaavalla (luento 2A).

Piirrä  $q(a)$  välillä  $a \in [0, 6]$  sen muodon hahmottamiseksi. Abel valitsee ennusteensa  $a$  siten, että  $q(a)$  on mahdollisimman pieni. Etsi tämä luku. Onko se jompikumpi luvuista  $\mu$  tai  $m$ ?

- (d) Myös Bertan työnantaja kannustaa osumaan lähelle oikeaa, mutta käyttää *lineaarista sakkoo*,  $|X - b|$ , jossain rahayksiköissä, kun  $X$  on havaittu lämpötila ja  $b$  on hänen ennusteensa. Bertan sakon odotusarvo on  $\ell(b) = E(|X - b|)$ . Kirjoita tämä integraalina ja sievennä se yksinkertaiseksi  $b$ :n polynomiksi. Piirrä funktio ja selvitä, minkä luvun Bertta valitsee, jotta  $\ell(b)$  minimoituu. Onko se jompikumpi luvuista  $\mu$  tai  $m$ ?

Vihje: Integraali kannattaa hajottaa kahteen osaan, integraaliksi välillä  $x \in [0, b)$  ja integraaliksi välillä  $x \in [b, 6]$ , jolloin pääset eroon itseisarvomerkkeistä. Huom: Bertan ongelma johtaa vähän hankalampaan integraaliin kuin Abelin, mutta laske integraalit sinnikkäästi auki.

- (e) Jos Abel ja Bertta antoivat eri lämpötilaennusteet, miksi? Koeta selittää arkijärjellä, miten eri sakkofunktiot vaikuttivat heidän toimintaansa.
- (f) (Vapaaehtoinen lisätehtävä, vaikeampi) Jos sait selville, että Abelin ja Bertan ennusteet vastaavat ennustejakauman tiettyjä kohtia, koeta todistaa että sama pätesi vaikka ennustejakauma olisi *mikä tahansa* tiheysfunktio.

Tämä on esimerkki sijoitteluongelmasta (engl. *facility location problem*, [https://en.wikipedia.org/wiki/Facility\\_location\\_problem](https://en.wikipedia.org/wiki/Facility_location_problem)). Sama matemaattinen muoto tulee vastaan esim. valittaessa jonkin palvelun (kauppa, kirjasto, tukkuvarasto) sijaintia, jos halutaan minimoida esim. käyttäjien keskimääräinen etäisyys kyseiseen palveluun, ja käyttäjien sijaintien jakauma tunnetaan. Keskimääräisen etäisyyden minimointi vastaa Bertan ongelmaa; keskimääräisen *neliöidyn* etäisyyden minimointi vastaa Abelin ongelmaa.

## Ratkaisu.

- (a) Funktio on kaikkialla epänegatiivinen. Tarkistetaan että lisäksi sen integraali on ykkönen:

$$\int_0^6 f(x)dx = \int_0^6 (x/18) dx = \int_{x=0}^6 \left( \frac{x^2}{36} \right) = 36/36 = 1.$$

- (b)

$$\mu = E(X) = \int_0^6 x f(x)dx = \int_0^6 (x^2/18)dx = \int_{x=0}^6 (x^3/54) = 4.$$

Mediaanin etsiminen:

$$P(X \leq m) = \int_0^m f(x)dx = \int_0^m (x/18) dx = \int_{x=0}^m (x^2/36) = m^2/36.$$

Merkitsemällä ko. lauseke yhtäsuureksi kuin  $\frac{1}{2}$  ja ratkaisemalla  $m$  saadaan  $m = \sqrt{18} \approx 4.2426$ .

- (c) Seurataan ensimmäistä vihjettä. Jos Abelin ennuste on  $a$ , niin sakon odotusarvo on

$$q(a) = E((X - a)^2) = E(X^2 - 2aX + a^2) = E(X^2) - 2a E(X) + a^2. \quad (*)$$



Tiedämme jo luvun  $E(X) = 4$ , ja tarvitsemme vielä luvun

$$E(X^2) = \int_0^6 x^2 f(x) dx = \int_0^6 (x^3/18) dx = \int_{x=0}^6 (x^4/72) = 18.$$

Siispä

$$q(a) = 18 - 2a \cdot 4 + a^2 = a^2 - 8a + 18.$$

Huomataan, että  $q$  on ylöspäin kaartuva paraabeli (2. derivaatta on positiivinen) ja sen minimikohta löytyy derivaatan

$$q'(a) = 2a - 8$$

nollakohdasta, joka on  $a = 4$ . Huomaamme, että Abelin valitsema ennuste on lämpötilan odotusarvo.

Huomataan lisäksi, että koska  $a = \mu$ , niin Abelin sakon odotusarvo on itse asiassa  $E((X - a)^2) = E((X - \mu)^2) = \text{Var}(X)$ .

(d) Jos Bertan ennuste on  $b$ , niin hänen sakkonsa odotusarvo on

$$\begin{aligned} \ell(b) &= E(|X - b|) = \int_0^6 |x - b| f(x) dx \\ &= \int_0^b (b - x) f(x) dx + \int_b^6 (x - b) f(x) dx \\ &= \int_0^b (bx/18 - x^2/18) dx + \int_b^6 (x^2/18 - bx/18) dx \\ &= \int_0^b (bx^2/36 - x^3/54) + \int_b^6 (x^3/54 - bx^2/36) \\ &= b^3/54 - b + 4. \end{aligned}$$

Tämä on kolmannen asteen polynomi. Tarkastelemalla sen muotoa (kuvasta tai 2. ja 3. derivaatasta) huomaamme, että  $\ell(b)$  kaartuu ylöspäin kaikkialla välillä  $[0, 6]$ . Sen derivaatalla

$$\ell'(b) = b^2/18 - 1$$

on kyseisellä välillä nollakohta  $b = \sqrt{18} \approx 4.2426$ , ja tämä on siis funktion minimikohta. Bertta valitsee ennustekseen lämpötilan mediaanin  $m$ .

(e) Sekä Abel että Bertta pyrkivät välttämään suuria ennustevirheitä, mutta koska Abelin ennusteen virhe *neliöidään*, hänellä on vielä Berttaakin voimakkaampi motiivi niiden välttämiseen. Jos Abel asettaisi ennusteensa esim. lähelle välin loppupistettä (6 astetta), ja lämpötila sattuisikin lähelle välin alkupistettä (0 astetta), niin Abelin sakko olisi aika suuri. Välttääkseen tällaisen suuren sakon mahdollisuuden hän asettaa ennusteensa hiukan alemmaksi kuin Bertta.

Toisin sanoen neliöllinen sakko johtaa varomaan suuria virheitä vielä voimakkaammin kuin lineaarinen sakko. Kumpi sitten on "parempi" sakkofunktio, riippuu täysin siitä millaisia ennusteita yhtiö haluaa. Sakkofunktioita voisi muotoilla muunkinlaisia, esim. sen perusteella, millaiset todelliset kustannukset syntyvät minkäkinlaisesta virheestä.

- (f) (Vapaaehtoinen lisätehtävä) Sama ilmiö tosiaan tapahtuisi muillakin jakaumilla: jos halutaan minimoida neliöllinen sakko, kannattaa valita odotusarvo, ja jos halutaan minimoida lineaarinen sakko, kannattaa valita mediaani.

Odotusarvon tapaus on yllättävän helppo todistaa. Huomataan, että yllä yhtälössä (\*) funktiossa  $q(a)$  on kolme termiä, mutta ensimmäinen termi  $E(X^2)$  ei mitenkään riipu  $a$ :sta, joten se katoaa derivoitaessa. Funktion derivaatta onkin

$$q'(a) = 0 - 2E(X) + 2a$$

ja tämä on nolla täsmälleen silloin kun  $a = E(X)$ . Tässä ei mitenkään käytetty tietoa siitä, mikä jakauma  $X$ :llä on, joten tulos pätee kaikilla jakaumilla.

Mediaanin kohdalla todistus on hiukan pidempi (löytyy esim. googlaamalla “median minimizes expected absolute deviation”).