

II

MATEMAATILINE STATISTIKA, TEOREETILISED JAOTUSED, PARAMEETRITE HINDAMINE

Antud peatükk püüab anda lühiülevaate matemaatilise statistika olemusest, teoreetilistest jaotustest, jaotuste peamistest parameetritest ja nende hindamismeetoditest.

2.1 JUHUSLIK SUURUS

Statistika jagatakse sageli kaheks – **kirjeldavaks** ehk **lihtsaks statistikaks** ning **matemaatiliseks statistikaks**. Neist esimese all mõistetakse andmete kokkuvõtlikku esitamist olulisemate tendentside tuvastamiseks, teine on aga mõeldud teaduslike järelduste tegemiseks empiiriliste (vaatlustest, katsetest, mõõtmistest pärinevate) andmete (**valimi**) põhjal.

Matemaatilise statistika meetodid baseeruvad tõenäosusteoorial ja on defineeritud juhusliku katse tulemuse, nn **juhusliku suuruse**, kaudu. Et juhuslik katse on iga katse, mille täpset tulemust pole võimalik ette näha, on juhuslikuks katseks ka näiteks kõikvõimalikud mõõtmised ning seega on järgevast ehk lihtsam aru saada, kui mõista juhusliku suuruse all mingit tunnust. Tunnuste mõõtmistulemusi võib siis tõlgendada juhuslike suuruste realiseerunud väärtustena. Tähistatakse juhuslikke suurusi traditsiooniliselt suurte tähtedega X, Y, Z jmt.

Teoreetiline matemaatiline statistika tegeleb teoreetiliste juhuslike suuruste (üldkogumit iseloomustavate lõpmatul hulgal objektidel mõõdetud tunnuste) ja nende funktsioonide omaduste uurimisega tuginedes juhuslike suuruste realiseerumisele (valimil mõõdetud tunnuste väärtustele).

2.2 OLULISEMAD JUHUSLIKKE SUURUSI ISELOOMUSTAVAD PARAMEETRID

2.2.1 Keskväärtus

Juhusliku suuruse **keskväärtuseks** (ooteväärtuseks e oodatavaks väärtuseks, inglise keeles *expected value, expectation*) nimetatakse tema lõpmatu hulga väärtuste keskmist. Juhusliku suuruse X keskväärtust tähistatakse $E(X)$, sageli kasutatakse ka tähte μ .

Omadused

1. $E(k) = k$, k on konstant.
2. $E(kX) = k E(X)$.
3. $E(X_1 + X_2) = E(X_1) + E(X_2)$.
4. Kui X_1 ja X_2 on sõltumatud, siis $E(X_1 X_2) = E(X_1) E(X_2)$.
5. Kui $\mathbf{X} = (X_1 \ X_2 \ \dots \ X_n)^T$, siis $E(\mathbf{X}) = (E(X_1) \ E(X_2) \ \dots \ E(X_n))^T$.

2.2.2 Dispersioon

Juhusliku suuruse X **dispersioon** $\text{var}(X)$ (ka $V(X)$, $D(X)$, σ_X^2 , inglise keeles *variance*) on defineeritud seosega

$$\text{var}(X) = E[X - E(X)]^2 = E(X^2) - E(X)E(X). \quad (2.1)$$

Omadused

1. $\text{var}(k) = 0$, k on konstant.
2. $\text{var}(kX) = k^2 \text{var}(X)$.
3. Kui X_1 ja X_2 on sõltumatud, siis $\text{var}(X_1 + X_2) = \text{var}(X_1) + \text{var}(X_2)$.

Ruutjuurt juhuliku suuruse dispersioonist nimetatakse **standardhälbeks** – $\sqrt{\text{var}(X)} = \sigma_X$.

2.2.3 Kovariatsioon

Juhuslike suuruste X_1 ja X_2 vaheline **kovariatsioon** $\text{cov}(X_1, X_2)$ (ka σ_{X_1, X_2} , inglise keeles *covariance*) kirjeldab nende juhulike suuruste vahelist lineaarset statistilist sõltuvust ja on defineeritud seosega

$$\text{cov}(X_1, X_2) = E(X_1 X_2) - E(X_1)E(X_2). \quad (2.2)$$

Omadused

1. $\text{cov}(X_1, X_2) = \text{cov}(X_2, X_1)$.
2. $\text{cov}(X, X) = \text{var}(X)$.
3. $\text{cov}(aX_1, bX_2) = ab \text{cov}(X_1, X_2)$, a, b on konstandid.
4. $\text{var}(X_1 \pm X_2) = \text{var}(X_1) + \text{var}(X_2) \pm 2 \text{cov}(X_1, X_2)$.
5. Kui X_1 ja X_2 on sõltumatud, siis $\text{cov}(X_1, X_2) = 0$ (sellest ja eelmisest omadusest järeldub dispersiooni 3. omadus).

Enam kui kahe juhuliku suuruse korral on nende dispersioone ja kovariatsioone mugav esitada maatriksina: juhulike suuruste vektori $\mathbf{X} = (X_1 \ X_2 \ \dots \ X_n)^T$ **dispersioonimaatriksiks** (kovariatsioonimaatriksiks) nimetatakse maatriksit

$$\begin{aligned} \text{var}(\mathbf{X}) &= E(\mathbf{X} - E\mathbf{X})(\mathbf{X} - E\mathbf{X})^T \\ &= \begin{pmatrix} \text{var}(X_1) & \text{cov}(X_1, X_2) & \dots & \text{cov}(X_1, X_n) \\ \text{cov}(X_1, X_2) & \text{var}(X_2) & \dots & \text{cov}(X_2, X_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(X_1, X_n) & \text{cov}(X_2, X_n) & \dots & \text{var}(X_n) \end{pmatrix} = \begin{pmatrix} \sigma_{X_1}^2 & \sigma_{X_1 X_2} & \dots & \sigma_{X_1 X_n} \\ \sigma_{X_1 X_2} & \sigma_{X_2}^2 & \dots & \sigma_{X_2 X_n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{X_1 X_n} & \sigma_{X_2 X_n} & \dots & \sigma_{X_n}^2 \end{pmatrix}, \end{aligned} \quad (2.3)$$

kus peadiagonaalil paiknevad juhulike suuruste dispersioonid ja väljaspool peadiagonaali juhulike suuruste vahelised kovariatsioonid.

Dispersioonimaatriksi omadused on analoogsed kovariatsiooni ja dispersiooni omadustega, väikesed erinevused tulenevad vaid sellest, et maatriksite korrutamisel ei kehti kõik üksikelementide korrutamise reeglid. Näiteks dispersioon konstantse maatriksiga \mathbf{A} korrutatud juhulike suuruste vektorist \mathbf{X} avaldub kujul $\text{var}(\mathbf{AX}) = \mathbf{A} \text{var}(\mathbf{X}) \mathbf{A}^T$.

2.2.4 Korrelatsioon

Jagades juhulike suuruste vahelise kovariatsiooni läbi nende standardhälvete korrutisega saame normeeritud kovariatsiooni, mida nimetatakse **korrelatsioonikordajaks**:

$$r(X_1, X_2) = \frac{\text{cov}(X_1, X_2)}{\sqrt{\text{var}(X_1) \text{var}(X_2)}}. \quad (2.4)$$

Omadused

1. $r(X_1, X_2) = r(X_2, X_1)$.
2. $r(X, X) = 1$.
3. Kui X_1 ja X_2 on sõltumatud, siis $r(X_1, X_2) = 0$.
4. $|r(X_1, X_2)| \leq 1$.

Sarnaselt kovariatsioonimaatriksiga defineeritakse juhulike suuruste vektori $\mathbf{X} = (X_1 \ X_2 \ \dots \ X_n)^T$ **korrelatsioonimaatriks** $r(\mathbf{X})$ elementidega r_{ij} , $i, j = 1, \dots, n$, kus

$$r_{ij} = r(X_i, X_j) = \text{cov}(X_i, X_j) / \sqrt{\text{var}(X_i) \text{var}(X_j)} \quad \text{ja} \quad r_{ii} = r(X_i, X_i) = 1.$$

Nii kovariatsiooni- kui ka korrelatsioonimaatriksid on sümmeetrilised maatriksid.

2.2.5 Lineaarne regressioon

Proгноosimaks juhusliku suuruse Y käitumist tingimusel, et sellega lineaarses sõltuvuses olev teine juhuslik suurus X omandab mingi kindla väärtuse, kasutatakse lineaarset regressiooniseost kujul

$$E(Y) = a + bX .$$

Võrrandi vasakul poolel olev suurus $E(Y)$ näitab, et tulemuseks on juhusliku suuruse Y oodatav (keskmine) väärtus juhusliku suuruse X mingi väärtuse korral.

Regressioonivõrrandi kordajad a ja b hinnatakse seostest

$$\hat{b} = \frac{\text{cov}(X, Y)}{\text{var}(X)} \quad \text{ja} \quad \hat{a} = EY - \hat{b}EX . \quad (2.5)$$

Kasuks tuleb teadmine, et kordaja b väljendab juhusliku suuruse Y muutumise suurust juhusliku suuruse X muutumisel ühe ühiku võrra.

2.3 TEOREETILISED JAOTUSED

Kõikvõimalikke erinevaid juhuslikke suurusi (tunnuseid) on lõpmatu palju ja ei ole mõeldav tuletada iga uue tunnuse puhul uusi eeskirju, iseloomustamaks seda tunnust juhuslike mõõdetud väärtuste (valimi) alusel. Seetõttu jagatakse juhuslikud suurused gruppidesse vastavalt nende tekkemehhanismile e olemusele ning kirjeldatakse iga gruppi parameetritest sõltuva eeskirjaga, mille abil on võimalik leida vastavasse gruppi kuuluvate juhuslike suuruste (tunnuste) väärtuste esinemise tõenäosused. Selliseid eeskirju nimetatakse **teoreetilisteks jaotusteks**. Teoreetilised jaotused on aluseks statistiliste hüpoteeside kontrollimisel, sageli ka parameetrite väärtuste ja nende hinnangute usaldusväärsuse hindamisel.

Jaotusi on kaht tüüpi – diskreetsed ja pidevad jaotused. **Diskreetne jaotus** esitatakse tõenäosusfunktsiooniga

$$p(k) = P(X=k)$$

või jaotustabeliga $\{p(k), k\}$, mis defineerib täpselt ära iga üksiku väärtuse esinemise tõenäosuse (k on juhusliku suuruse võimalik väärtus).

Pidev jaotus esitatakse tihedusfunktsiooniga

$$f(x) = dF(x)/dx,$$

mille abil on defineeritud juhusliku suuruse iga väärtuse mingisse fikseeritud vahemikku (a, b) sattumise tõenäosus:

$$P(a < X < b) = \int_a^b f(x)dx = F(b) - F(a),$$

$F(x) = P(X \leq x)$ on jaotusfunktsiooni väärtus kohal x .

Diskreetsed jaotused, mis tekivad millegi kokkulugemisel, subjektiivsel hindamisel või mingi nähtuse toimumise või mittetoimumise fikseerimisel, vajavad analüüsil siinses kursuses käsitletavate meetodite edasiarendusi – nn üldistatud lineaarseid mudelid – mistõttu on järgnevalt tutvustatud vaid traditsioonilises üldiste lineaarsete mudelite teoorias ja geneetiliste parameetrite hindamisel kasutatud leidvaid pidevaid jaotusi (ausalt öeldes neid ehk selle lühikursuse raames vaja ei lähegi, aga kui tekib soov või vajadus peale aretusväärtuste hindamise tutvuda ka dispersioonikomponentide ja päritavuskoeffitsientide hindamisega, on vaja teadmisi normaaljaotusest, ning kui huvi pakub hinnangute täpsus ja hüpoteeside kontrollimine mingi parameetri nullist erinemise kohta, kulub ära ettekujutus χ^2 -, F - ja t -jaotusest).

2.3.1 Normaaljaotus

Normaaljaotusega juhuslikku suurust X keskväertusega $E(X) = \mu$ ja dispersiooniga $\text{var}(X) = \sigma^2$ tähistatakse $X \sim N(\mu, \sigma^2)$ ja tema tihedusfunktsioon esitub valemiga

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}. \quad (2.6)$$

Keskväertuse ja dispersiooni, kui normaaljaotuse ainukeste parameetrite läbi, on tihedusfunktsioon iga juhusliku suuruse X väärtuse tarvis määratud.

Tähtsaks omaduseks on, et normaaljaotusega juhuslike suuruste lineaarkombinatsioon on samuti normaaljaotusega (muutuvad vaid parameetrite väärtused). Sagedasimaks lineaarteisenduseks on standardiseerimine

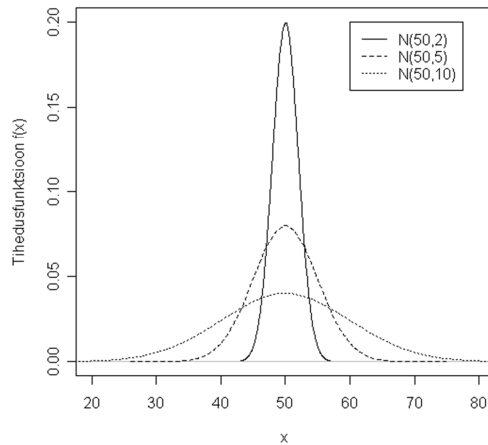
$$Z = \frac{X - \mu}{\sigma} \sim N(0,1),$$

kus $N(0,1)$ on standardne normaaljaotus, mille jaotusfunktsiooni $\Phi(x)$ väärtused on tabuleeritud (vt tabel 2.1). Seejuures kehtivad seosed $F(x) = \Phi\left(\frac{x-\mu}{\sigma}\right)$ ja $\Phi(-x) = 1 - \Phi(x)$.

Tabel 2.1. Standardse normaaljaotuse enamkasutatavad jaotusfunktsiooni väärtused; $\Phi(x) = P(X \leq x)$, kus $X \sim N(0,1)$.

$\Phi(x)$	0,005	0,025	0,05	0,5	0,95	0,975	0,995
x	-2,58	-1,96	-1,64	0	1,64	1,96	2,58

Näide. Normaalkaotusega $N(50, \sigma^2)$ on näiteks vere kogus indiviidi 50 ml vereproovis, kus σ^2 iseloomustab proovivõtmise täpsust.



Joonis 2.1. Näiteid normaaljaotuse $N(50, \sigma)$ tihedusfunktsiooni graafikuist erinevate parameetri σ väärtuste korral.

2.3.2 χ^2 -jaotus

Sõltumatute standardse normaaljaotusega juhuslike suuruste X_1, \dots, X_n ruutude summa on χ^2 -jaotusega vabadusastmete arvuga n . Sümbolite kaudu defineeritult: kui X_1, \dots, X_n on sõltumatud juhuslikud suurused, kus $X_i \sim N(0,1)$, $i = 1, \dots, n$, siis

$$\sum_{i=1}^n X_i^2 \sim \chi^2(n) \quad \text{ja} \quad \sum_{i=1}^n (X_i - \bar{X})^2 \sim \chi^2(n-1),$$

kus $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$.

Kui $X \sim \chi^2(n)$, siis $E(X) = n$ ja $\text{var}(X) = 2n$.

2.3.3 F-jaotus

Kui juhuslikud suurused $U \sim \chi^2(k_1)$ ja $V \sim \chi^2(k_2)$ ning U ja V on sõltumatud, siis on juhuslik suurus Z F-jaotusega:

$$Z = \frac{U/k_1}{V/k_2} \sim F(k_1, k_2).$$

Jaotuse parameetrid k_1 ja k_2 on positiivsed täisarvud, mida nimetatakse F-jaotuse vabadusastmeteks (inglise keeles *degrees of freedom*).

2.3.4 Student'i t-jaotus

Kui juhuslik suurus $X \sim N(0,1)$ ja juhuslik suurus $Y \sim \chi^2(n)$, kusjuures X ja Y on sõltumatud, siis

$$Z = \frac{X}{\sqrt{Y/n}} \sim t(n),$$

ehk Z on t-jaotusega vabadusastmete arvuga n .

Matemaatilise statistika rakenduste tarvis on oluline tulemus, et kui $X_i \sim N(\mu, \sigma^2)$, $i = 1, \dots, n$, siis

$$\frac{\bar{X} - \mu}{s/\sqrt{n}} \sim t(n-1),$$

kus $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ ja $s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2}$.

2.4 PARAMEETRITE HINDAMINE

Kui me ka teame, mis teoreetilisele jaotusele meie poolt uuritav suurus oma olemuselt vastab, ei tea me siiski selle jaotuse parameetreid. Viimased tuleb hinnata tuginedes andmetele (valimile). Vastav arvutusvalem on lihtsamal juhul leitav intuiitiivselt (nö mõistlikkuse printsiibist lähtudes – näiteks on loomulik võtta keskväärtuse hinnanguks valimi andmete aritmeetiline keskmine). Keerulisemal juhul ei pruugi aga intuitsioonist piisata või kui ka piisab, siis peab seda intuitsiooni kontrollima. Matemaatiline statistika pakub mitmeid võimalusi parameetrite hindamisvalemite teoreetiliseks konstrueerimiseks. Levinuimad meetodid on suurima tõepära meetod ja vähimruutude meetod.

2.4.1 Suurima tõepära meetod (*maximum likelihood method*, ML-meetod)

Suurima tõepära meetodit kasutatakse siis, kui teoreetiline jaotus on teada ja hinnatav parameeter kujutab enesest selle teoreetilise jaotuse tihedus- või tõenäosusfunktsiooni parameetrit (argumenti). Hinnanguks on siis loomulik valida see parameetri väärtus, mis realiseerunud juhul (st uuritavate andmete korral) kõige paremini sobib ehk teisisõnu on antud valimi jaoks tõepäraseim väärtus. Et eelduse kohaselt sõltus hinnatavast parameetrist ka üldkogumi jaotus, siis on ka hinnangule vastav jaotus tõepäraseim antud valimi jaoks.

Teoreetilise definitsioonina sõnastades: parameetri θ suurima tõepära hinnanguks nimetatakse väärtust $\hat{\theta}$, mille korral tõepärafunktsioon $L(\theta)$ saavutab maksimaalse väärtuse oma parameeterruumis,

$$L(\theta) = \begin{cases} f(x_1; \theta) \cdot f(x_2; \theta) \cdot \dots \cdot f(x_n; \theta), & \text{pideval juhul,} \\ p(x_1; \theta) \cdot p(x_2; \theta) \cdot \dots \cdot p(x_n; \theta), & \text{diskreetsel juhul.} \end{cases}$$

Tõepärafunktsioon kujutab enesest sama valemit, mis tihedusfunktsioongi. Erinevus seisneb selles, et tõepärafunktsioonis loetakse vastupidiselt teoreetilistele jaotustele fikseerituks andmete osa (meil on ju olemas mingid realiseerunud väärtused) ja juhuslikuks, kirjeldamist vajavaks, parameetrite osa. Juhul, kui tihedus- või tõenäosusfunktsiooni abil defineeritud teoreetiline jaotus vastab tegelikkusele, on parameetrite suurima tõepära hinnangud täpseimad.

Näide. Olgu meil vaatluse all dialleelne lookus alleelide a ja A sagedustega populatsioonis vastavalt p ja $1 - p$. Olgu eelnevalt teada, et alleeli a sagedus on kas $\frac{1}{2}$ või $\frac{1}{4}$, st $p \in \{\frac{1}{2}; \frac{1}{4}\}$. Olgu meil kaks vaatlust: $x_1 = 'a'$ ja $x_2 = 'a'$. Kumb on tõepärasem hinnang p -le, kas $\frac{1}{2}$ või $\frac{1}{4}$?

Tõepärafunktsioon: $L(p) = P(X=x_1) \times P(X=x_2) = p^2$, millest $L(\frac{1}{2}) = \frac{1}{4}$, $L(\frac{1}{4}) = \frac{1}{16}$. Kuna $L(\frac{1}{2}) > L(\frac{1}{4})$, siis $\hat{p} = \frac{1}{2}$ on suurima tõepära hinnang p -le.

2.4.2 Vähimruutude meetod (*least square method*, LS-meetod)

Vähimruutude meetod ei eelda mingi tihedus- või tõenäosusfunktsiooni kasutamist, mistõttu on selle abil saadavad hinnangud sageli lihtsamal kujul võrreldes teiste hindamismeetoditega. Nagu suurima tõepära meetod, püüab vähimruutude meetodki valida hinnanguks seda parameetri väärtust, mis realiseerunud juhul (uuritavate andmete korral) kõige paremini sobib. Ainult sobivus on defineeritud pisut teisel kujul – parim hinnang on selline, mille korral ruuterinevus realiseerunud väärtuste ja parameetri hinnangule vastavate väärtuste vahel on minimaalne. Vähimruutude meetodil on tuletatud näiteks lineaarse regressioonivõrrandi parameetrite hinnangufunktsioonid (2.5).

2.4.3 Hinnangute omadused

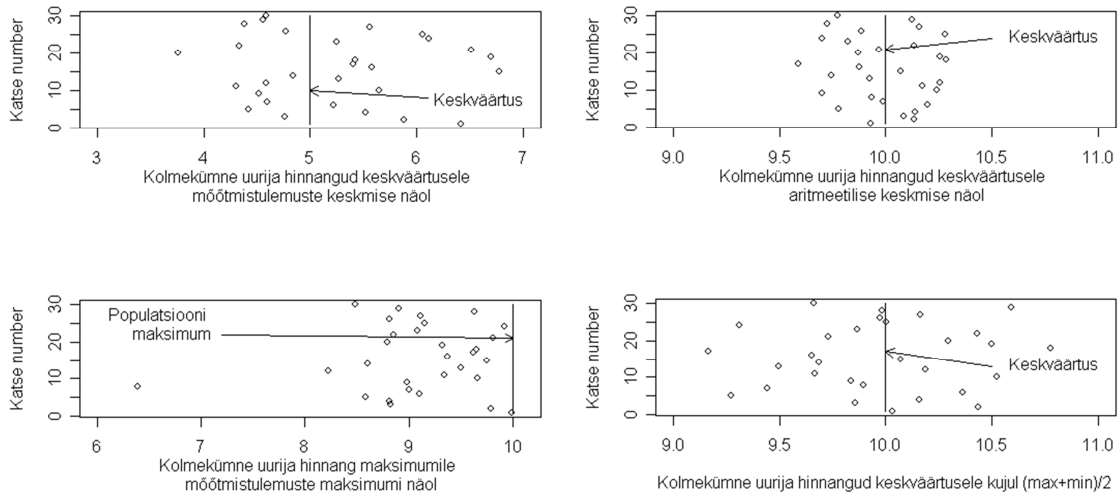
Parameetri θ hinnangut nimetatakse **nihketa hinnanguks** (*unbiased*), kui $E(\hat{\theta}) = \theta$; ehk hinnang on „keskmiselt õige“, puudub süstemaatiline viga (vt ka joonis 2.2).

Parameetri θ hinnangut $\hat{\theta}$ nimetatakse **efektiivseks hinnanguks**, kui $\text{var}(\hat{\theta})$ on vähim kõigi parameetri θ nihketa hinnangute dispersioonide hulgas; ehk – efektiivne hinnang on täpseim hinnang.

Näide. Valimi (juhusliku suuruse realiseerunud väärtuste) keskmine $\bar{x} = \frac{1}{n} \sum x_i$ on nihketa hinnang populatsiooni (juhusliku suuruse) keskvaärtusele $E(X)$:

$$E(\bar{x}) = E\left(\frac{1}{n} \sum_{i=1}^n x_i\right) = \frac{1}{n} \sum_{i=1}^n \underbrace{E(x_i)}_{E(X)} = \frac{1}{n} n E(X) = E(X).$$

Samuti on \bar{x} populatsiooni keskmise efektiivne hinnang (vt ka joonis 2.2).



Joonis 2.2. Näide nihketa ja nihkega hinnangust.

Joonis 2.3. Näide täpsemast (efektiivsest) ja vähemtäpsemast hindamismeetodist

2.4.4 Hinnangu standardviga

Et andmete alusel leitud parameetri θ hinnang $\hat{\theta}$ on juhuslik suurus, siis eksisteerib tal ka dispersioon $\text{var}(\hat{\theta})$. Viimane on aga jällegi tundmatu üldkogumi parameeter. Seega, et saada tegelikkuses aimu oma andmete alusel leitud parameetri hinnangu täpsusest, tuleb andmetest hinnata ka hinnangu dispersioon, millest reeglina parema mõistetavuse huvides võetakse veel ruutjuur (et saada varieeruvuse hinnangut samal skaalal parameetri endaga).

Hinnangu standardhälbe hinnangut nimetatakse hinnangu **standardveaks**:

$$se(\hat{\theta}) = \sqrt{\text{vâr}(\hat{\theta})}.$$

Näide. Keskvaärtuse $E(X) = \mu$ hinnangu $\hat{\mu} = \bar{x}$ dispersiooni hinnang on $\text{vâr}(\hat{\mu}) = s^2/n$ ja standardviga on

$$se(\hat{\mu}) = \frac{s}{\sqrt{n}},$$

kus s^2 on valimi dispersioon.