



Biomeetria

Biomeetria olemus, andmed, tunnuste tüübid, kirjeldav statistika




Biomeetria

Biomeetria (biostatistika) – matemaatiliste meetodite (tõenäosusteooria ja matemaatilise statistika) kasutamine bioloogiliste objektide uurimisel.


Õppeaine peamine eesmärk on anda baasteadmised andmeanalüüsist:

- 1) statistikaalase põhiterminoloogia tundmine,
- 2) oskus valida ja rakendada MS Excelis uurimisprobleemiga ja andmete tüübiga sobivaid analüüsimeetodeid,
- 3) oskus korrektselt kirjeldada ja interpreteerida andmeanalüüsi tulemusi.

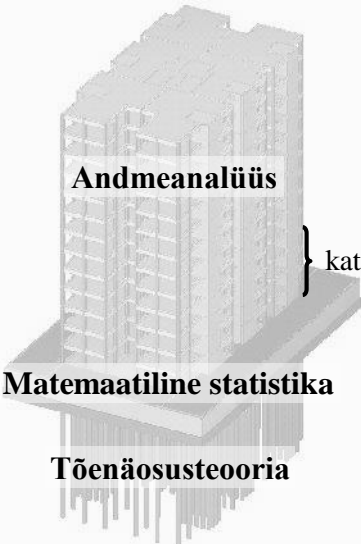


Statistika ja tõenäosusteooria loomakasvatuses ja veterinaarmeditsiinis

- ✓ **andmetöötlus** – loomade jõudlusomaduste, sigivusomaduste ja neid mõjutavate faktorite uurimine, haiguste põhjuste ja ravi mõju uurimine jne;
- ✓ **teaduslik kirjandus, ajakirjandus** – arusaamine esitatud tulemustest (mõisted, tähistused jne), kriitiline hinnang ajakirjanduses ilmuvate artiklite kohta;
- ✓ **populatsioonigeneetika ja evolutsiooniteooria** – aitab mõista geneetilise informatsiooni edasikandumist ajas ja ruumis (aretus);
- ✓ **modelleerimiseksperimendid, infotöötlus** – uuritakse haiguse levimust (epidemiologia), sööda lõhustuvust, mingi keemilise elemendi segunemist jne.



Andmeanalüüsi olemus




Andmeanalüüs teeb teaduslikke järeldusi reaalsete (vaatlustest, katsetest, mõõtmistest pärinevate) andmete põhjal, valides rakendatavad statistikameetodid nii, et need võimalikult hästi andmetega sobiksid.

} katseplaan; andmed ja nende esitus;
kirjeldav statistika;
valim *versus* populatsioon

Matemaatiline statistika tegeleb teoreetiliste andmete $\mathbf{X} = (X_1, \dots, X_n)$ ja nende funktsioonide $T(\mathbf{X})$ (e statistikute) tõenäosuslike omaduste uurimisega ning statistiliste otsustuste tegemisega.

Tõenäosusteooria




Objekt

Objekt on uurimisalune ühik, üksikindiviid (näiteks lehm, talu, põllulapp, firma, inimene, punkt metsas või järvel).

Ka samade andmete puhul võib uurimisobjekti valikuks olla mitu erinevat võimalust. Näiteks: 2 pesakonda kutsikaid – ühes 2, teises 6 kutsikat

	Objekt – kutsikas			Objekt – pesakond	
	Kutsikas	Psk nr	Psk suurus	Psk nr	Psk suurus
Haiguse levimust uurides võivad uurimispunktideks olla kas lehmad või laudad:	1	1	6	1	6
	2	1	6	2	2
	3	1	6		
“5% vaadeldud lehmadest põdesid aasta jooksul uuritavat haigust”	4	1	6		
	5	1	6		
vs	6	1	6		
	7	2	2		
“80% vaadeldud lautadest esines uuritavat haigust aasta jooksul”.	8	2	2		
	Keskmine psk suurus: 5			Keskmine psk suurus: 4	

Lehmade piimajõudlust uurides võib objektideks valida näiteks lüpsikorrad või lüpsipäevad või hoopis laktatsiooni, kusjuures uuritava tunnuse väärtuste stabiilsus võib märgatavalt sõltuda meie valikust (näiteks lüpsikorral lüpsitud piimakoguste varieeruvus on ilmselt suurem võrreldes päevalüpside varieeruvusega).



Tunnus

Tunnus on objekti iseloomustav näitaja, mida põhimõtteliselt on võimalik mõõta või vaadelda.

Näiteks:
päevane piimaand, tõug ja vanus lehmi uurides,

talusid uurides talu aastane sissetulek, töötajate arv, põllumaa pindala ja kaugus lähimast linnast jne.




Statistiline andmestik

Objekt-tunnus-maatriks – tabel, kus iga veerg kujutab ühte tunnust ja iga rida ühte objekti.

Näide. Ühest tallist koguti andmed müügis olevate hobuste tõu, värvuse ja hinna kohta. Saadud andmemaatriks (-tabel) on järgmine:

Tõug	Värvus	Hind
tori	raudjas	9000
trakeen	hall	26000
tori	raudjas	16000
eesti hobune	kõrb	9000
eesti hobune	raudjas	12000
trakeen	kõrb	35000
tori	kõrb	17000
eesti hobune	raudjas	14000
tori	must	21000
eesti hobune	kõrb	19000

Objekt-tunnus-maatriks



Tunnuste tüübid


Arvulised e. kvantitatiivsed tunnused [*numerical*]

Diskreetse [*discrete*] **tunnuse** väärtused saavad olla vaid täisarvulised, peaaegu alati on need tekkinud millegi loendamisel.
Näiteks pesakonna suurus, terade arv viljapeas, laktatsiooni number, ...

Pideva [*continuous*] **tunnuse** võimalike väärtuste arv lõpmatu ja iga kahe võimaliku väärtuse vahele mahub veel vähemalt üks pideva tunnuse võimalik väärtus; pideva tunnuse väärtused saadakse enamasti millegi otsesel mõõtmisel.
Näiteks piimatoodang, villa pikkus, esmapoegimise iga, saagikus, pH, ...

Soovitused:

- ✓ kõik tunnuse väärtused peaksid olema mõõdetud sama täpsusega,
- ✓ sama tunnuse väärtuste puhul tuleks kasutada samu ühikuid
(lehma 1 toodang 6300 (kg), lehma 2 toodang 7,9 (tonni) – keskmine toodang 3153,95!?).



Tunnuste tüübid

Mittearvulised e. kvalitatiivsed tunnused [*categorical*]

Järjestustunnuse [*ordinal*] väärtuste vahel on võimalik objektiivne järjestus (hinnangud etteantud skaalal jm).
Näiteks haridus (alg- / kesk- / kõrgharidus), poegimiskeskus, hinnang mulla niiskusele (väga kuiv / kuiv / paras / niiske / liigniiske), hinnang pulli välimusele (niru / normaalne / kaunis), ...

Probleemiks võimalikud subjektiivsed hinnangud (milline pull on kaunis?)!

Nominaalsed tunnused [*nominal*] on mittearvulised tunnused, mille väärtuste vahel ei ole sisulist järjestust.
Näiteks tõug, värvus, farm, kasvukoht, ...

Binaarsed (dihhotoomsed) tunnused on kahe väärtusega järjestustunnused.
Näiteks sugu, haige/terve, tiinestus/ei tiinestunud.

Tunnuste tüübid



Näide. Uuriti Lõuna-Eestis asuvaid talusid, kogutud andmed on esitatud järgnevas tabelis.

<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>
suurus (ha)	müügitulu	peamine tegevusala	põllumaa kvaliteet	talupere suurus
25	340000	1	1	5
15	220000	2	2	4
44	700000	1	3	4
12	500000	3	3	2
20	1200000	2	1	2

Tunnused *A* ja *B* on pidevad, *C* nominaalne (kasutatud kodeering: 1-karjakasvatus; 2-viljakasvatus; 3-turism), *D* on järjestustunnus (kasutatud kodeering: 1-väga hea; 2-keskmine; 3-kehv) ning tunnus *E* on diskreetne tunnus.

Tunnuste kodeerimine



Kodeerimine – sõnaliste vastusevariantide arvudega asendamine.

Näiteks tunnuse “arvamus valitsusest” väärtuste sisestamisel võime vastsevariandi “valitsus on hea” asemel sisestada numברי “1”, vastusevariandi “valitsus on kesk-pärane” asemel numברי “2” ja vastusevariandi “valitsus on saast” asemel numברי “3”

✓ **Järjestustunnuste kodeerimisel tuleb jälgida, et koodid säilitaksid väärtuste sisulise järjestuse.**

Vastuse variandid	Vastuse variantide koodid			
hea	1	1	3	1
halb	2	3	1	-1
ei tea	3	2	2	0

✓ Binaarse tunnuse kodeerimisel on eelistatavaim nõ 0/1-kodeering (huvipakkuv sündmus ei toimunud / toimus).

✓ Nominaaltunnuseid ei ole enamasti vaja arvuliseks kodeerida, ja kui kodeerida, siis koodid sisulist tähendust ei oma (loogiline oleks näiteks järjestada väärtused tähestiku järjekorras).

Statistiline andmestik – märkusi ja soovitusi



Ka hästi planeeritud uurimuse korral võib juhtuda, et kõigi objektide korral ei ole teada kõigi tunnuste väärtusi ja andmestik jääb lünklikuks.

Puuduv väärtus peab olema tähistatud nii, nagu ei tähistata andmestikus midagi muud.

Plaanides andmeid analüüsida standardse statistikatarkvara (SAS, R, Statistica, ...) või mõne tabelarvutussüsteemi (MS Excel, Open Office, ...) abil,

on mõistlik **jätta puuduvale väärtusele vastav lahter tühjaks**.

Arvu 0 puuduva väärtuse tähisena võiks vältida.

Andmeanalüüsi tüübid



✓ **kirjeldav statistika** – andmete kokkuvõtlik/ülevaatlik esitamine:

- arvkarakteristikud,
- sagedustabelid,
- joonised.

✓ **analüüsiv statistika** – andmete põhjal üldiste järelduste ja otsustuste tegemine:

- parameetrite hindamine,
- hüpoteeside kontroll,
- mudelite konstrueerimine.

Kirjeldav statistika – arvarakteristikud



- andmestiku suurus (valimi maht) – n
- (aritmeetiline) keskmine [*average, mean*] – $\bar{x} = \sum_{i=1}^n x_i / n$
- mediaan (nn 50%-punkt) [*median*]
- mood [*mode*] – enim esinev (suurima sagedusega) väärtus

Konkreetsete mõõdetud väärtuste tähistamiseks kasutatakse väikeseid tähti ja soovides täpsustada objekti, kellel/millel see väärtus on mõõdetud, esitakse objekti number alaindeksis:

x_3 on tunnuse X väärtus 3. objektil (näiteks 3. lehmäl).

Kirjeldav statistika – arvarakteristikud



- andmestiku suurus (valimi maht) – n
- (aritmeetiline) keskmine [*average, mean*] – $\bar{x} = \sum_{i=1}^n x_i / n$
- mediaan (nn 50%-punkt) [*median*]
- mood [*mode*] – enim esinev (suurima sagedusega) väärtus


Näide. Uuringu all olnud 5-l haigestunud loomal määrati haiguse peiteajaks vastavalt 8, 16, 12, 60 ja 14 päeva (üks uuritud loomadest oli ilmselt geneetiliselt erinev või siis sai juba mingit muud, haiguse avaldumist pärssivat ravi). Haiguse keskmine peiteaeg on

$$\bar{x} = \frac{8+16+12+60+14}{5} = \frac{110}{5} = 22 \text{ päeva.}$$

Peiteaeg, millest pooltel loomadel avaldus haigus varem ja pooltel hiljem, on leitav kui kasvavalt järjestatud peiteaegade keskmine väärtus e mediaan:

$$8, 12, \underline{14}, 16, 60 \\ = med$$

Kirjeldav statistika – arvarakteristikud



- standardhälve [*standard deviation*] – $s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$
- dispersioon [(*sample*) *variance*] – s^2
- standardviga [*standard error*] – $se = s/\sqrt{n}$
- variatsioonikordaja [*coefficient of variation*] – $v = \frac{s}{\bar{x}} \times 100\%$
(omab mõtet üksnes positiivsete väärtustega tunnuste korral!)

Näide. Uuriti 5 metsiku ja 4 puhtatõulise laborihiire reaktsiooni ärritajale.
Tulemuseks saadi järgmised väärtused:
metsikud hiired – 15, 45, 30, 10, 25; labori hiired – 20, 25, 30, 25.


Keskmisel reaktsioonid kummagi grupi jaoks on

$$\bar{x}_m = \frac{15+45+30+10+25}{5} = \frac{125}{5} = 25, \quad \bar{x}_l = \frac{20+25+30+25}{4} = \frac{100}{4} = 25.$$

$$s_m = \sqrt{\frac{(15-25)^2 + (45-25)^2 + (30-25)^2 + (10-25)^2 + (25-25)^2}{5-1}} = \sqrt{\frac{750}{4}} = \sqrt{187,5} \approx 13,69;$$

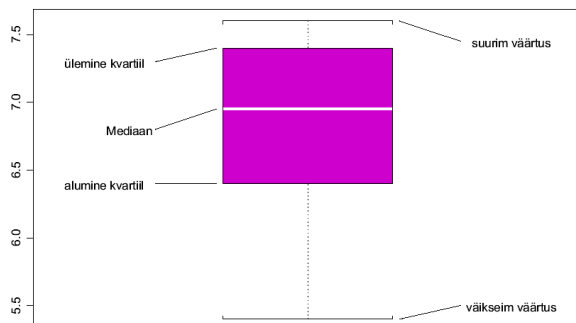
$$s_l = \sqrt{\frac{(20-25)^2 + (25-25)^2 + (30-25)^2 + (25-25)^2}{4-1}} = \sqrt{\frac{50}{3}} \approx \sqrt{16,67} \approx 4,08.$$

Kirjeldav statistika – arvarakteristikud

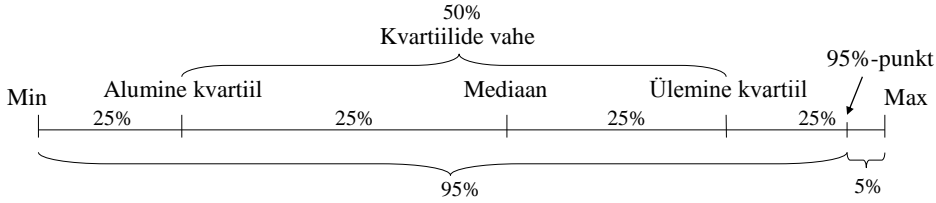


Karp-vurrud diagramm

- kvartiilid, kvartiilide vahe [*lower, upper quartile*]
- detšiilid, protsentiilid/**kvantiilid**
- min, max



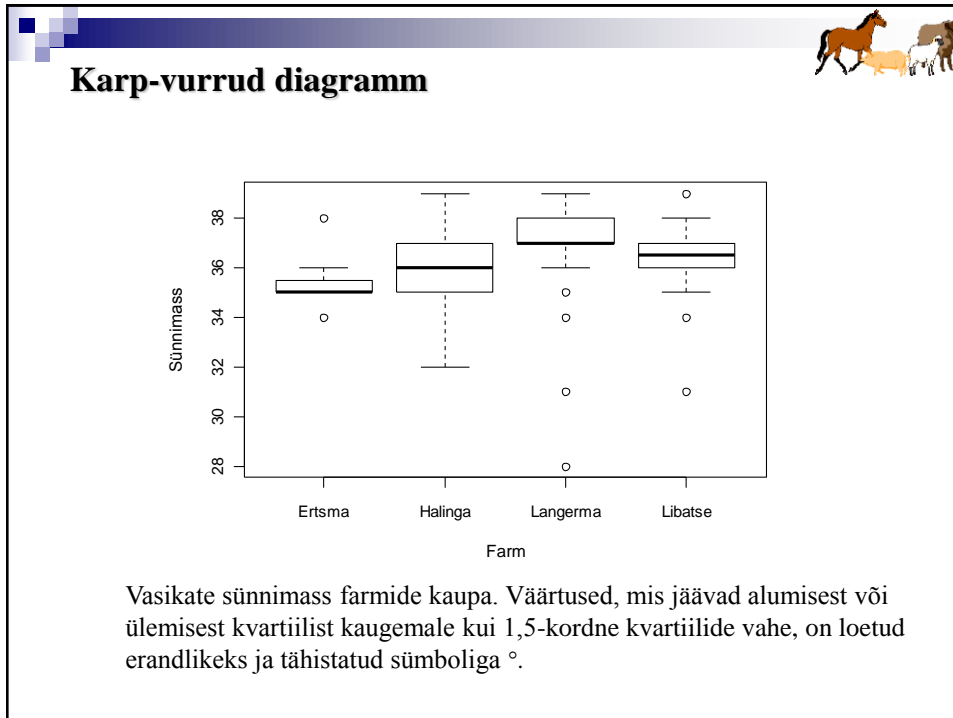
50%
Kvartiilide vahe



Min Alumine kvartiil Mediaan Ülemine kvartiil 95%-punkt Max

25% 25% 25% 25%

95% 5%



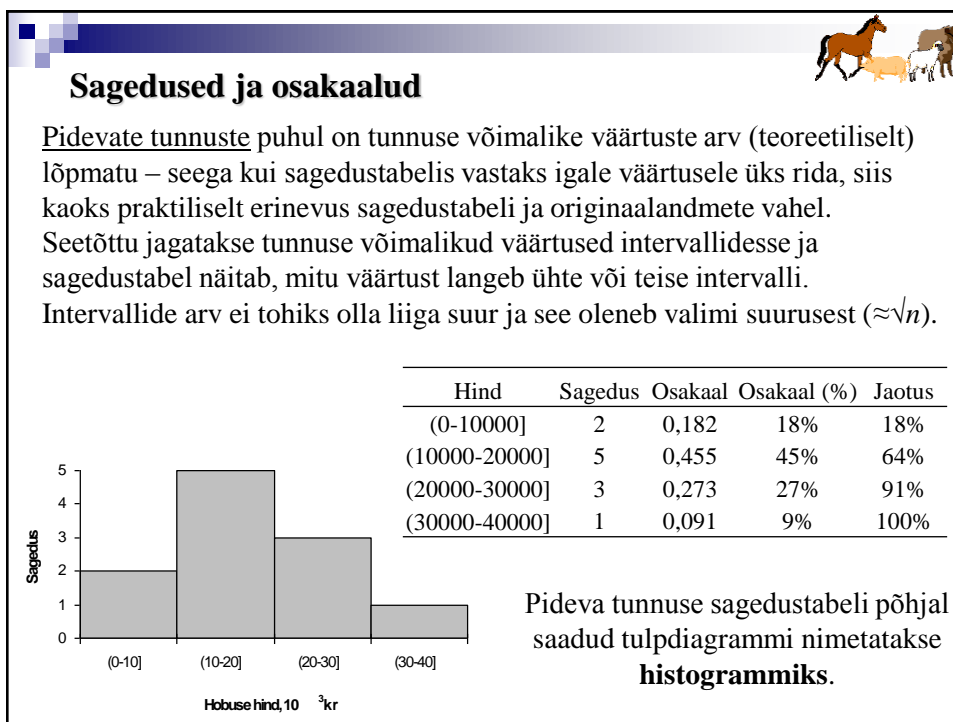
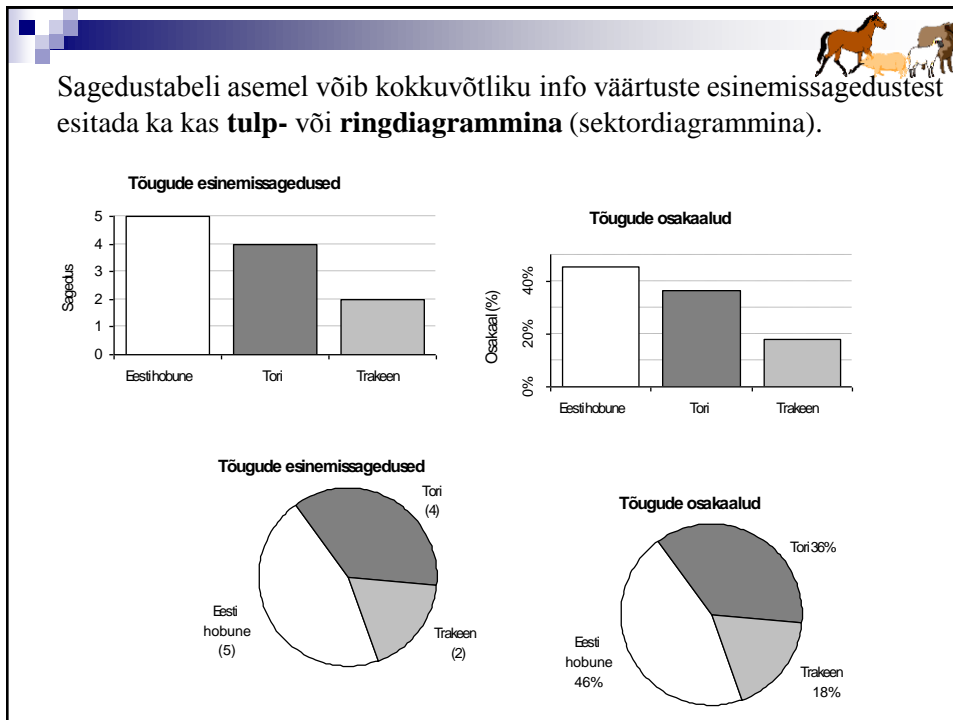
Kirjeldav statistika – sagedused ja osakaalud

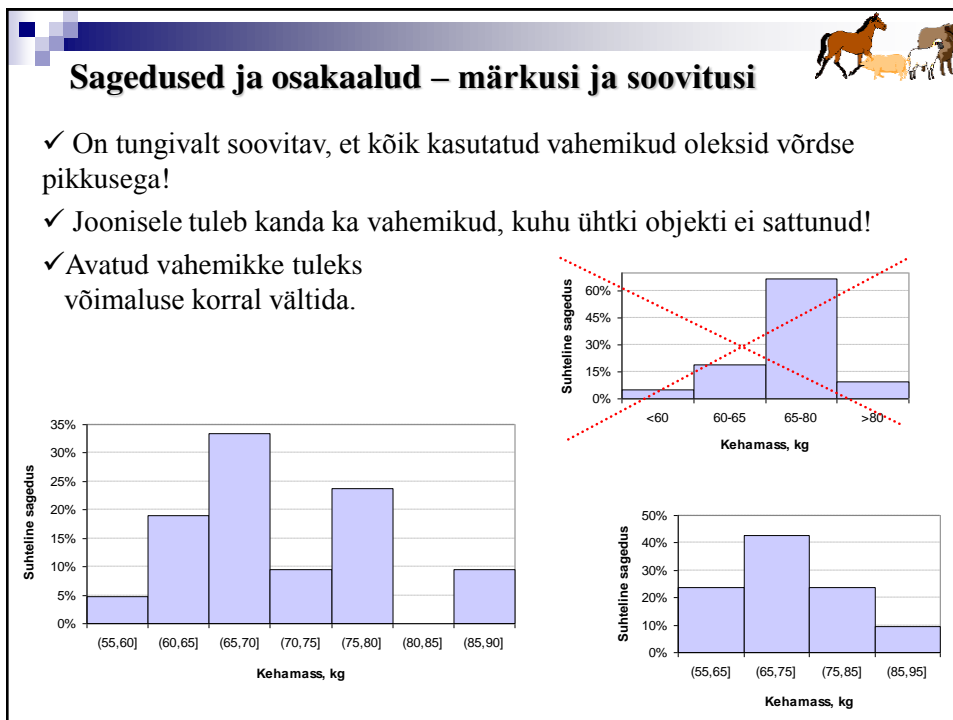
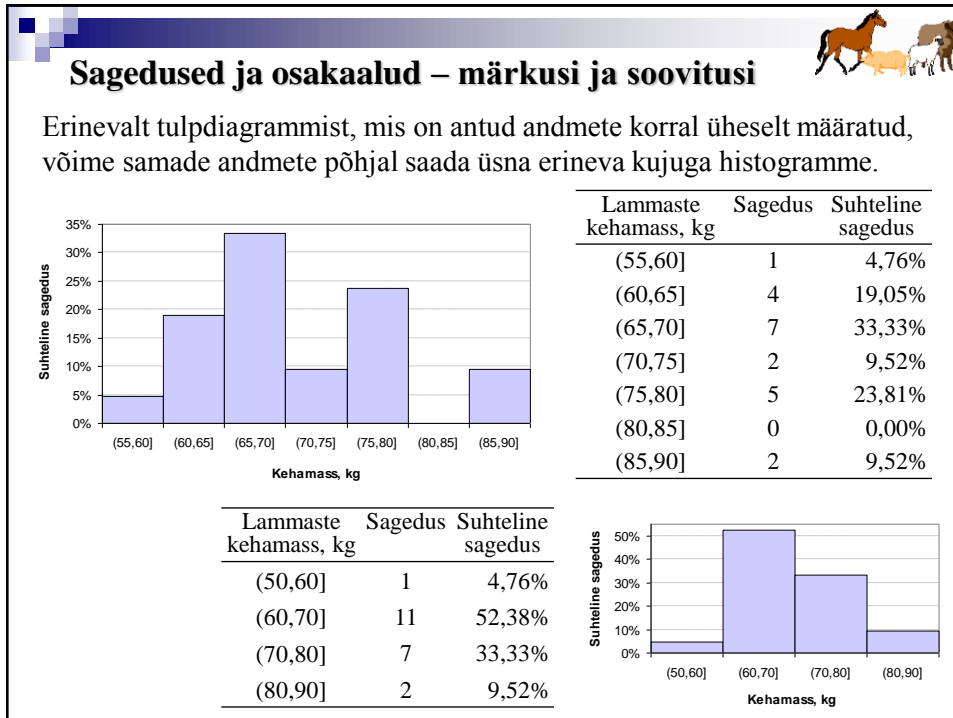
Mittearvuliste või diskreetsete tunnuste (erinevate väärtuste arv suhteliselt väike) ülevaatlikuks kirjeldamiseks on lihtne lugeda kokku, mitu korda iga erinevat väärtust esineb ja kirjutada saadud arvud tabeli kujul. Väärtuse esinemiste arvu nimetatakse tema **sageduseks**.

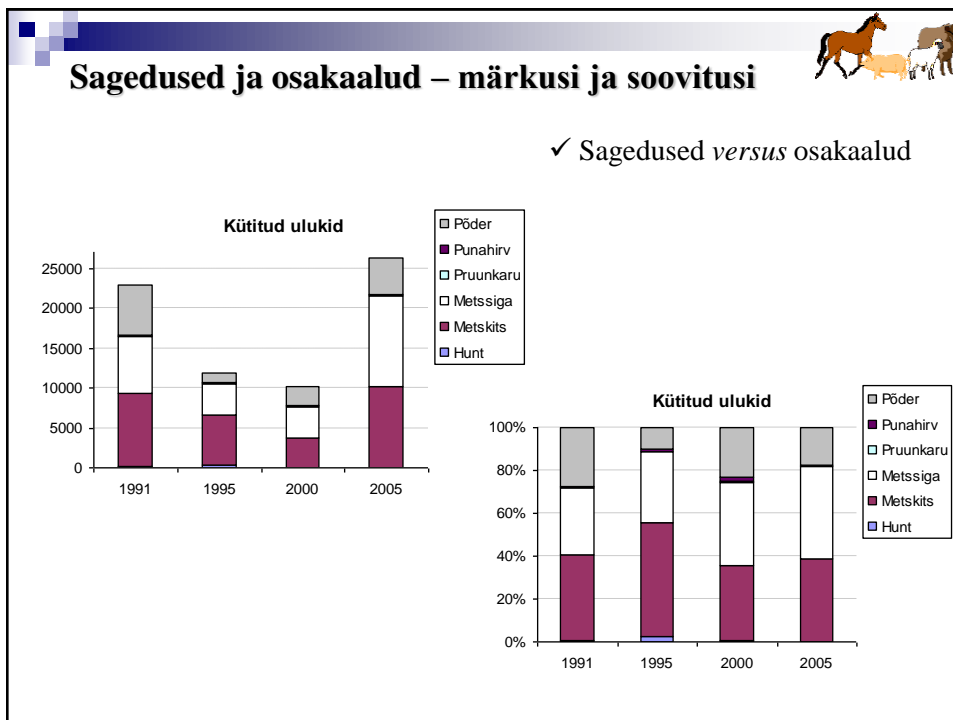
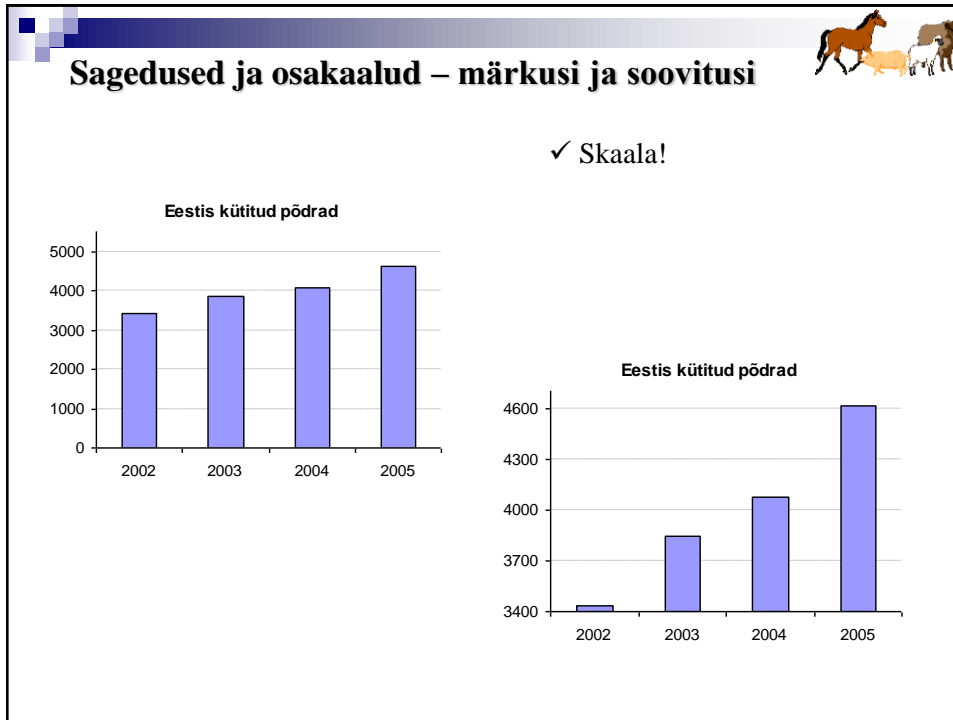
Tihti leitakse lisaks iga väärtuse (protsentuaalne) **osakaal** valimis, mida nimetatakse ka **suhteliseks sageduseks**.

Tõug	Värvus	Hind
tori	raudjas	9000
trakeen	hall	26000
tori	raudjas	16000
eesti hobune	kõrb	9000
eesti hobune	raudjas	12000
trakeen	kõrb	35000
tori	kõrb	17000
eesti hobune	raudjas	14000
tori	must	21000
eesti hobune	kõrb	19000

Tõug	Sagedus	Osakaal	Osakaal (%)
eesti hobune	5	0,455	45,5%
tori	4	0,364	36,4%
trakeen	2	0,182	18,2%





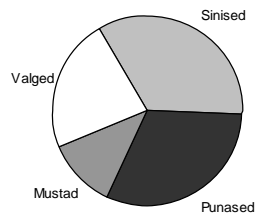


Sagedused ja osakaalud – märkusi ja soovitusi



✓ Vältida tuleks 3-mõõtmelisi graafikuid, eriti ringdiagramme.

Rahva eelistused poliitiliste erakondade osas



Erakond	Osakaal (%)
Sinised	34
Punased	31
Mustad	12
Valged	23

Rahva eelistused poliitiliste erakondade osas

