

Biomeetria praks 3

Illustreeritud (mittetäielik) tööjuhend

Eeltöö

1. Avage *MS Excel*'is oma kursuse ankeedivastuseid sisaldav andmestik,
2. lisage uus tööleht (*Insert / Lisa -> Worksheet / Tööleht*), nimetage see ümber leheküljeks 'Praks3' ja
3. kopeerige kogu 'Andmed'-lehel paiknev andmetabel lehekülje 'Praks3' ülemisse vasakusse nurka.

Ülesanne 1.

- Leidke andmetabeli alla (NB! Vähemalt üks tühi rida jätke vahele!) kõigi arvtunnuste kohta vaatluste arv (n), keskmine väärtus (\bar{x}), mediaan, standardhälve (s), standardviga (se), minimaalne ja maksimaalne väärtus, kasutades *Exceli* funktsioone.
- Lisage andmetabelisse uus tunnus nimega 'KMI' (kehamassiindeks) ja arvutage selle väärtused kõigile tudengitele valemist $KMI = \text{Kehamass, kg} / (\text{Pikkus, m})^2$.
Leidke eelnevalt nimetatud arvarakteristikute väärtused ka uuele tunnusele.

Tööjuhend

1. Jätke andmetabeli alla vähemalt üks tühi rida

(see on vajalik, et *Excel* mitmete operatsioonide teostamisel – näiteks andmete sorteerimisel või filtreerimisel, *Pivot Table*'i rakendamisel – ei tõlgendaks arvutatud keskmisi ja muid näitajaid andmetabeli osana)

ja kirjutage esimesse veergu leitavate arvarakteristikute nimed (siis on hiljem lihtsam aru saada, mida kuhugi arvutatud on).

66	N	170	60	58	3	K
67	N	170	73	57,5	4	sal
68						
69	Vaatluste arv					
70	Keskmine					
71	Mediaan					
72	Standardhälve					
73	Standardviga					
74	Min					
75	Max					

2. Arvutage kõigi arvarakteristikute väärtused tudengite pikkuse kohta, kasutades *Exceli* funktsioone.

a) Selleks võite valida vastava funktsiooni *Exceli* funktsioonide listist (vajalike funktsioonide nimed leiate järgmiselt leheküljelt punktist b):

The image shows a multi-step process in Microsoft Excel:

- Data Table:** A table with columns A-E and rows 65-75. Row 69 is highlighted with a dashed box, containing the text "Vaatluste arv" and the value "65".
- Insert Function Dialog:** A dialog box titled "Insert Function" is open. The "Or select a category" dropdown is set to "Statistical". The "COUNT" function is selected in the list. The description below reads: "COUNT(value1;value2;...) Counts the number of cells that contain numbers and numbers within the list of arguments."
- Function Arguments Dialog:** A dialog box titled "Function Arguments" for the COUNT function. The "Value1" field contains the range "B2:B67". The resulting value is shown as "65".
- Final Result:** The value "65" is entered into cell B69 of the spreadsheet.

- b) Teades funktsiooni nime ja süntaksit, võite trükkida vastava käsu ka kohe *Exceli* töölehe vastavasse lahtrisse.
(NB! Ärge unustage alustamast käsku võrdusmärgiga '='!)

Kõik need funktsioonid on rakendatavad ka eelmisel leheküljel esitatud viisil – valige ise, milline variant omale arusaadavam ja mugavam tundub.

Vaatluste arv	=COUNT(B2:B67)
Keskmine	=AVERAGE(B2:B67)
Mediaan	=MEDIAN(B2:B67)
Standardhälve	=STDEV(B2:B67)
Standardviga	
Min	=MIN(B2:B67)
Max	=MAX(B2:B67)

- c) Et *Excelis* puudub eraldi funktsioon standardvea leidmiseks, tuleb arvutused teostada, tuginedes standardvea arvutusvalemile $se = s/\sqrt{n}$ (st, et vastav valem tuleb ise sisestada):

	A	B	C
68			
69	Vaatluste arv	65	
70	Keskmine	174,877	
71	Mediaan	175	
72	Standardhälve	10,1608	
73	Standardviga	=B72/SQRT(B69)	
74	Min	153	
75	Max	197	

3. Rakendage samu funktsioone, arvutamaks soovitud arvkarakteristikute väärtused kõigi andmestikus sisalduvate arvtunnuste jaoks.
Ümardage keskmised, standardhälbed ja standardvead ühe kohani peale koma.

66	N	1/U	6U	56	3	K
67	N	170	73	57,5	4	sal
68						
69	Vaatluste arv	65				
70	Keskmine	174,877				
71	Mediaan	175				
72	Standardhälve	10,1608				
73	Standardviga	1,2603				
74	Min	153				
75	Max	197				

Vaatluste arv	65	64	41	66
Keskmine	174,877	69,46	55,676	3,54545455
Mediaan	175	70	56	3
Standardhälve	10,1608	13,8	3,185	0,66057826
Standardviga	1,2603	1,725	0,4974	0,08131156
Min	153	42	50	3
Max	197	100	66	5

Vaatluste arv	65	64	41	66
Keskmine	174,9	69,5	55,7	3,5
Mediaan	175	70	56	3
Standardhälve	10,2	13,8	3,2	0,7
Standardviga	1,3	1,7	0,5	0,1
Min	153	42	50	3
Max	197	100	66	5

7. Leidke vajalikud arvarakteristikute väärtused ka uuele tunnusele.

N	170	73	57,5	4	saksa keel	Hum	jah	ei	25,25952
Vaatluste arv	65	64	41	66					65
Keskmine	174,9	69,5	55,7	3,5					#DIV/0!
Mediaan	175	70	56	3					#DIV/0!
Standardhälve	10,2	13,8	3,2	0,7					#DIV/0!
Standardviga	1,3	1,7	0,5	0,1					#DIV/0!
Min	153	42	50	3					#DIV/0!
Max	197	100	66	5					#DIV/0!

Kopeeri / Copy -> Kleebi / Paste

- Milles probleem?

Probleem on andmetabeli reas, kus on puudu nii pikkus kui ka kehamass.

	A	B	C	D	E	F	G	H	I	J
1	SUGU	PIKKUS	MASS	PEA_P	MAT_HINNE	EBA_AINE	AINEKOOD	PUDER	ÖNNELIK	KMI
33	N	173	70		4	matemaatika	Reaal	ei	jah	23,38869
34	M				5	kehaline kasv.	Muu	jah	ei	#DIV/0!
35	N	175	70		3	matemaatika	Reaal	ei	jah	22,86714

Et Excel tõlgendab

puuduvaid väärtusi aritmeetilistes tehetes (aga mitte statistikafunktsioonides!) 0-dena ja kehamassiindeksi arvutusvalemis tuleb teostada jagamine pikkuse ruuduga, siis on puuduva pikkuse korral tulemuseks 0-ga jagamine, mida ei ole aga võimalik teostada. Sestap ka veateade #DIV/0!, mis just 0-ga jagamist tähendabki.

Ja kui selline veateadet sisaldav andmeblokk määrata argumentiks näiteks keskmist arvutavale funktsioonile AVERAGE, siis on ka viimase tulemuseks analoogne veateade.

- Mida teha? Antud juhul on lihtsaim variant **kustutada ära vigane väärtus kehamassiindeksi vastavast lahtrist:**

	A	B	C	D	E	F	G	H	I	J
1	SUGU	PIKKUS	MASS	PEA_P	MAT_HINNE	EBA_AINE	AINEKOOD	PUDER	ÖNNELIK	KMI
33	N	173	70		4	matemaatika	Reaal	ei	jah	23,38869
34	M				5	kehaline kasv.	Muu	jah	ei	#DIV/0!
35	N	175	70		3	matemaatika	Reaal	ei	jah	22,86714

Tulemus:

Vaatluste arv	65	64	41	66						65
Keskmine	174,9	69,5	55,7	3,5						22,2
Mediaan	175	70	56	3						22,41027
Standardhälve	10,2	13,8	3,2	0,7						4,2
Standardviga	1,3	1,7	0,5	0,1						0,5
Min	153	42	50	3						0
Max	197	100	66	5						30,7174

- Kõik korras? **Minimaalne kehamassiindeks on 0!** Põhjuseks see, et lisaks ühele juba kustutatud väärtusele on andmestikus veel üks rida, kus pikkus on küll teada, aga kehamassi mitte, mistõttu Excel käsitleb viimast 0-na ja saab ka kehamassiindeksi väärtuseks 0-i.

Õige oleks **kustutada ära ka see 0-ga võrduv kehamassiindeks.**

ELIK	KMI	ELIK	KMI
naa	21,92922	naa	21,92922
naa	0	naa	0
naa	22,06881	naa	22,06881

Tulemus:

	SUGU	PIKKUS	MASS	PEA_P	MAT_HINNE	EBA_AINE	AINEKOOD	PUDER	ÖNNELIK	KMI
Vaatluste arv		65	64	41	66					64
Keskmine		174,9	69,5	55,7	3,5					22,50
Mediaan		175	70	56	3					22,54
Standardhälve		10,2	13,8	3,2	0,7					3,17
Standardviga		1,3	1,7	0,5	0,1					0,40
Min		153	42	50	3					16,01
Max		197	100	66	5					30,72

Ülesanne 2.

- Arvutage tudengite pikkuse, massi ja peaumbermõdu kohta nii palju arvkarakteristikuid, kui protseduur *Descriptive Statistics (Tools / Tööriistad -> Data Analysis ...)* võimaldab.
- Leidke ka 90%, 95% või 99% usalduspiirid keskmistele väärtustele. Mida need usalduspiirid näitavad?

Tööjuhend

- Arvkarakteristikute arvutamine: *Tools / Tööriistad -> Data Analysis ... -> Descriptive Statistics*

	A	B	C	D	E	F	G	H	I	J	K	L	M
1	SUGU	PIKKUS	MASS	PEA_P	MAT	HINNE	EBA_AINE	AINEKOOD	PUDER	ÖNNELIK	KMI		
2	N	165	59		4	füüsika	Reaal	jah	jah	21,6713			
3	N	165	49		4	kirjandus	Hum	jah	jah	17,9982			
4	N	171	70	56	3	ajalugu	Hum	ei	jah	23,939			
5	M	186	75		4	eesti keel	Hum	jah	jah	21,6788			
6	M	182	73	56						22,0384			
7	N	172,5	60							20,1638			
8	N	167	75							26,8923			
9	N	175	64	50						20,898			
10	N	169	72	55						25,2092			
11	M	174	93	58	3	matemaatika	Reaal	jah	jah	30,7174			
12	M	191	88	60									
13	M	175	70	56									
14	N	157	57										
15	M	187	70	56									
16	M	187	56										
17	N	164	50										
18	M	185	80	56									
19	M	179	89	56									
20	N	177	79	53									
21	M	191	80	60									

Võimalik on analüüsida mitut tunnust korraga tingimusel, et nende väärtused paiknevad kõrvuti veergudes

Descriptive Statistics

Input

Input Range:

Grouped By: Columns Rows

Labels in first row

Output options

Output Range:

New Worksheet Ply:

New Workbook

Summary statistics

Confidence Level for Mean: %

Kth Largest:

Kth Smallest:

OK Cancel Help

Valik 'Labels in first row' peab olema märgitud, kui andmed on ette antud koos esimeses reas paikneva nimega.

Lisavalikute 'Summary statistics' jt kohta vt järgmine lk.

Väljundtabeli vasaku ülemise nurga asukoht

28	N	171	54	53,2									
29	N	178	80	56									
30	N	157	48										
31	N	180	67	61									
32	M	190	85										
33	N	173	70										
34													
35													
36													
37													
38													
39													
40	M	183	80		3	keemia	Reaal	jah	jah	22,6627			
45	M	177	75	56	4	eesti keel	Hum	jah	nii ja naa	23,9395			
46	N	165	46		3	eesti keel	Hum	jah	nii ja naa	16,8962			
47	M	176	80	52	4	saksa keel	Hum	nii ja naa	ei	25,8264			
48	M	185	100	58	3	matemaatika	Reaal	ei	jah	29,2184			
49	M	177	75		3	ajalugu	Hum	jah	jah	23,9395			
50	N	177	62	56	3	kehaline kasv.	Muu	ei	jah	19,79			
51	M	186	76	55	3	keemia	Reaal	jah	jah	21,9679			
52	M	178	68	57	4	vene keel	Hum	jah	jah	21,4619			
53	M	173	72	58	4	vene keel	Hum	nii ja naa	nii ja naa	24,0569			
54	N	165	59	56	3	vene keel	Hum	jah	nii ja naa	21,6713			
55	N	167	66		4	matemaatika	Reaal	jah	jah	23,6652			
56	N	177	85		4	keemia	Reaal	jah	jah	27,1314			
57	N	161	44	56	3	ajalugu	Hum	ei	jah	16,9747			
58	N	159	50		4	füüsika	Reaal	jah	jah	19,7777			
59	N	175	51		3	geograafia	Hum	jah	nii ja naa	16,6531			
60	N	158	50	55	4	vene keel	Hum	ei	jah	20,0288			
61	N	167	62,5	55	3	füüsika	Reaal	jah	jah	22,4103			
62	M	186	95	57	3	matemaatika	Reaal	ei	jah	27,4598			
63	N	170	60	55	3	füüsika	Reaal	ei	jah	20,7612			
64	M	197	85	50	3	matemaatika	Reaal	ei	jah	21,9021			
65	N	153	65	57	4	ajalugu	Hum	jah	jah	27,7671			
66	N	170	60	58	3	keemia	Reaal	jah	jah	20,7612			
67	N	170	73	57,5	4	saksa keel	Hum	jah	ei	25,2595			

Selgitus protseduuri *Descriptive Statistics* lisavalikutest eelmisel lehel:

- valiku *Summary statistics* tulemusena arvutab *Excel* kaheteistkümne põhilise arvarakteristiku väärtused;
- valiku *Confidence Level for Mean: 95%* tulemusena arvutatakse suurus, mis tuleb keskmisele juurde liita või lahutada, saamaks ülemist ja alumist usalduspiiri; vaikimisi kasutatava 95% asemele võib ise trükkida mõne teise arvu (näiteks 90 või 99);
- valikute *Kth Largest* ja *Kth Smallest* tulemusena väljastatakse järjekorranumbriga *K* väärtus vastavalt suuremate ja väiksemate väärtuste poolt lugedes; *K* = 1 korral on tulemuseks maksimaalne ja minimaalne väärtus, et aga need suurused sisalduvad ka valiku *Summary statistics* väljundis, on antud juhul mõistlik tellida näiteks suuruselt järgmised väärtused (siis *K* = 2).

• Tulemus:

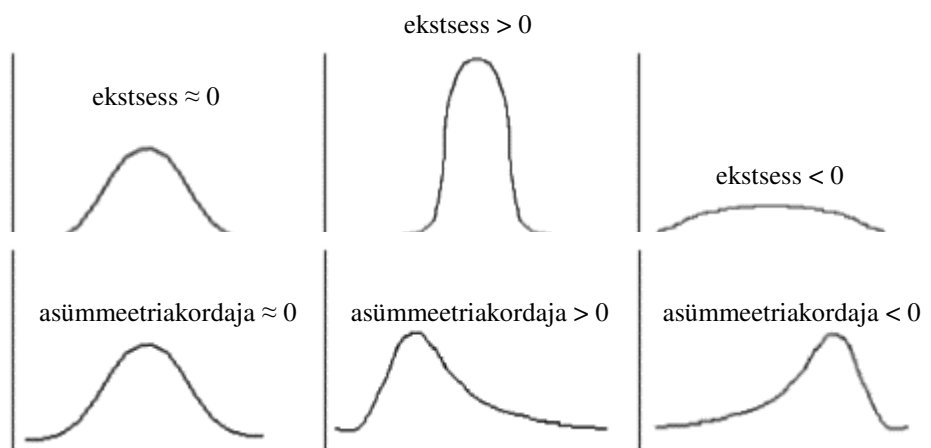
PIKKUS		MASS		PEA_P			
Mean	174,8769	Mean	69,46094	Mean	55,67561	} Valiku Summary statistics tulemus	Keskmine
Standard Error	1,260296	Standard Error	1,724923	Standard Error	0,497417		Standardviga
Median	175	Median	70	Median	56		Mediaan
Mode	177	Mode	70	Mode	56		Mood
Standard Deviation	10,16083	Standard Deviation	13,79939	Standard Deviation	3,185026		Standardhälve
Sample Variance	103,2424	Sample Variance	190,4231	Sample Variance	10,14439		Dispersioon
Kurtosis	-0,59575	Kurtosis	-0,63602	Kurtosis	1,876824		Ekstsess e järsakuskordaja
Skewness	-0,02697	Skewness	0,004988	Skewness	0,419495		Asümmeetriakordaja
Range	44	Range	58	Range	16		Ulatus = Max - Min
Minimum	153	Minimum	42	Minimum	50		
Maximum	197	Maximum	100	Maximum	66		
Sum	11367	Sum	4445,5	Sum	2282,7		
Count	65	Count	64	Count	41		Vaatluste arv
Largest(2)	195	Largest(2)	95	Largest(2)	61		
Smallest(2)	157	Smallest(2)	44	Smallest(2)	50		
Confidence Level(95,0%)	2,51773	Confidence Level(95,0%)	3,446984	Confidence Level(95,0%)	1,005318		

• Lisalugemine – uuritava tunnuse jaotuse kuju iseloomustamine

Enamustest protseduuri *Descriptive Statistics* väljundis sisalduvatest arvarakteristikutest on ennegi juttu olnud. Siiski on siin ka kaks uut suurust, mida kasutatakse peamiselt uuritava tunnuse jaotuse kuju iseloomustamiseks – need suurused on **ekstsess e järsakuskordaja** (ingl *kurtosis*) ja **asümmeetriakordaja** (ingl *skewness*). Sellest, mida need karakteristikud mõeldavad, annavad parema ettekujutuse järgnevad joonised:

Jaotuse märkimisväärsest erinevusest normaaljaotusest on mõtet rääkida siis, kui ükskõik kumb neist kordajatest omandab absoluutväärtuselt 1-st suurema väärtuse ...

Eriti palju neid kordajaid siiski ei kasutata.



- **Jaotuse sümmeetrilisuse üle otsustamisel kasutatakse sageli (asümmeetriakordaja asemel) keskmise ja mediaani võrdlust.**

Nimelt, kuna aritmeetiline keskmine on tundlik erandlike väärtuste suhtes, siis vihjab

$\bar{x} > med$ sellele, et jaotuse kuju on parempoolse ebasümmeetriaga (leiduvad üksikud teistest palju suuremad väärtused, ja seega asümmeetriakordaja > 0),

$\bar{x} < med$ aga sellele, et jaotuse kuju on vasakpoolse ebasümmeetriaga (leiduvad üksikud teistest palju väiksemad väärtused, ja seega asümmeetriakordaja < 0).

- **Vaata, kas kirjeldatud seos keskmise ja mediaani erinevuse ning asümmeetriakordaja väärtuse vahel peab paika ka teie kursuse tudengite pikkuste, masside ja peaümberrõõmõõdude korral.**

2. Leidke 90%, 95% või 99% usalduspiirid keskmistele väärtustele. Mida need näitavad?

Kuna *Excel* ise usalduspiire välja ei arvuta, tuleb need enesest leida.

Selleks võib protseduuri *Descriptive Statistics* väljundtabelit täiendada kahe reaga, kuhu tuleks selguse huvides ka kirja panna, mida uued arvutatavad suurused enesest kujutavad.

	L	M	N
1		PIKKUS	
2			
3		Mean	174,8769
4		Standard Error	1,260296
5		Median	175
6		Mode	177
7		Standard Deviation	10,16083
8		Sample Variance	103,2424
9		Kurtosis	-0,59575
10		Skewness	-0,02697
11		Range	44
12		Minimum	153
13		Maximum	197
14		Sum	11367
15		Count	65
16		Largest(2)	195
17		Smallest(2)	157
18		Confidence Level(95,0%)	2,51773
19			
20		Alumine 95% usalduspiir	=N3-N18
21		Ülemine 95% usalduspiir	=N3+N18

Usalduspiirid keskmisele leitakse valemist

$$\bar{x} \pm t_{1-\alpha/2, n-1} \frac{s}{\sqrt{n}}$$

Excel väljastab toodud valemi mõlemad liidetavad, mille alusel on lihtne mõlemad usalduspiirid välja arvutada.

Alumine 95% usalduspiir	172,3592
Ülemine 95% usalduspiir	177,3947

Seega, tõlgendades antud andmestikku kui valimit kogu teie kursusest, võib väita, et loomakasvatussaaduste tootmise esimese kursuse tudengite keskmine pikkus jääb 95% tõenäosusega vahemikku 172,4 cm kuni 177,4 cm. St, et mõttes ära **kõigi** esmakursuslaste pikkused ja arvutades keskmise, peaks saadud tegelik keskmine 95% tõenäosusega jääma leitud piiridesse.

- Kui keegi leidis 95% usalduspiiride asemel 90% või 99% usalduspiirid, siis need peaksid tulema vastavalt (172,8; 177,0) ja (171,5; 178,2). **Miks on 90% usaldusintervall kitsam?**
- **Arvutage usaldusintervall ka keskmisele massile ja keskmisele peaümberrõõmõõdule ning püüdke neist vähemalt ühe kohta sõnastada lõppjärgeldus!!**