

# Biomeetria

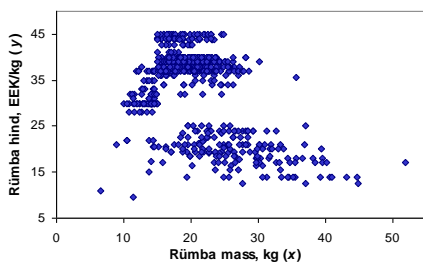
## Kahe arvtunnuse ühine käitumine, regressioonanalüüs

### Lineaarne regressioonanalüüs



#### Millal kasutada ja mida näitab?

Kasutatakse prognoosimaks ühe arvtunnuse väärtusi teis(t)e järgi.



Regressioonivõrrand:  $\hat{Y} = a + bX$

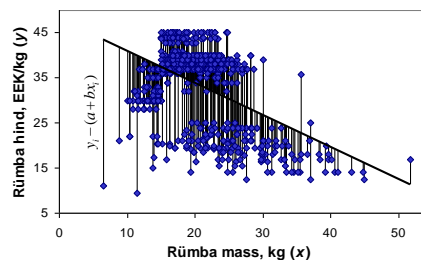
Tunnust  $Y$  nimetatakse funktsioon- ja tunnust  $X$  argumenttunnuseks.

Näide.

Rümbe 1 kg hind =  $a + b \cdot$  Rümbe mass

Regressioonivõrrandi parameetrid  $a$  ja  $b$  hinnatakse vähimruutude meetodil, st et minimeeritakse prognoosi jäägid:

$$\sum_{i=1}^n [y_i - (a + bx_i)]^2 \Rightarrow \min$$





## Lineaarne regressioonanalüüs

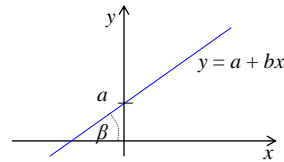
Regressioonivõrrandi (ehk regressioonimudeli) parameetrite  $a$  ja  $b$  vähimruutude hinnangud:

$$b = \frac{\sum_{i=1}^n (x_i y_i - \bar{x}\bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad a = \bar{y} - b\bar{x}$$

Regressioonivõrrandi kordajate geomeetriline tähendus:

vabaliige  $a$  märgib kohta, kus regressioonisirge lõikab  $y$ -telge, ning regressioonikordaja  $b$  iseloomustab regressioonisirge ja  $x$ -telje vahelist nurka (matemaatilisemalt väljendudes  $\tan(\beta) = b$ , kus  $\beta$  on sirge tõusunurk).

Sisulise tähenduse kohaselt näitab regressioonikordaja  $b$ , kui mitme ühiku võrra muutub funktsioontunnuse väärtus, kui argumenttunnus muutub 1 ühiku võrra.



### Näide.

$$\text{Rümba 1 kg hind} = a + b \cdot \text{Rümba mass} = 48,18 + (-0,711) \cdot \text{Rümba mass}$$

Seega kaasneb lamba lihakeha massi suurenemisega 1 kg võrra 0,71-kroonine hinnalangus 1 kg liha eest.




## Lineaarne regressioonanalüüs – jäägid

Näide.

$$\begin{aligned} \text{Rümba 1 kg hind} &= a + b \cdot \text{Rümba mass} \\ &= 48,18 + (-0,711) \cdot \text{Rümba mass} \end{aligned}$$

Vaatluse jrk nr ( $i$ )	Tegelik rümba 1 kg hind ( $y_i$ )	Prognoositud rümba 1 kg hind ( $a + b \cdot x_i$ )	Prognoosi jäägid $y_i - (a + b \cdot x_i)$
1	39	31,250	7,750
2	39	37,367	1,633
3	40	37,083	2,917
4	39	30,824	8,176
5	39	33,171	5,829
6	37	34,593	2,407
7	39	32,175	6,825
8	39	32,744	6,256
9	37	31,250	5,750
10	34	31,677	2,323
11	37	35,304	1,696
12	33	37,936	-4,936
13	33	39,216	-6,216
14	37	33,597	3,403
15	39	32,175	6,825



## Regressioonimudeli sobivus


Determinatsioonikordaja  $R^2$  ütleb, kui suure osa uuritava tunnuse varieeruvusest mudel ära kirjeldab,  $0 \leq R^2 \leq 1$ . Mida suurem, seda parem!

Leitakse kui mudelile vastava hajuvuskomponendi  $SS_1 = \sum_{i=1}^n [(a + bx_i) - \bar{y}]^2$  ja uuritava tunnuse koguhajuvust kirjeldava hälvete ruutude summa  $SS = \sum_{i=1}^n (y_i - \bar{y})^2$  jagatis:  $R^2 = SS_1/SS$ .

Mudeli standardviga  $SE$  on mudeli prognoosijäägi standardhälbe hinnang, mis näitab keskmist prognoosiviga. Mida väiksem, seda parem!

Hüpoteeside kontroll

- Hüpotees mudeli, kui terviku kohta (võrreldakse konstrueeritud mudeli ja nn konstantse mudeli  $y=a$  jääkide varieeruvust):
  - $H_0$  : mudel ei ole parem võrreldes konstantse mudeliga,
  - $H_1$  : mudel on parem võrreldes konstantse mudeliga.
- Hüpoteesid mudeli parameetrite kohta – kontrollitakse väidet iga parameetri nullist erinemise kohta:
  - $H_0 : a = 0$        $H_0 : b = 0$
  - $H_1 : a \neq 0$       $H_1 : b \neq 0$



## Regressioonimudeli sobivus

**Näide.**  
Rümba 1 kg hind =  $a + b \cdot \text{Rümba mass} = 48,18 + (-0,711) \cdot \text{Rümba mass}$

MS Exceli protseduuri *Regression* väljund:

SUMMARY OUTPUT

Mitmene korrelatsioonikordaja – mõõdab uuritava tunnuse ja tema prognoositud väärtuste vahelist korrelatsiooni. Mida suurem, seda parem!

Regression Statistics	
Multiple R	0,4730
R Square	0,2237
Adjusted R Square	0,2226
Standard Error	7,8450
Observations	686

Determinatsioonikordaja  $R^2$  ja selle väikeste valimite tarvis kohandatud [adjusted] väärtus

Mudeli standardviga


$H_0$  : mudel ei ole parem võrreldes konstantse mudeliga  
 $H_1$  : mudel on parem võrreldes konstantse mudeliga

ANOVA					
	df	SS	MS	F	Significance F
Regression	1	12130,21573	12130,2157	197,0999	1,57702E-39
Residual	684	42095,74764	61,5435		
Total	685	54225,96337			

	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%
Intercept	48,1782	1,0843	44,4341	8,2819E-204	46,0493	50,3070
R_mass	-0,7113	0,0507	-14,0392	1,5770E-39	-0,8107	-0,6118

Mudeli parameetrite hinnangud

Hüpoteeside kontrolli mudeli iga parameetri kohta:  
 $H_0 : a = 0$   
 $H_1 : a \neq 0$   
 $H_0 : b = 0$   
 $H_1 : b \neq 0$



## Regressioonanalüüsi eeldused

Regressioonimudeli parameetrite hindamine ei eelda tunnuste jaotumist vastavalt normaaljaotuse seaduspäradele!

Mudeli täpsuse ja statistilise olulisuse hindamiseks peavad:

- mudeli (prognoosi)jäägid olema ligikaudu normaaljaotusega (kontrollimiseks histogramm, tõenäosuspaber);
- mudeli jäägid olema ühtlase varieeruvusega (hajuvusdiagramm).

Ükskõik kumma eelduse rikutuse korral ei pruugi mudeli kohta käivate hüpoteeside kontrollimisel arvatavate teststatistikute jaotusseadused kehtida, mistõttu ei pruugi õiged olla ka otsustused mudeli sobivuse ja rakendatavuse üle.

Eelkõige teise eelduse paikapidamatus võib vihjata mittesobivale mudelile (vale matemaatiline funktsioon, mõni arvestamata jäänud argument vmt).

