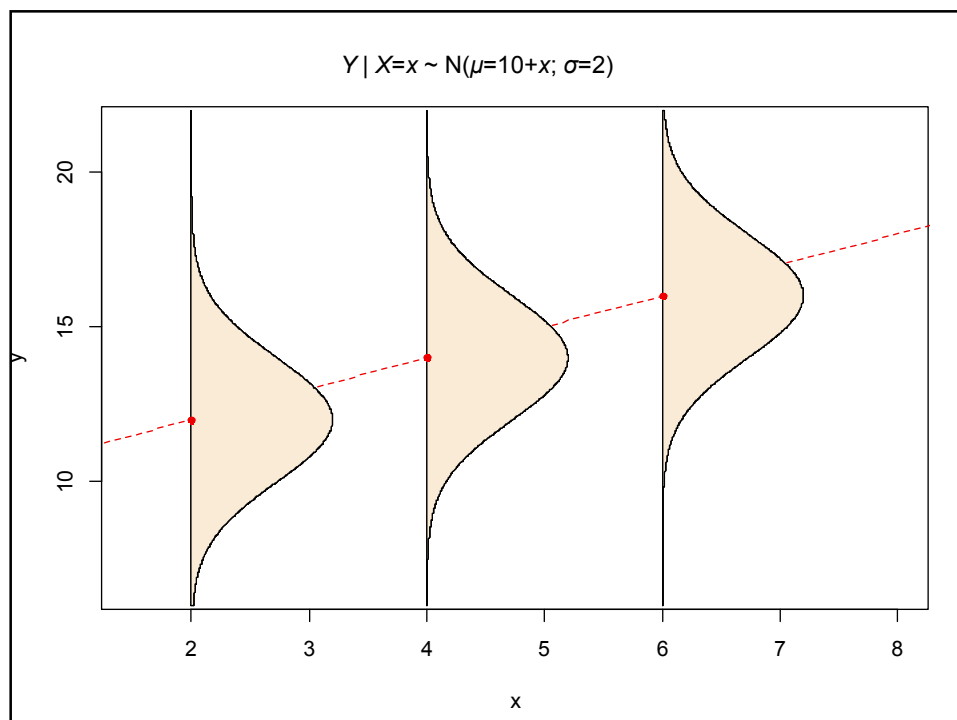


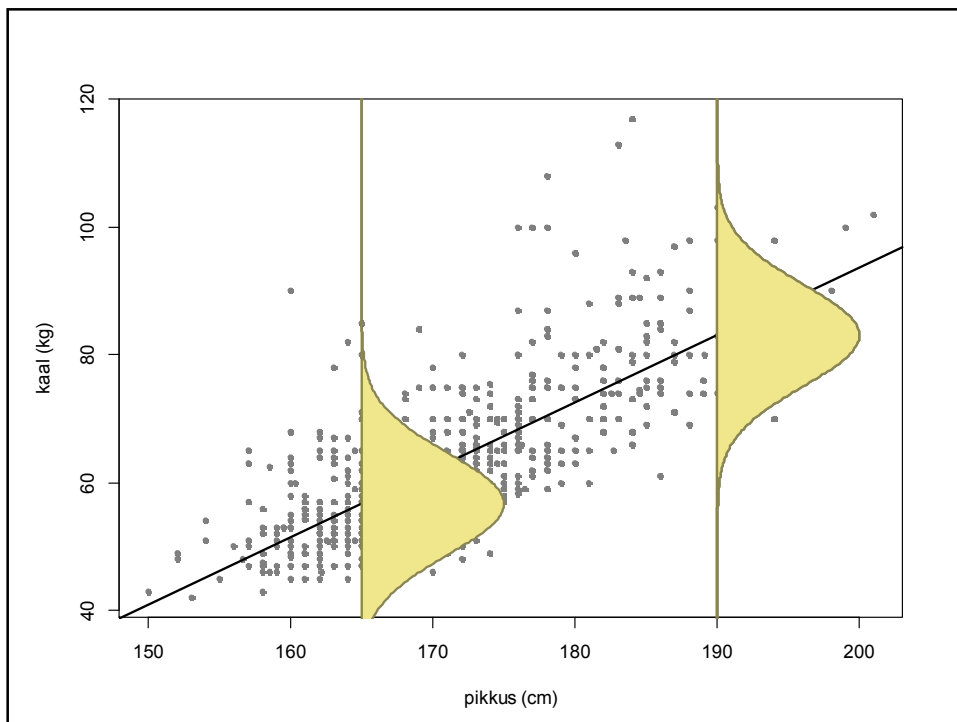
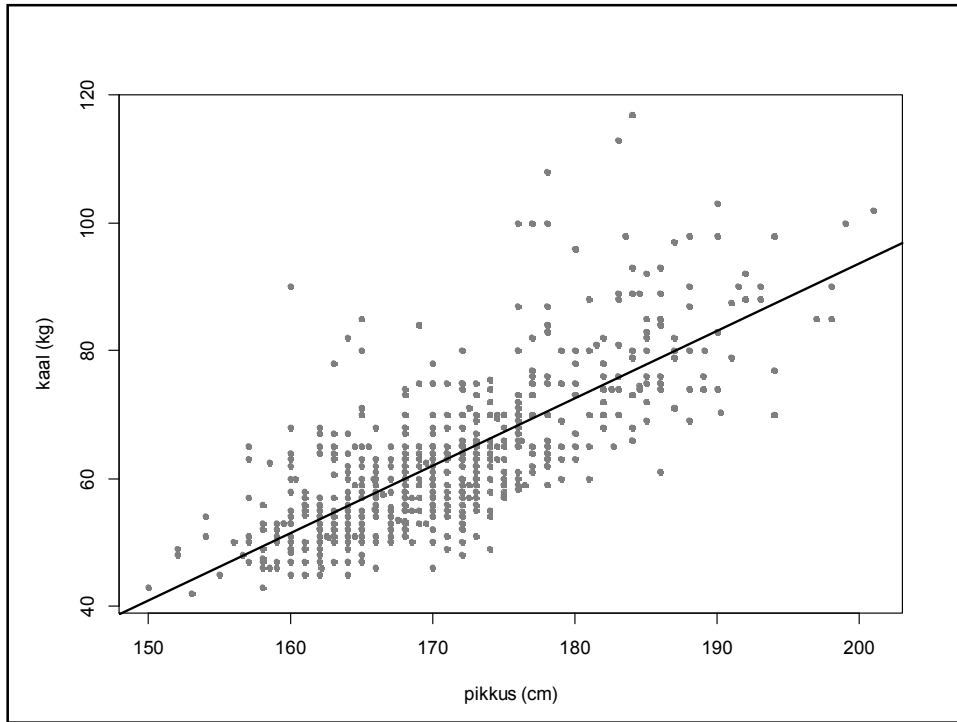
Biomeetria  
2. loeng

## Lihtne lineaarne regressioon

mudeli hindamisest; usaldusintervall;  
prognosiintervall; determinatsioonikordaja; ...

Märt Möls  
martm@ut.ee





## Regressioonanalüüsi mudel(id)

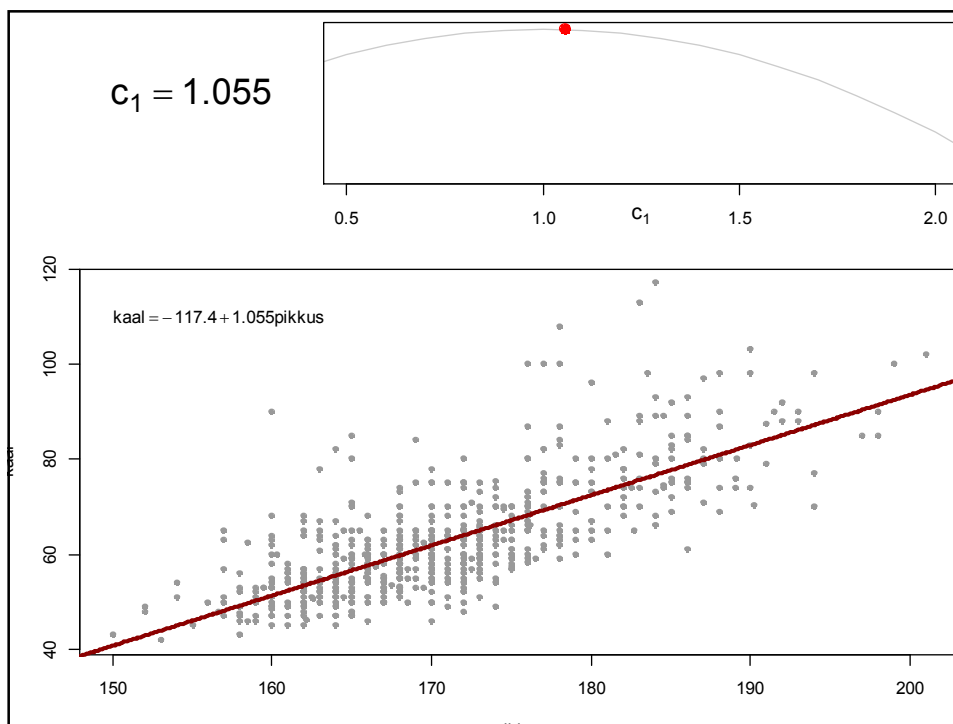
$$Y | X=x \sim N(\mu=c_0+c_1x; \sigma=\sigma_\varepsilon)$$

$$Y \sim N(\mu=c_0+c_1x; \sigma=\sigma_\varepsilon)$$

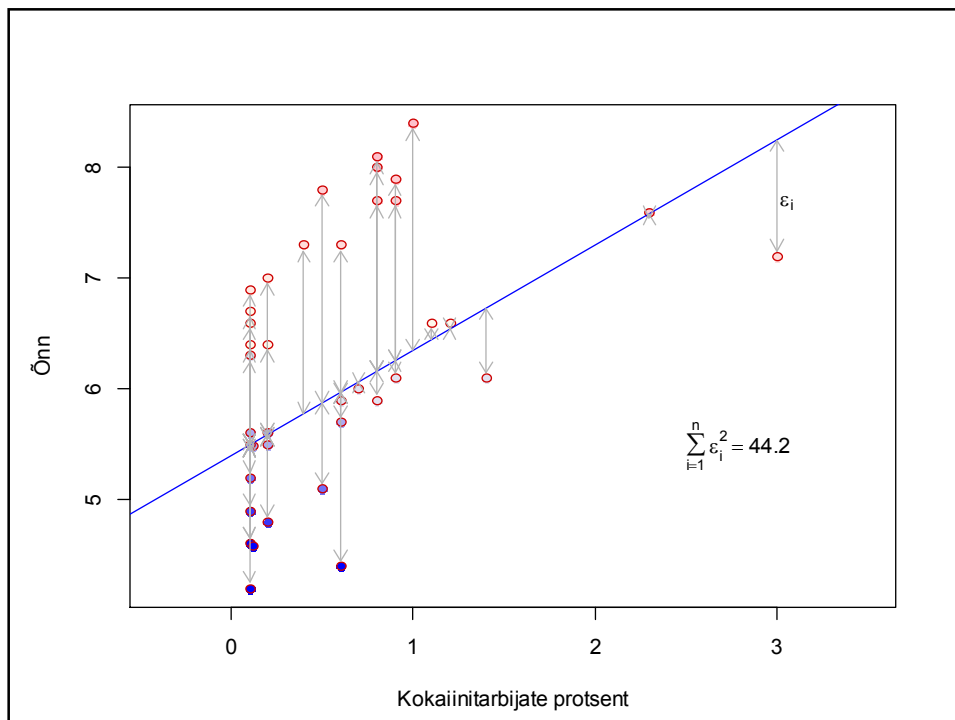
$$Y = c_0 + c_1x + \varepsilon; \quad \varepsilon \sim N(0; \sigma_\varepsilon)$$

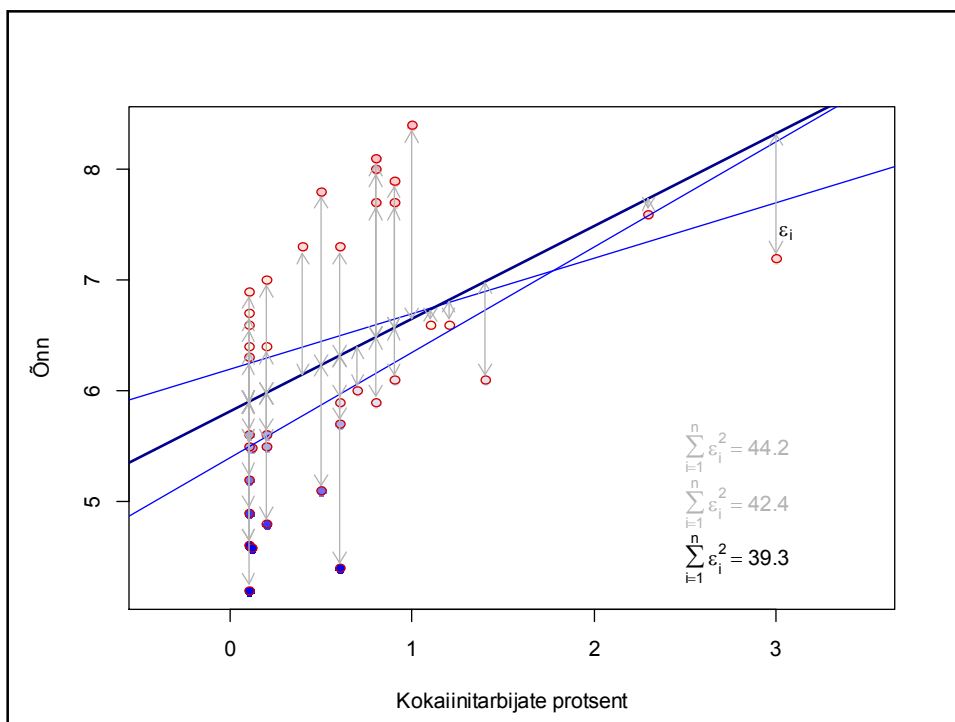
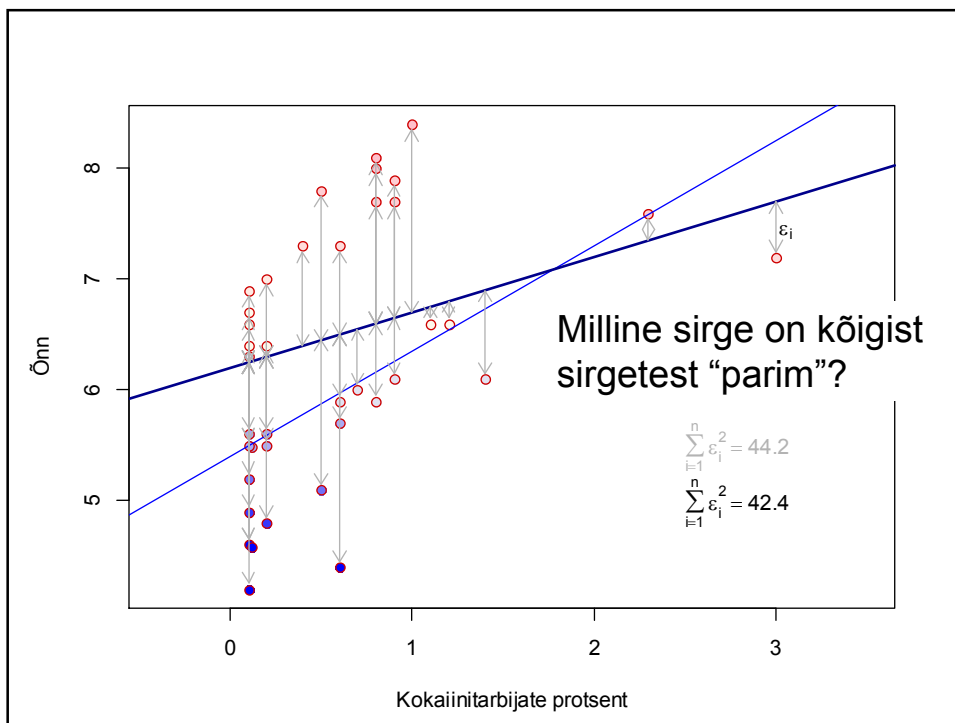
$$EY = c_0 + c_1x$$

$c_0$ ,  $c_1$  väärtused hinnatakse suurima tõepära meetodi abil



## Alternatiivne interpretatsioon





## Sama sirge...

- $Y \sim N(\mu=c_0+c_1x; \sigma=\sigma_\varepsilon)$ , sõltumatud vaatlused.  
otsime  $c_0$  ja  $c_1$  väärtuseid mille puhul nähtud andmete saamise tõenäosus oleks kõige suurem;
- Otsime sirget  $y=c_0+c_1x$ , mille puhul prognoosivigade ruutude summa oleks minimaalne;
- $EY=c_0+c_1x$ ,  $DY=\sigma_\varepsilon$ , vaatlused sõltumatud.  
Otsime kõige täpsemat (lineaarset) nihketa hinnangut  $c_0+c_1x$  -le

## Kas (lineaarne) seos $Y$ ja $X$ vahel eksisteerib?

Kas

$$Y \sim N(\mu=c_0+c_1x; \sigma=\sigma_\varepsilon)$$

või

$$Y \sim N(\mu=c_0; \sigma=\sigma_\varepsilon) ?$$

Mida prognoosime      Mille abil prognoosime

```

> mudel=lm(kaal~pikkus)
> summary(mudel)

[...]
Coefficients:
      Estimate Std. Error t value Pr(>|t|)
(Intercept) -117.36675    6.33019  -18.54 <2e-16 ***
pikkus       1.05474     0.03695   28.54 <2e-16 ***
[...]

```

Kas  $c_0$  võib olla 0?      (arrow pointing to Pr(>|t|) for Intercept)

Kas  $c_1$  võib olla 0?      (arrow pointing to Pr(>|t|) for pikkus)

Mida prognoosime      Mille abil prognoosime

```

> mudel=lm(kaal~pikkus)
> summary(mudel)

[...]
Coefficients:
      Estimate Std. Error t value Pr(>|t|)
(Intercept) -117.36675    6.33019  -18.54 <2e-16 ***
pikkus       1.05474     0.03695   28.54 <2e-16 ***
[...]

```

Kas  $c_0$  võib olla 0?      (arrow pointing to Pr(>|t|) for Intercept)

Kas  $c_1$  võib olla 0?      (arrow pointing to Pr(>|t|) for pikkus)

```

> confint(mudel)

              2.5 %      97.5 %
(Intercept) -129.7967014 -104.936795
pikkus       0.9821798    1.127295

```

95% usaldusintervallid      (arrows pointing to the 2.5% and 97.5% columns)

```

> predict(mudel, newdata=data.frame(pikkus=170))
      1
61.93859      170cm pikkuste inimeste kaalude keskvärtus (hinnang)

> predict(mudel, newdata=data.frame(pikkus=170),
  interval="confidence")
      fit      lwr      upr      Kui täpselt me teame 170cm pikkuste
1 61.93859 61.3231 62.55408 inimeste kaalude keskvärtust?

> predict(mudel, newdata=data.frame(pikkus=170),
  interval="prediction")
      fit      lwr      upr      Kui täpselt me teame 170cm
1 61.93859 46.29655 77.58063 pikkuse inimese (Jaani) kaalu?

> install.packages("gmodels")
> library(gmodels)
> estimable(mudel, c(1,170))
      Estimate Std. Error t value  DF Pr(>|t|)
(1 170) 61.93859  0.3134512 197.602 654      0

> estimable(mudel, c(1,170), conf.int=0.95)
      Estimate Std. Error t value  DF Pr(>|t|) Lower.CI Upper.CI
(1 170) 61.93859  0.3134512 197.602 654      0 61.3231 62.55408

```

## Estimable -käsust

Võimalik kasutada keerukamate mudelite juures;

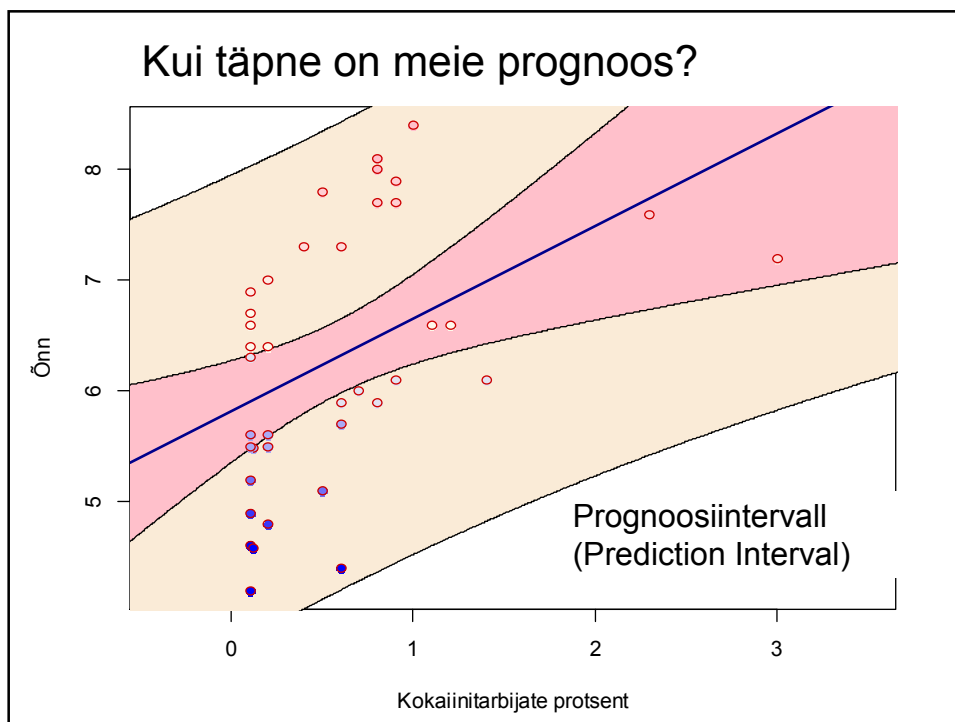
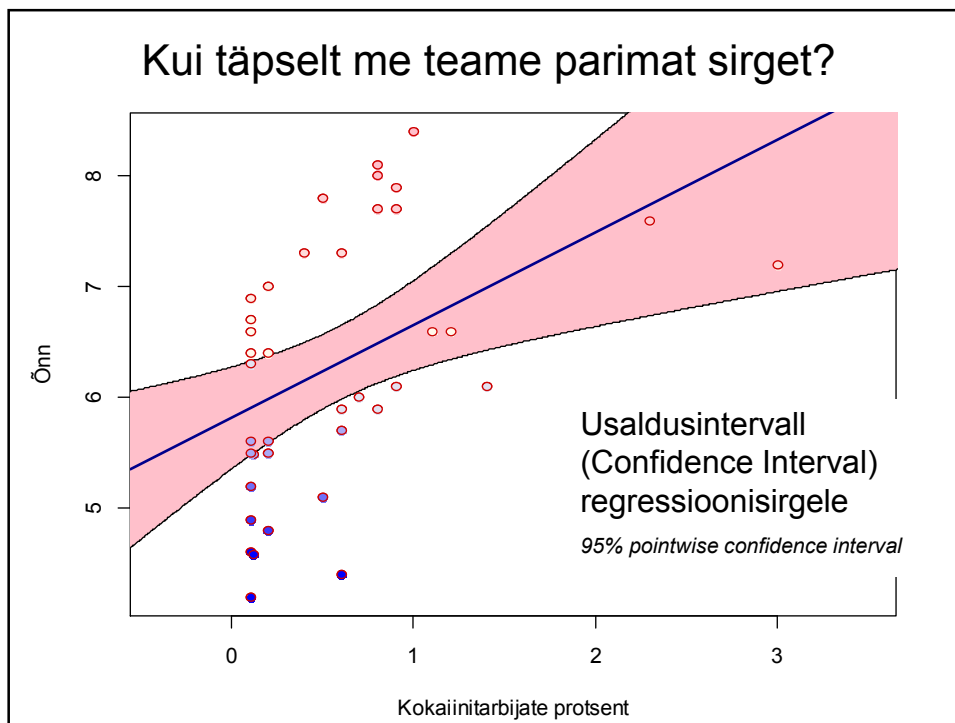
Võimalik kasutada keerukamate küsimuste jaoks (milline on 170cm ja 160cm pikkuste tudengite keskmiste kaalude erinevus?)

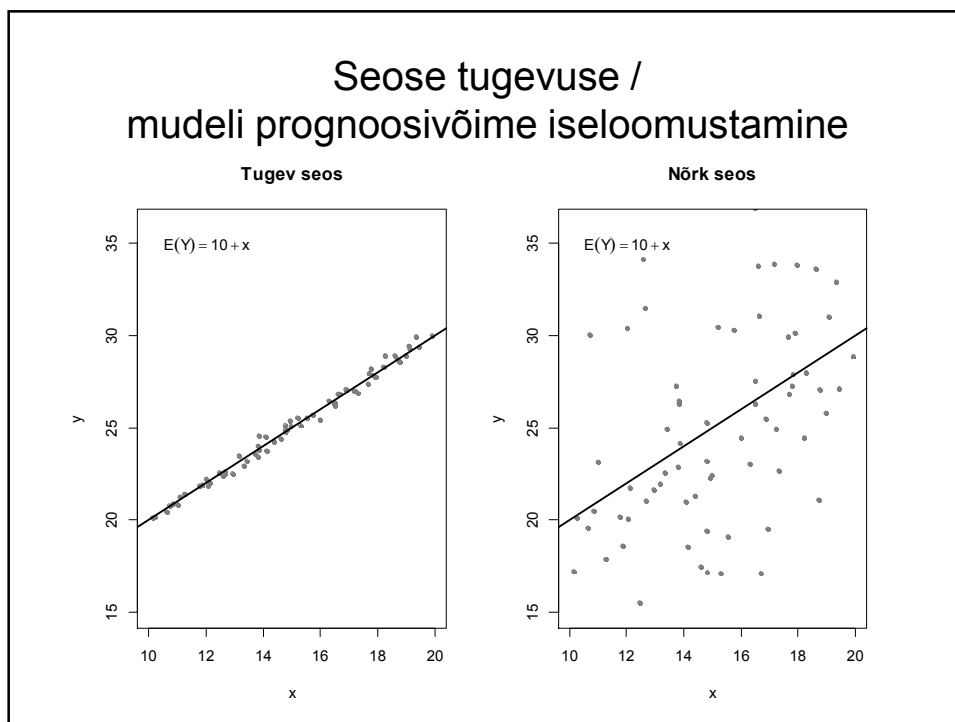
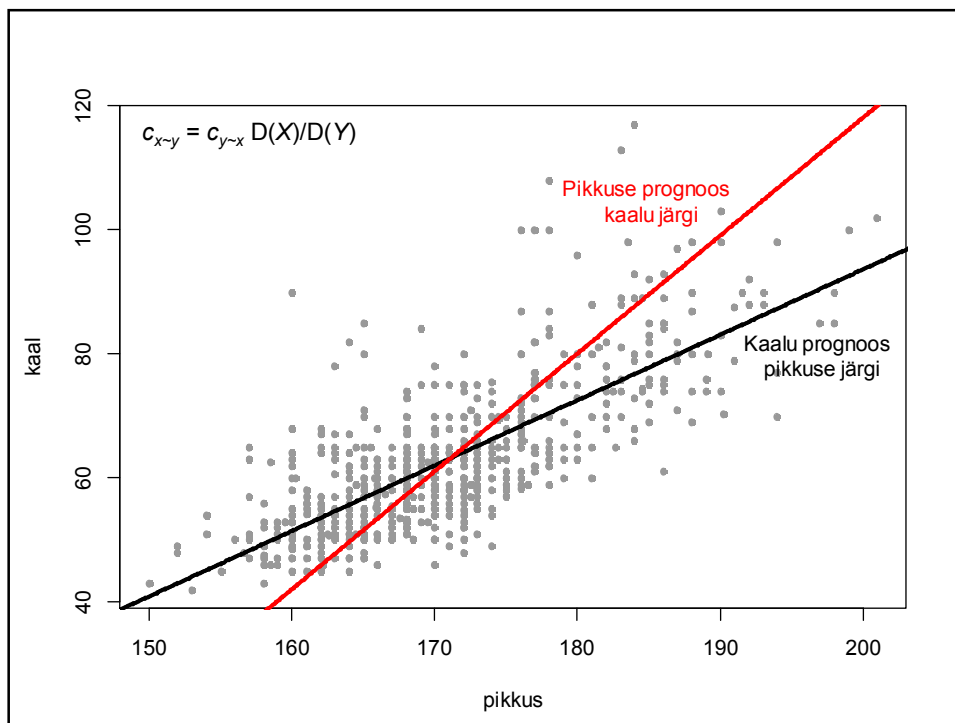
```

> estimable(mudel, c(1-1,170-160), conf.int=0.95)
      Estimate Std. Error t value  DF Pr(>|t|) Lower.CI Upper.CI
(0 10) 10.54737  0.369513 28.54398 654      0 9.821798 11.27295

```







```

> summary(mudel)

```

Tegelik – Prognoos  
(tegelik kaal – pikkuse järgi prognoositud kaal)

```

Residuals:
  Min       1Q   Median       3Q      Max
-17.814  -5.361  -1.220   3.499  40.295

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept) -117.36675    6.33019  -18.54  <2e-16 ***
pikkus       1.05474    0.03695   28.54  <2e-16 ***

Residual standard error: 7.96 on 654 degrees of freedom
(1 observation deleted due to missingness)
Multiple R-squared:  0.5547,    Adjusted R-squared:  0.554
F-statistic: 814.8 on 1 and 654 DF,  p-value: < 2.2e-16

```

```

Residual standard error: 7.96

> predict(mudel, data.frame(pikkus=170))
  1
61.93859

> predict(mudel, data.frame(pikkus=170))-2*7.96
  1
46.01859

> predict(mudel, data.frame(pikkus=170))+2*7.96
  1
77.85859

> predict(mudel, data.frame(pikkus=170),
          interval="prediction")
      fit      lwr      upr
1 61.93859 46.29655 77.58063

```

## Determinatsioonikordaja

Mudel 1:  $Y = c_0 + \varepsilon_1$

Mudel 2:  $Y = c_0 + c_1 X + \varepsilon_2$

$D\varepsilon_1$  - iseloomustab lihtsama mudeli prognooside täpsust

$D\varepsilon_2$  - iseloomustab keerukama mudeli prognoositäpsust

$D\varepsilon_1 - D\varepsilon_2$  - täpsuse suurenemine tänu x-tunnuse teadmisele

$(D\varepsilon_1 - D\varepsilon_2)/D\varepsilon_1$  - suhteline võit prognoositäpsuses

$$R^2_{\text{tegelik}} = (D\varepsilon_1 - D\varepsilon_2)/D\varepsilon_1$$

$$Y = c_0 + \varepsilon_1$$

$$D(Y) = D(c_0 + \varepsilon_1) = D(\varepsilon_1)$$

$$R^2_{\text{tegelik}} = (DY - D\varepsilon_2)/DY$$

$R^2$  ja  $R^2_{\text{adjusted}}$  (kohandatud  $R^2$ )

$$R^2_{\text{tegelik}} = (DY - D\varepsilon_2) / DY$$

Vaja hinnata.  $DY$  hindamisega probleeme pole,  $D\varepsilon_2$  hindamisega küll.

Kasutades nihkega ("vale") hinnangut  $D\varepsilon_2$ -le, saame determinatsioonikordaja  $R^2$ .

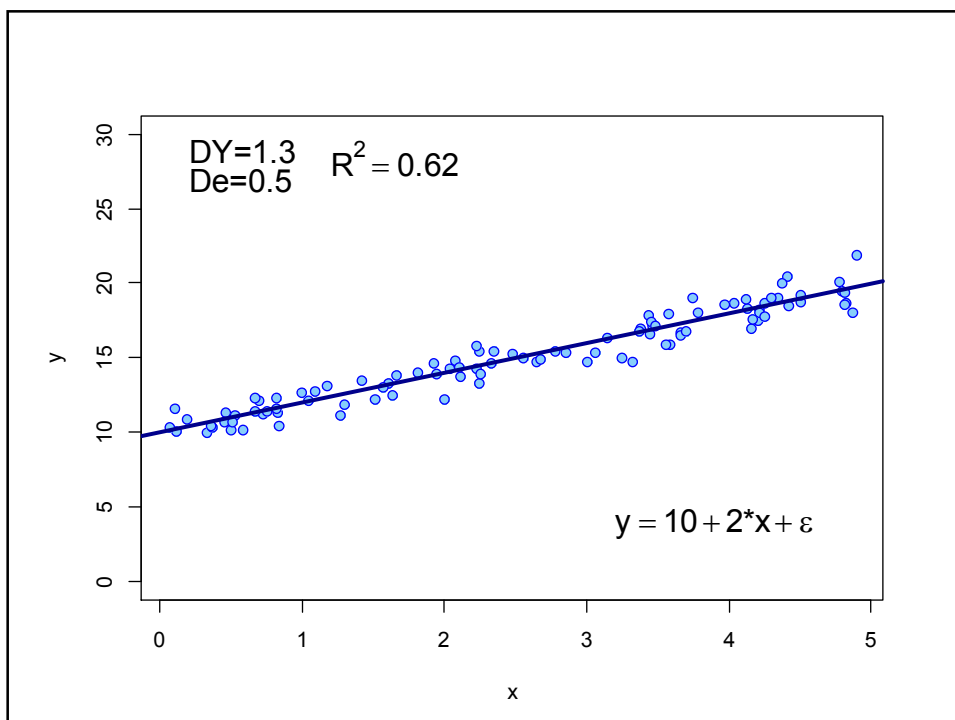
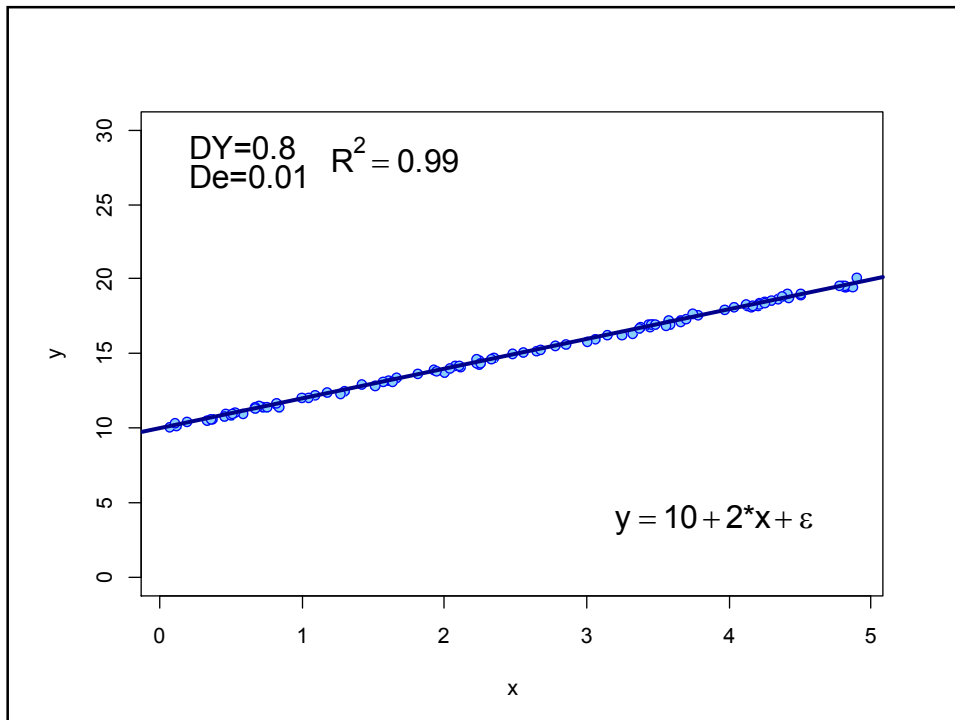
Kasutades nihketa ("õiget") hinnangut  $D\varepsilon_2$ -le, saame kohandatud determinatsioonikordaja  $R^2_{\text{adjusted}}$ .

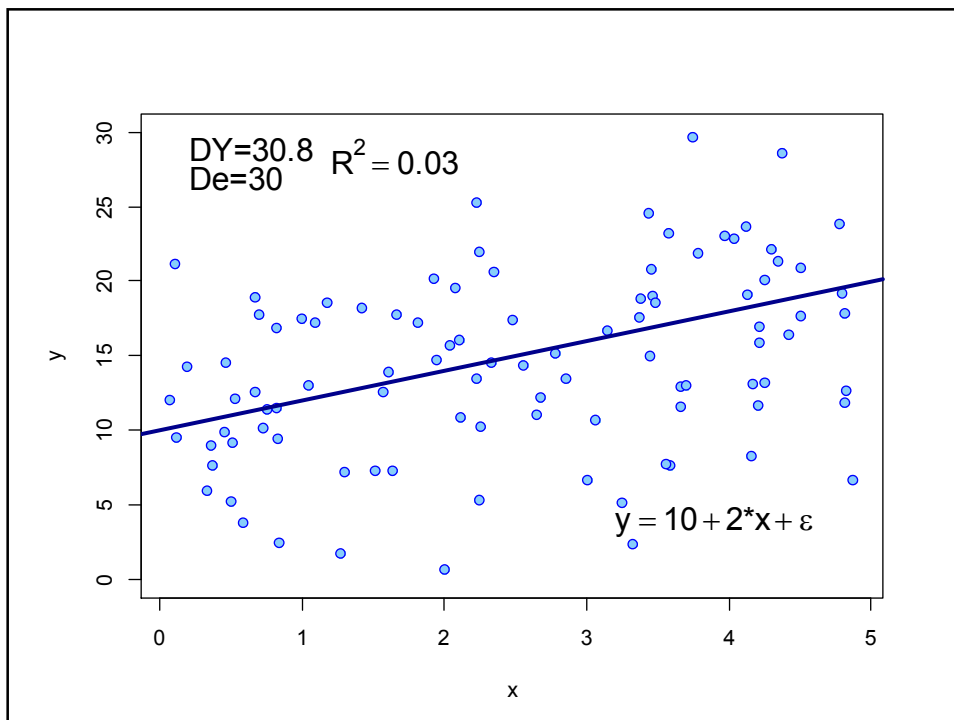
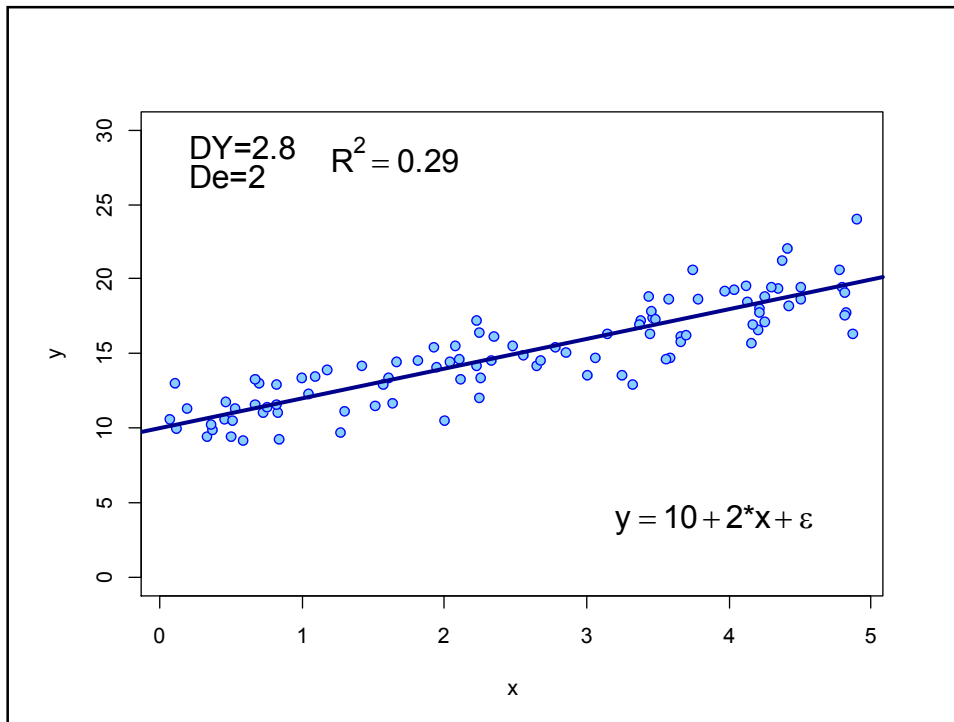
```
> summary(lm(kaal~pikkus))
[...]
```

Residual standard error: **7.96** on 654 degrees of freedom  
(4 observations deleted due to missingness)

Multiple R-squared: 0.5547, Adjusted R-squared: **0.554**

```
> Dkaal=var(kaal, na.rm=TRUE)
> (Dkaal-7.96**2)/Dkaal
[1] 0.5540274
```

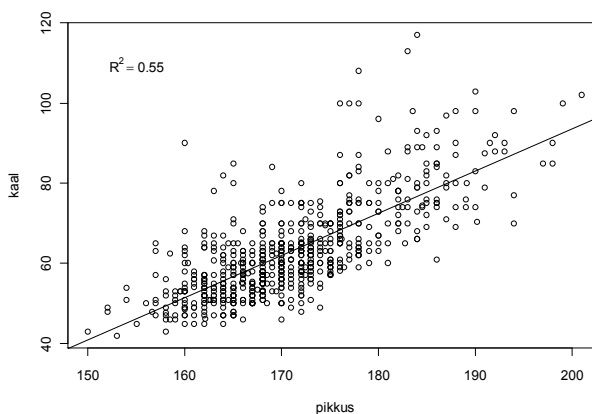




## Determinatsioonikordajaga manipuleerimine

Determinatsioonikordaja on eelkõige interpreteeritav siis, kui uuritav valim on tõepoolest juhuslik valim mingist populatsioonist. Kui uuritavad andmed on kogutud eksperimenteerides (olukorras, kus me ise otsustame, millised saavad olema  $X$ -tunnuse väärtused) on determinatsioonikordaja  $R^2$  teatavates piirides eksperimentaatori enda valida/otsustada.

### Kaalu prognoosimine pikkuse järgi. Juhuslik valim: $R^2=0,55$





## Kuidas suurendada või vähendada determinatsioonikordajat?

$$KAAL = c_0 + c_1 PIKKUS + e.$$

Determinatsioonikordaja

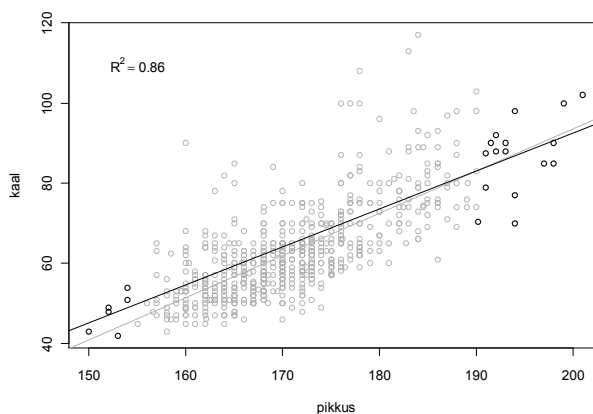
$$R^2 = 1 - D(e)/D(KAAL).$$

Vaja oleks kas muuta  $D(e)$  või  $D(KAAL)$  väärtust. Antud näites on lihtsam muuta  $D(KAAL)$  väärtust:

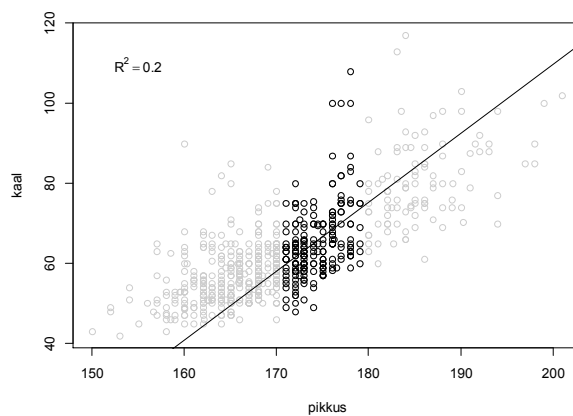
$$D(KAAL) = c_1^2 D(PIKKUS) + D(e)$$

Valides uuringusse väga erinevate pikkustega tudengeid saame suure  $R^2$  väärtuse, valides sarnase pikkusega tudengeid saame väikese  $R^2$  väärtuse.

## Pikkus < 155cm või pikkus > 190cm $R^2 = 0,86$



Pikkus > 170cm ja pikkus < 180cm  
 $R^2=0,2$



```
> mudelMees=lm(kaal~pikkus, data=kokku2[sugu==2,])
> summary(mudelMees)
Multiple R-squared:  0.3704,    Adjusted R-squared:  0.3661

> mudelNaine=lm(kaal~pikkus, data=kokku2[sugu==1,])
> summary(mudelNaine)
Multiple R-squared:  0.2968,    Adjusted R-squared:  0.2954

> mudel=lm(kaal~pikkus, data=kokku2)
> summary(mudel)
Multiple R-squared:  0.5547,    Adjusted R-squared:  0.554
```