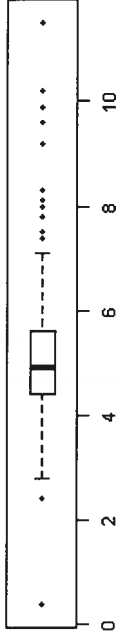


```
attach(ex02.26beer)
```

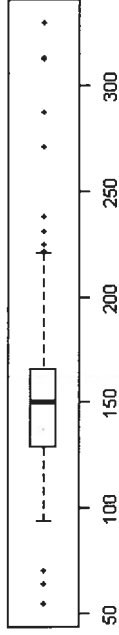
```
par(mfrow=c(3,1))  
boxplot(PercentAlcohol, horizontal=TRUE, main="PercentAlcohol")  
boxplot(Calories, horizontal=TRUE, main="Calories")  
boxplot(Carbohydrates, horizontal=TRUE, main="Carbohydrates")
```

PercentAlcohol



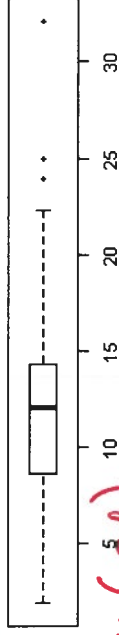
$$\bar{X} = 5.2, \text{ sd} = 1.4$$

Calories



$$\bar{X} = 154, \text{ sd} = 44$$

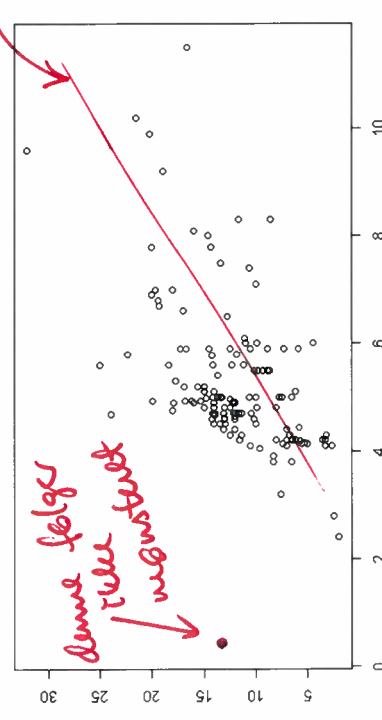
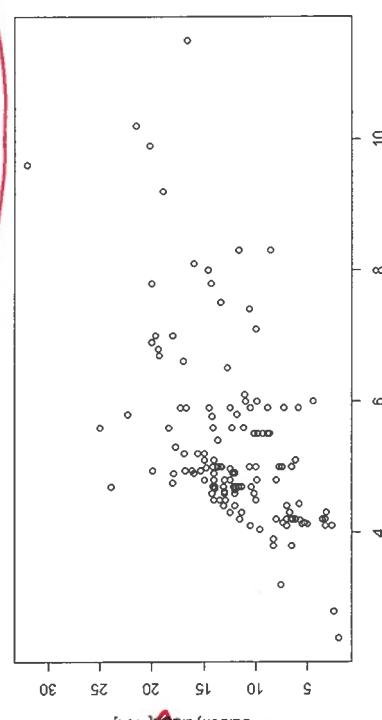
Carbohydrates



$$\bar{X} = 12.0, \text{ sd} = 4.9$$

Ingen av disse er helt symmetriske, men ingen er heller helt skjeve. Ser at det er outliers i alle tre variabler, og regner derfor med at sd er noe større enn det vi vil forventet ved n-ford. Velger likevel å oppsummere cha mean(sd)

```
> c(summary(PercentAlcohol), sd(PercentAlcohol))  
  Min.  1st Qu.  Median      Mean      3rd Qu.      Max.      sd  
 0.40000  4.40000  4.90000  5.229000  5.600000  11.500000  1.428737  
> c(summary(Calories), sd(Calories))  
  Min.  1st Qu.  Median      Mean      3rd Qu.      Max.      sd  
 55.00000 129.00000 150.00000 154.10000 166.00000 330.00000 44.49254  
> c(summary(Carbohydrates), sd(Carbohydrates))  
  Min.  1st Qu.  Median      Mean      3rd Qu.      Max.      sd  
 1.900000  8.600000  12.000000  11.960000  14.300000  32.100000  4.905587
```

<pre> which.min(PercentAlcohol) &gt; which.min(PercentAlcohol) [1] 104 ← tall nr 104 i vektoren &gt; PercentAlcohol[104] [1] 0.4 &gt; PercentAlcohol[which.min(PercentAlcohol)] [1] 0.4 </pre>	<pre> which.min(PercentAlcohol) which.min(PercentAlcohol) [1] 104 ← tall nr 104 i vektoren PercentAlcohol[104] [1] 0.4 PercentAlcohol[which.min(PercentAlcohol)] [1] 0.4 </pre>
<pre> plot(PercentAlcohol, Carbohydrates) </pre>  <p>denne følger ikke sammenheng håndtegnet</p>	<pre> plot(PercentAlcohol[-104], Carbohydrates[-104]) </pre>  <p>samme plott, men uten outliers.</p>

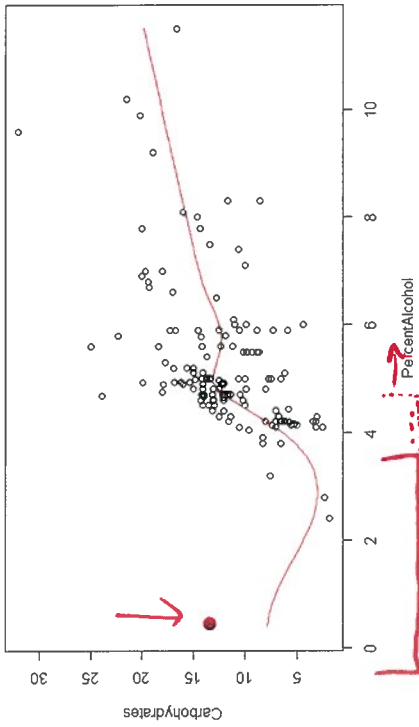
generelt: høye verdier av % alkohol → høye karbohydrat-tall, og } positiv sammenheng  
 lave ————— i. ————— lave

Foresetter en positiv korrelasjon.

Ingen klare ikke-lineære smh, som f.eks  etc.

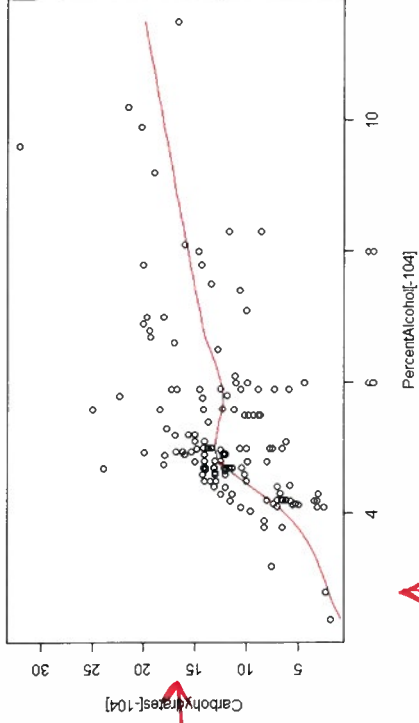
↳ en lineær smh fanger opp det generelle i sammenhengene.

```
lines(loess.smooth(PercentAlcohol, Carbohydrates), col="red")
```



uten  
outlier

```
lines(loess.smooth(PercentAlcohol[-104], Carbohydrates[-104]))
```



A

lokalt tilpasset regresjonslinje: tilpasse linja til et gittende område, for å se hva data forteller om sammenhengen. Viktig å ikke legge stor vekt på endene. Der er det lite data, og den lokale tilpassingen blir usikker.

se relativt lineært ut, hvis vi se litt stort på det. Vi har i hvert fall ingen (kjemiske) teorier som kan forklare hvorfor det evt skulle være en "kul" i midten, og da vil vi heller ikke tilpasse parametere til det kurve-

<pre> cor(PercentAlcohol, Carbohydrates) cor.test(PercentAlcohol, Carbohydrates) cor(PercentAlcohol, Carbohydrates)^2 &gt; cor(PercentAlcohol, Carbohydrates) [1] 0.5210898 &gt; cor.test(PercentAlcohol, Carbohydrates)  Pearson's product-moment correlation  data: PercentAlcohol and Carbohydrates t = 7.5023, df = 151, p-value = 5.005e-12 alternative hypothesis: true correlation is not equal to 0 95 percent confidence interval:  0.3950790 0.6278535 sample estimates:  cor 0.5210898  &gt; cor(PercentAlcohol, Carbohydrates)^2 [1] 0.2715346 </pre>	<pre> cor(PercentAlcohol[-104], Carbohydrates[-104]) cor.test(PercentAlcohol[-104], Carbohydrates[-104]) cor(PercentAlcohol[-104], Carbohydrates[-104])^2 &gt; cor(PercentAlcohol[-104], Carbohydrates[-104]) [1] 0.5484863 &gt; cor.test(PercentAlcohol[-104], Carbohydrates[-104])  Pearson's product-moment correlation  data: PercentAlcohol[-104] and Carbohydrates[-104] t = 8.0338, df = 150, p-value = 2.576e-13 alternative hypothesis: true correlation is not equal to 0 95 percent confidence interval:  0.4265300 0.6508548 sample estimates:  cor 0.5484863  &gt; cor(PercentAlcohol[-104], Carbohydrates[-104])^2 [1] 0.3008373 </pre>
---	---

Korrelasjon 0.52 med outliers, 0.55 uten. Positiv korrelasjon som forventet, og den øker nå outlier er borte.

Hypotesetest: Parameter: Korrelasjonen i populasjonen,  $\rho$

$H_0$ : Ingen (null) Sammenheng;  $\rho = 0$

$H_1$ : Sammenheng mellom alkohol% og karbohydrater;  $\rho \neq 0$

Velges tosidig hypotese delvis av konservasjon & fordi det er konservativt, delvis fordi jeg vet svært lite om dette og kan derfor ikke utelukke både negativ og positive sammenheng.

Kausale trengs mer sikkerhet for å få høyere % alkohol  $\rightarrow$  positiv smh. Tosidig kausale brukes mest sikkerhet når det lages alkohol  $\rightarrow$  negativ smh.

$$\rho = \frac{E[(x - \mu_x)(y - \mu_y)]}{\sigma_x \sigma_y}$$

Bruker  $r$ , altså korrelasjonen i utvalget, som testobservator.

Velger signifikansnivå 0.01 denne gang.

Hva betyr signifikansnivå?

Det maksimale akseptable nivå for type I-feil. Vi ønsker altså at

$$P(\text{type I-feil}) = P(\text{Forkaste } H_0 \mid H_0 \text{ sann}) = P\left(\begin{array}{l} \text{"r er veldig stor"} \\ \text{ell. veldig liten"} \\ \text{(høy positiv korrelasjon)} \\ \text{(høy negativ korrelasjon)} \end{array} \mid \text{tross } \rho = 0\right) \leq 0.01$$

For å finne ut hvor stor  $P$  (type I-feil) vi ville fått dersom vi hadde forkastet  $H_0$  med de observasjonene vi har nå, beregner vi  $p$ -verdien:

$p$ -verdi = "signifikanssannsynligheten" = den sannsynligheten for type I-feil vi måtte ha godtatt dersom vi valgte å forkaste  $H_0$  på bakgrunn av det vi har observert

$$\rightarrow P(\text{minst like ekstreme observasjoner} \mid H_0) \\ = P(\text{som det vi har observert})$$

$$= P(|r| > 0.52 \mid \rho = 0) \stackrel{\uparrow}{=} 0.00000000005005$$

OK, dette var lite sannsynlig under  $H_0$ . Forkast  $H_0$ .

Fra R-utskrift

<pre>lm(Carbohydrates~PercentAlcohol)</pre>	<pre>lm(Carbohydrates~PercentAlcohol)</pre>
<pre>&gt; lm(Carbohydrates~PercentAlcohol)</pre>	<pre>&gt; lm(Carbohydrates[-104]~PercentAlcohol[-104])</pre>
<pre>Call: lm(formula = Carbohydrates ~ PercentAlcohol)</pre>	<pre>Call: lm(formula = Carbohydrates[-104] ~ PercentAlcohol[-104])</pre>
<pre>Coefficients: (Intercept) 2.605</pre>	<pre>Coefficients: (Intercept) 1.650</pre>

En annen måte å tallfeste den lineære sammenhengen på, er å beregne den linja som passer best til data. La  $y = \text{Carbohydrates}$ ,  $x = \text{PercentAlcohol}$

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon$$

$\beta_0$  angir hvor linja skjærer y-aksen (intercept)

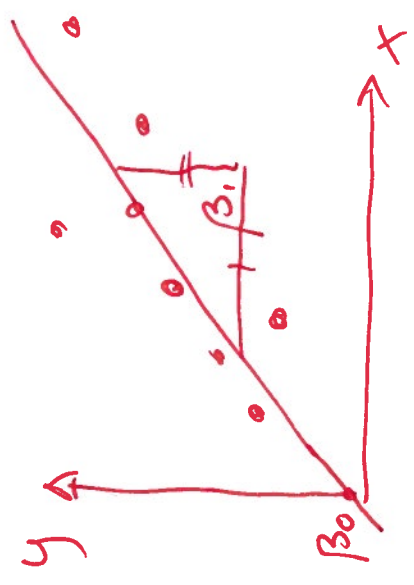
$\beta_1$  angir stigningstallet til linja, altså hvor mange enheter vi forventer at  $y$  skal øke når  $x$  øker med én enhet

Her er  $\hat{\beta}_0 = 2.605$   
 $\hat{\beta}_1 = 1.789$

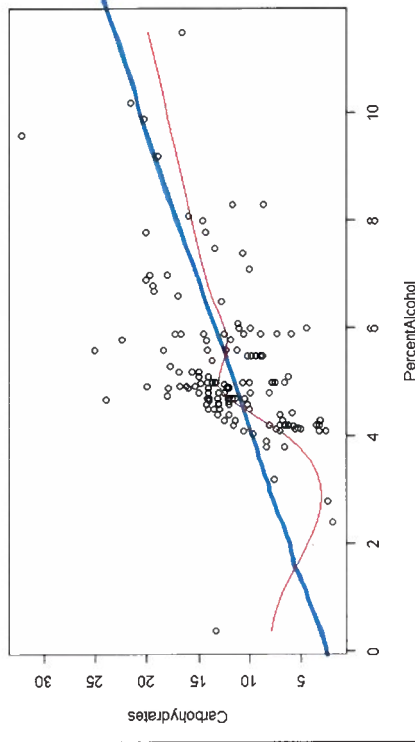
→ det endrer seg hvis vi fjerner outlieren. Nye estimat:

$$\hat{\beta}_0 = 1.650$$

$$\hat{\beta}_1 = 1.958$$



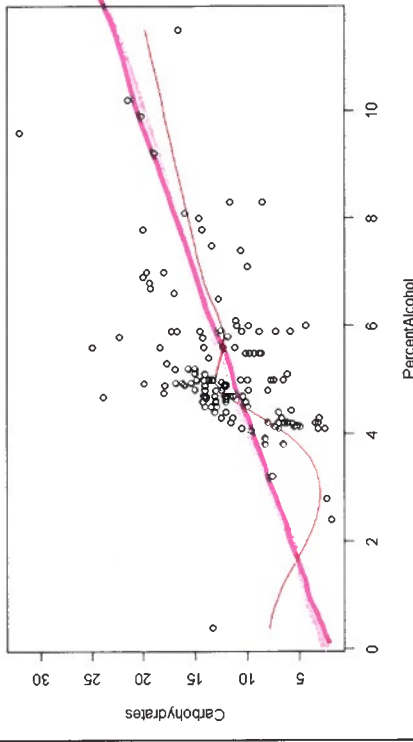
```
?abline(2.605, 1.789, col="blue")
```



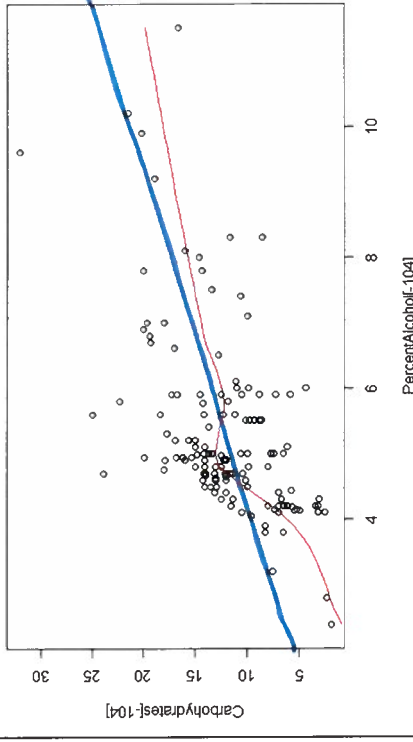
```
# abline(2.605, 1.789, col="blue")
```

Disse estimatene måtte byttes ut med 1.650, 1.958. Velger heller den generelle måten:

```
abline(lm(Carbohydrates~PercentAlcohol), col="pink")
```



```
abline(lm(Carbohydrates[-104]~PercentAlcohol[-104]), col="blue")
```



Tegnes inn denne linja på to ulike måter i R

```

summary(lm(Carbohydrates~PercentAlcohol))
summary(lm(Carbohydrates[-104]~PercentAlcohol[-104]))
summary(lm(Carbohydrates[-104]~PercentAlcohol[-104]))
Call:
lm(formula = Carbohydrates ~ PercentAlcohol)

Residuals:
    Min       1Q   Median       3Q      Max
-8.940 -3.145  0.586  2.628 12.904

Coefficients:
(Intercept) 2.6049
PercentAlcohol 1.7892

Estimate Std. Error t value Pr(>|t|)
(Intercept) 2.6049    1.2924    2.016  0.0456 *
PercentAlcohol 1.7892    0.2385    7.502  5e-12 ***

Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 4.201 on 151 degrees of freedom
Multiple R-squared:  0.2715,    Adjusted R-squared:  0.2667
F-statistic: 56.29 on 1 and 151 DF,  p-value: 5.005e-12

Call:
lm(formula = carbohydrates[-104] ~ PercentAlcohol[-104])

Residuals:
    Min       1Q   Median       3Q      Max
-9.3034 -3.0670  0.7311  2.6830 13.0661

Coefficients:
(Intercept) 1.6495
PercentAlcohol[-104] 1.9583

Estimate Std. Error t value Pr(>|t|)
(Intercept) 1.6495    1.3253    1.245  0.215
PercentAlcohol[-104] 1.9583    0.2438    8.034  2.58e-13 ***

Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 4.128 on 150 degrees of freedom
Multiple R-squared:  0.3008,    Adjusted R-squared:  0.2962
F-statistic: 64.54 on 1 and 150 DF,  p-value: 2.576e-13

```

uten outliers

Men  $\hat{\beta}_0$  og  $\hat{\beta}_1$  er som sagt estimater, og har usikkerhed knyttet til sig.  
 Vi kan både gøre hypotesetester og beregne konfidensintervaller.  
 La oss ta hypotesetest først:

$H_0$ : Ingen (null) sammenheng mellom alkohol% og karbohydrater :  $\beta_1 = 0$   
 $H_1$ : Det er en \_\_\_\_\_ :  $\beta_1 \neq 0$

Nøye tilfidshypoteser av samme grunn som tidligere.  
 (Obs: R tester for om  $\beta_0 = 0$  også, men det er ikke det som er interessant her.)



P-verdi =  $P(\hat{\beta}_1 \text{ minst like ekstremt } | \beta_1 = 0)$   
verdi for  $\hat{\beta}_1$  som det vi har observert

$$= P(|\hat{\beta}_1| > 1.7892 \mid \beta_1 = 0) \stackrel{\text{Forkast } H_0}{=} 0.0000000000005005$$

↑  
i følge R-utskrift

Eksempel på at selv om effekt målet ( $r$  og  $\hat{\beta}_1$ ) er forskjellig, så er metode basert på samme antakelser og samme matematikk, og gir like konklusjoner.

Konfidensintervall for  $\beta_1$ :

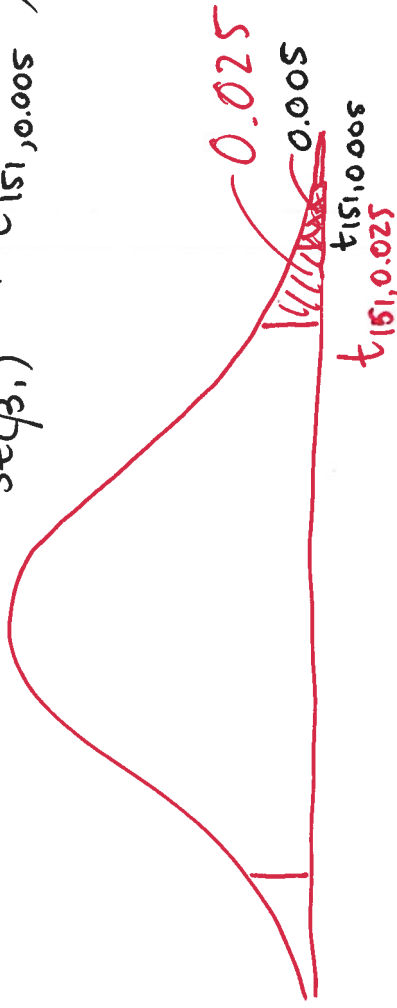
$$\text{Ta utgangspunkt i at } \frac{\hat{\beta}_1 - \beta_1}{\text{std. Error}(\hat{\beta}_1)} \sim T_{n-2}$$

$$\text{Her } n = 153, \text{ og } df = 151$$

Da er

$$P(-t_{151, 0.025} < \frac{\hat{\beta}_1 - \beta_1}{SE(\hat{\beta}_1)} < t_{151, 0.025}) = 0.95$$

$$P(-t_{151, 0.005} < \frac{\hat{\beta}_1 - \beta_1}{SE(\hat{\beta}_1)} < t_{151, 0.005}) = 0.99$$



velger denne

vi finner i table D :  $t_{100, 0.025} = 1.984$

eg  $t_{1000, 0.025} = 1.962$

$t_{100, 0.005} = 2.626$

$t_{1000, 0.005} = 2.581$

$$P(-1.984 < \frac{\hat{\beta}_1 - \beta_1}{SE(\hat{\beta}_1)} < 1.984) = 0.95$$

-2.626

2.686

0.99

$$P(\hat{\beta}_1 - 1.984 \cdot SE(\hat{\beta}_1) < \beta_1 < \hat{\beta}_1 + 1.984 \cdot SE(\hat{\beta}_1)) = 0.95$$

-2.686

+ 2.686

0.99

99% KI for  $\beta_1$  :

$$1.7892 \pm 2.686 \cdot 0.2385$$

$$\rightarrow \underline{[1.15, 2.43]}$$

$H_0$ -verdien  
for  $\beta_1$   
↓

Intervallt inneholder ikke 0

→ Forkast  $H_0$  på nivå 0.01

95% KI for  $\beta_1$  :

$$1.7892 \pm 1.984 \cdot 0.2385$$

$$\rightarrow \underline{[1.32, 2.26]}$$

Intervallt inneholder ikke 0

→ Forkast  $H_0$  på nivå 0.01