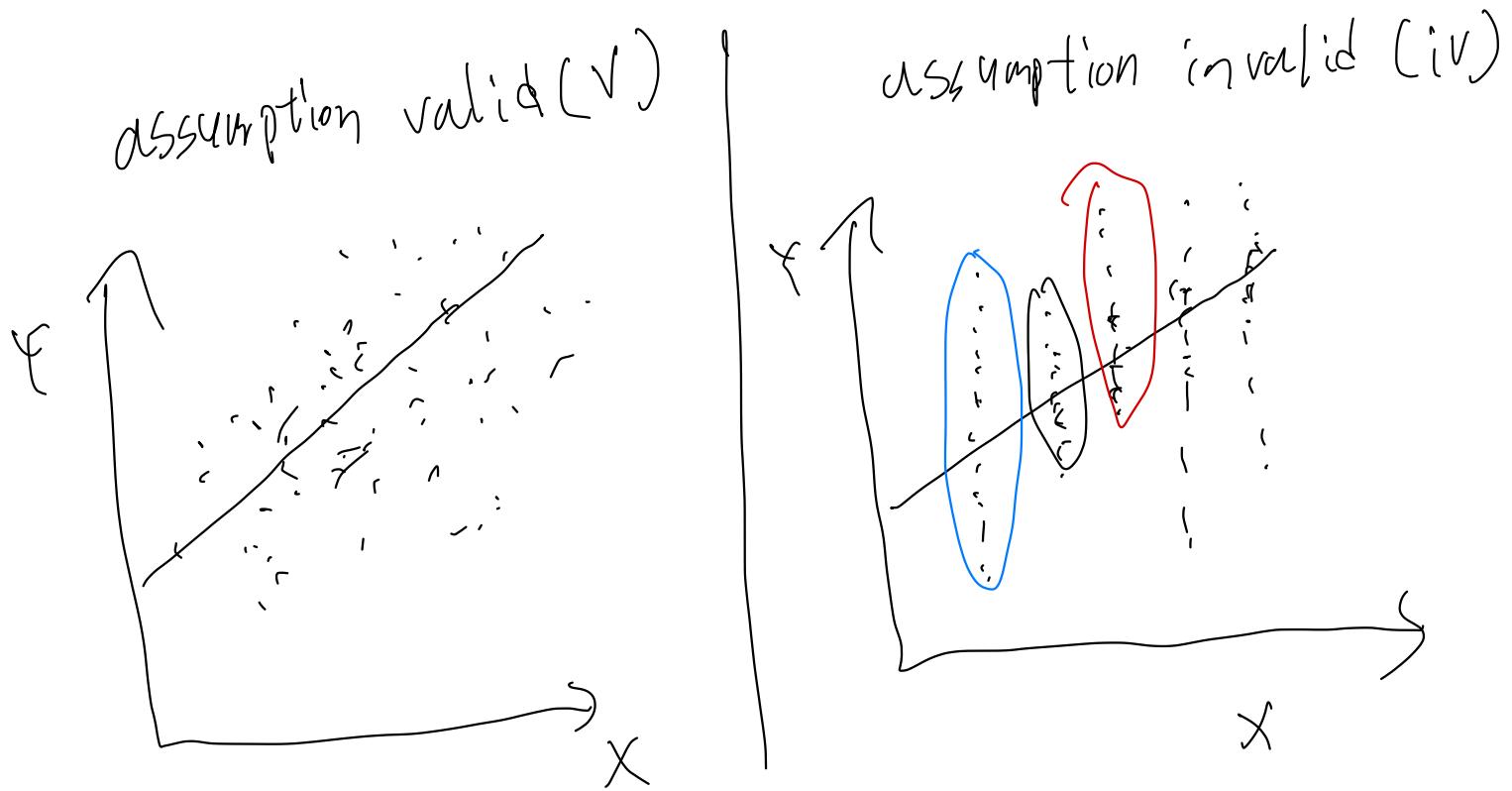


2주차 분석학(W3)

2020 3 31 월

01주차

(b)



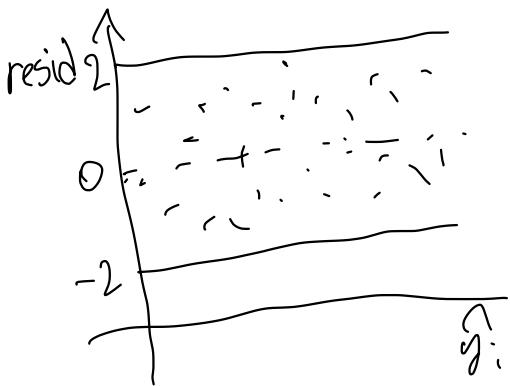
In graph (v) Shows that observations are independent
there's no specific pattern on graph. \rightarrow that's valid case.

In graph (iv) shows that observations are not independent
each others. They formed like some group. That kind of
pattern can seen like medicine test with (treatment group
with control group). \rightarrow that is invalid case.

linear regression is conditional given observation.
if independent assumption violated, that can be
not effective.

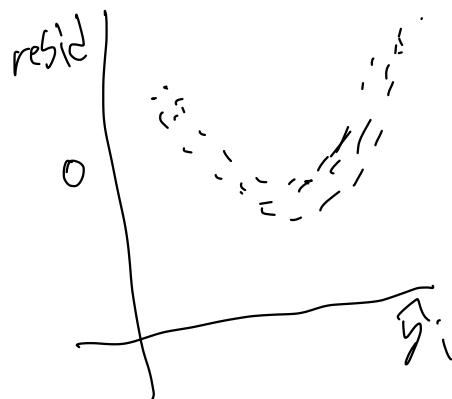
1. (C)

valid (V)



invalid

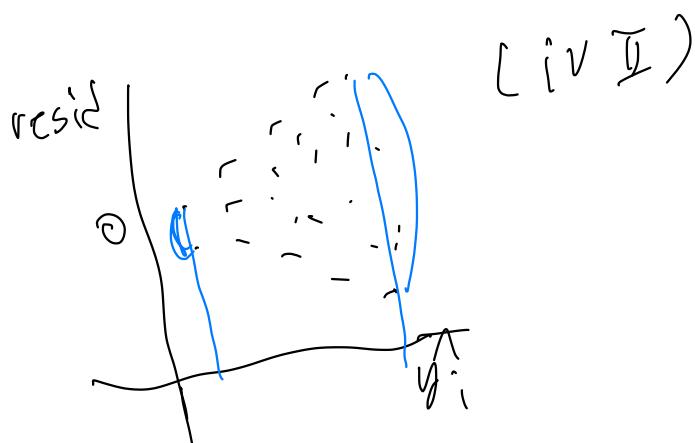
(IV I)



valid case

in (V) residual \sim standard normal dist, so it will plotted randomly around 0, and almost within $[-2, 2]$, \rightarrow that is constant variance case.

But invalid case like (IV II) that has hetero variance. And invalid case like (IV I) that shows non linear and not constant variance,



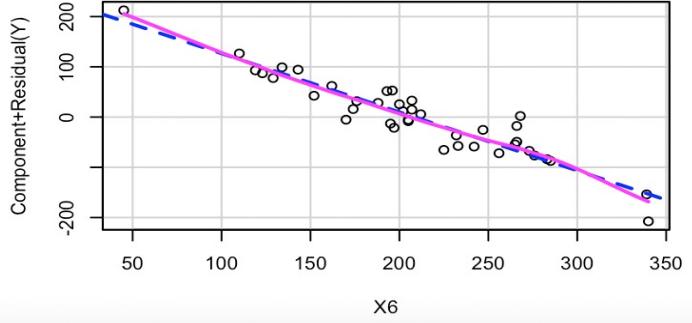
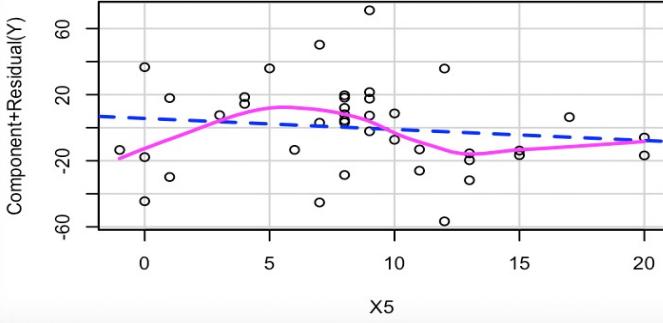
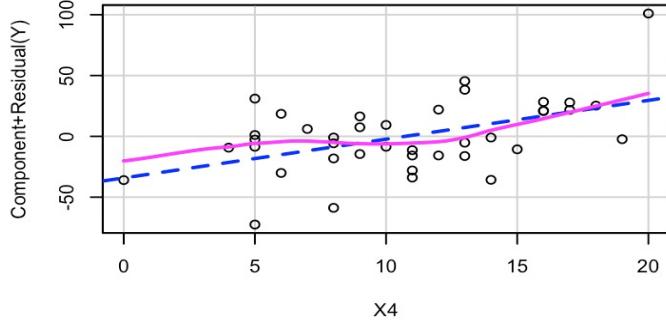
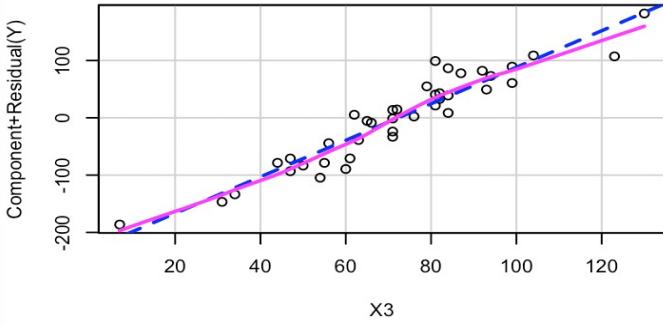
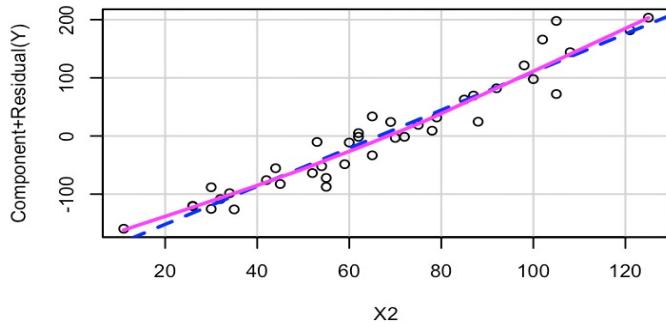
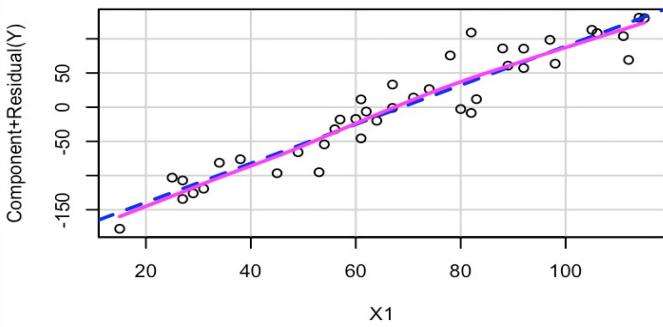
2.

(1)

R 4.1.2 · ~/ ➔

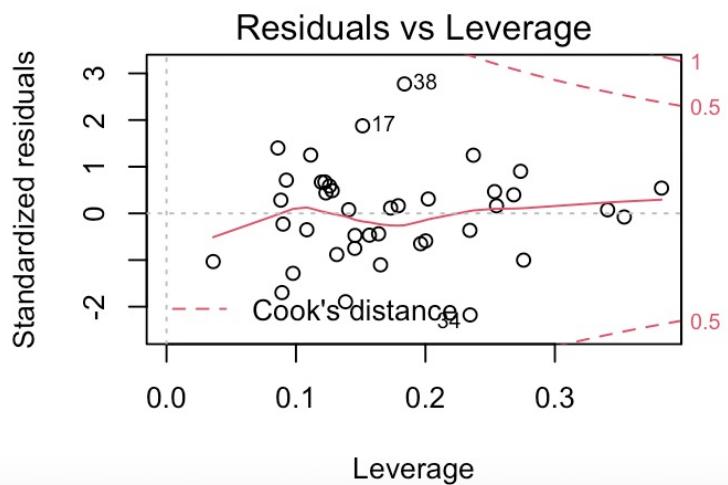
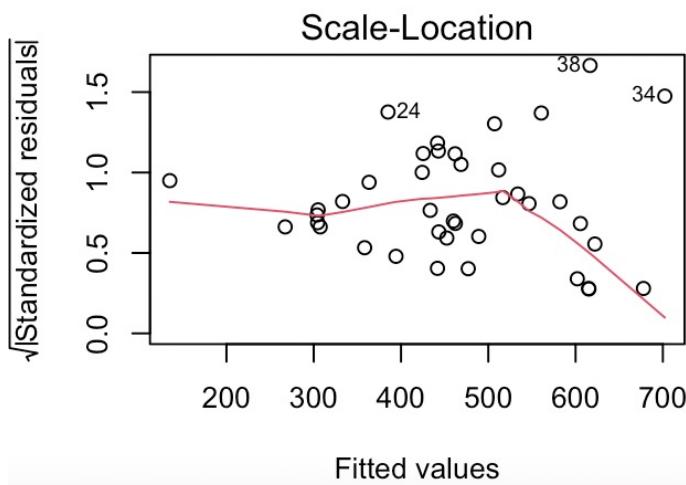
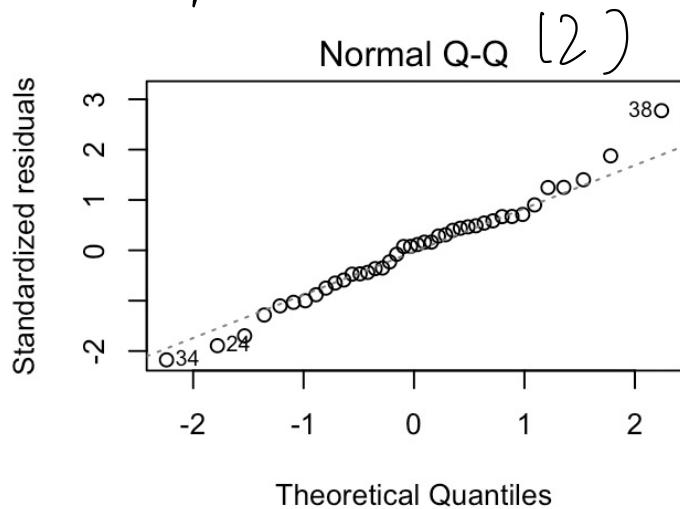
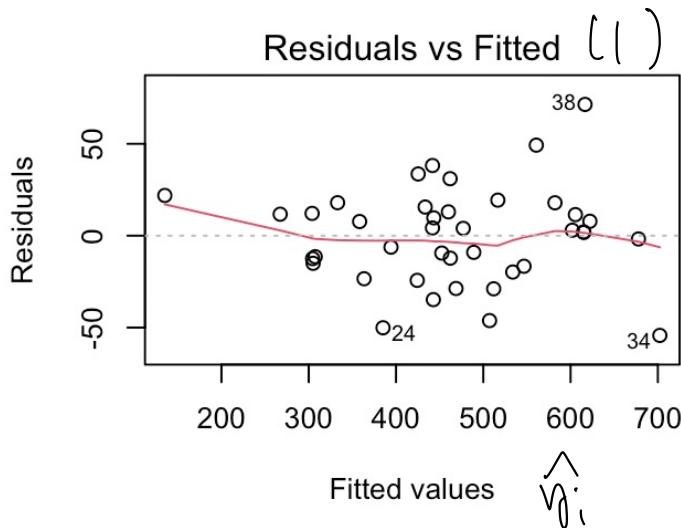
	Y	X1	X2	X3	X4	X5	X6
1	443	49	79	76	8	15	205
2	290	27	70	31	6	6	129
3	676	115	92	130	0	9	339
4	536	92	62	92	5	8	247
5	481	67	42	94	16	3	202
6	296	31	54	34	14	11	119

Component + Residual Plots



linearity assumption seems fine

Let's check error distribution assumption.



(2) and shapiro test shows normality assumption would fine.

(1) shows some non constant variance

Constant Variance assumption might break.

$DW = 2.30$ thus over 2.0 and

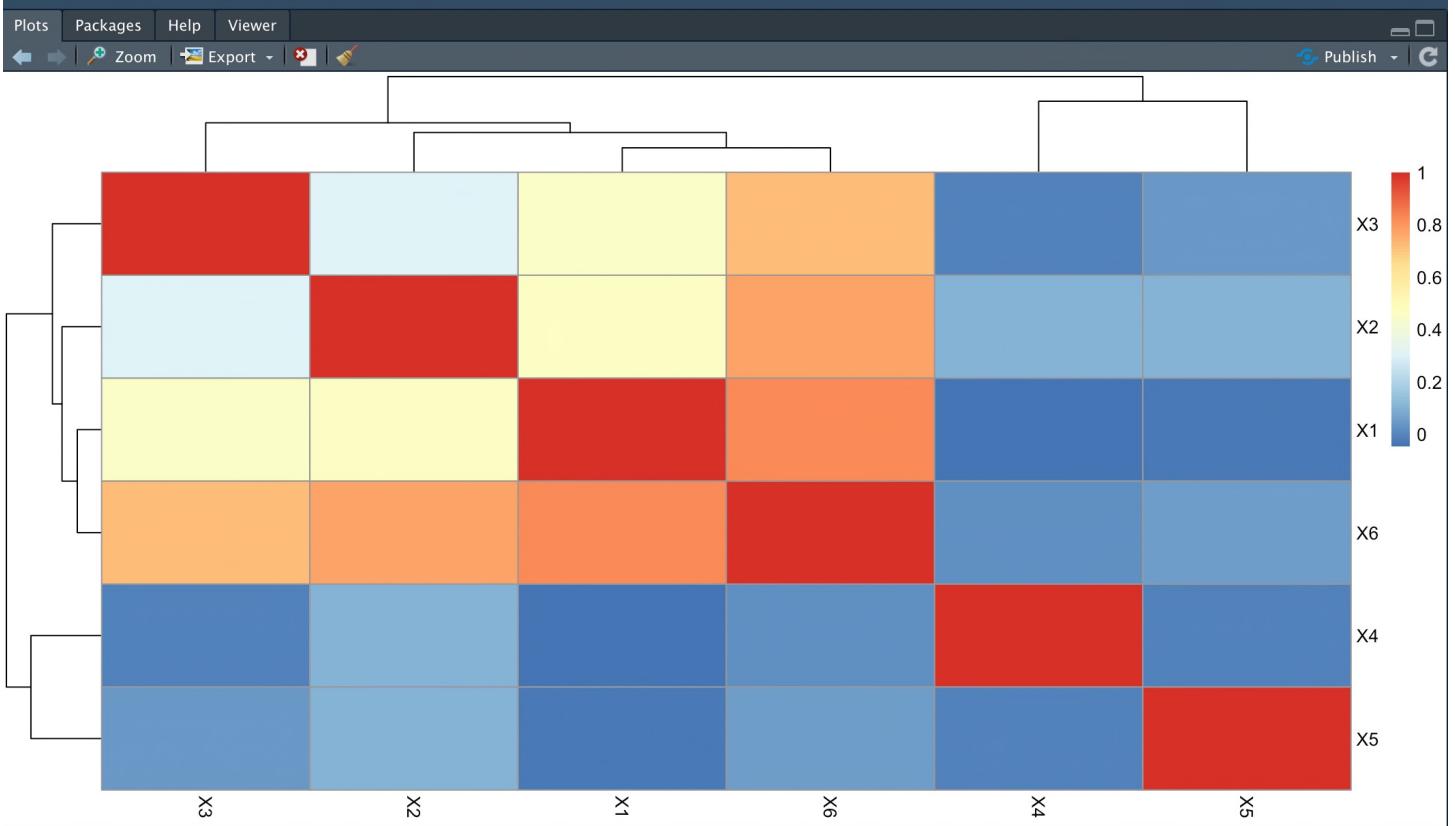
(3) shows some kinds of non-independency.
Next page, I'll check predictor variables assumptions

```
> shapiro.test(lm$residuals)
Shapiro-Wilk normality test

data: lm$residuals
W = 0.98425, p-value = 0.8407

> durbinWatsonTest(residuals(lm))
[1] 2.305141
```

```
> library(pheatmap)
> cor_matrix <- cor(df[,2:7])
> pheatmap(cor_matrix, legend = T)
>
```



Well check correlation matrix with independent variables.

And this plot shows some kinds of multicollinearity.

So, I think as error distribution assumption,

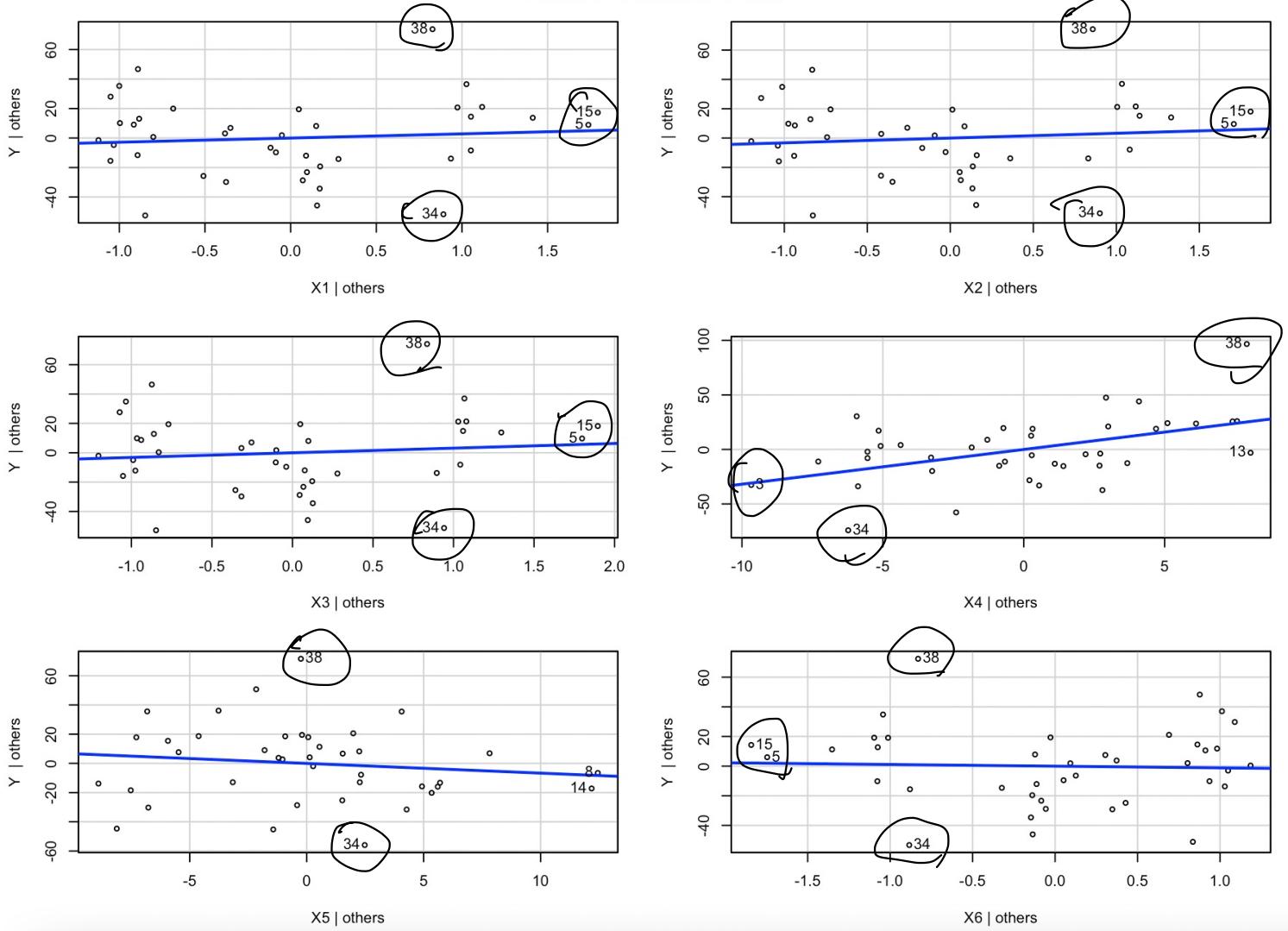
constant variance assumption were violated.
Independency

And about predictor variables, it shows some

high correlation, so linearly independent assumption in predictor variables was violated. And it seems like have multicollinearity.

(2) Unusual observations.

Added-Variable Plots



(3, 5, 15, 34, 38) points are suspicious.

(3, 15 → leverage points and outliers.

(34, 38 → influential points and outliers.

And I'll check the details in next pages.

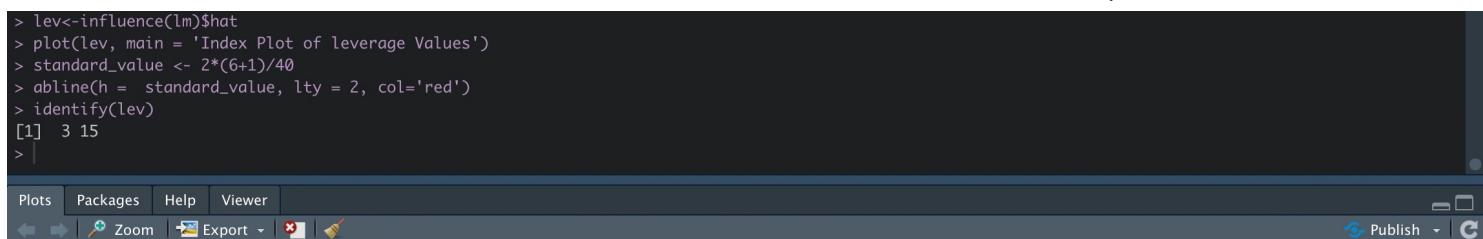
Check leverage point

We can use $P_{ii} = i^{\text{th}} \text{ diagonal element of } P = X(X^T)^{-1}X$ average

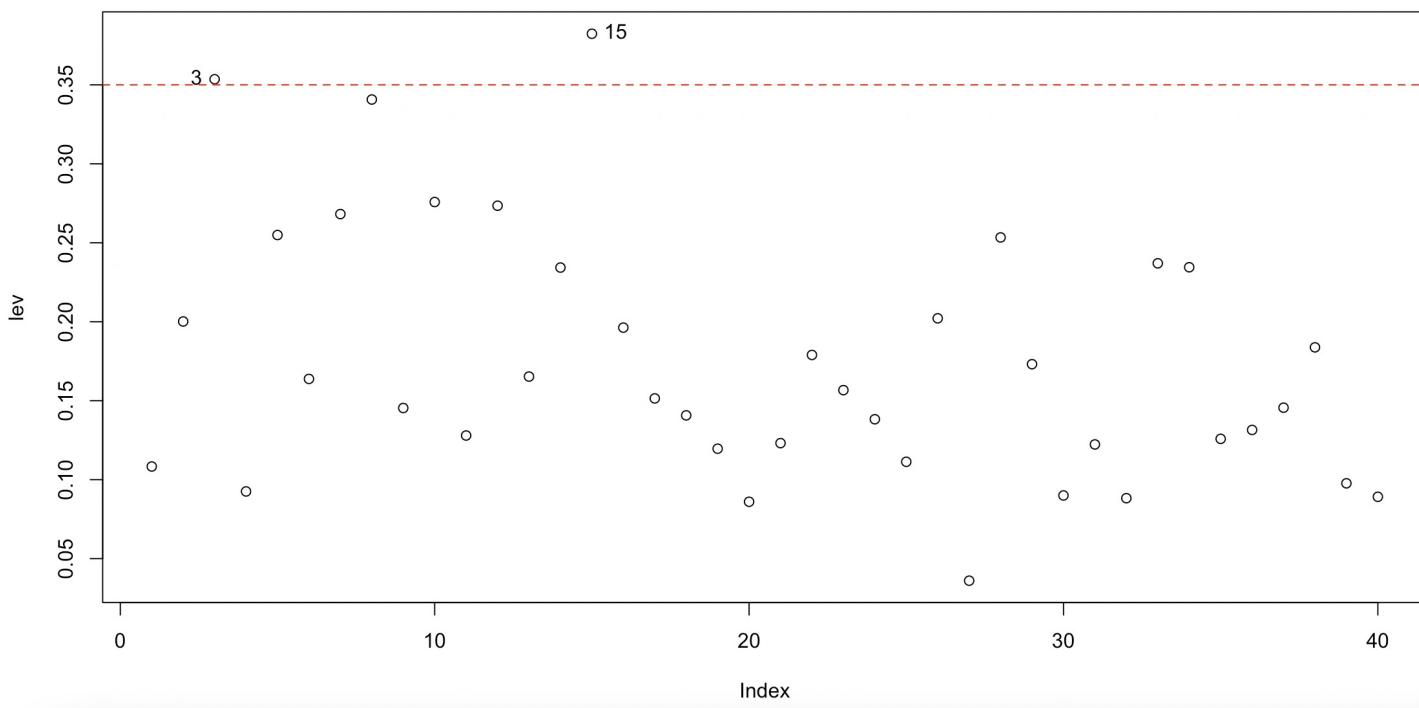
$$n \leq 40, P = 6$$

↳ and compare with $2\left(\frac{p+1}{n}\right)$

Standard value $= 2\left(\frac{p+1}{n}\right)$, we can compare $P_{ii} > 2\left(\frac{p+1}{n}\right)$ and if that's true \rightarrow that points are high leverage points.



Index Plot of leverage Values

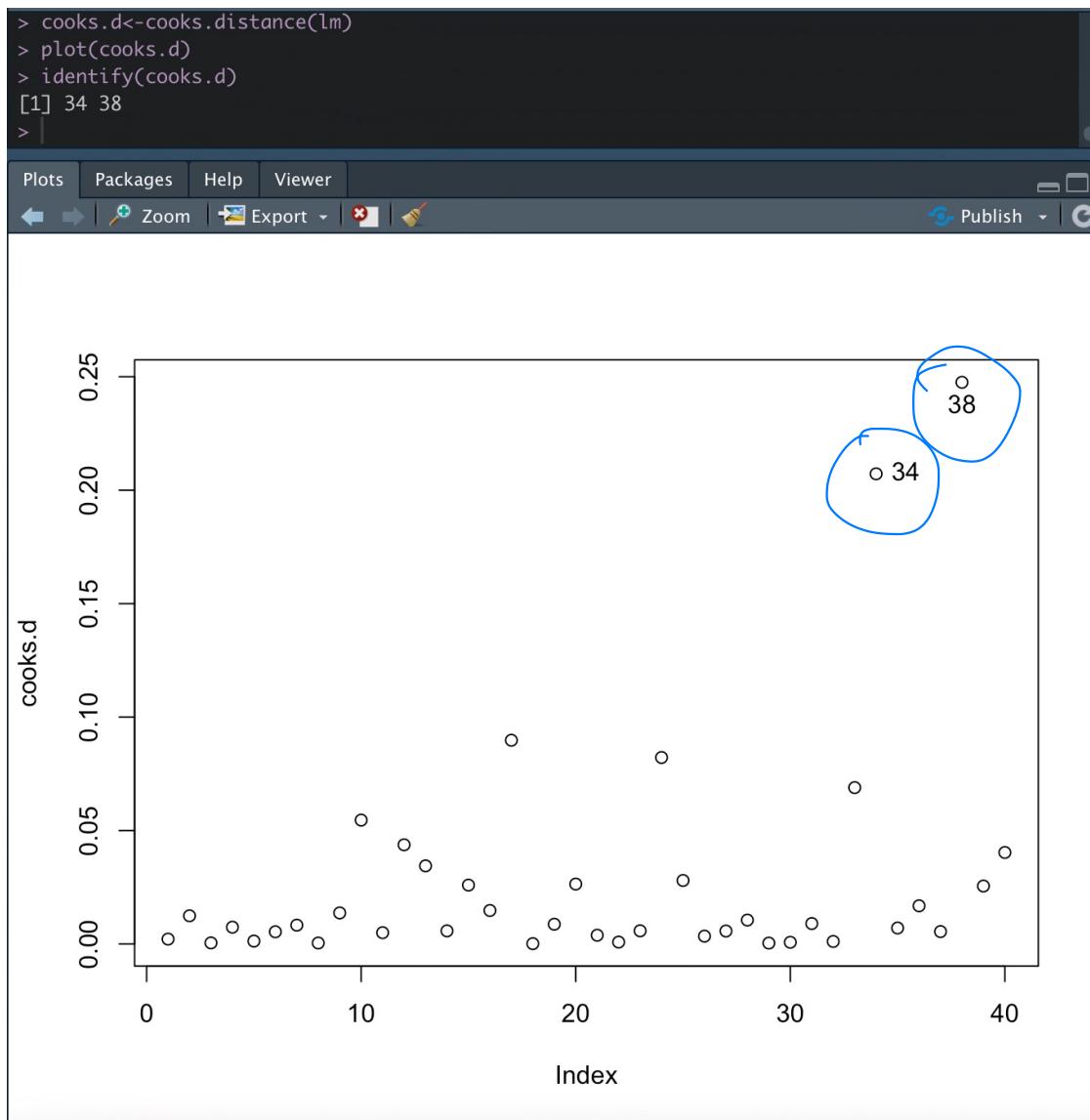


3, 15 points are $P_{ii} > 2\left(\frac{p+1}{n}\right) = 0.35$ (above red line)
3, 15 points are leverage points.

check influential points \rightarrow we can use Cook's distance or DFFITS
 if $C_i \geq F(p=1, n-p-1; 0.5)$ or $|C_i| \geq 1$
 then i^{th} is influential.
 $p=6, n=40$
 $F(7, 33, 0.5) =$

```
> qf(0.5, 7, 33)
[1] 0.9252857
```

\rightarrow we just need to check $C_i \geq 0.9252857$



There's no points that $C_i \geq F(7, 33, 0.5)$
 but we have to check 34, 38 points since
 it's far from other points. And it might be influential points.

outliers

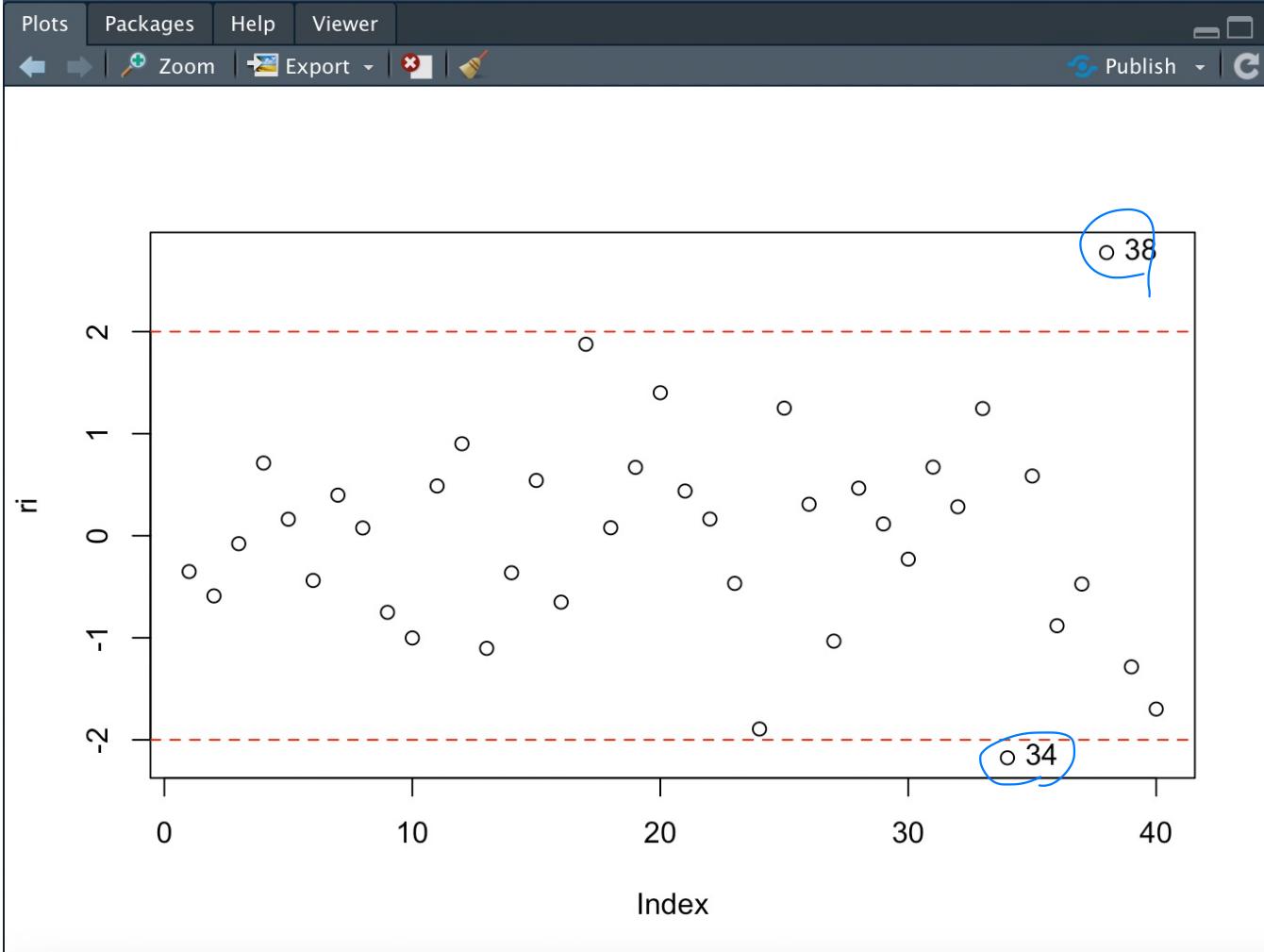
→ x-axis → leverage points.

→ y-axis → standardized residual (standardized)

we can check outliers.

in (2) we found leverage points (3rd, 15th)

```
> ri<-rstandard(lm); plot(ri)
> abline(h = c(-2, 2), col='red', lty=2)
> identify(ri)
[1] 34 38
>
```



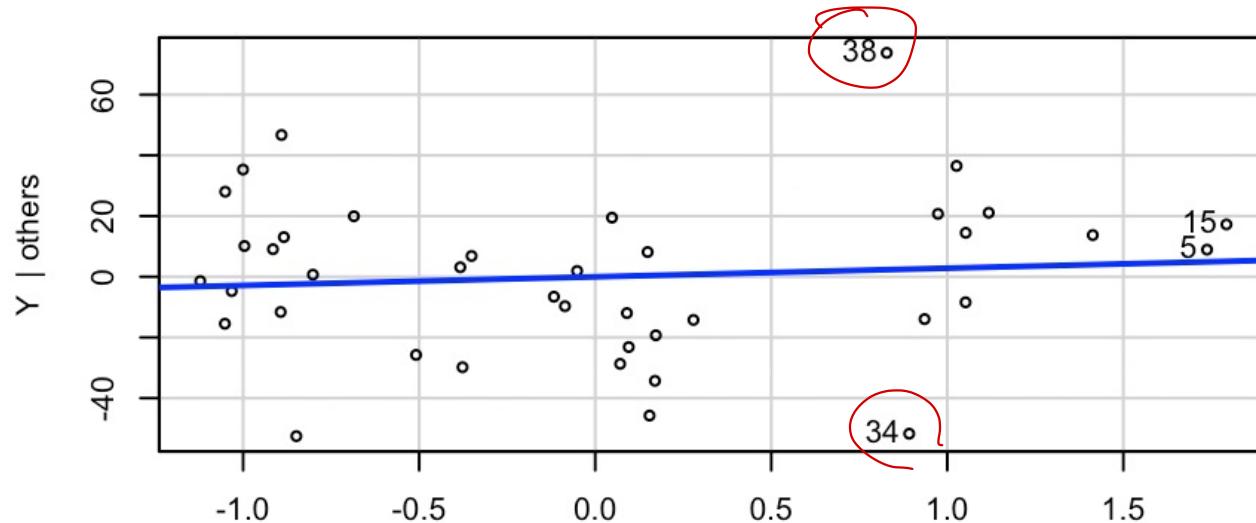
We can use standardized residual-index plot to find outliers. And 34, 38 seems outliers. Since residual $\sim N(0, 1)$, those points are out of (-2, 2). ∴ (x-axis, 3, 15 points) (34, 38)
y-axis, 34, 38 points

We can conclude outliers (3, 15, 34, 38) points.

(3)

A-V plot for X_1

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \beta_5 X_5 + \beta_6 X_6 + \epsilon$$

 $X_1 | \text{others}$

> avPlots(lm, id.n = 1, layout = c(1,1))

Hit <Return> to see next plot: avPlots(lm, id.n = 1, layout = c(1,1))

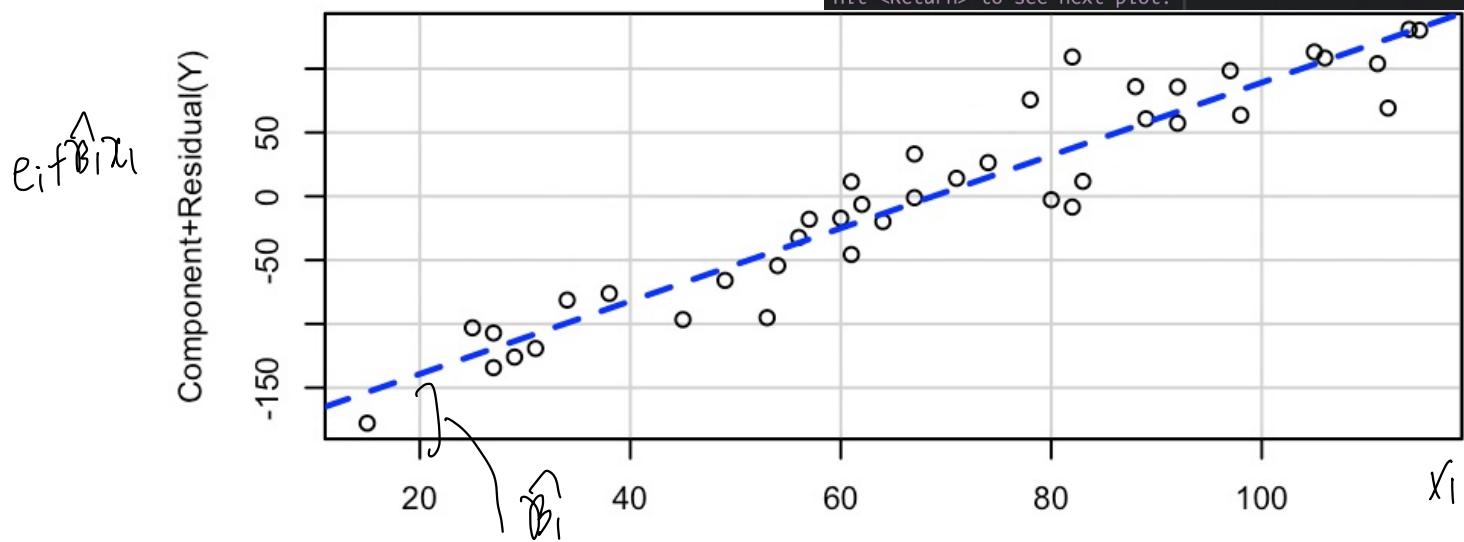
Hit <Return> to see next plot:

R+C plots for X_1

> crPlots(lm, smooth=F, layout = c(1,1))

Hit <Return> to see next plot: crPlots(lm, smooth=F, layout = c(1,1))

Hit <Return> to see next plot:



A-V plot shows, remove (X_2, X_3, \dots, X_6) 's effects for X_1 and y and display linearity. And we can check X_1 can explain Y 's linearity.

R+C plot also shows linearity with X_1 and Y . Slope might be β_1 , and points seems good fitted with line. We can sure it's linearity.