

# Recall: ggplot()- layered grammar of graphics

key tool in Data  
Visualisation

```
ggplot(data = <DATA>) +  
<GEOM_FUNCTION>(  
  mapping = aes(<MAPPINGS>),  
  stat = <STAT>,  
  position = <POSITION>  
) +  
<COORDINATE_FUNCTION> +  
<FACET_FUNCTION>
```

Coordinate chart

Scatter plot

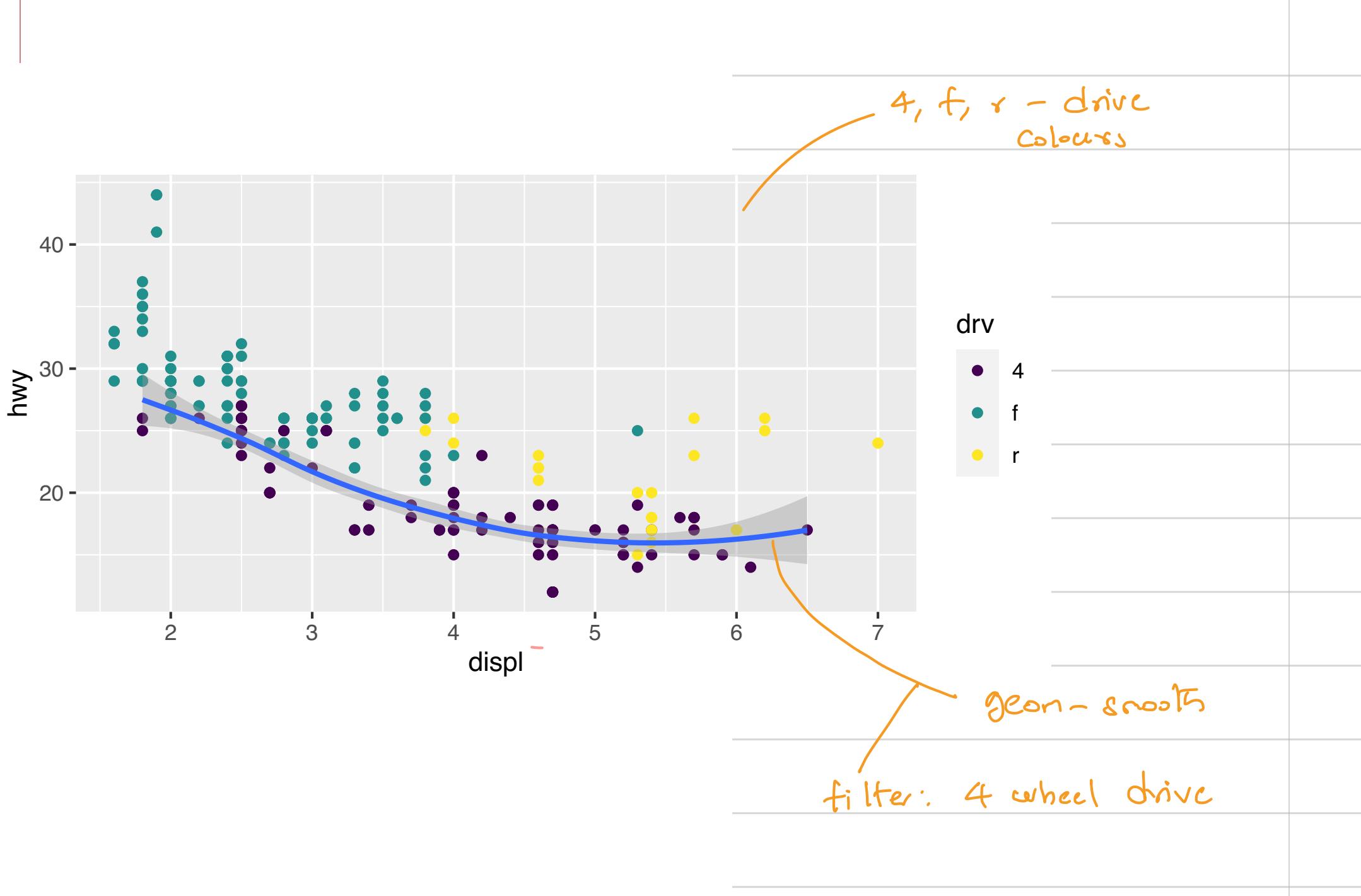
Smooth line

bar chart

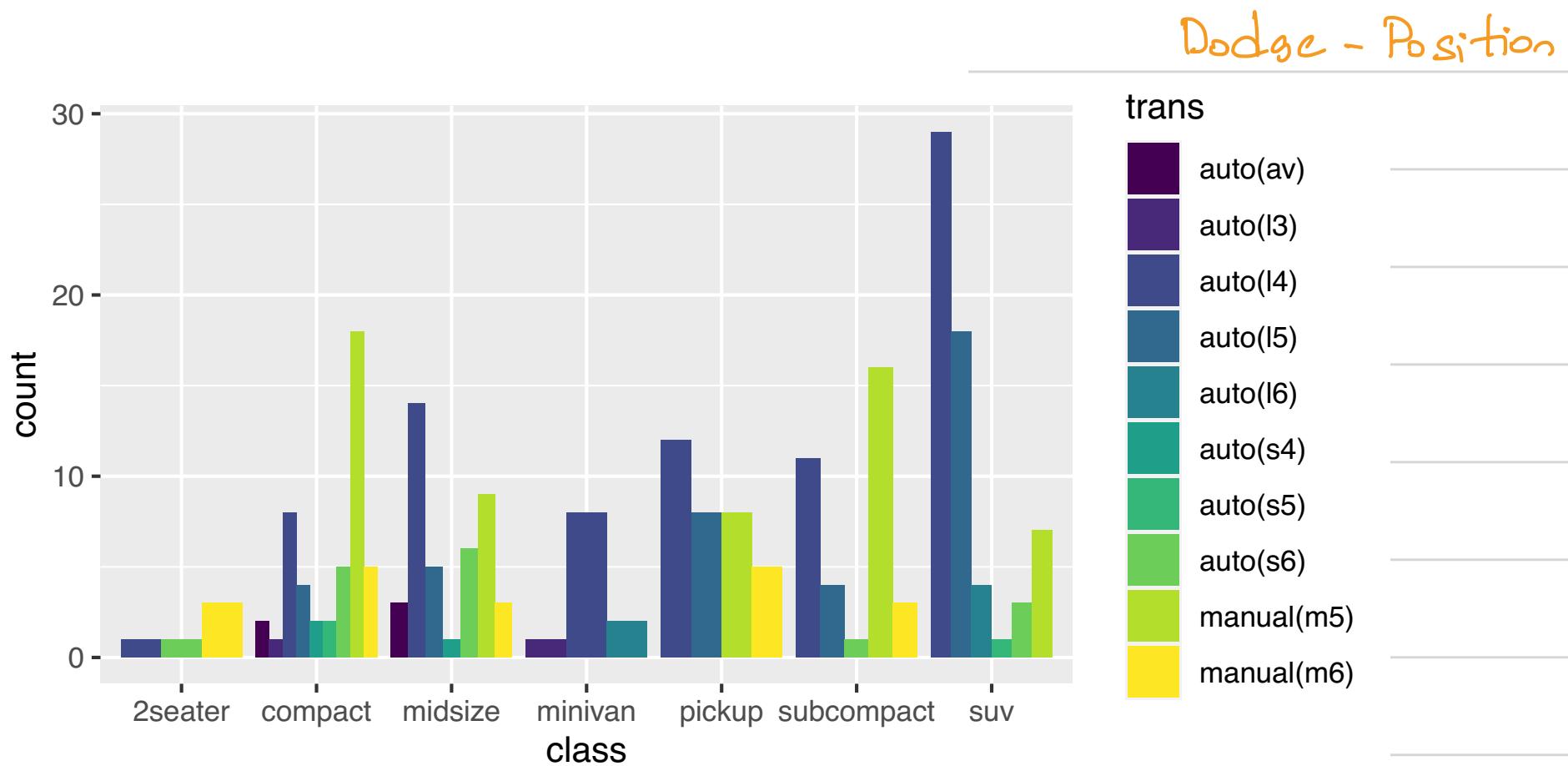
polar coordinates

Faceting.

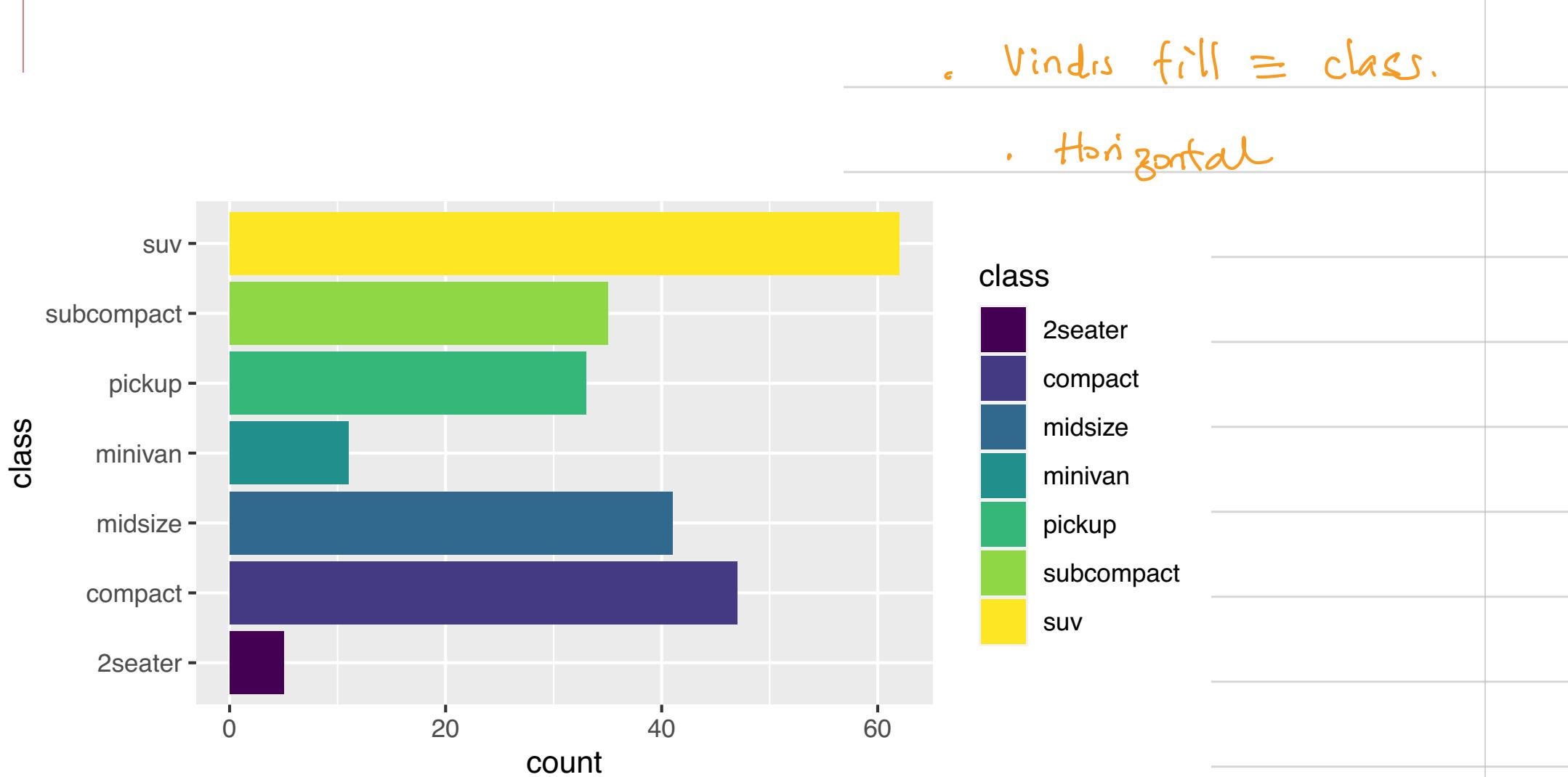
# ggplot()-geoms



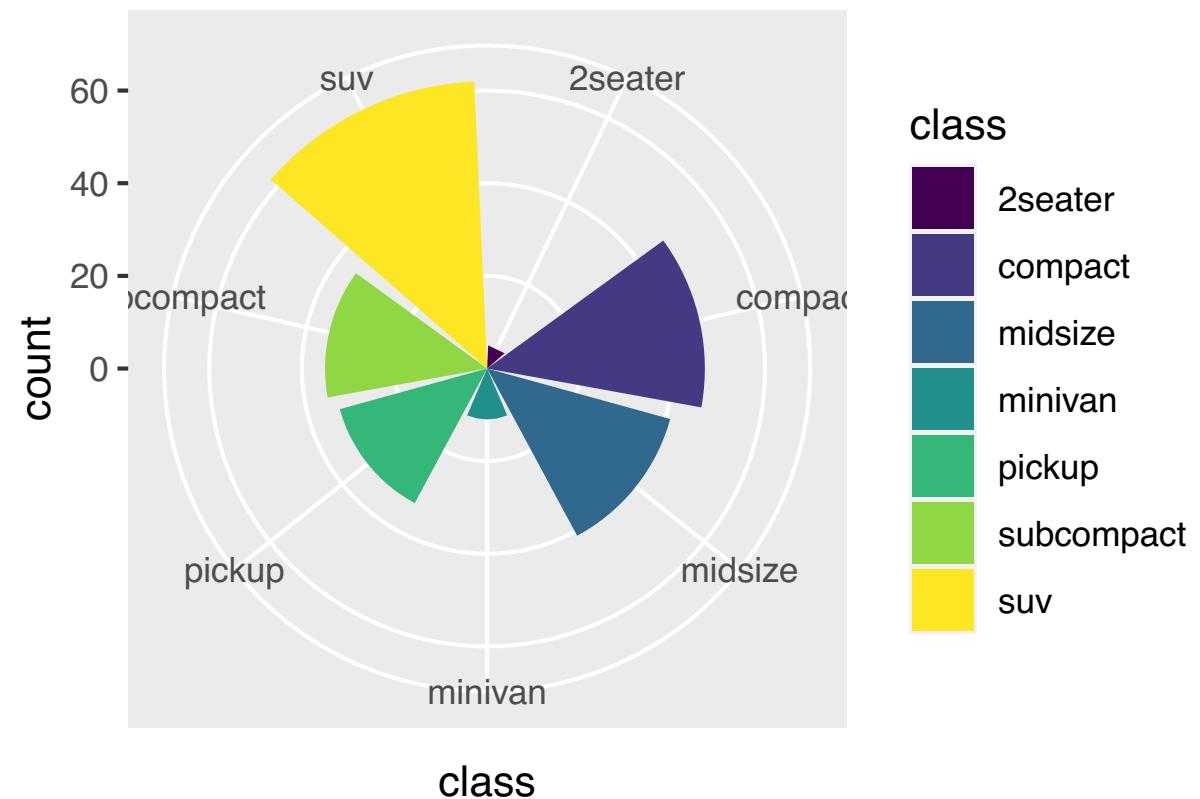
# ggplot()-barchart



# ggplot()-barchart



# ggplot()-barchart



Polar Coordinates.

## Model: Aims and Methods

- No statistical Analysis in today's classes

- Datasets & week 2 R Code

Website: <https://www.isibang.ac.in/~rathreya/Teaching/LSM/>

- while working with the R-Code
- Please set correct working directory

- introduce

- linear model
- (Exploratory)

how it works ?

Please download

the data sets

&

R-Code.

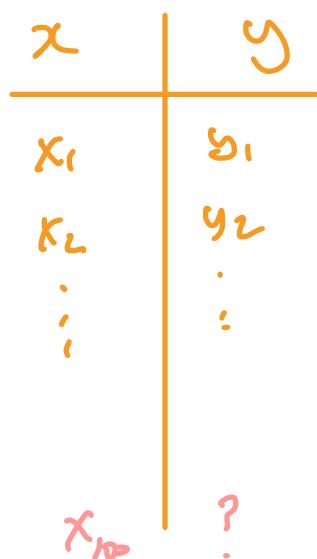
(Place R-Code and  
data set in the  
same directory)

Aim:- To find relationship between  $y$  and  $x$ .

Given: Data Set

$x$  - independent  
 $y$  - dependent on  $x$

} variable



Caution:-

All models are wrong but some are useful.

- $PV = RT$  - model for ideal gases
  - good approximation to understand the system

Goal (today):

- low-dimensional

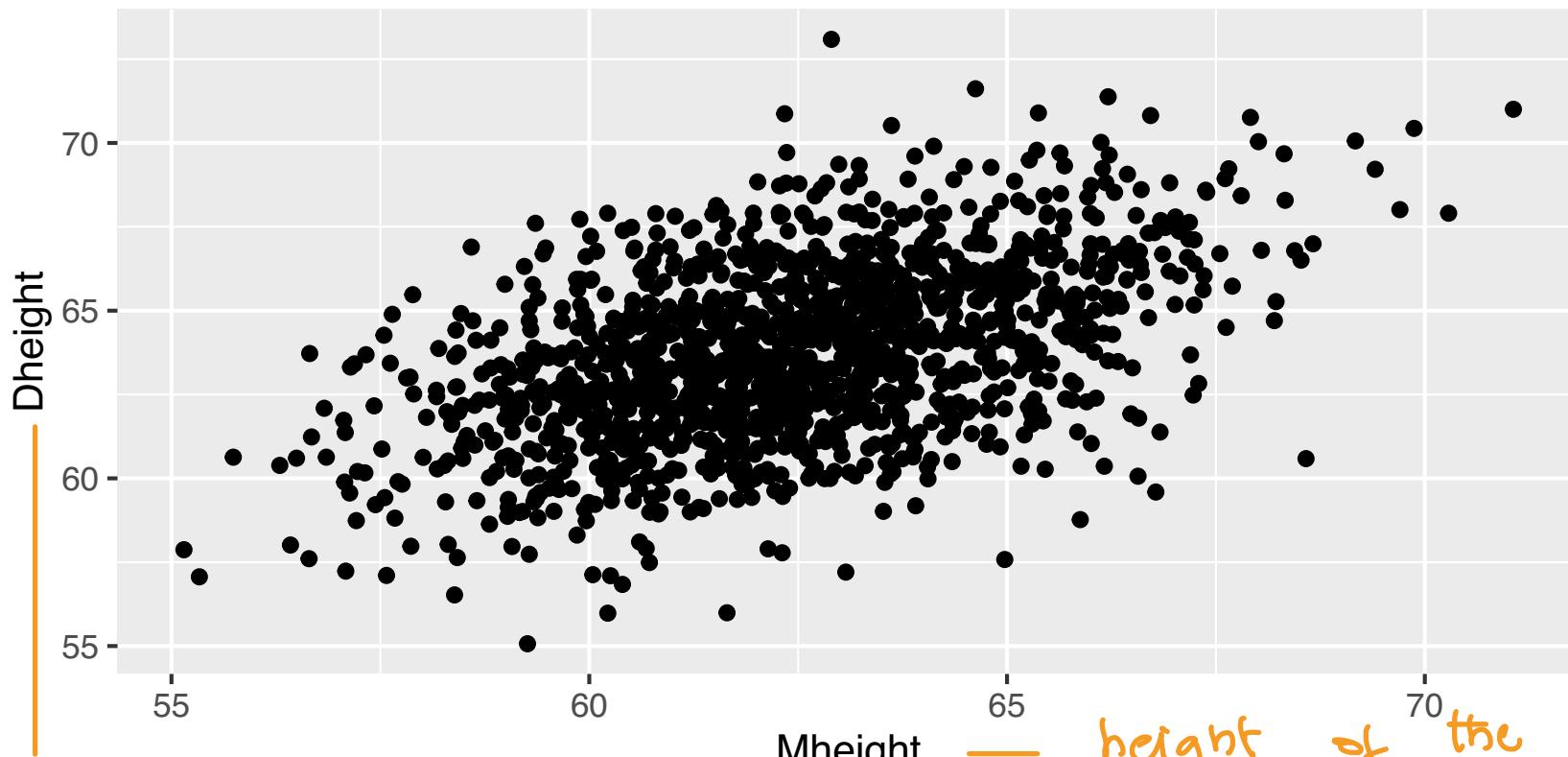
summary of the data set.

- "patterns" in data

- Prediction ?

# Scatterplots and Regression

1893 - 98 : E.S. Pearson collected ~ 1500 odd  
(heights.txt)



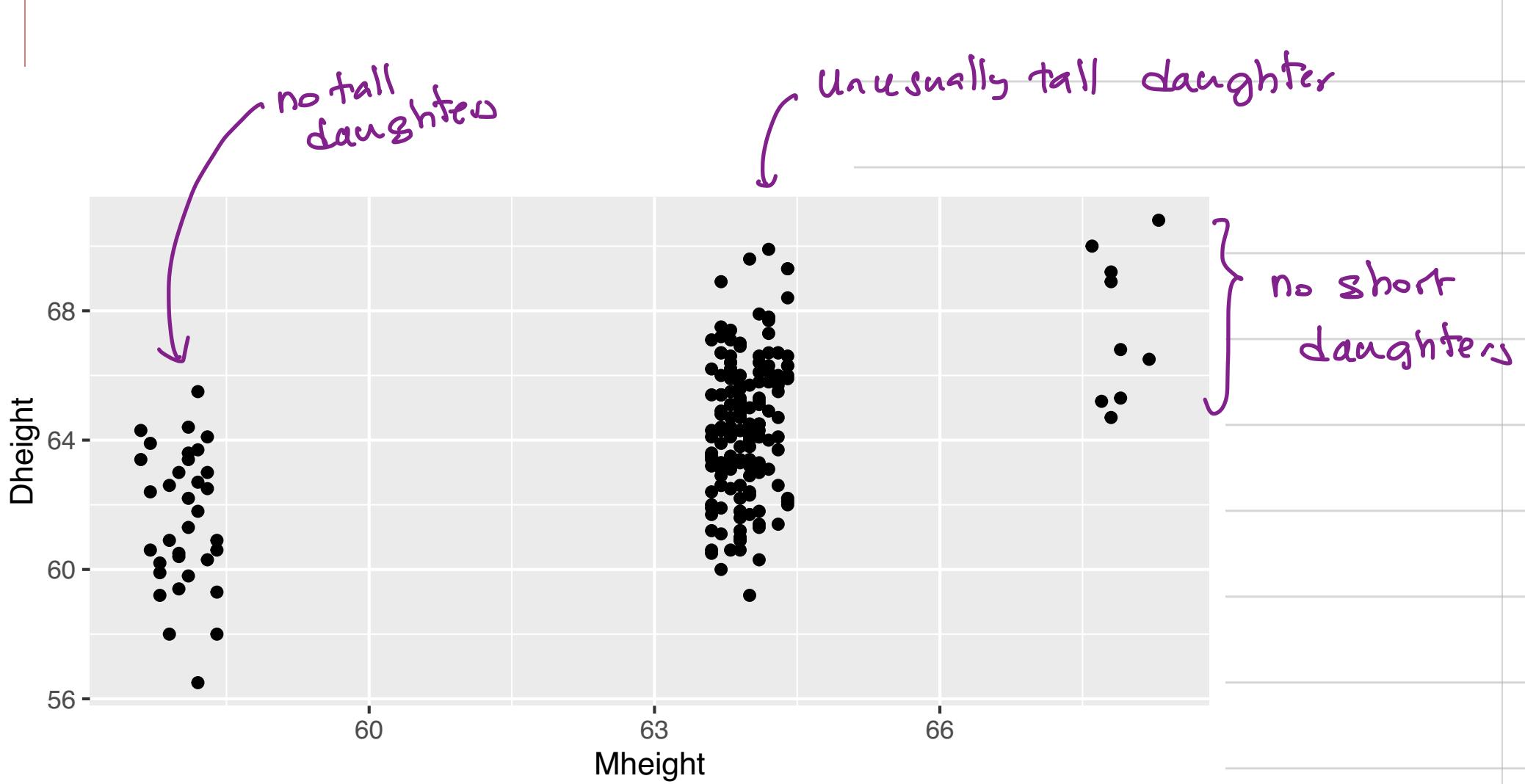
jitter plot  
moves  
each  
data point  
by a small  
uniform  
S.V.

height of one  
adult daughter  
(age > 18)

Question :- How does the daughter inherit  
object of study the height of the mother?

height of the mother  
(age < 65)

## Scatterplots and Regression



identify outliers.

- Develop a technique to use a linear model

$\{(x_i, y_i) : 1 \leq i \leq n\}$  - Data

Find:  $a_1, a_2$ :

$$y = a_2 x + a_1$$

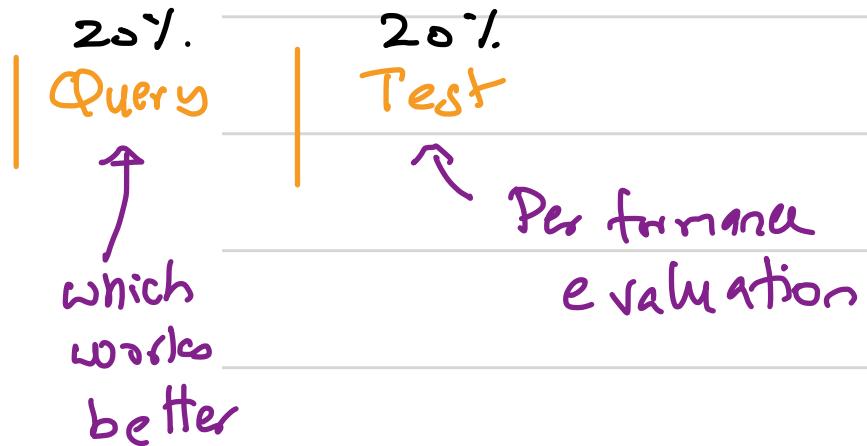
is a good approximation

In practice:

- what are hypotheses on the data?

one wants to test for

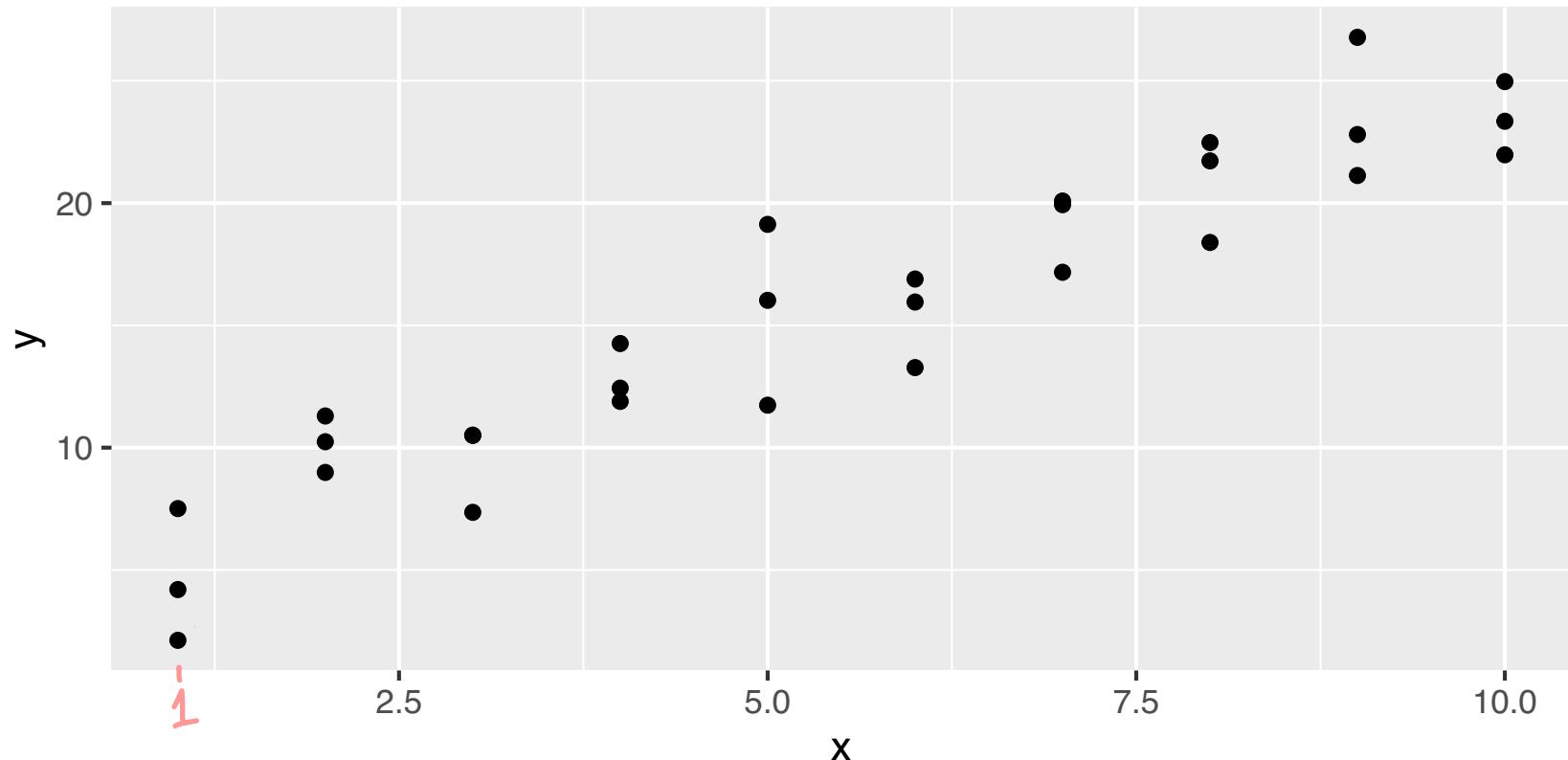
- Training  
60%  
1 or 2 models



- . work with simulated data set
  - package called "model r"
- . **sim1** - data set that we will use.

# Linear Models

- Sim1 - data set



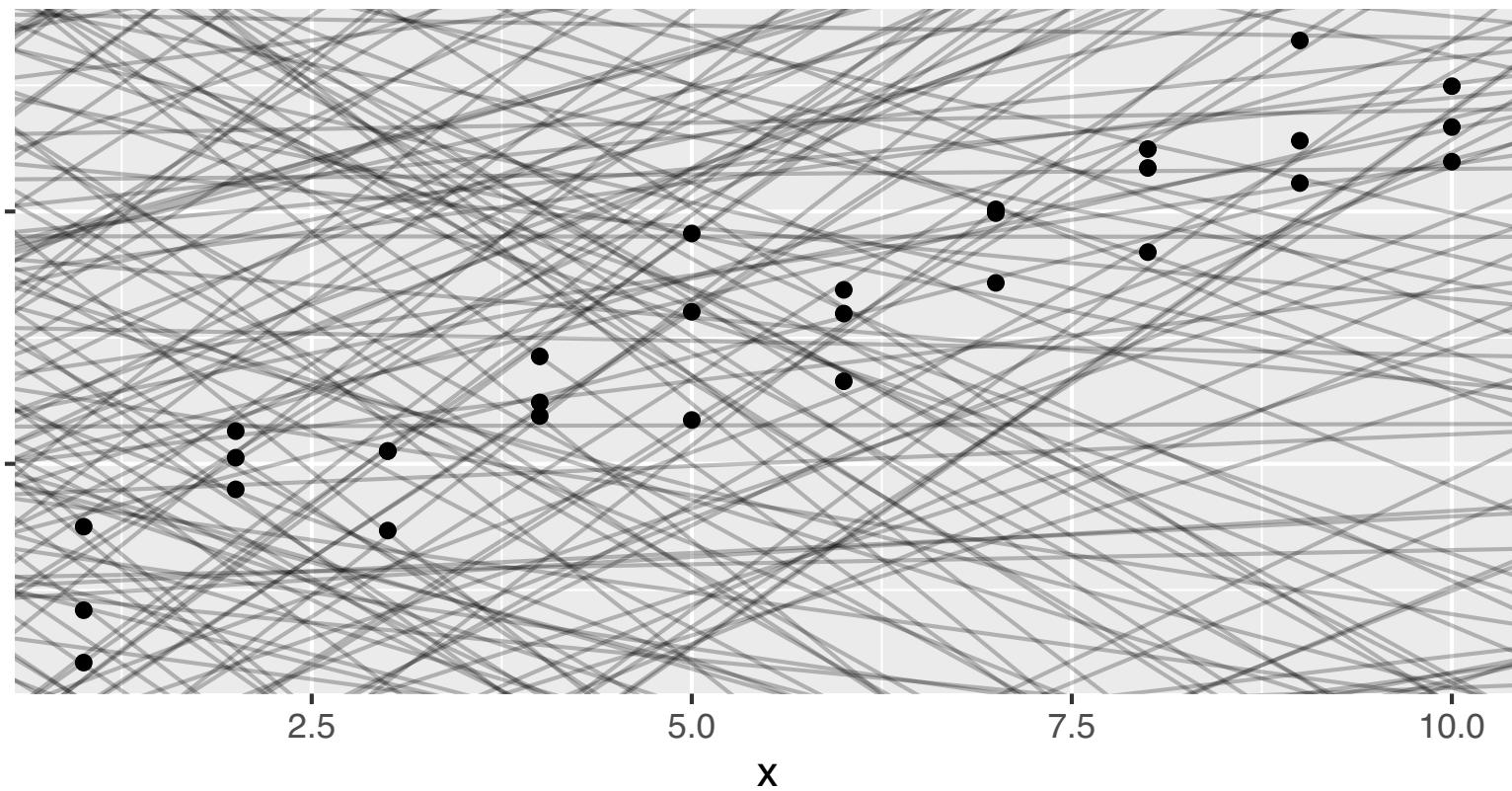
• There is  
a  
pattern

- "linear relationship"

# Linear Models

$a_1 \equiv \text{intercept} - 250$  uniformly chosen r.v. in  $[-20, 40]$

$a_2 \equiv \text{slope} - 250$  uniformly chosen r.v. in  $[-5, 5]$



} Plotted  
the  
250

lines

&

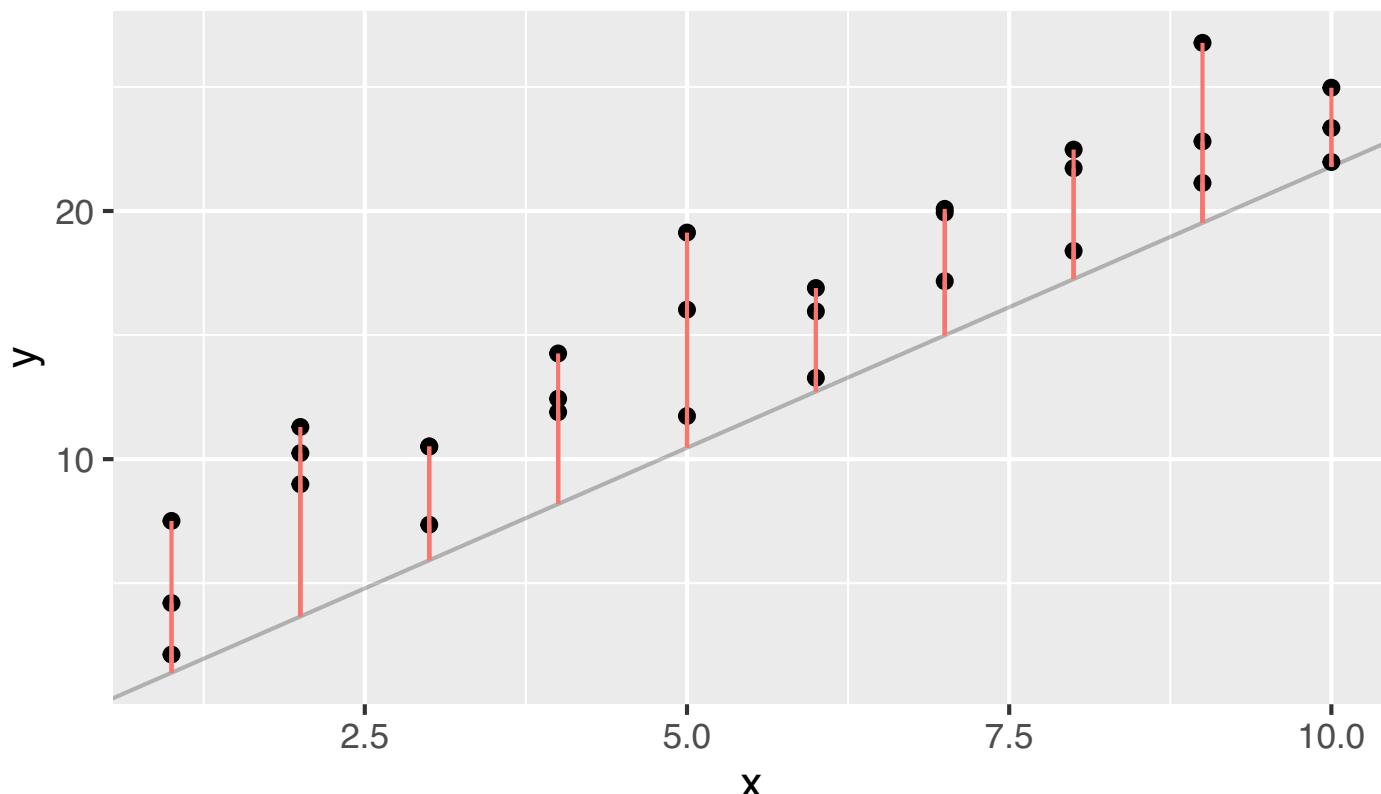
data  
set

250 linear models

# Linear Models : Understand Errors

- Chosen  $l30^{th}$  line
- plotted the line

[1] 130



colour  
— resid

Residuals  
} are the  
distance of  
the points  
from the  
line.

Criteria for good :- minimize the "sum of distances" (X)

The "sum of distances" (X)  
The "sum of (distances)" (Third)

The "sum of (distances)<sup>2</sup>" ✓ least square"

# Linear Models : find the best possible line

```
> model1 <- function(a, data) {  
+   a[1] + data$x * a[2]  
+ }  
  
> #measure distance  
  
> measure_distance <- function(mod, data) {  
+   diff <- data$y - model1(mod, data)  
+   sqrt(mean(diff ^ 2))  
+ }  
  
> best <- optim(c(0, 0), measure_distance, data = sim1)  
> best$par  
[1] 4.222248 2.051204  
  
> #> [1] 4.222248 2.051204  
  
> measure_distance(c(best$par[1], best$par[2]), sim1)  
[1] 2.128181
```

Data:  $\{(x_i, y_i) : 1 \leq i \leq n\}$

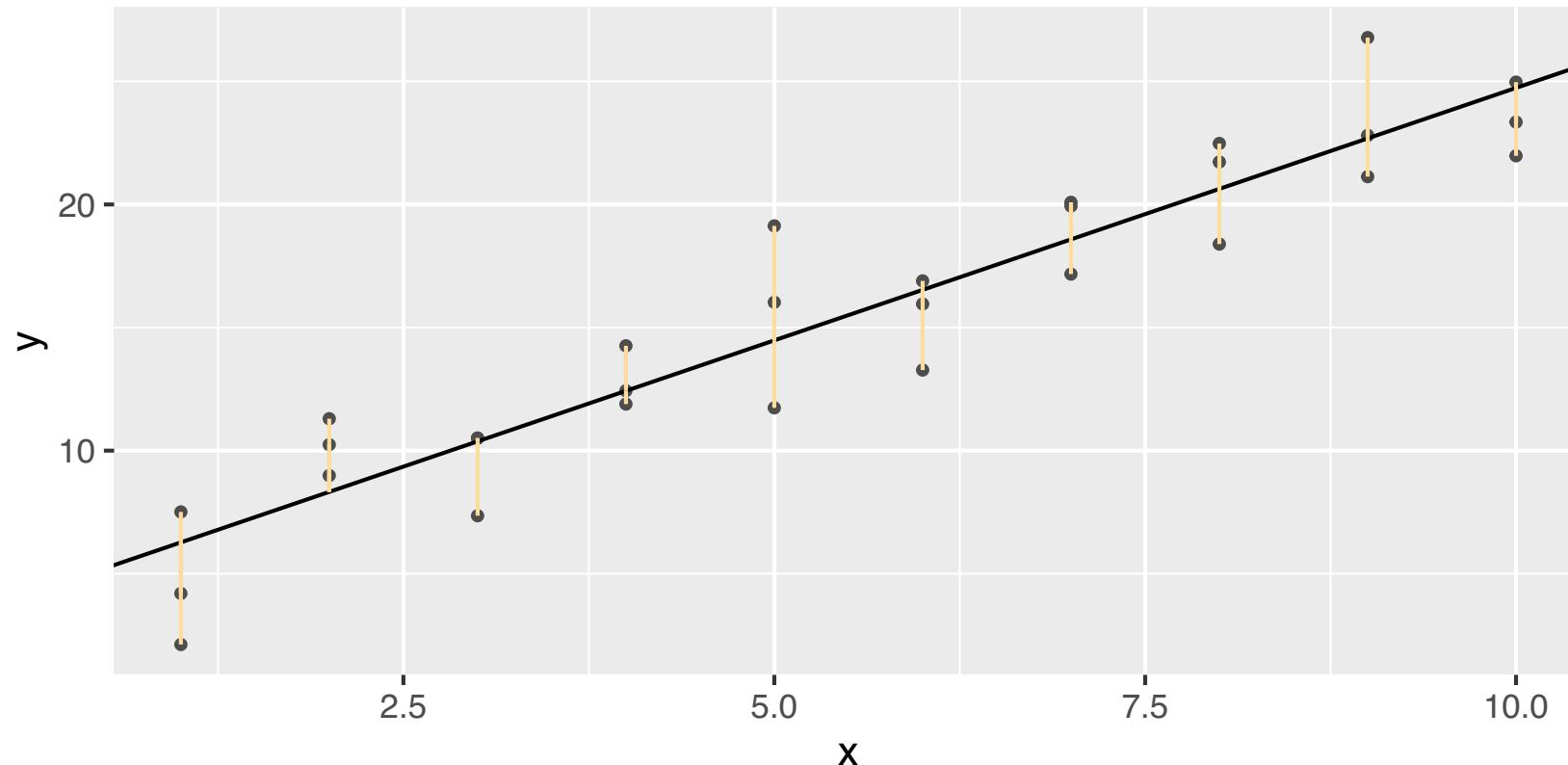
$$\text{measure-distance}(a_1, a_2) = \left[ \frac{1}{n} \sum_{i=1}^n [y_i - (a_1 + a_2 x_i)]^2 \right]^{\frac{1}{2}}$$

Find the  $(a_1, a_2)$  that minimizes:

measure-distance  $(a_1, a_2)$

## Linear Models : best line

Plotted the line : intercept:- 4.22 , Slope:- 2.05



Seems to provide a good approximation

. How to judge the model?

/

. Is this the best possible solution under the model?

- . work with simulated data set
  - package called "model r"
- .  $\text{Sim1}$  - data set that we will use.

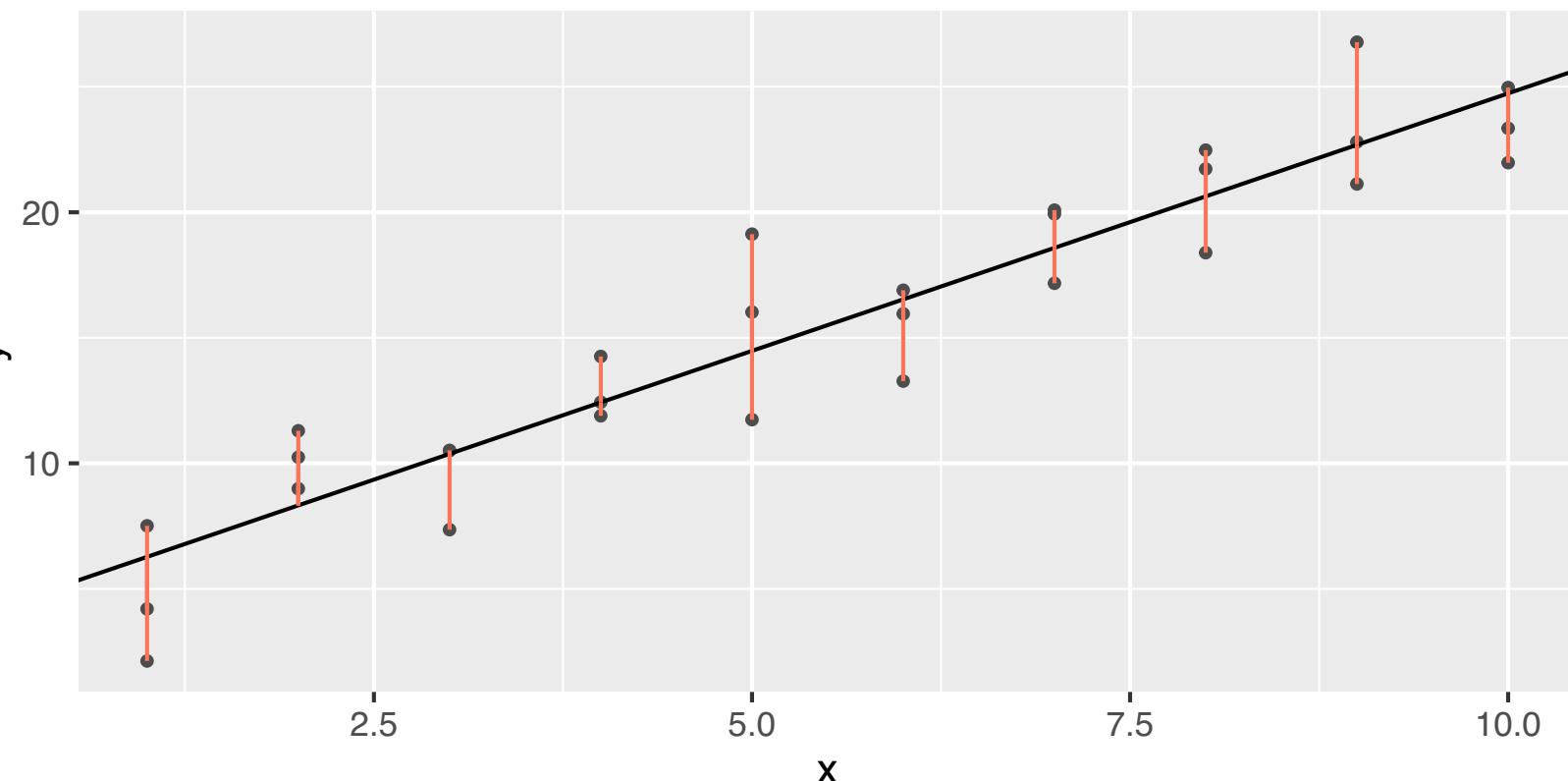
- .  $\{(x_i, y_i) . 1 \leq i \leq n\} = \text{Data}$
- .  $f(a_1, a_2) = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - (a_1 + a_2 x_i))^2}$
- . Found  $a_1, a_2$  that minimised  $f$ 
  - least square line

To do  
thus  
entire  
procedure

R- in built function = "lm"

## Linear Models : best line

(Intercept) x  
4.220822 2.051533



Using lm function

- check week 2 v code

to get precise

Syntax.

Forbes data :- atmospheric pressure & boiling point of water

~ 1857

Could this ←

- replace the fragile Barometer

• collected data

Also in Scotland

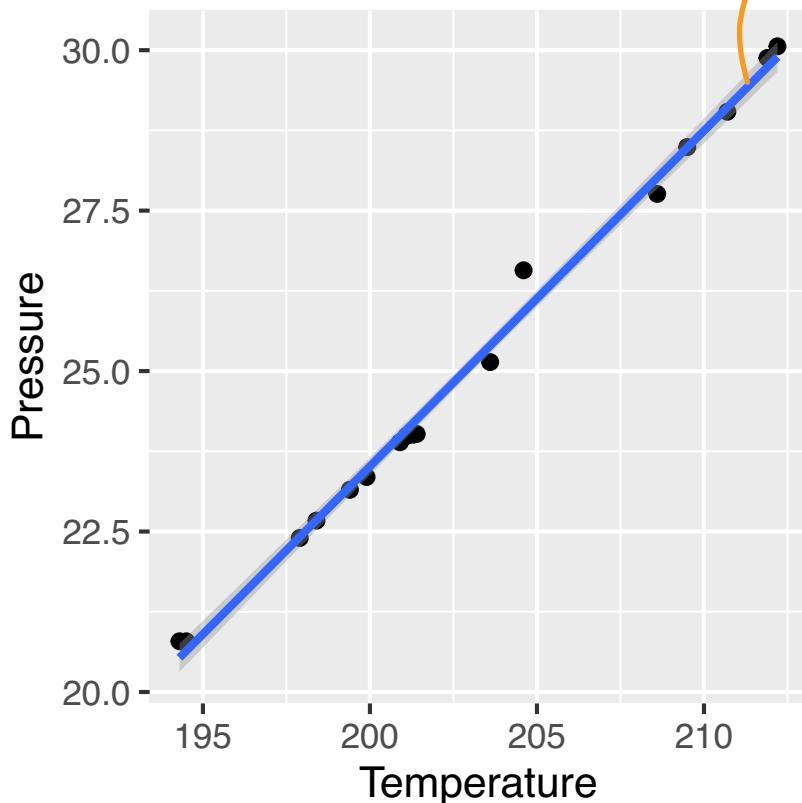
• 17 data points

• Plot the data point

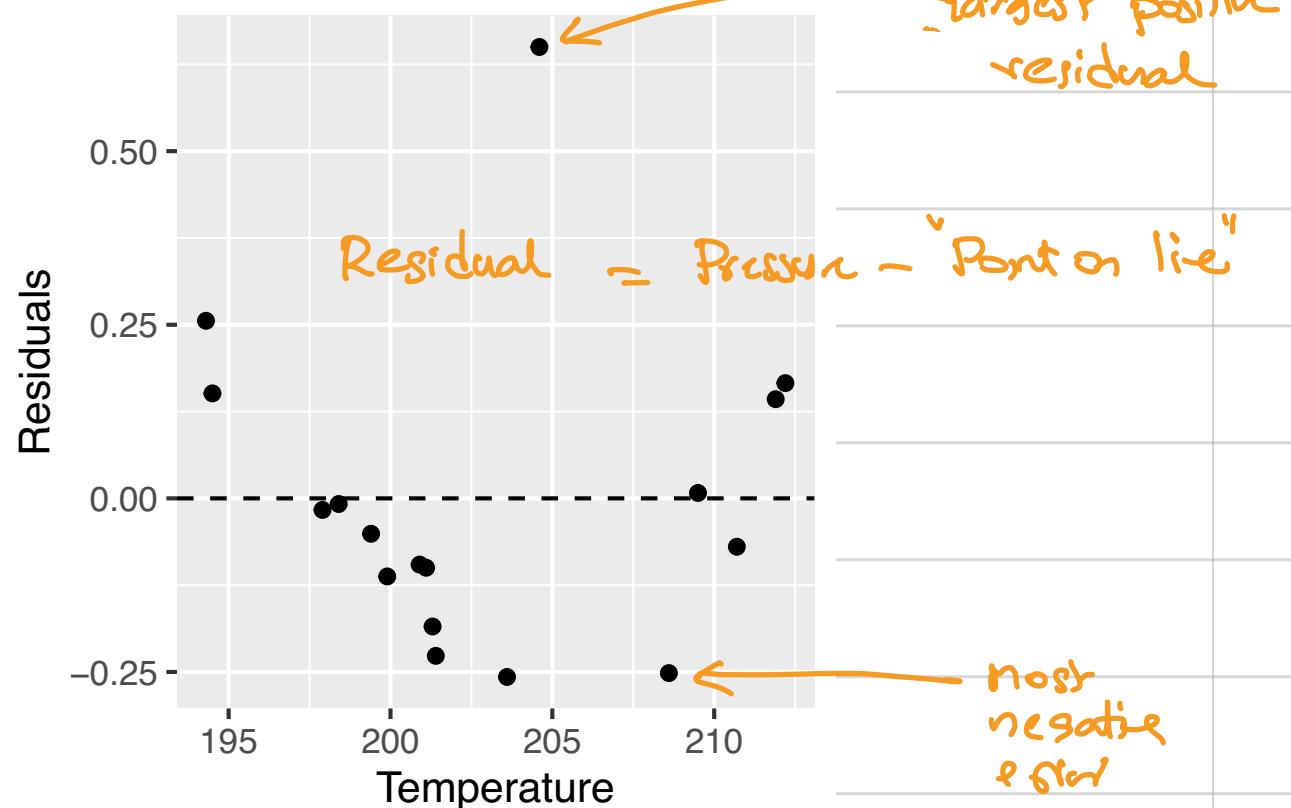
- ln function to find  
the best line

# Scatterplots and Regression

Forbes data

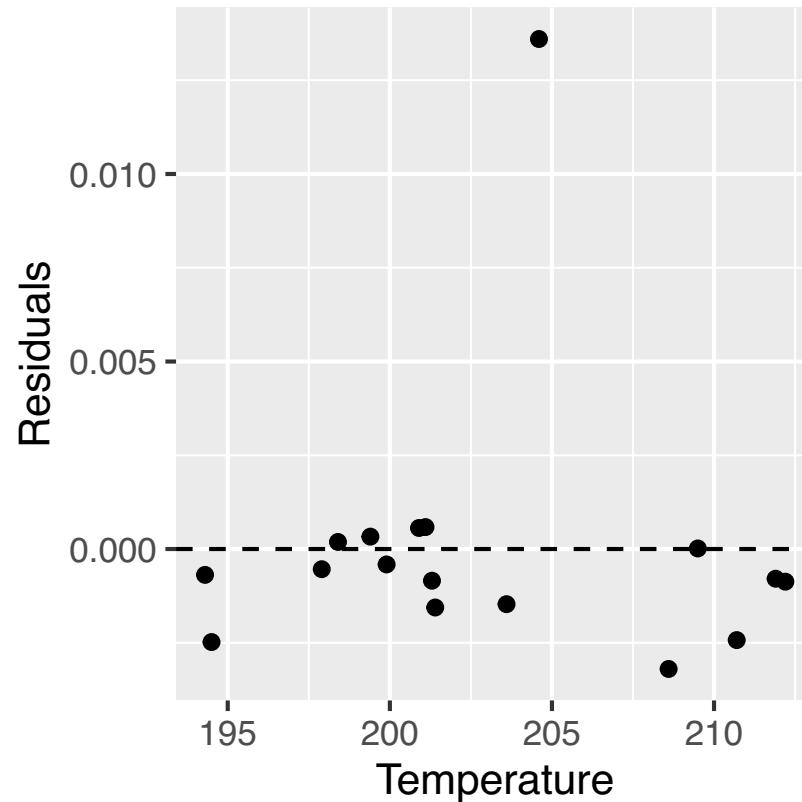
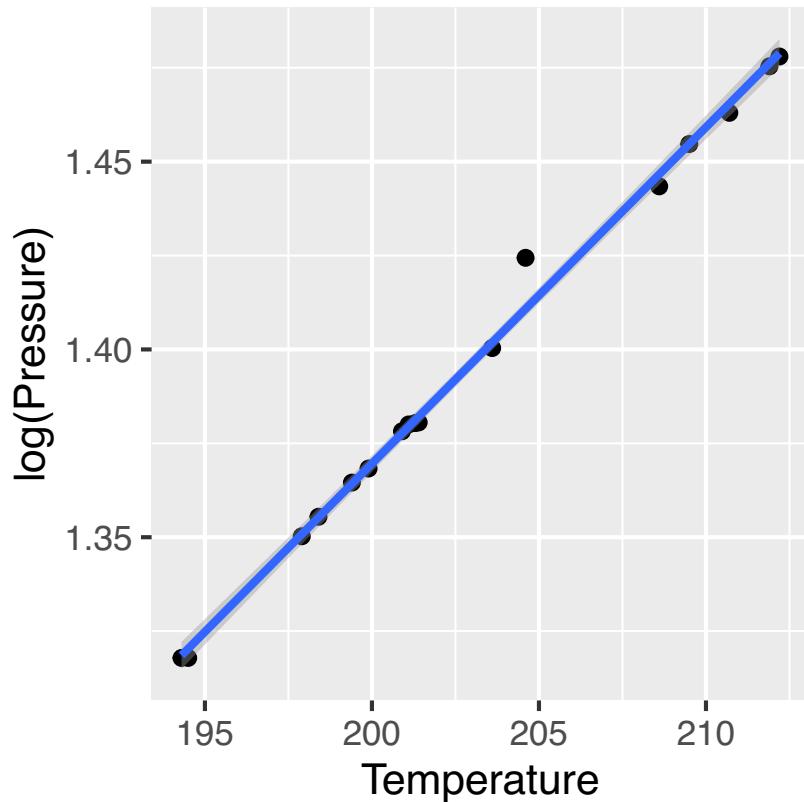


linear - best least square line



# Scatterplots and Regression

Forbes : physical Theory  $\sim \log(\text{Pressure}) = \text{linearly}$   
related to temperature



- length at Age for smallmouth Bass (fish dataset)

mm

└ Lake in Minnesota U.S.A.  
~ 1991

fish have annular  
ring like trees

← Age:  $\leq 7$   
→ (use to compute)

- graph = cannot use the line to predict age  
from length.

- Predicting the weather

Sept - Dec ..... → Predict Jan - June  
Snowfall ?

93 years data

- early & late are not related!

