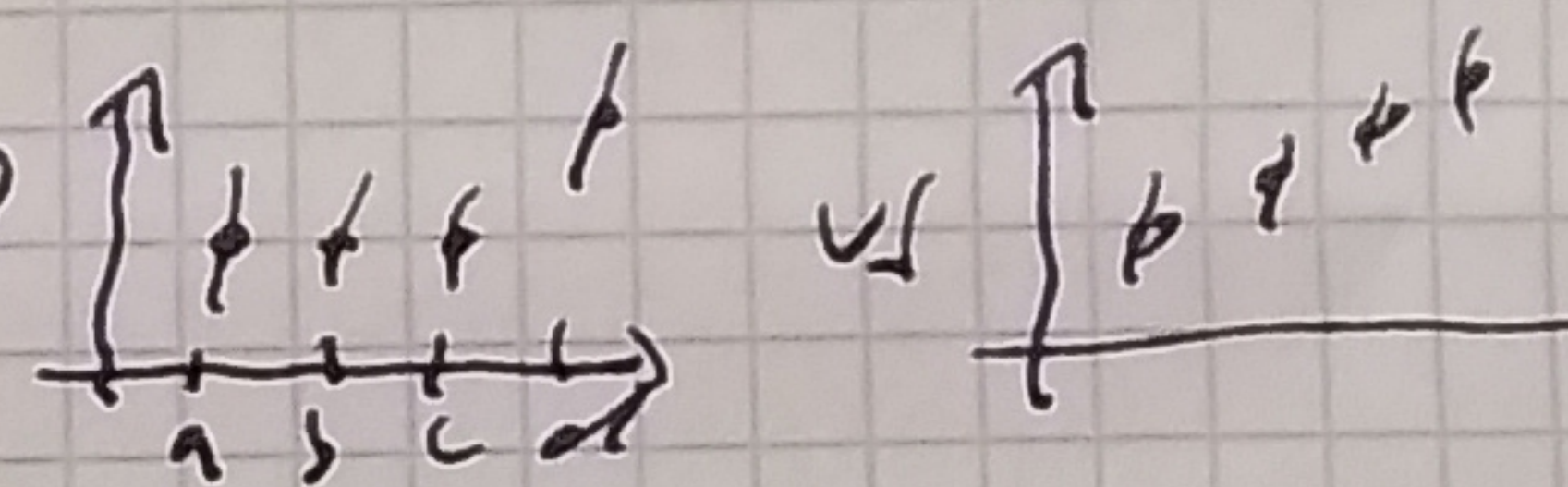
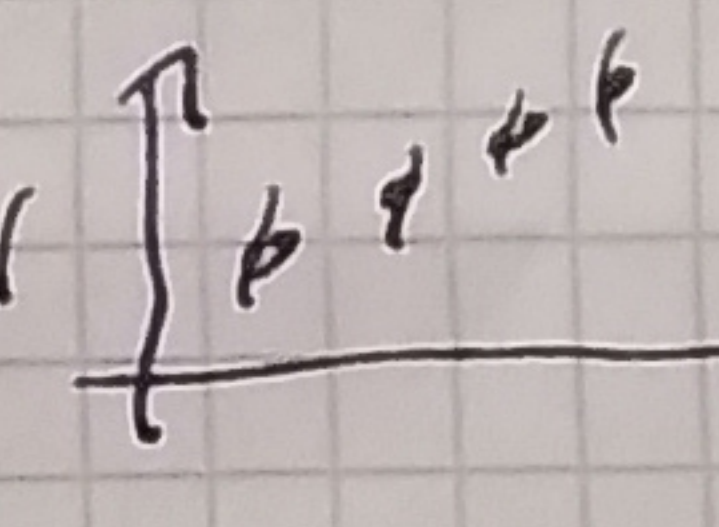
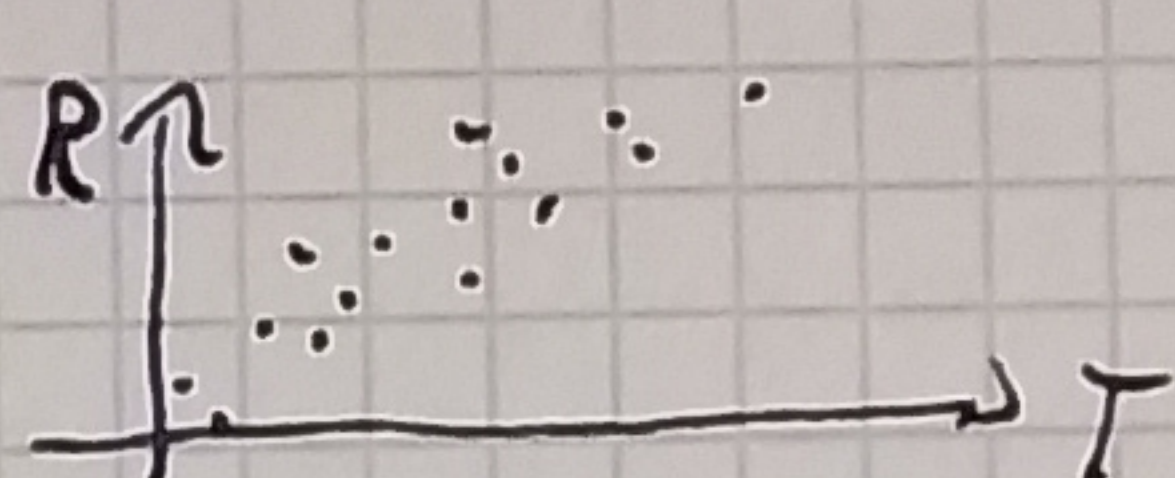


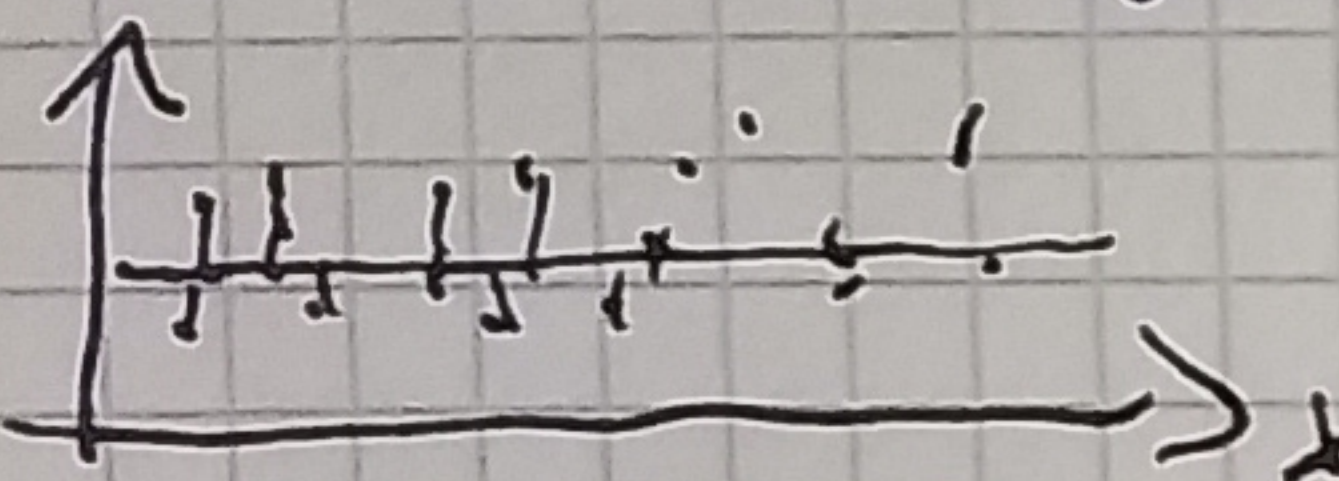
10. Linear Regression

Ho fails in ANOVA - what now?  vs 
We have (y_i, x_i) instead of groups 

⇒ Regression.

10.1 Constant Function

We have data (y_i, x_i) and want to fit it with a function $\overset{\text{MODEL:}}{g(x)} = a$
 $\hat{=}$ a constant

 Our "estimator" for $\hat{y} = \hat{a}$

We define the distance to the estimator $r_i = y_i - \hat{y}(x_i)$

~~residual~~ $\hat{=}$ residual $= y_i - \hat{a}$

Minimize square deviation from model:

$$SSE = \sum_{i=1}^n (y_i - \hat{y})^2 = \sum_{i=1}^n (y_i - \hat{a})^2$$

$$\frac{\partial SSE(a)}{\partial a} = \sum_{i=1}^n 2(y_i - a) \cdot (-1) \stackrel{!}{=} 0$$

$$= -2 \sum_{i=1}^n y_i + 2na \stackrel{!}{=} 0$$

$$= -2n\bar{y} + 2na \stackrel{!}{=} 0 \Rightarrow \hat{a} = \bar{y}$$

\hat{a} is an extremum check $\frac{\partial^2 SSE(a)}{\partial a^2} = +2n > 0$ always...

⇒ $\hat{a} = \bar{y}$ minimizes the residual ⇒ "best fit"

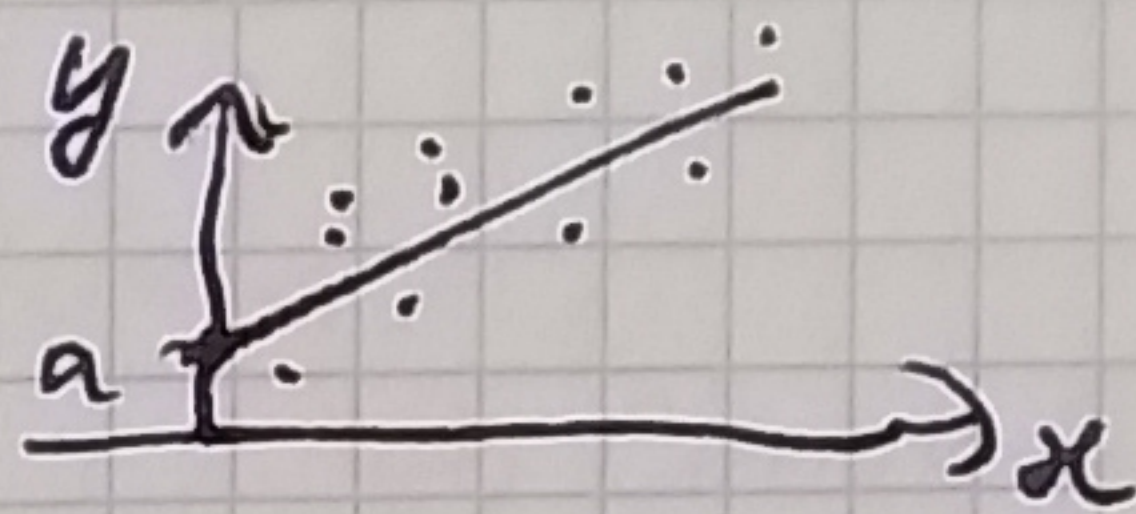
→ property of the arithmetic mean.

10.2 Linear function

$$\hat{y} = a + bx \rightarrow \text{find } \hat{a} \text{ and } \hat{b} \text{ that minimize the residual}^2 / SSE$$

\uparrow intercept \uparrow slope

$$r_i = y_i - (\hat{a} + \hat{b}x_i)$$



$$\frac{\partial SSE(a, b)}{\partial a} = \frac{\partial}{\partial a} \left(\sum_i (y_i - a - bx_i)^2 \right)$$

$$= 2 \sum_i (y_i - a - bx_i) (-1) \stackrel{!}{=} 0$$

~~$$= -2(n\bar{y} - na - nb\bar{x}) \stackrel{!}{=} 0$$~~

$$\boxed{\hat{a} = \bar{y} - \hat{b}\bar{x}} \quad \text{I}$$

$$\text{check } \frac{\partial^2 SSE}{\partial a^2} = 2n > 0$$

$$\frac{\partial SSE}{\partial b} = 2 \sum_i (y_i - a - bx_i) (-x_i) \stackrel{!}{=} 0$$

$$\Rightarrow \sum_i (x_i y_i) - na\bar{x} - b \sum_i x_i^2 \stackrel{!}{=} 0$$

$$\text{I} \Rightarrow \sum_i x_i y_i - n(\bar{y} - b\bar{x})\bar{x} - b \sum_i x_i^2 \stackrel{!}{=} 0$$

$$\Rightarrow \boxed{\hat{b} = \frac{\sum_i (x_i y_i) - n\bar{x}\bar{y}}{\sum_i x_i^2 - n\bar{x}^2} = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2}}$$

$$\text{check } \frac{\partial^2 SSE}{\partial b^2} = 2 \sum_i x_i^2 \geq 0$$

$$\text{check: } \frac{\partial^2 SSE}{\partial b^2} = 2 \sum_i x_i^2 \geq 0$$

10.3 ANOVA for (one - slope) linear regression:

	ν	SS	MS	F
Model	1	$SS_M = \sum_i (\hat{y}(x_i) - \bar{y})^2$	SS_M / ν_M	MS_M / MSE
Error/ Residual	$n-2$	$SS_E = \sum_i (y_i - \hat{y}(x_i))^2$	SS_E / ν_E	
Total	$n-1$	$SS_T = \sum_i (y_i - \bar{y})^2$		

F test: if $F_{measured} > F_{var, e} (P)$ ^{$\alpha = 0.05$}

=> reject H_0 : "all data does not depend on x "

(all x "groups" are the same)

=> conclude there is some correlation.

Goodness of fit vs

$$R^2 = \frac{SS_M}{SS_T} = 1 - \frac{SS_E}{SS_T}$$

$0 \leq R^2 \leq 1$
 ↑ meaningles ↑ perfect fit

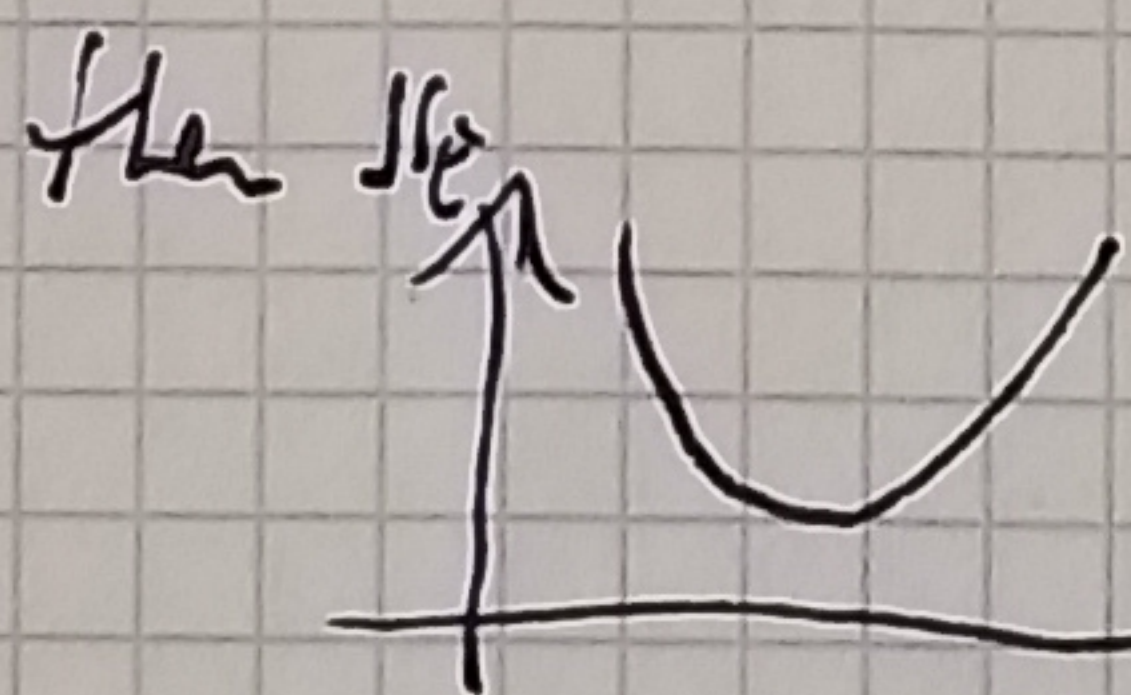
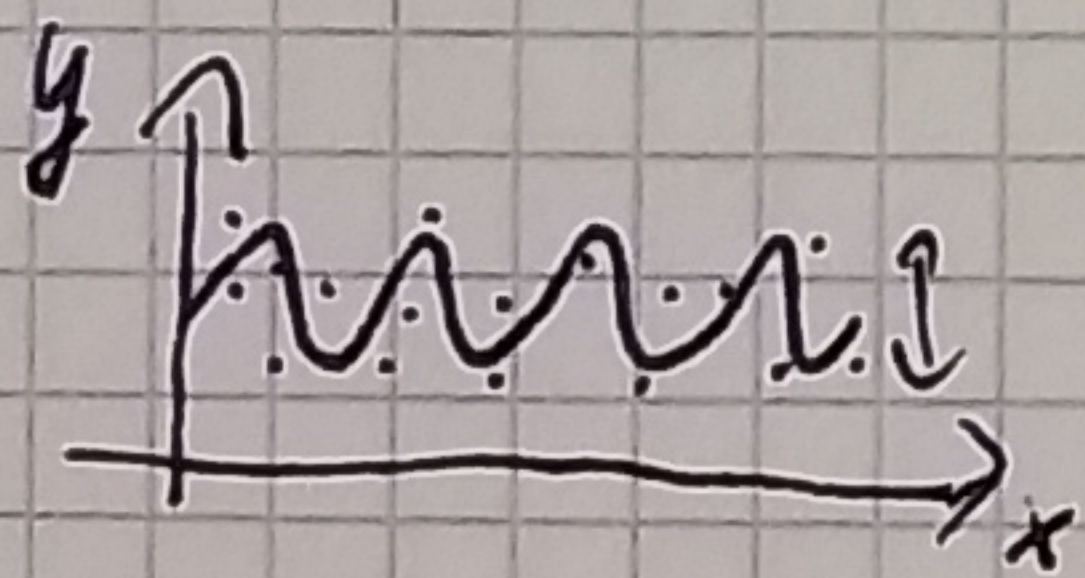
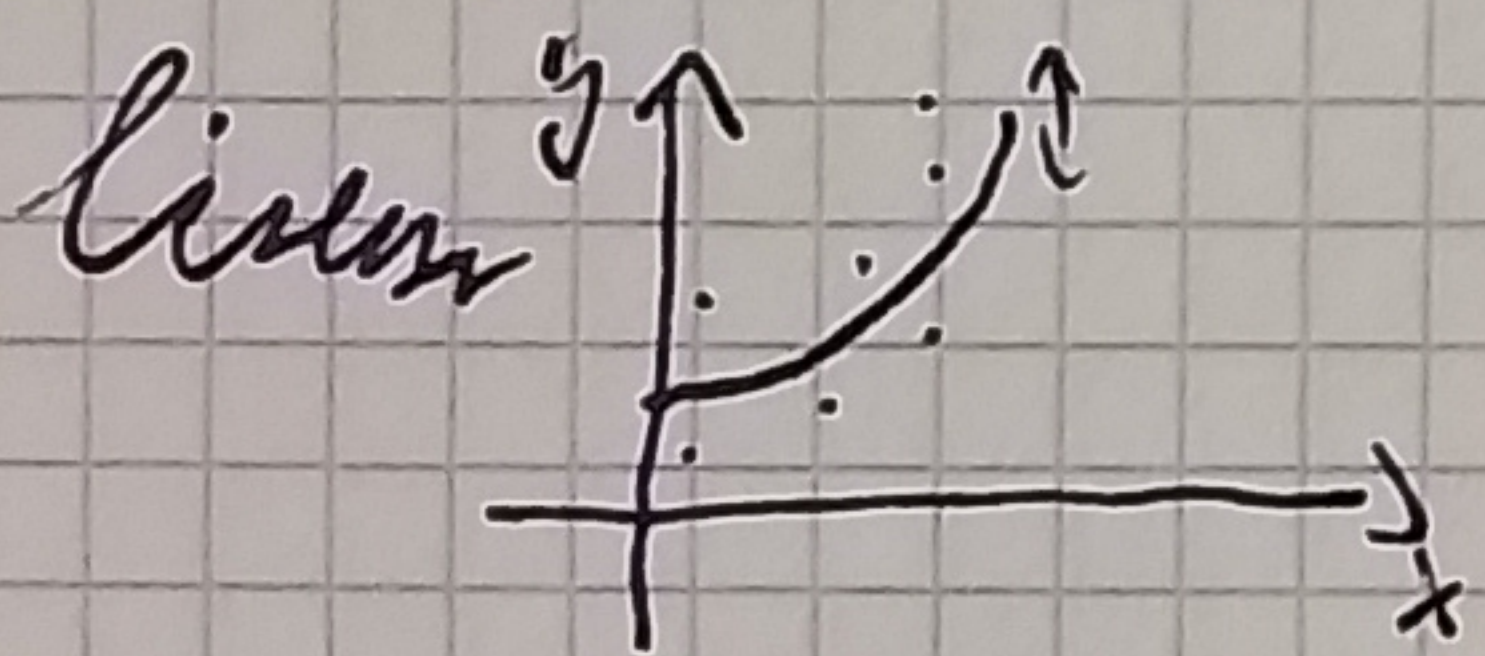
→ slides

10.4 What is "linear"?

→ the parameters a, b .

$y = a + bx$ ✓ $y = a + b \underbrace{x^2}$ ✓ $y = a + b \underbrace{\sin(x)}$ ✓

$y = a + \sin(bx)$ ~~linear~~ → Taylor expansion $y = a + b_1 x + b_2 x^2 + \dots$
 → more parameters.



local optimum = global optimum.

→ deterministic

