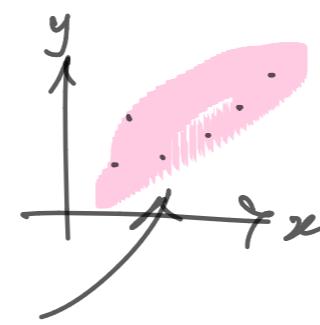
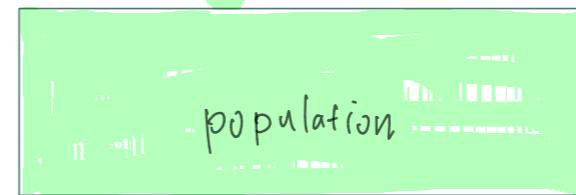


Bivariate data:  $(x, y)$

$\nearrow$  independent variable  
 $\nwarrow$  dependent variable  
(explanatory) (regressor)



Module 2: Modelling Data



Correlation Coefficient ( $r$ )

measures linear association =  
how tightly points are  
clustered around a line

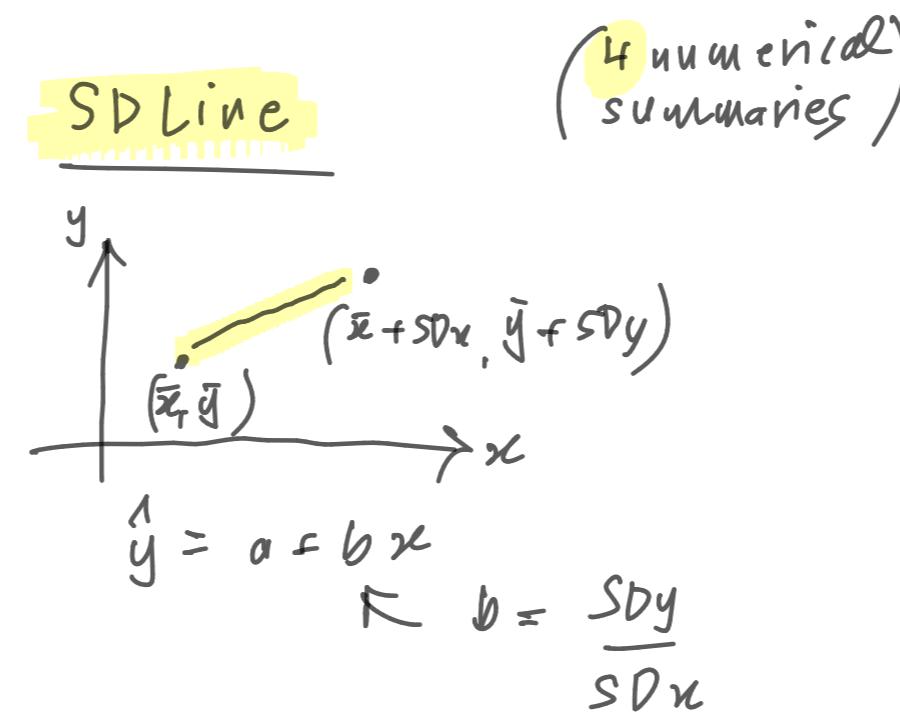
$$r_{\text{pop}} = \frac{1}{n} \sum_{i=1}^n \frac{x_i - \bar{x}}{SD_x} \cdot \frac{y_i - \bar{y}}{SD_y}$$

$= r_{\text{sample}}$

$$-1 \leq r \leq 1$$

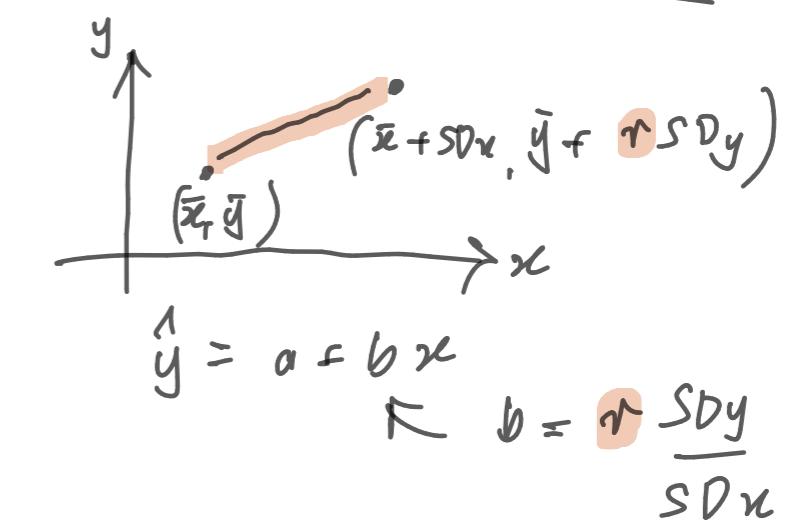
classic mistakes with  $r$   
eg  $r=0 \neq$  no association

### SD Line



(4 numerical summaries)

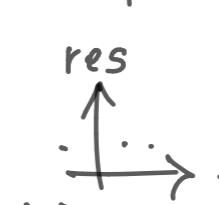
### Regression Line



(5 numerical)  
Summaries

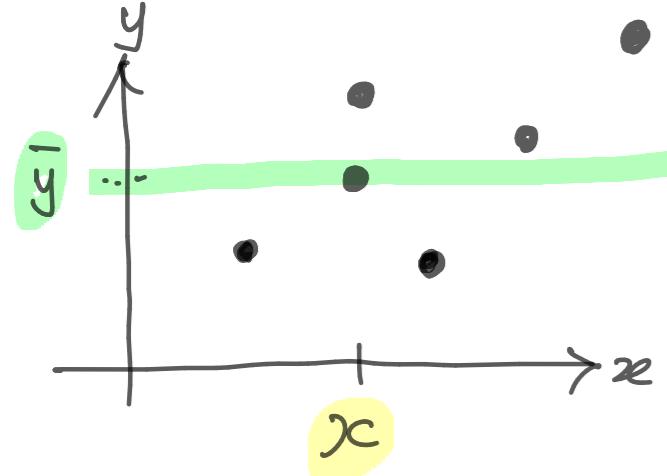
$$\bar{x} \bar{y} SD_x SD_y r$$

(for  $r > 0$ )

Model	Variable	When?	What?	Diagnostics	Predictions
Linear	2 Aulant	$V_2 \rightarrow V_1$ looks linear	$\begin{array}{l} \textcircled{1} \text{ linear regression line } \\ \hat{y} = a + b x \end{array}$ $\begin{array}{l} \textcircled{2} \text{ correlation coefficient } \\ -1 \leq r \leq 1 \end{array}$	$\begin{array}{l} \textcircled{1} \text{ Does the scatter plot look linear? } \\ \textcircled{2} \text{ Does the residual plot look random? } \end{array}$  	$\begin{array}{l} \textcircled{1} \text{ Baseline} \\ \textcircled{2} \text{ strip} \\ \textcircled{3} \text{ Regression line} \\ \textcircled{4} \text{ Percentile Rules} \\ \textcircled{5} \text{ Strips } \sim \end{array}$

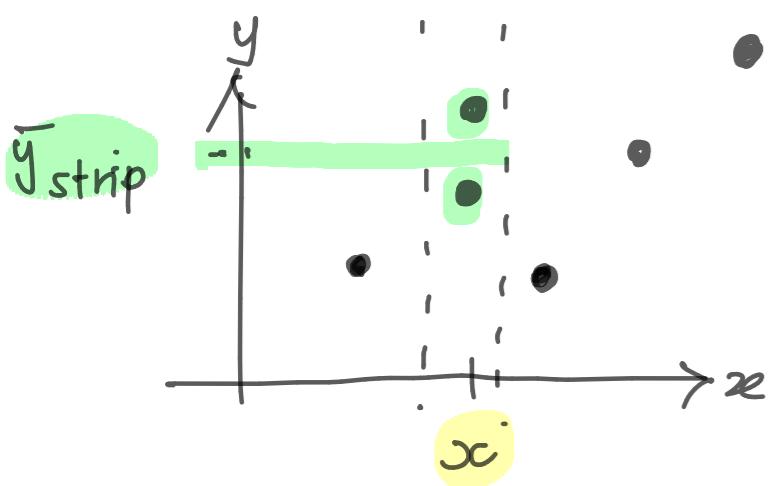
# Predictions based on the linear Model

## ① Base Line



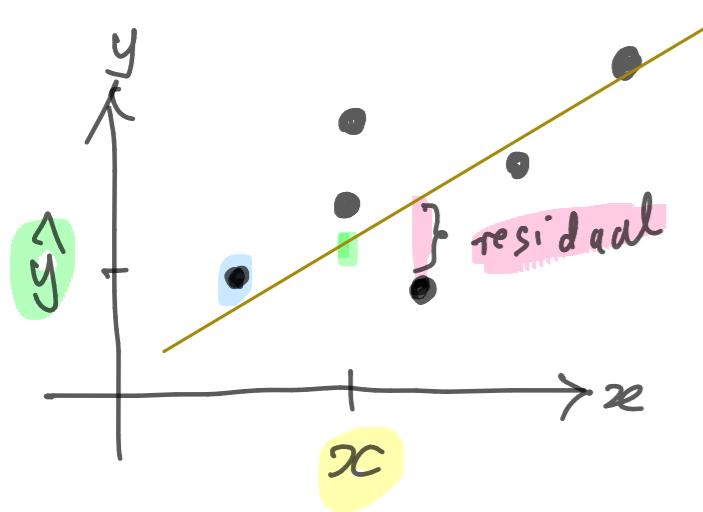
for any  $x$ ,  
we predict  $\bar{y}$

## ② Strip

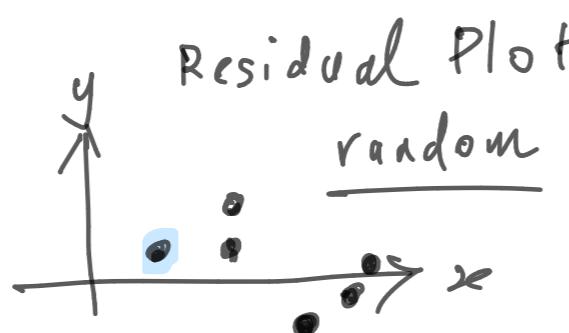


we take the  
average of  
the  $y$ s in the  
 $x$  strip

## ③ Regression Line



$\hat{y}$  = the point  
corresponding to  
 $x$  on the line



$$\begin{aligned}\text{RMS Error} &= \text{"SD of Line"} = \text{RMS of (gaps from the)} \\ &= \sqrt{1 - r^2} \text{ SD}_y \\ &\Rightarrow \text{SD}_y \quad (\text{Baseline})\end{aligned}$$

(DV15)

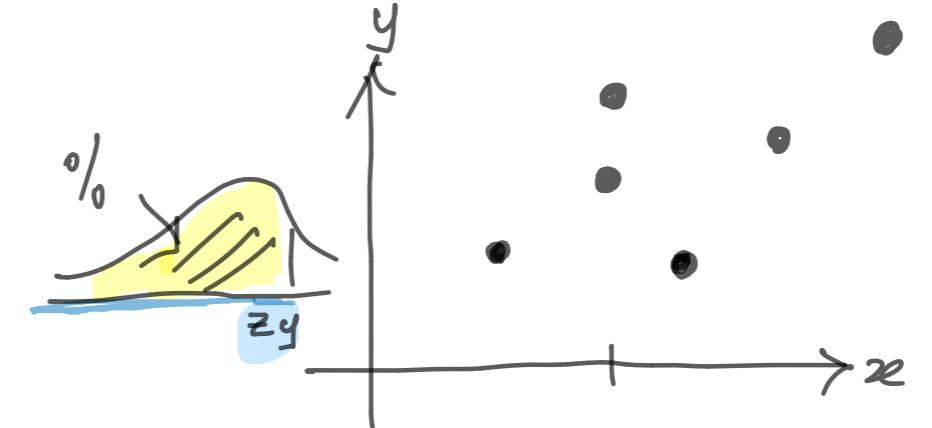
## Steps

1. Convert % in  $x$  direction  
to  $z$  score  
 $z_x = qnorm(\%)$

2. Calculate  $z_y = r z_x$

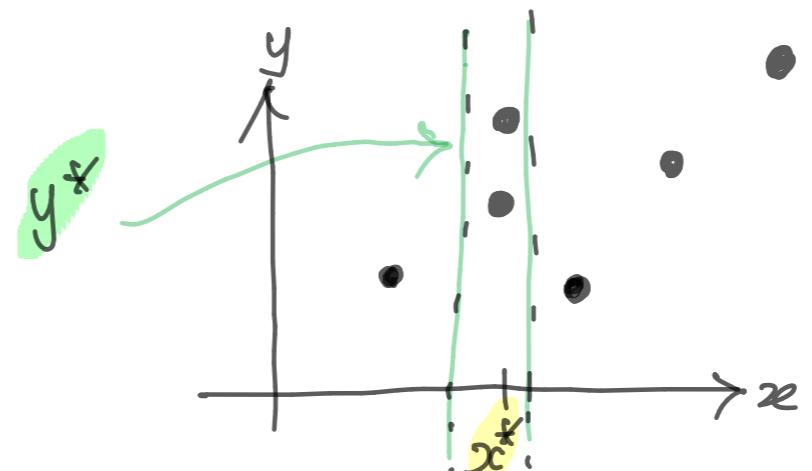
3. Convert  $z$  score in  $y$   
direction to %  
 $pnorm(z_y)$

## ④ Percentile Ranks



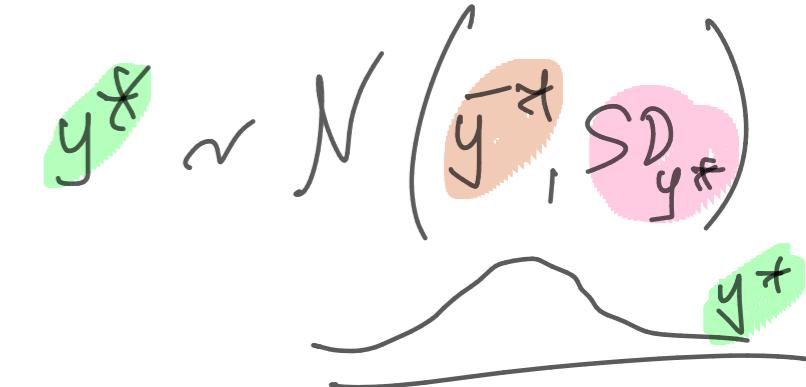
Summary:  $pnorm(r * qnorm(\%))$

## ⑤ Strips (homoscedastic)



$$\bar{y}^* = \bar{y} + r z_{x^*} \text{SD}_y$$

$$\text{SD}_{y^*} = \text{RMSE}$$



(DV16)