

8 Regression Analysis

Purpose: To predict values of Y

Model:

$$Y \approx b_0 + b_1 x_1 + b_2 x_2 + \dots + b_p x_p$$



dependent
variable



independent variables
(Do not mistake for

stochastically independent

variables)

variables)

8.1 Multiple linear Regression

Y, X_1, X_2, \dots, X_p scale leveled variables

Example

Miles - Per - Gallon .sav

$Y = \text{MPG} = \text{miles per gallon}$

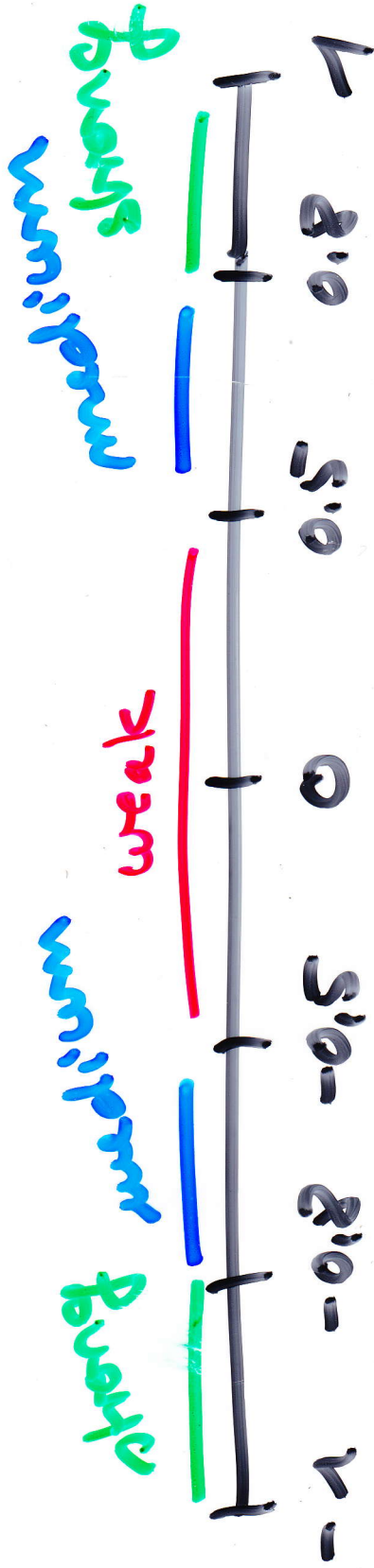
$X_1 = \text{weight (in pounds)}$

$X_2 = \text{horsepower}$

8.1.1 One Independent Variable

$r(\text{MPG, weight}) = -0.825$ stronger than

$r(\text{MPG, horsepower}) = -0.788$ 



Model:

$$\boxed{\text{MPG}} \approx \beta_0 + \beta_1 \cdot$$

$\boxed{\text{weight}}$

a) Scatter plot without regression line



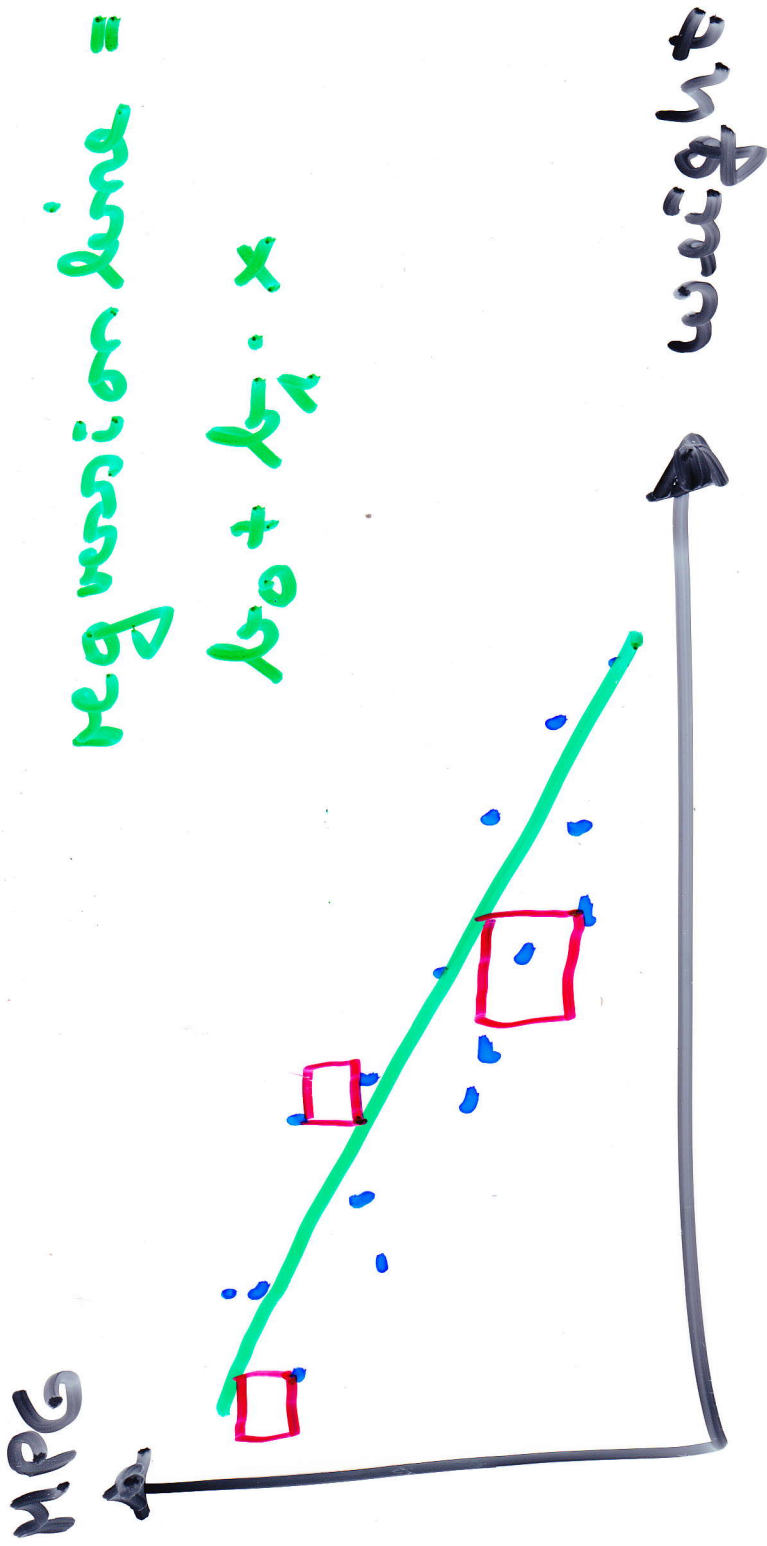
SPSS commands §. 1. 10

with increasing weight the miles per gallon are decreasing

b)

regression line =

$$b_0 + b_1 \cdot x$$



Method of the least squares = the sum of the squared distances of the points

from the line is minimum

$$\hat{b}_0 = 57.797$$

$$\hat{b}_1 = -0.011$$

Scatter plot with regression line

→ regression coefficient if the weight

increases by one pound the miles

per gallon will decrease by 0.011
miles

→ if a car has a weight of 0 pounds

the car will drive 57.797 miles per

gallon = nonsense!

c) linear regression model

commands 8.1.10

d) predict values

4500 pounds weight

Miles per gallon = ?

$$57.797 - 0.011 \cdot 4500 = \text{predicted}$$

miles due to the model = 10.04 miles

$$= \text{PRE} = \text{predicted values}$$

e) 2250 pounds weight

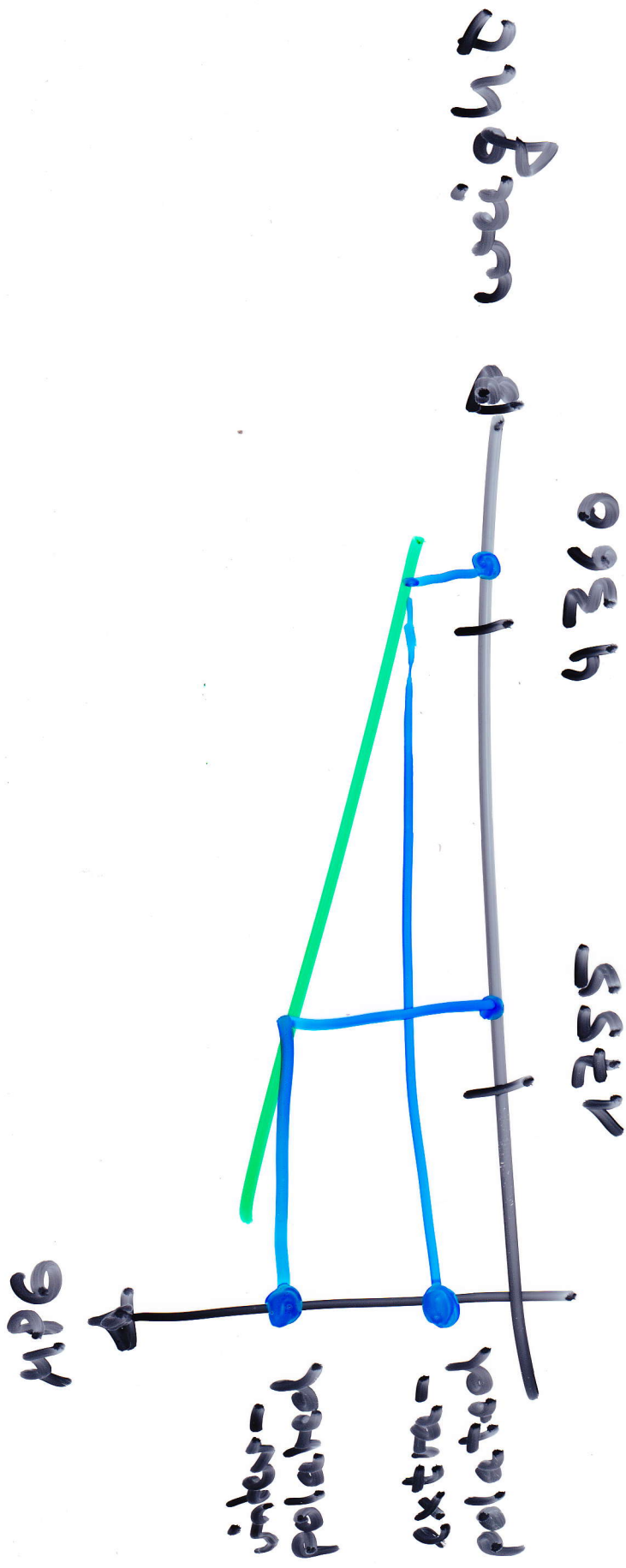
miles per gallon = ?

$$57.797 - 0.11 \cdot 2250 = 33.92 \text{ miles}$$

f) Are the predicted values reliable?

$$x_{\min} = 1755 \text{ weight}$$

$$x_{\max} = 4360 \text{ weight}$$



33.92 ~~10.04~~ miles is an *interpolated value*

because 2250 ∈ [1755; 4360]

10.04 miles is an *extrapolated value*

because 4500 ∉ [1755; 4360]

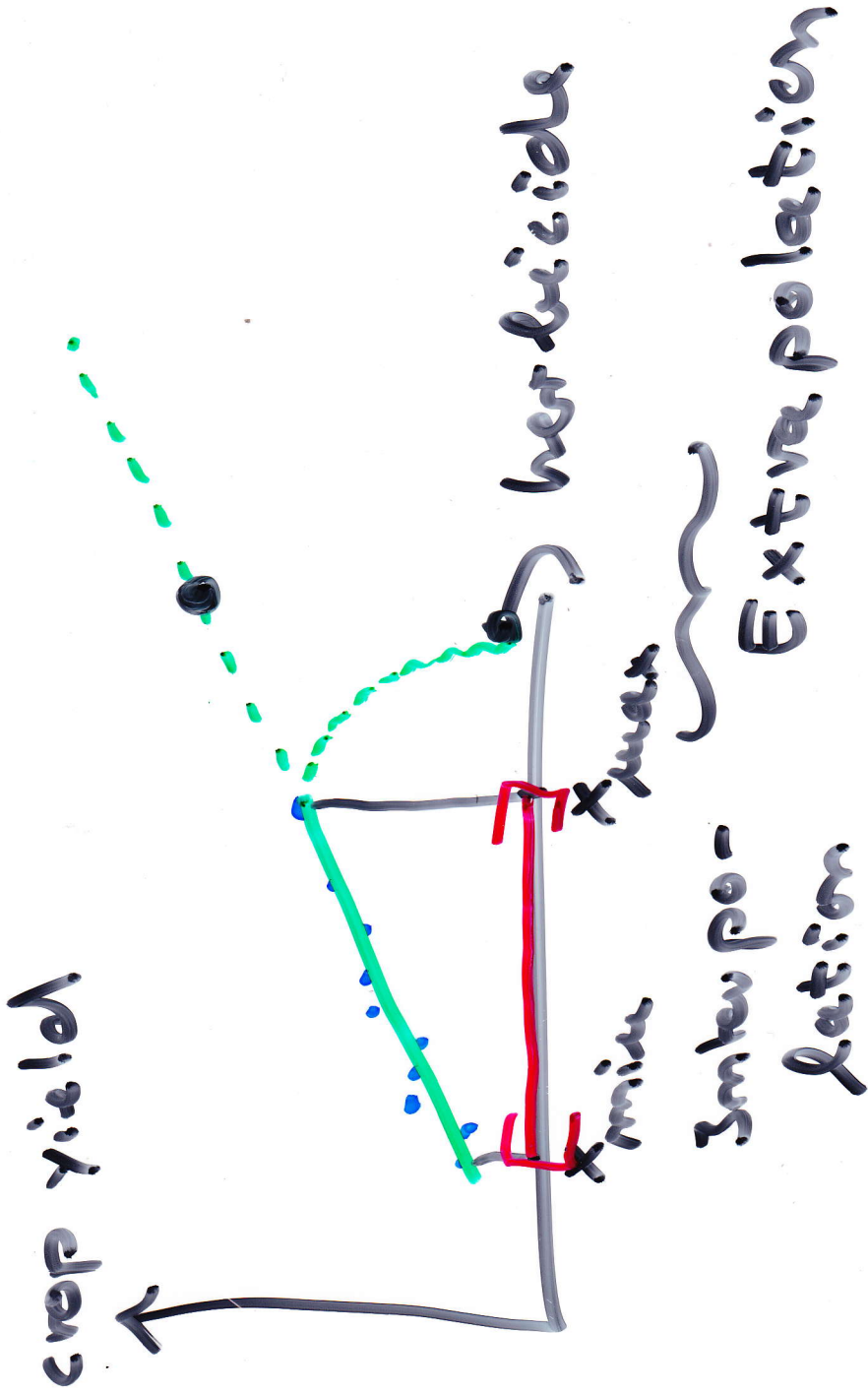
		Correlation	
		weak	medium
intervention	no	no	yes
	no	no	no
extremepolation	no	no	no
	no	no	no

The predicted value of 33.92 miles is reliable, because this value is an

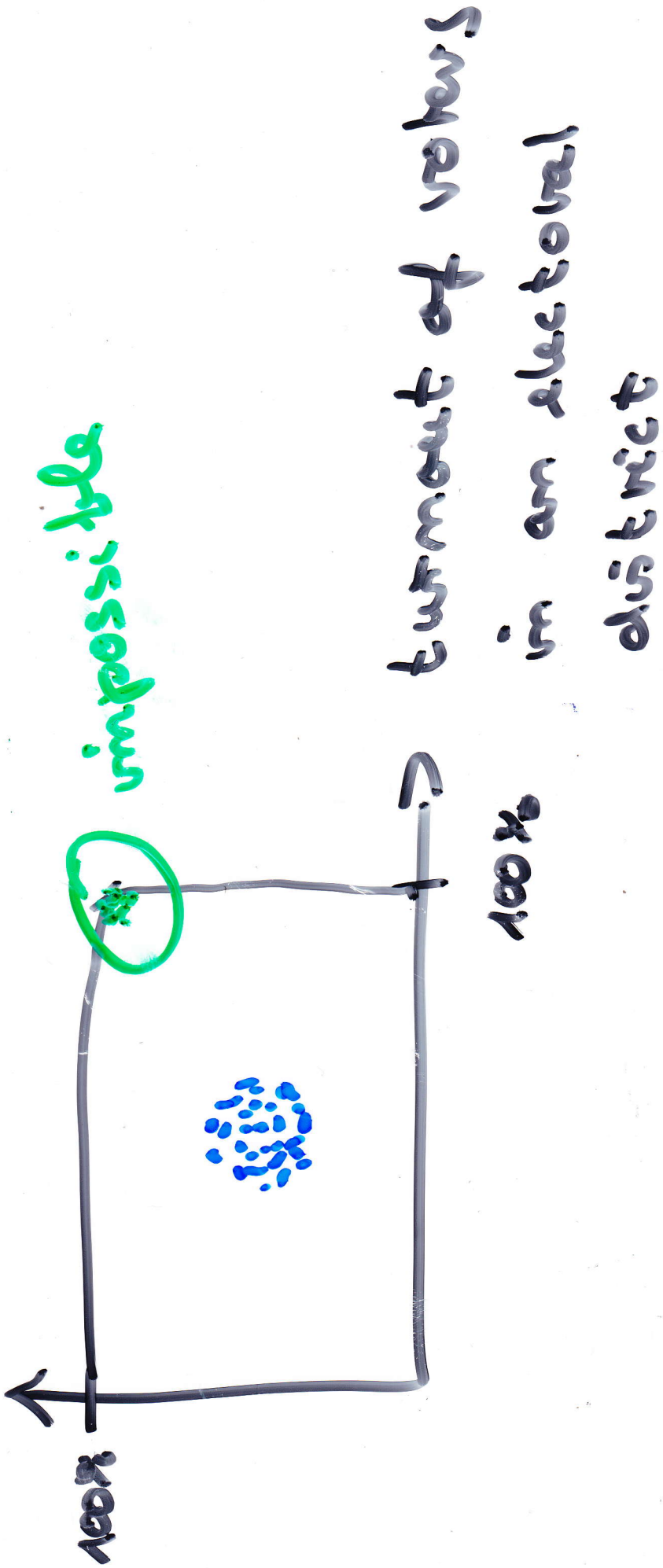
interpolated value and the correlation of -0.825 is strong.

The predicted value of 10.04 miles is not reliable, because this value is an extrapolated value.

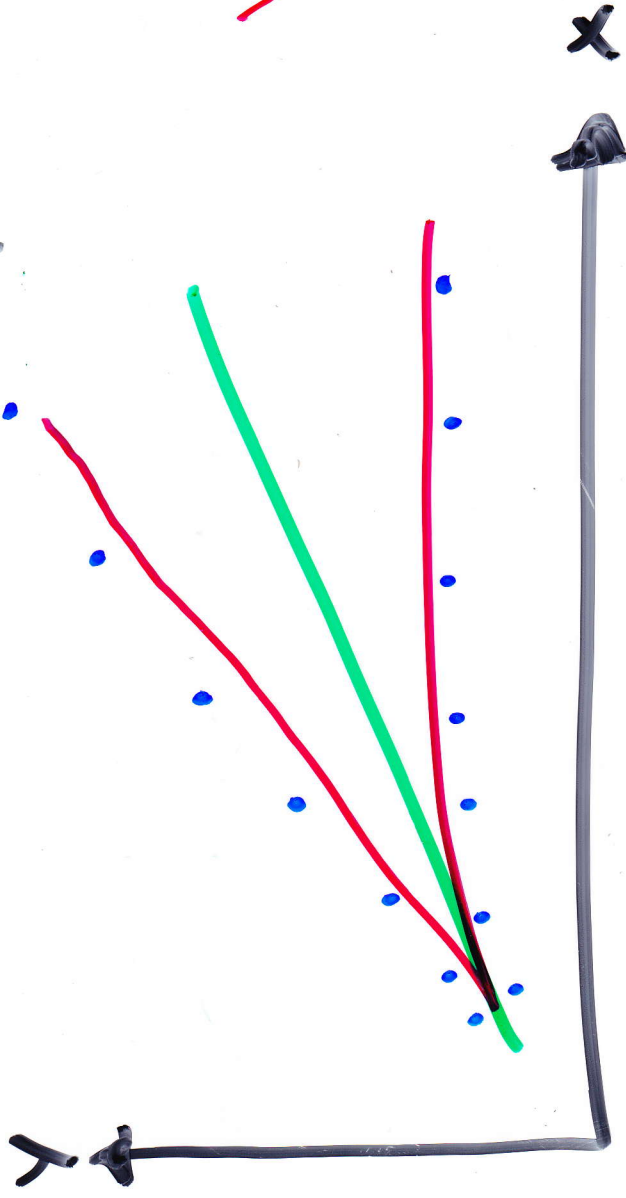
Example



Example: Election fraud
percentage of votes of the winner



Heteroscedasticity



V-formation

heteroscedasticity

In case of heteroscedasticity the values of b_0 and b_1 are inappropriate if they were calculated due to

the method of least squares.

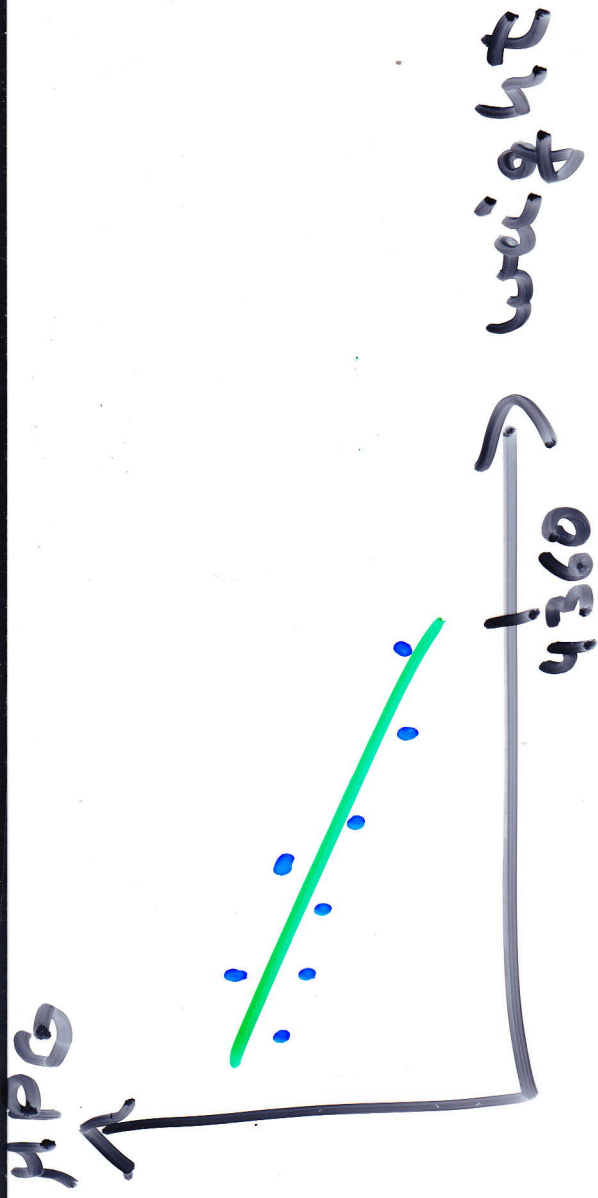


homoscedasticity

Example

Miles - Per - Gallon. sav

$$\boxed{\text{MPG}} \approx \beta_0 + \beta_1 \cdot \boxed{\text{weight}}$$



homoscedasticity

LEV = leverage points

Exercise 8.9

Leverage value of car with 4360 pounds weight = 0.13011 largest effect

the shape of the regression line

8.1.2 Two or more independent variables

Example

Miles - Per - Gallon . sav

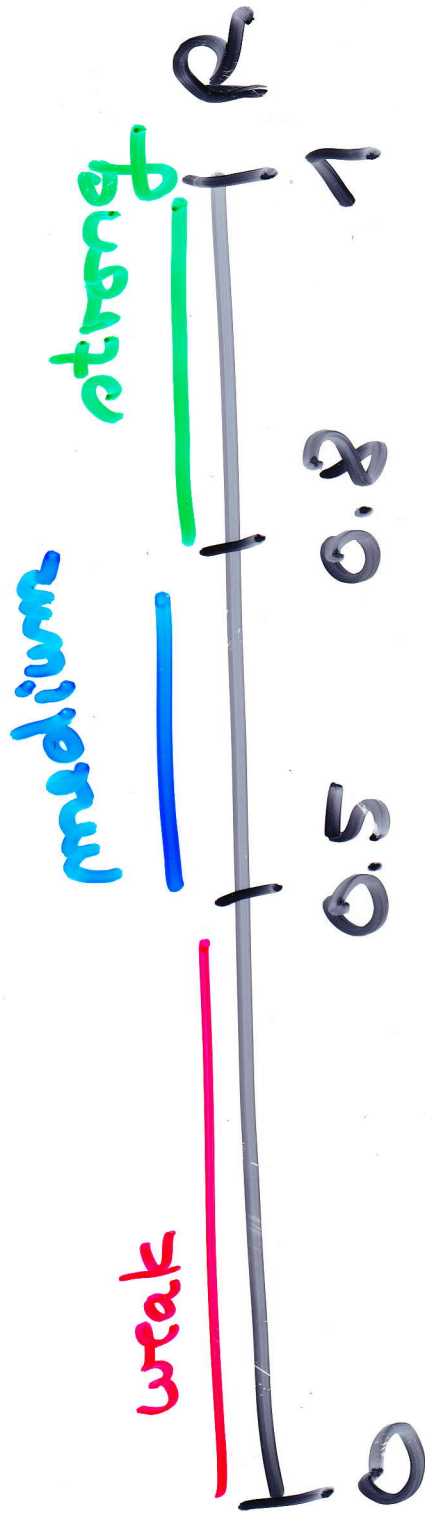
$$\boxed{\text{MPG}} = b_0 + b_1 \cdot \boxed{\text{weight}} + b_2 \cdot \boxed{\text{Horse-power}}$$

$$r(\text{MPG}, \text{weight}) = -0.825$$

$$r(\text{MPG}, \text{Horsepower}) = -0.788 \quad \left. \begin{array}{l} \text{columns,} \\ \text{zero-order} \end{array} \right\}$$

8.1.3 Multiple correlation coefficient

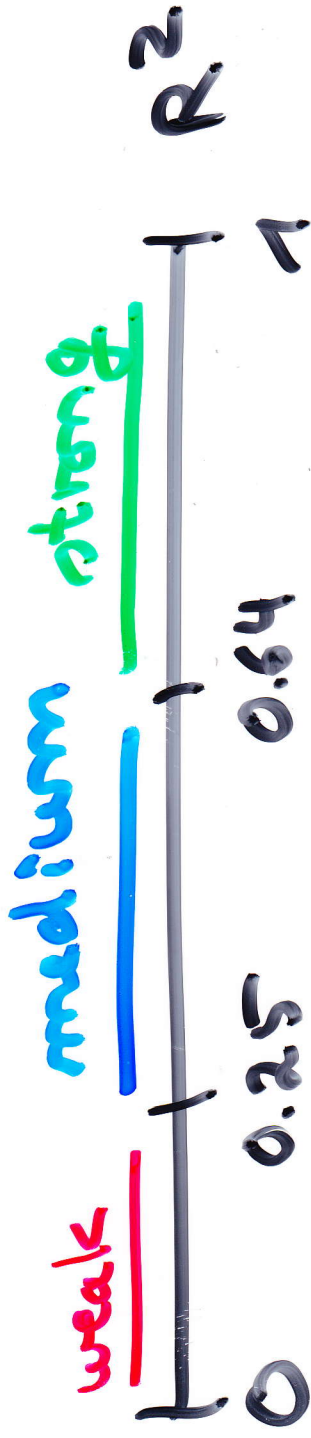
$$R \in [0; 1]$$



$$R(\text{YPG, waight, Horsepower}) = 0.866$$

strong correlation

R^2 = coefficient of determination $\in [0; 1]$



$R^2 = 0.866^2 = 0.749$ strong correlation

8.1.4 Adjusted R-Square

1. Model

$$\boxed{MPG} \approx b_0 + b_1 \cdot \boxed{\text{weight}} \quad R_a^2 = 0.674$$

2. Model

$$\boxed{MPG} \approx \beta_0 + \beta_1 \cdot \boxed{\text{weight}} + \beta_2 \cdot \boxed{\text{horsepower}}$$

$$R_a^2 = 0.739$$

Aim of the linear regression model: only
a few independent variables should be
in the model

Problem: If we add another independent variable into the model, the values

of R will never decrease.

R_a^2 is an indicator whether it is a good or bad idea, to add another independent variable into the model:

If the value of R_a^2 is increasing, the additional variable should remain in the model. If the value of R_a^2 is decreasing, the additional variable

should be cancelled from the model.

The increase of R_a^2 indicates that the noise power should be added into the model.

8.1.5 Multi collinearity

multi collinearity = if two or more independent variables are depending on each other

multiple correlation coefficient

of the independent variables:

$$R \geq 0.95$$

$$R^2 \geq 0.9025 \approx 0.9 \quad | \cdot (-1)$$

$$-R^2 \leq -0.9 \quad | +1$$

$$1 - R^2 \leq 0.1 \quad | \text{reciprocal}$$

$$\frac{1}{1 - R^2} \geq \frac{1}{0.1} = 10$$

variance inflation factor = VIF

multi collinearity \Leftrightarrow VIF ≥ 10

In case of multi collinearity the estimators of the regression coefficients are incorrect.

VIF = 2.224 < 10

no multi collinearity of weight and horsepower

8.1.6 Forecasting

$$\hat{y}_0 = 58.157$$

$$\hat{y}_1 = -0.007$$

$$\hat{y}_2 = -0.118$$

2500 pounds

80 horse power

MPG = ?

$$58.157 - 0.007 \cdot 2500 - 0.118 \cdot 80 = 34.6$$

MPG = 34.6
2500
80

miles

Is the predicted value reliable?

31.6 is an unrounded value and

$R = 0.866$ strong correlation, thus the

predicted value of 31.6 miles is reliable.

$\beta_2 = -0.118$ if the horsepower in-

creases by one unit and the weight is

unchanged the miles per gallon will

decrease by 0.118 miles.

8.1.7 Heteroscedasticity

Model with only one independent variable



heteroscedasticity

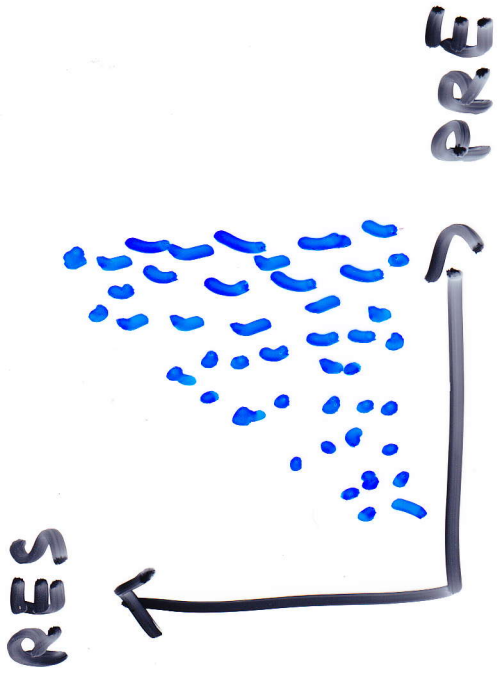


homoscedasticity

Model with two or more independent variables

$RES = \text{residuals} = \text{observed value} - \text{minimum}$

predicted value



(V-form pattern
or a parabolic)

heteroscedasticity

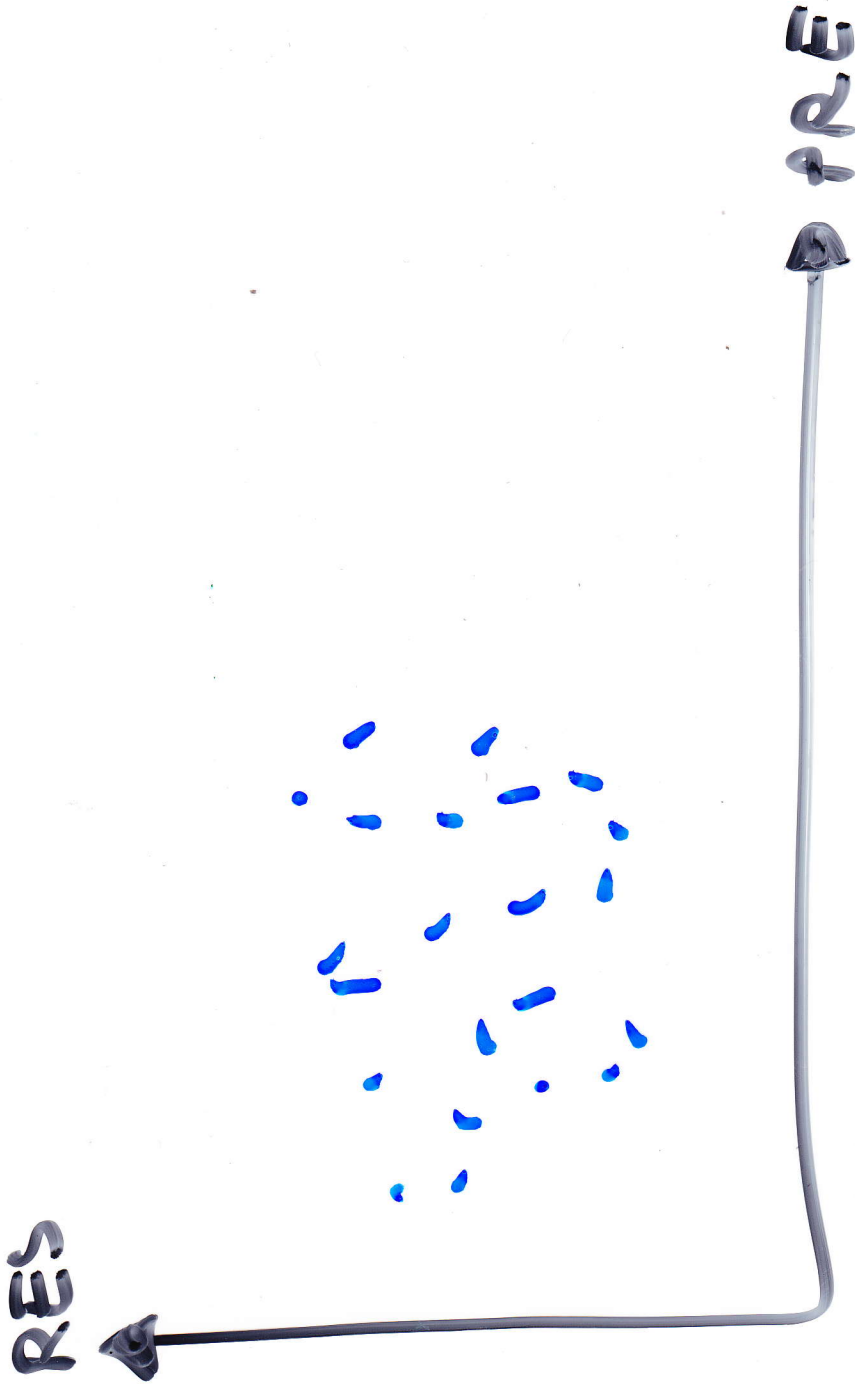
Example

Miles - Per - Gallon . sav



(points are plotted
by random)

homoscedasticity



homoscedasticity

Home work : 6.1 and 6.2