

Day 1

1.1 Linear Models

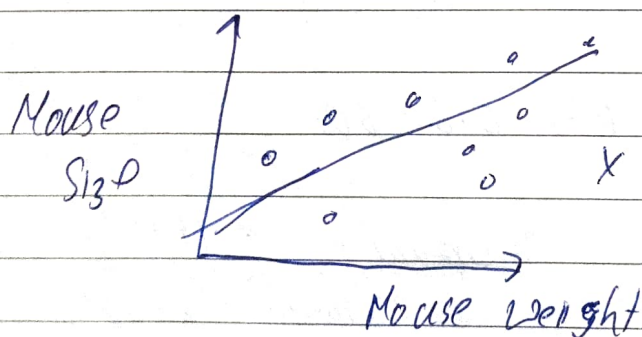
Linear Regression

the Main Ideas!

- 1) Use Least-squares to fit a line to the data
- 2) Calculate R^2
- 3) Calculate a p-value for R^2

Distance from a line to data is called residual

First draw a line then find the residuals and square it and have the sum



$$y = 0.1 + 0.78x$$

↓ ↓
intercept slope

Calculating R^2 gives the probability of getting the guess right or not

To get R^2

- 1) Calculate avg mouse size
 - 2) Draw a line in avg
 - 3) Sum up the squared residuals
- This is called $SS(\text{mean})$

$$SS(\text{mean}) = (\text{data} - \text{mean})^2$$

$$\text{Variation around the mean} = \frac{(\text{data} - \text{mean})^2}{n}$$

4) $SS(\text{fit})$ Sum of squares around the least-square fit basically the line

$$SS(\text{fit}) = (\text{data} - \text{line})^2$$

$$R^2 = \frac{\text{Var}(\text{mean}) - \text{Var}(\text{fit})}{\text{Var}(\text{mean})}$$

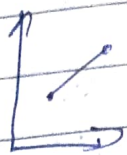
$$= \frac{11.1 - 4.4}{11.1}$$

$$R^2 = 0.6 = 60\%$$

$$R^2 = \frac{SS(\text{mean}) - SS(\text{fit})}{SS(\text{mean})}$$

Mouse weight explains 60% of mouse size

When there is only 2 points

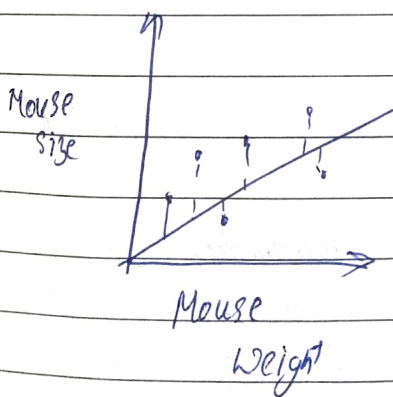


$SS(\text{mean}) = \text{a value}$ but $SS(\text{fit}) = 0$
So $R^2 = 100\%$.

Which means if you have only two pts the $R^2 = 100\%$
So to ~~get~~ determine the r squared value is statistically significant. We need a p value

$$R^2 = \frac{\text{Var}(\text{mean}) - \text{Var}(\text{fit})}{\text{Var}(\text{mean})}$$

$$= \frac{\text{Var}(\text{mouse size}) - \text{var}(\text{After taking mouse weights into account})}{\text{var}(\text{mouse size})}$$



$F = \frac{\text{The variation in mouse size explained by weight}}{\text{The variation in mouse size not explained by weight}}$

$$R^2 = \frac{SS(\text{mean}) - SS(\text{fit})}{SS(\text{mean})}$$

$$F = \frac{SS(\text{mean}) - SS(\text{fit})}{SS(\text{fit})} \cdot \frac{(P_{\text{fit}} - P_{\text{mean}})}{(n - P_{\text{fit}})}$$

P_{fit} is the number of parameters in the fit line

The line is given by $y = y_{\text{intercept}} + \text{slope}$
2 para

$$P_{\text{fit}} = 2$$

P_{mean} is the no of parameters in ^{line} mean value
 $y = y_{\text{intercept}}$
 $P_{\text{mean}} = 1$

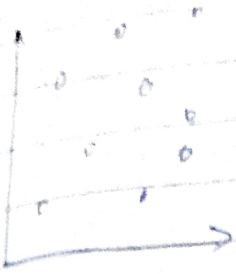
Why divide $SS(\text{fit})$ by $n - P_{\text{fit}}$ instead of just n ?

Intuitively, the more parameters you have in your equation, the more data you need to estimate them. For example, you only need two pts to estimate a line, but you need 3 points to estimate a plane

if the fit is good

$$F = \frac{\text{large no}}{\text{small no}}$$

the F is a large value



$$F = \frac{SS(\text{mean}) - SS(\text{fit})}{\text{error } P(\text{fit})} \quad \frac{SS(\text{fit})}{(n - P(\text{fit}))}$$

$$F = 2$$

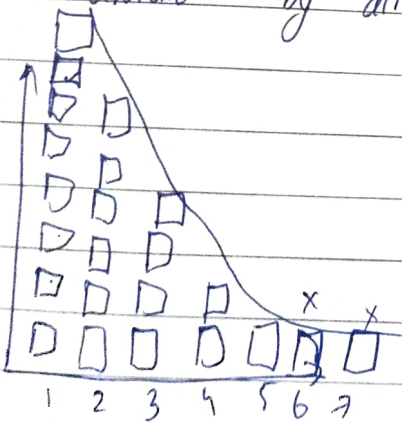
then generate another data

then you get new value of $F = 3$

The plot all the F 's in histogram

then find the F value with the original data
lets say $F = 6$

the p value is number of more extreme values divided by all the values



The decrease of distribution

lets take a p value with $(n - P(\text{fit})) = 10$ $(P(\text{mean} - P(\text{fit})) = 1$
 $(n - P(\text{fit})) = 6$

Here the 10 gets taper off faster

this means the p-value will be smaller when there are more parameters sample relating to the no of parameters in the fit eqn

Linear regression

1) Quantifying the relationship in the data (this is R^2)

1) This need to be large

2) Determining how reliable that relationship is (this is the p-value that we calculate with F) This need to be small

Mh

import pandas as pd

pd.read_csv(path)

Building your Model

Steps

Define : What type of model

Fit : Capture patterns from provided data

Predict

Evaluate

Melbourne_model = DecisionTreeRegressor (random_state=1)

melbourne_model.fit(x, y)