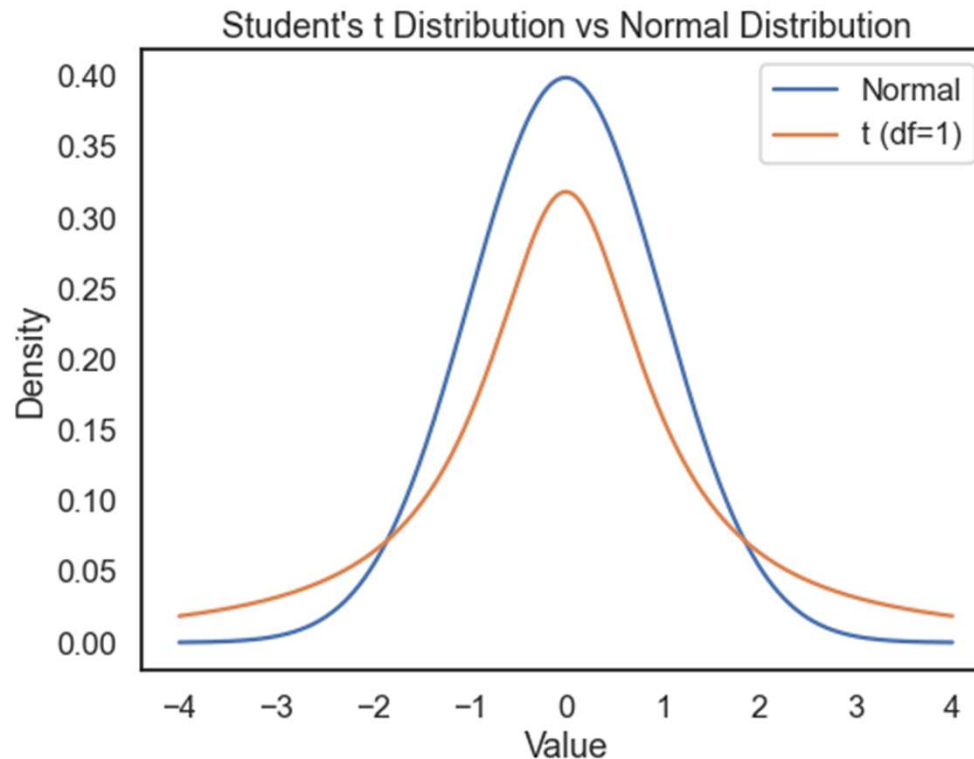


Linear Regression

The Student's t distribution vs. the Normal distribution

Student's t distribution by English statistician William Gosset under the pseudonym "Student".
Used when the sample size is small (<30)

Normal distribution by Carl Friedrich Gauß



Student t distribution can be applied with small sample size with less than 30

| Interval | area |
|------------------------------|------|
| $[\mu-1\sigma, \mu+1\sigma]$ | 68.3 |
| $[\mu-2\sigma, \mu+2\sigma]$ | 95.4 |
| $[\mu-3\sigma, \mu+3\sigma]$ | 99.7 |

only applicable when having enough samples

Correlation

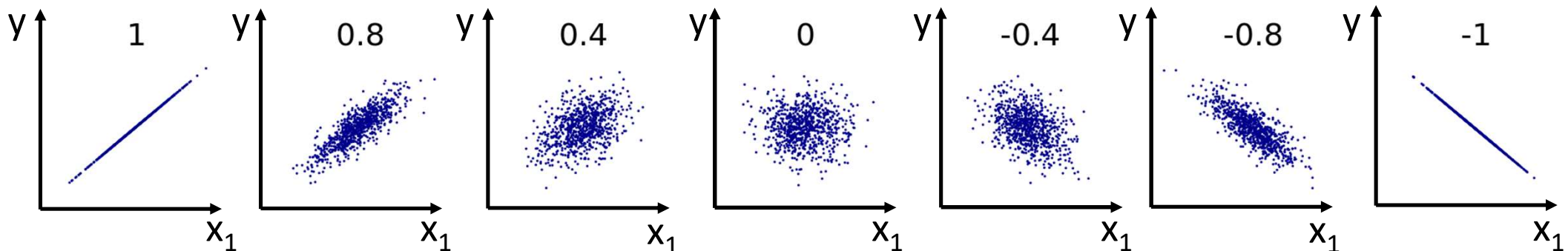
In statistics, correlation or dependence is any statistical relationship between two random variables (aka. bivariate data), e.g. x_1 and x_2 or x_1 and y , ...

The **Pearson correlation measure** is defined as

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} = \frac{\text{cov}(X, Y)}{\sigma_x \sigma_y}$$

deviation of x to mean of x and y and mean of y
covariance
std deviation

Examples for different r :

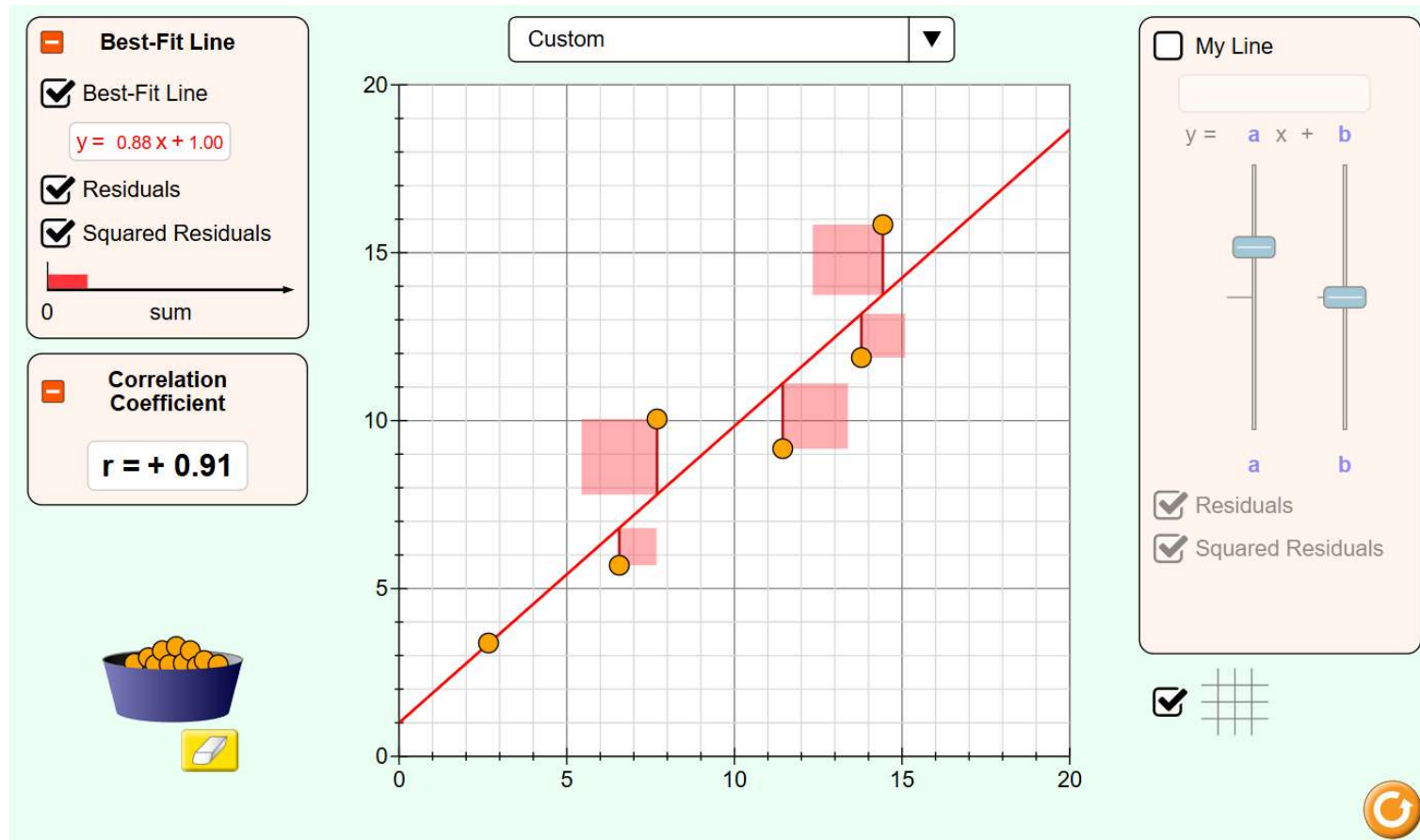


when $r = 1$, correlation of x to y is positively dependent

no dependence between x and y

the bigger x , the smaller negatively correlated

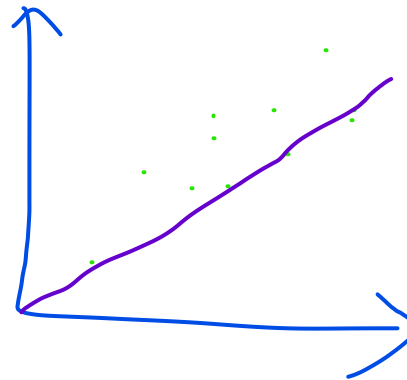
Least squares linear regression - Animation



https://phet.colorado.edu/sims/html/least-squares-regression/latest/least-squares-regression_en.html

Linear Regression


- A simple approach to supervised learning
- It assumes that the dependence of Y on X is linear
- Conceptually and practically very useful
- Which questions we are trying to answer by doing this linear regression
 - + is there relationship between certain feature & the response
 - + what feature are useful in predicting the response
 - + How accurately can we extrapolate
 - + The relationship actually linear



$y = +b + m.x$ - linear function

Linear regression using a single predictor X

Let's $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$ be the prediction for Y based on the i th value of x with $\hat{\beta}_0$ as intercept & $\hat{\beta}_1$ as slope. Then $e_i = y_i - \hat{y}_i$ represents the i th residual


residual

We then define the residual sum of squares (RSS) as:

$$RSS = e_1^2 + e_2^2 + \dots + e_n^2$$

the least squares approach choose $\hat{\beta}_0$ and $\hat{\beta}_1$ to minimize the RSS. the minimizing values can be shown to be.

$\hat{\beta}_1 = \frac{\text{Sum}(\text{deviation of } x \text{ and } y)}{\text{Sum}(\text{deviation of } x)^2}$ - how to find the slope

$\hat{\beta}_0 = \text{sample } y - \hat{\beta}_1(\text{sample } x)$

=> find \bar{y} and \bar{x} (sample means)

Assessing the accuracy of the coefficient estimates

The standard error of an estimate reflects how it varies under repeated sampling.
We have:

Std error can be used to compute confidence intervals. a 95% confidence interval is defined a range of values such that with 95% probability, the range will contain the true unknown value of the parameter.

It has the form: - screenshot

This means, there is approx 95% chance that the interval β_1 will contain the true value β_1

Quiz question

Suppose we accidentally doubled our data, leading to observations and error terms identical in pairs. If we ignored this, our standard error calculations would be as if we had a sample of size $2n$, when in fact we have only n samples.

What happens in that case to the confidence interval $\beta_1 \in [\hat{\beta}_1 - 2 \cdot SE(\hat{\beta}_1), \hat{\beta}_1 + 2 \cdot SE(\hat{\beta}_1)]$

with $SE(\hat{\beta}_1)^2 = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$

- A) It narrows by a factor of 2
- ☒ B) It narrows by a factor of $\sqrt{2}$
- C) It remains unchanged
- D) It widens by a factor of 2



Hypothesis testing

- + std Error can be used to perform hypothesis tests on the coefficients:

- Null hypothesis H_0 : there is no relationship between X and Y, this correspond to $H_0: \beta_1 = 0$
- Alternative hypothesis H_A : there is some relationship between X & Y, this correspond to $H_A: \beta_1$ is different from 0

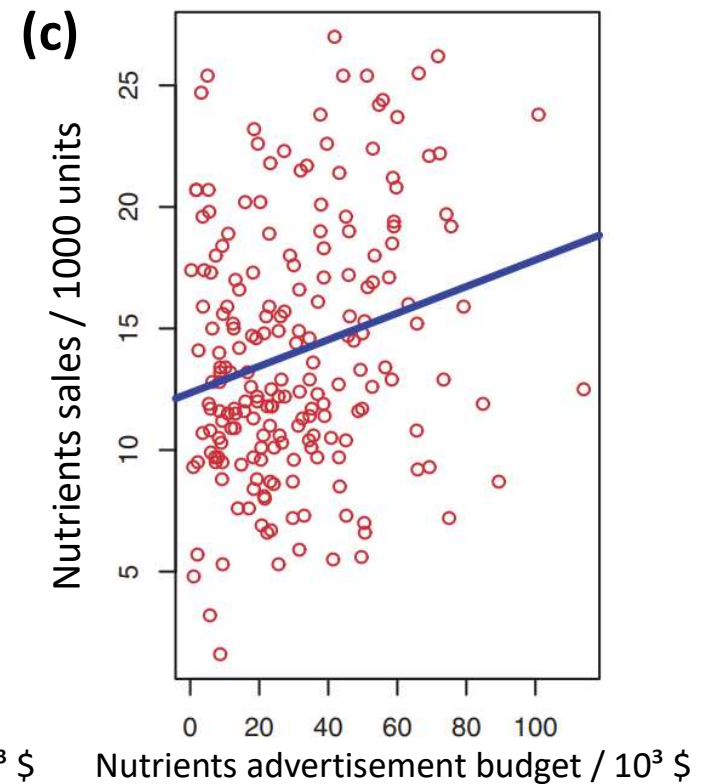
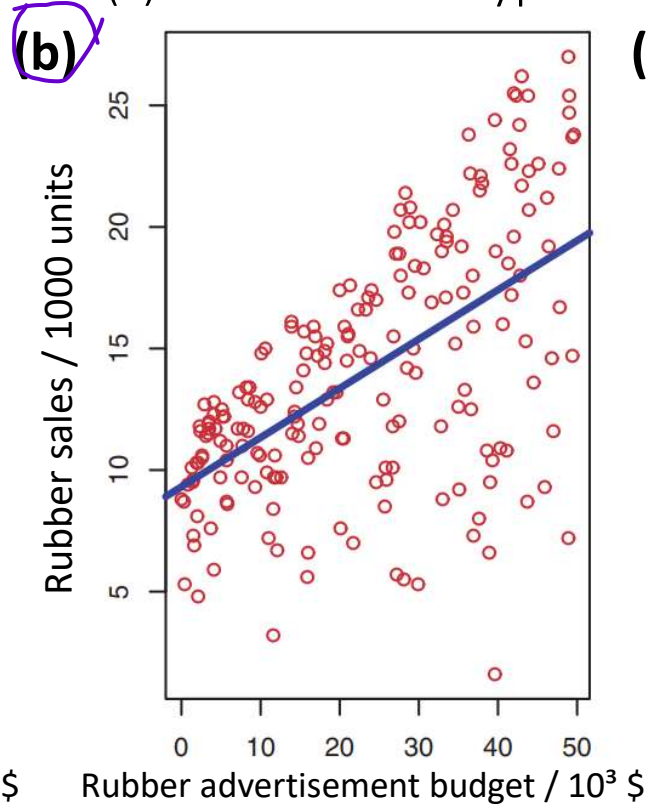
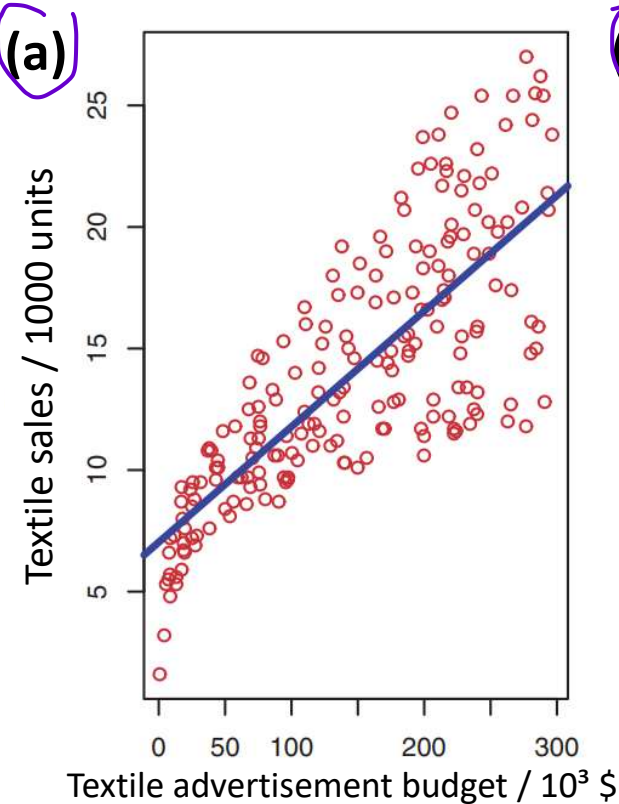
To test the Null hypothesis we compute a t- statistics, given by: $(\hat{\beta}_1 - \beta_1) / SE(\hat{\beta}_1) = \hat{\beta}_1 / SE(\hat{\beta}_1) \Rightarrow t$ measures the number of stds that $\hat{\beta}_1$ is away from 0

we reject $H_0 : \beta_1 = 0$ when : t is large (usually > 10) & p (the probability of observing any value equal to |t| or larger) value is small (< 0.05)

- + F is large: finds out if at least one predictor is useful $F = ((\text{Total sample square} - \text{RSS})/p) / \text{RSS} / (n-p-1)$
- + Z is large: measures accuracy estimate by $\hat{\beta}_1 / SE(\hat{\beta}_1)$

Quiz question

A chemical company spends a certain budget to advertise for textiles, rubber, and nutrients in 200 different markets. In which case(s) will the Null-Hypothesis $H_0: \beta_1 = 0$ be rejected?



p-value: <0.0001
t-statistics: 32.8
Std. error: 0.0014

<0.0001
21.9
0.0086

0.8599
-0.2
0.0059

p smaller than 0.05 and t greater than 10

Quiz question

The following correlation matrix is provided:

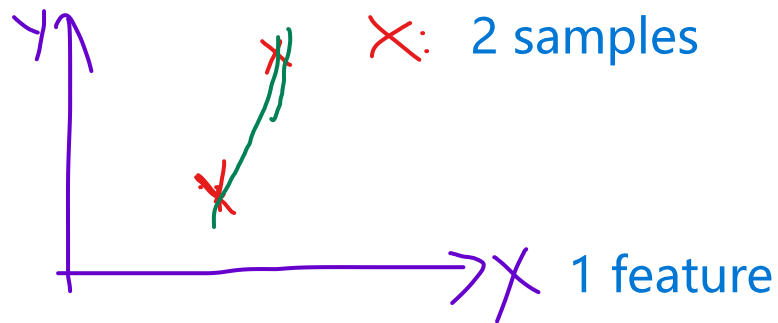
| | Textile ad. | Rubber ad. | Nutrients ad. | Sales |
|---------------|-------------|------------|---------------|-------|
| Textile ad. | 1 | 0.05 | 0.06 | 0.78 |
| Rubber ad. | | 1 | 0.35 | 0.58 |
| Nutrients ad. | | | 1 | 0.23 |
| Sales | | | | 1 |

Which advertisement seems to have the lowest success in terms of increasing sales?

- a) Textile ad.
- b) Rubber ad.
- c) Nutrients ad.

Generalization of least squares linear regression

A supervised parametric model for p (#feature) $<$ n (#sample)



$f(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n + \text{epsilon}$ (irreducible error) with β_0 as intercept and β_1 to β_n as slope. β_i are also called the coefficients or parameters. and epsilon is the irreducible error.

RSS = Residual Sum of square = $\sum_{i=1..n} (y_i - \bar{y})^2$ with $(y_i - \bar{y})$ being the i th residual

TAA = Total sum of squares = $(\sum_{i=1 \dots n} (y_i - \bar{y}))^2$

Assessing the overall model accuracy:

R - squared = fraction of variance explained

$R^2 = 1 - \text{RSS}/\text{TSS}$

r^2 (correlation) = R^2 in case of linear regression

R^2 belong to $[0;1]$. The higher R^2 the better is the quality of the prediction

Quiz question

Which of the following are potential problems of linear regression?

- a) The response-predictor relationship is non-linear
check the trend in residual plots to find out $e_i = y_i - y_i \text{ predicted}$
- b) Heteroscedasticity (Non-constant variance in error terms)
For example. time series data
- c) Outliers and High-leverage points
outliers y is far from the value predicted by the model
high leverage point uses leverage statistics h
- d) Collinearity
spot it by calculating the variance inflation factor

Collinearity

collinearity is the linear relationship between two predictors

- for eg: temp colinear with humidity
- prediction of sickness according to body temp, or temp impacts the state of body sickness

To detect collinearity:

- + look at the correlation matrix of the predictors (if the value is high so the strong correlation is)
- + if collinearity exists between 3 or 4 variables, but not single variables then it is called multi-collinearity: investigate this using variance inflation factor (VIF)
when $VIF > 5$ signals problems with collinearity

for $VIF (\beta_i^2) = 1 / (1 - R^2_{x|y})$. R^2 from a regression of X_j onto all of the other predictors

Solutions:

- Drop one of the variables
- Combine the correlated variables into one single predictor. for eg: by averaging them

Adding feature interaction to the model

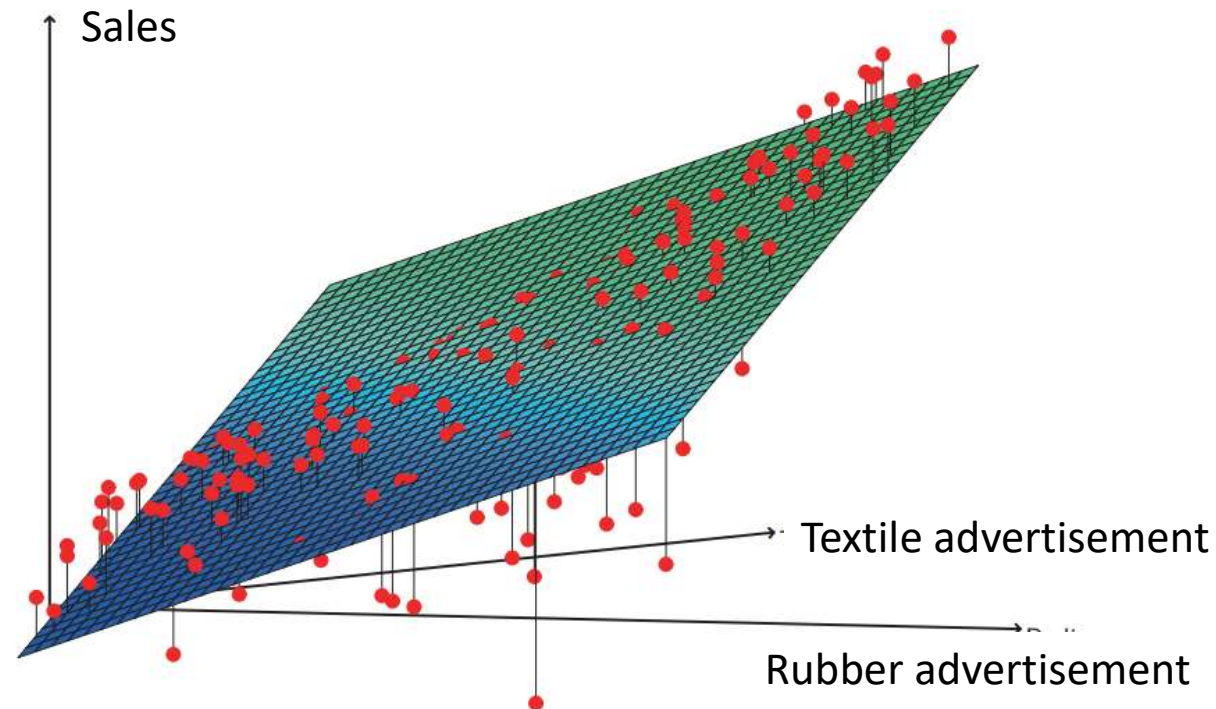
Figure: Introduction to Statistical Learning, T. Hastie et al. 2021

When advertising is split between rubber and textile then the true sales is underestimated. There seems to be a synergy effect, aka. an interaction, between textile and rubber advertisement. Instead of using the model:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon$$

We add an interaction term:

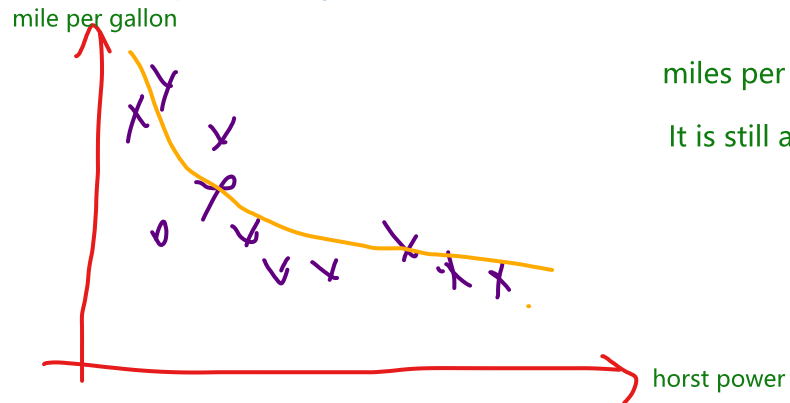
$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + \varepsilon$$



| | Coefficient | Std. error | t-statistic | p-value |
|----------------|-------------|------------|-------------|-------------------|
| Intercept | 6.7502 | 0.248 | 27.2 | <10 ⁻⁴ |
| Textile | 0.0191 | 0.002 | 13.7 | <10 ⁻⁴ |
| Rubber | 0.0289 | 0.009 | 3.2 | 0.001 |
| Textile·Rubber | 0.0011 | 0.001 | 21.7 | <10 ⁻⁴ |

Making the model non-linear

Simple way to add non-linear effects is to: (1) square the features



$$\text{miles per gallon} = \beta_0 + \beta_1 \cdot \text{horsepower} + \beta_2 \cdot \text{horsepower}^2 + \epsilon$$

It is still a linear model, what we have done is $x = \text{horse power}$, $x_2 = \text{horse power square}$

Incorporating qualitative features into the model

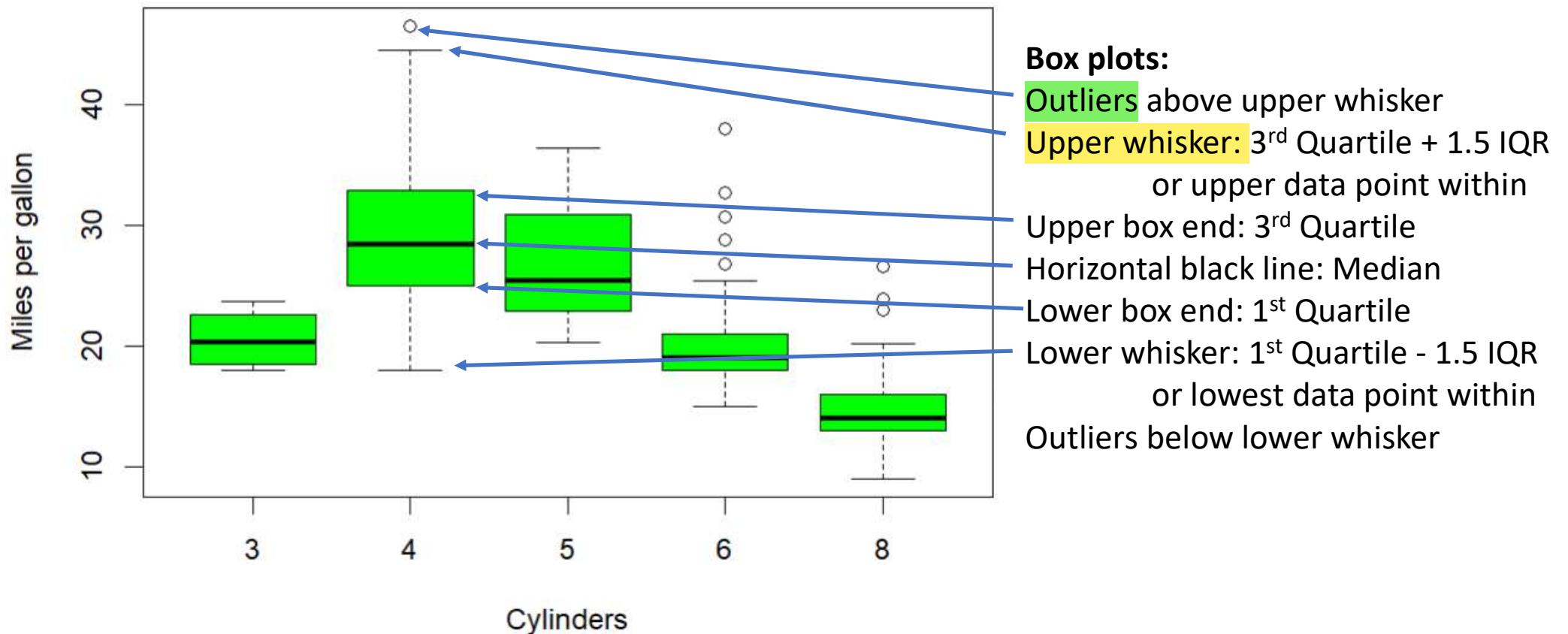
example: We want to incorporate into the model whether the car runs on 'diesel, electricity, or Biofuel, ...
For this: we create "dummy variables" and modify the model:

$x_{i1} = 1$: if the car runs on diesel; OR 0 if the car has no diesel engine

$x_{i2} = 1$: if car has electric engine; OR 0 if the car has no electric engine

Therefore, we $y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \epsilon_i$

Box plots, aka. whisker plots



Box plots for a Normal Gaussian distribution

