

Regression Modelling (Baayen)

$lm1 = lm(var1 \sim var2)$

↑
dependent
variable,
to be
modelled

↑
predictor

lm - linear model

Slope and intercept are estimated using
Least Square Regression Line

↓
by minimizing the squared vertical distance
between the data points and the line

$coef(lm1)$ - returns the model's coefficients

$abline(lm1)$ - will draw the line

Correlation - degree to which the data
points cluster around the regression line

The degree of correlation - correlation
coefficient

ρ - correlation coeff for a population

r - correlation coeff for a sample

$$-1 \leq \rho \leq 1$$

$$-1 \leq p \leq 1$$

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

↑
perfect
negative
correlation

↑
perfect
positive
correlation

this measure ~~can be useful~~
can be used to think is how useful is
to fit a straight line to the data

r^2 (or R^2) - measure for evaluating how much
of the scatter is accounted for

or R^2 quantifies the proportion of the variance
in the data that is captured and
explained by the regression model

When we try to fit a model to a data set,
the goal is to be able to predict values
of the dependent variables

⚡ ⚡

The better the prediction, the higher r .
And ~~the~~ with many ~~predictors~~ predictors,
little variance is explained by the model

↳ in this situation R^2 is close to 0.

summary(lm) returns a summary of the lm

- residuals
- coefficients
 - intercept, slope
 - each coeff. comes with 3 other numbers:
 - standard error
 - t-value
 - p-value

p-value tells us if the coefficient is ~~signifi-~~ significantly different from zero (and hence, ~~potentially useful~~.)

if coefficient is zero, there are no relations at all between the predictor and the dependant variable

t-test is used to ascertain if p-coeff is significantly different from zero, and, hence, potentially useful.

$$\underline{t\text{-value}} = p\text{-coeff} / \text{st. error}$$

st. error ^{measure on} - how we're sure about the estimate of the coefficient.

the smaller the standard error, the smaller is the confidence interval around the estimate \rightarrow

`names(summary(lm1))` ← names

`summary(lms)$coef`

`summary(lms)$coef[, 3]`

Residual standard error - measure of how unsuccessful the model is

the better model, the smaller its residual standard error

multiple R-squared - r^2 , squared correlation coefficient

We get r by square root of R^2

another way of calculating r :

`cor(a, b)`

`cor.test(a, b)`

↓

also tests if it's significantly different from 0.
and lists 95% confidence intervals

F-value - a test if the linear model as a whole succeeds in explaining a significant portion of the variance.

we also can describe quadratic terms with the

$$\ln 1 - \ln(a \sim b + I(c^2))$$

linear means that the relationship
can be expressed as sum
(linear combination)

Basic rules:

- visualize
- beware of outliers
- straight lines aren't always the best
- keep it simple

For factors.

$$\ln(x \sim \text{some factor})$$

↓
factor, not numerical value

↓
gets converted into numerical
vectors (one or more)

let's assume s.f. has 2 levels

↓
so it's converted into one numerical
vector with 1 and 0.

However, intercept and slope need special interpretation

Dummy Coding - way to deal with factor variables

one level is signed out as the Default or Reference level, which is contrasted to others

in this case intercept would ~~mean~~ represent the group mean for the default level

Class = c("animal", "plant")
↑
default level

so intercept would be mean for animals

the second value (Classplant) represents the contrast (i.e. difference) between the group mean of the plants and of the animals

In other words, mean of plants level is $[(\text{intercept}) + \text{classplant}]$

t-value tells that the adjustment is statistically significant

(i.e. the means differ significantly)

lm can be applied to ^{factors with} more than 2 levels

anova - analysis of variance

to explain some of the variation ~~in thing~~
within groups



reports F value
which reports that there are (not)
significant differences in the
means

but it doesn't tell us what differences
are involved
(need to run "summary")



it will show all variables
except the default one, so the
difference is with this variable
(this default row is (Intercept))

Tukey's Honestly Significant Difference

TukeyHSD in R

for multiple comparisons
for detecting significant differences

$\text{aov}(\text{breaks} \sim \text{tension}, \text{data} = \text{warpbreaks})$

↑
special for analysis of variance

(output is exactly the same as if we applied anova to lm)

$\text{tc} = \text{TukeyHSD}(\text{warpbreaks.aov})$

	diff	lwr	upr	p adj
M-L	-10	-19.55	-0.44	0.03
H-L	-14.72	-24.28	-5.16	0.00
H-M	-4.72	-14.28	-4.83	0.46

↑
this table lists the differences in the means

the lower and the upper end points of the confidence intervals

adjusted p-value

$\text{plot}(\text{tc})$ will draw these differences

(those that intersect --- dotted line --- are not significantly different)

numerical predictor

factor as predictor

analysis of covariance -
with both numeric and factors

In R, \ln used for all these analyses:
regression, variance, covariance

all are built on the same fundamental principles

$$\ln 1 = \ln(a \sim b \times c)$$

↑
factor

Let's consider the following model

$$\ln I = \ln(\text{mean Size Rating}_i \sim \text{mean Familiarity}_i \times \text{Class} + I(\text{mean Familiarity}^2))$$

↑
↑

factor
numeric

Coefficients

	Est
• (Intercept)	4,93
• mean Familiarity	-0,63
• \bar{R}^2 (mean Familiarity ²)	0,11
• Class plant	-1,01
• mean Familiarity : Class Plant	-0,21

Summarizes how the preceding coefficients should be modified in order to make them more precisely for the nouns that fall into "plant" category

(mean Size Rating)

(Class factor)

The coefficient "Classplant" tells us that we should subtract ~~-0,21~~ -1,01 from the intercept in order to obtain the (modified) mean for the plants

The final coefficient "mean Familiarity : Classplant" tells us that the coefficient "mean Familiarity" should be decreased by -0,21 in order to make it precise for plants

The last coefficient - interaction between mean Familiarity and Class

mean Familiarity \times Class

interaction

\Downarrow the same is

mean Familiarity \times Class +
mean Familiarity \div Class

mean Familiarity \div Class

interaction of the predictor to its left
and right

What the interaction tells us is that
the linear coefficient of mean Familiarity
has to be adjusted downwards when
dealing with plants rather than with
animals

(categories
in
factor)

• For animals the coefficient is
-0.63

• For plants, we add the coef.
for the interaction of
mean Familiarity by Class

to this coef: $-0.63 - 0.212 = -0.84$

\uparrow
i.e. the linear term of mean Familiarity
differs for plants and
animals