

ENVX2001 Lab 07 - Regression model development

ENVX2001 Applied Statistical Methods
Semester 1, 2026

Learning outcomes

In this lab, you will work towards achieving learning outcomes [L03](#), and [L05](#).

Lab objectives

In this lab, we will:

- ☐ Identify best predictors for model - Exercise 1
- ☐ Fit model and check assumptions - Exercise 1
- ☐ Interpret model output - Exercise 1



Tip

Please work on this exercise by creating your own R Markdown file.

Preparation

- ☐ Install or update the performance package

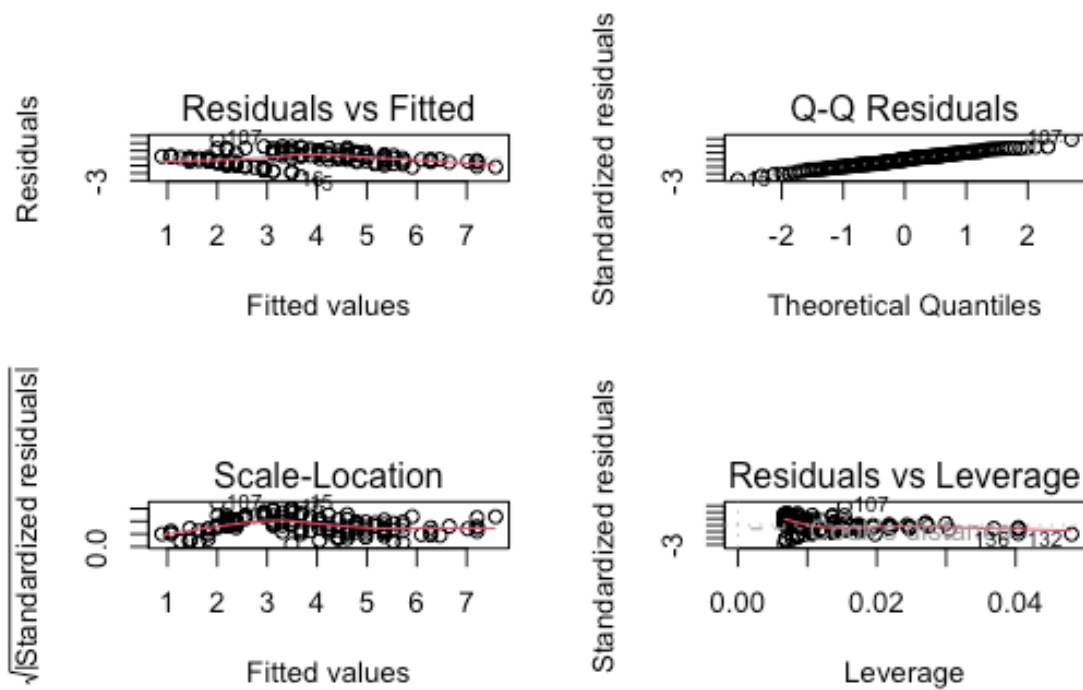
```
CODE
#install.packages("performance")
library(performance)
```

This package is really good for checking your models. For this lab, we will focus on the `check_model()` function, which gives us nice pretty diagnostic plots for models:

plot()

```
CODE
```

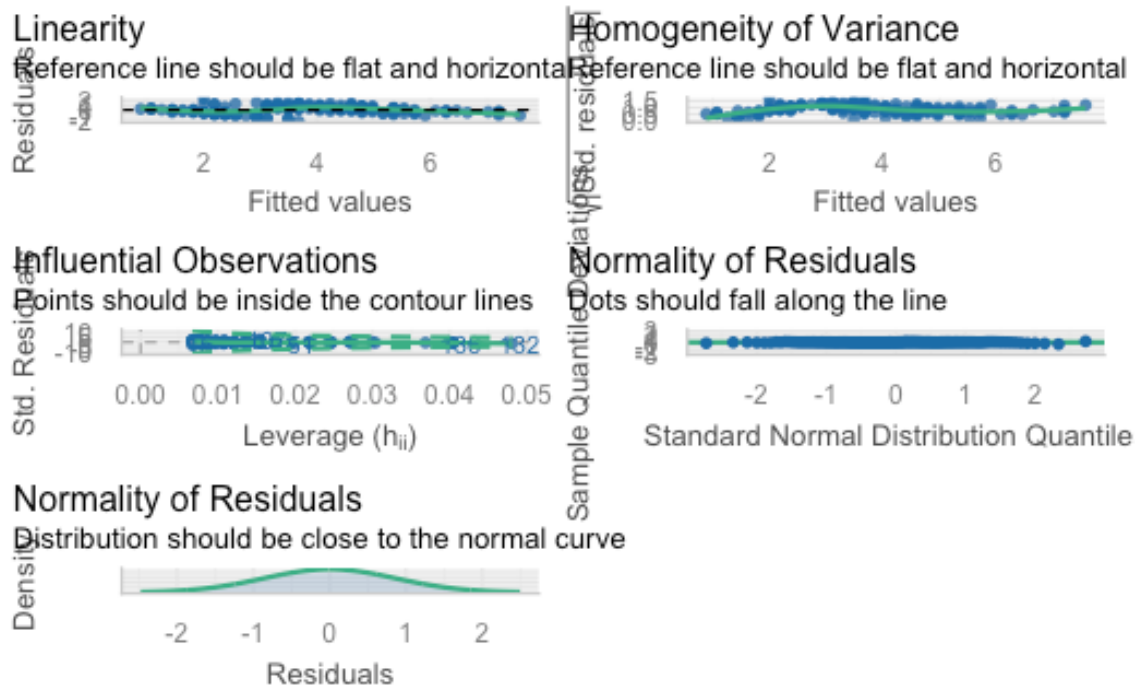
```
par(mfrow=c(2,2))
plot(iris_lm)
```



```
CODE
par(mfrow=c(1,1))
```

check_model() from performance

```
CODE
library(performance)
check_model(iris_lm)
```



Exercise 1: Modelling bird abundance

We will now use the transformed data in `loyn` for this exercise. If you have not already figured out how to perform the transformation, or if something is wrong, you may use the `loyn` tab in the `m1r.xlsx` MS Excel document. Alternatively, the code to convert the data is below.



Tip

This is the same data we used in the walkthrough exercise

CODE

```
# Load library if needed
library(readxl)
# reset the data import just in case it has been modified
loyn <- read_xlsx("data/m1r.xlsx", "Loyn")
# make transformations

loyn <- loyn %>%
  mutate(
    L10AREA = log10(AREA),
```

```

      L10DIST = log10(DIST),
      L10LDIST = log10(LDIST)
    )

# check
glimpse(loyn)

```

```

OUTPUT
Rows: 56
Columns: 10
$ ABUND <dbl> 5.3, 2.0, 1.5, 17.1, 13.8, 14.1, 3.8, 2.2, 3.3, 3.0, 27.6, 1.0...
$ AREA <dbl> 0.1, 0.5, 0.5, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 2.0, 2.0, 2.0...
$ YR.ISOL <dbl> 1968, 1920, 1900, 1966, 1918, 1965, 1955, 1920, 1965, 1900, 1.0...
$ DIST <dbl> 39, 234, 104, 66, 246, 234, 467, 284, 156, 311, 66, 93, 39, 4...
$ LDIST <dbl> 39, 234, 311, 66, 246, 285, 467, 1829, 156, 571, 332, 93, 39, 4...
$ GRAZE <dbl> 2, 5, 5, 3, 5, 3, 5, 5, 4, 5, 3, 5, 2, 1, 5, 5, 3, 3, 3, 2, 2...
$ ALT <dbl> 160, 60, 140, 160, 140, 130, 90, 60, 130, 130, 210, 160, 210, 160...
$ L10AREA <dbl> -1.00000000, -0.30103000, -0.30103000, 0.00000000, 0.00000000, 0.0...
$ L10DIST <dbl> 1.591065, 2.369216, 2.017033, 1.819544, 2.390935, 2.369216, 2...
$ L10LDIST <dbl> 1.591065, 2.369216, 2.492760, 1.819544, 2.390935, 2.454845, 2...

```

Best single predictor?

Question 1

Obtain the correlation between ABUND and all of the predictor variables using `cor()`. Based on these, what would you expect to be the best single predictor of ABUND?

```

CODE
cor(loyn)

```

Answer

```

CODE
cor(loyn)

```

```

OUTPUT
      ABUND      AREA      YR.ISOL      DIST      LDIST
ABUND  1.00000000  0.255970206  0.503357741  0.2361125  0.08715258
AREA   0.25597021  1.000000000 -0.001494192  0.1083429  0.03458035
YR.ISOL 0.50335774 -0.001494192  1.000000000  0.1132175 -0.08331686
DIST    0.23611248  0.108342870  0.113217524  1.0000000  0.31717234
LDIST    0.08715258  0.034580346 -0.083316857  0.3171723  1.00000000
GRAZE   -0.68251138 -0.310402417 -0.635567104 -0.2558418 -0.02800944
ALT      0.38583617  0.387753885  0.232715406 -0.1101125 -0.30602220
L10AREA  0.74003580  0.584651024  0.278414517  0.3047850  0.33680642
L10DIST  0.12672333  0.163054319 -0.019572228  0.8233190  0.29365797
L10LDIST 0.11812448  0.101607829 -0.161116108  0.4968169  0.82059568
      GRAZE      ALT      L10AREA      L10DIST      L10LDIST
ABUND  -0.68251138  0.3858362  0.7400358  0.12672333  0.11812448
AREA    -0.31040242  0.3877539  0.5846510  0.16305432  0.10160783
YR.ISOL -0.63556710  0.2327154  0.2784145 -0.01957223 -0.16111611

```

DIST	-0.25584182	-0.1101125	0.3047850	0.82331904	0.49681692
LDIST	-0.02800944	-0.3060222	0.3368064	0.29365797	0.82059568
GRAZE	1.00000000	-0.4071671	-0.5590886	-0.14263922	-0.03399082
ALT	-0.40716705	1.00000000	0.2751428	-0.21900701	-0.27404380
L10AREA	-0.55908864	0.2751428	1.00000000	0.30216662	0.38247952
L10DIST	-0.14263922	-0.2190070	0.3021666	1.00000000	0.60386637
L10LDIST	-0.03399082	-0.2740438	0.3824795	0.60386637	1.00000000

The best single predictor would be L10AREA as this has the highest r ($r = 0.74$)

Assumptions and interpretation

Question 2

Use multiple linear regression to see whether ABUND can be predicted from L10AREA and GRAZE. Are the assumptions met? Is there a significant relationship? *Note: we are using these 2 predictors as they have the largest absolute correlations. Use `lm()` and specify the model as `ABUND ~ L10AREA + GRAZE`.*

```
CODE
lm.mod1 <- lm(ABUND ~ GRAZE + L10AREA, data = loyn)

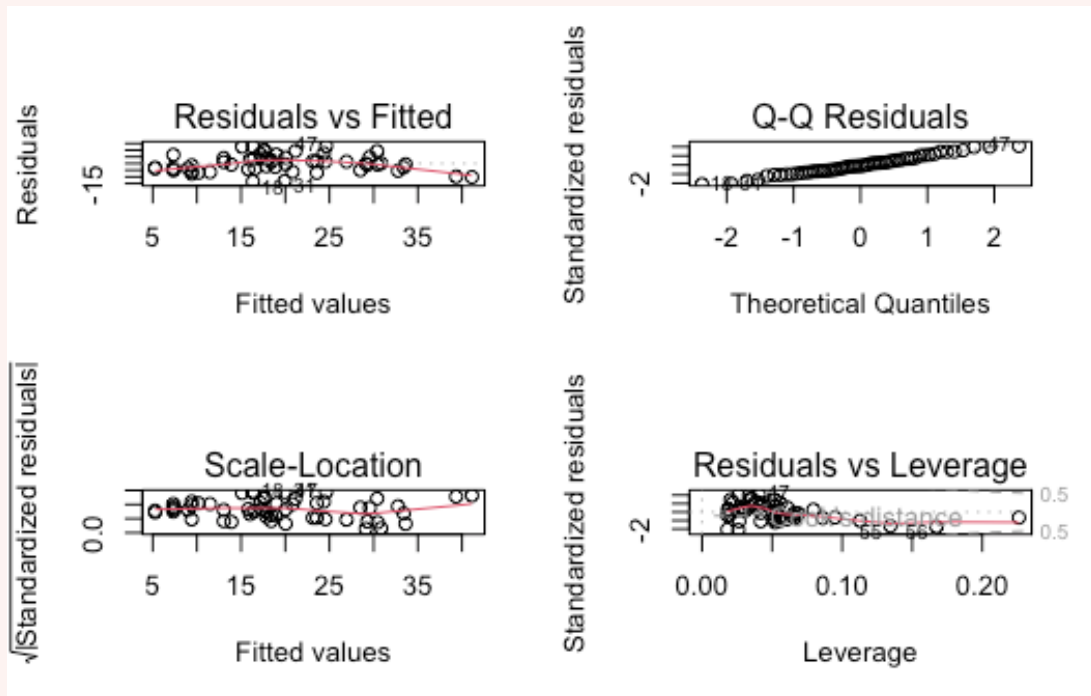
par(mfrow = c(2, 2))
plot(lm.mod1)
par(mfrow = c(1, 1))

summary(lm.mod1)
```

Answer

```
CODE
lm.mod1 <- lm(ABUND ~ GRAZE + L10AREA, data = loyn)

par(mfrow = c(2, 2))
plot(lm.mod1)
```



```
CODE
par(mfrow = c(1, 1))

summary(lm.mod1)
```

OUTPUT

```
Call:
lm(formula = ABUND ~ GRAZE + L10AREA, data = loyn)

Residuals:
    Min       1Q   Median       3Q      Max
-13.4296  -4.3186  -0.6323   4.1273  13.0739

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  21.6029     3.0917   6.987 4.73e-09 ***
GRAZE        -2.8535     0.7125  -4.005 0.000195 ***
L10AREA       6.8901     1.2900   5.341 1.98e-06 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 6.444 on 53 degrees of freedom
Multiple R-squared:  0.6527,    Adjusted R-squared:  0.6396
F-statistic: 49.81 on 2 and 53 DF,  p-value: 6.723e-13
```

This is a significant model as both b_1 and b_2 are significant and the model is significant.

The residuals look reasonable. They are approximately normally distributed (both right hand plots), but possibly the variance is not totally constant and there are possibly a few values with high leverage (left hand plots).

Question 3

How good is the model based on the (i) r^2 (ii) adjusted r^2 ? Use `summary()`.

CODE

```
summary(lm.mod1)$r.squared  
summary(lm.mod1)$adj.r.squared
```

Answer

The Adjusted r^2 is lower than the r^2 , but we would opt for the adjusted r^2 as it takes the number of predictors into account. Overall the model is ok, explaining 64.0% of variation in Abundance.

Question 4

Which variable(s) has the most significant effect(s)? *(Refer specifically to the t probabilities in the table of predictors and their estimated parameters or coefficients in the output of `summary()`).* Interpret the p-values in terms of dropping predictor variables.

Answer

Both L10AREA and GRAZE are highly significant, L10AREA is the most significant. In terms of effect, a 1 unit change in GRAZE results in a -2.9 decrease in abundance (with L10AREA remaining constant), while a 1 unit change in L10AREA, (therefore a 10 unit change in AREA) results in a 6.9 increase in abundance (GRAZE holding constant).

Question 5

Repeat the multiple regression, but this time include YRS.ISOL as a predictor variable (it has the 3rd largest absolute correlation). This will allow you to assess the effect of YRS.ISOL with the other variables taken into account.

Answer

CODE

```
lm.mod2 <- lm(ABUND ~ GRAZE + L10AREA + YR.ISOL, data=loyn)
```

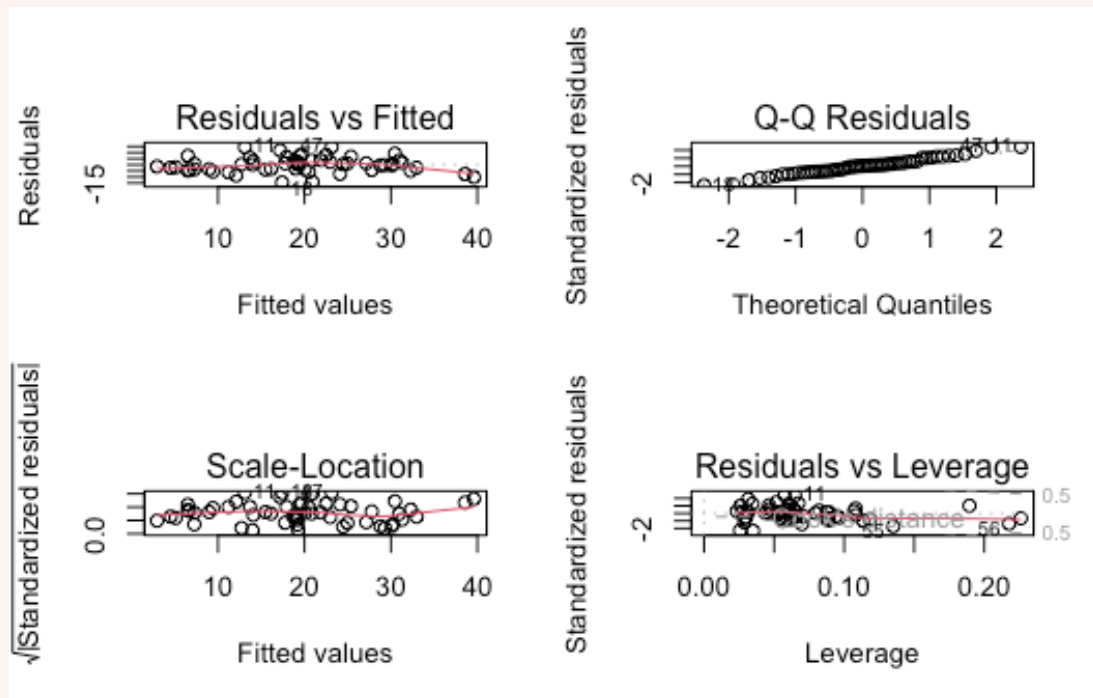
Question 6

Check assumptions, do the residuals look ok? If you are happy with the assumptions, you can proceed to interpret the model output.

Answer

CODE

```
par(mfrow=c(2,2))  
plot(lm.mod2)
```



CODE

```
par(mfrow=c(1,1))
```

The residuals look OK, but YR.ISOL is borderline significant ($p = 0.0768$).

Question 7

Compare the r^2 and adjusted r^2 values with those you calculated for the 2 predictor model, Which is the better model? Why?

CODE


```
summary(lm.mod2)
```

Answer

Both of these are greater than for model in step 3, so this is a better model.

Exercise 2: California streamflow

The following dataset contains 43 years of annual precipitation measurements (in mm) taken at (originally) 6 sites in the Owens Valley in California. I have reduced this to three variables labelled lake_sabrina (Lake Sabrina), pine_creek (Big Pine Creek), rock_creek (Rock Creek), and the dependent variable stream runoff volume (measured in ML/year) at a site near Bishop, California (labelled runoff_volume).

Note the variables have already been log-transformed to increase normality of the residuals in the regressions.

Start with a full model and manually remove the variables one at a time, checking every time whether removal of a variable actually improves the model.

CODE

```
# read in the data
stream_data <- read_xlsx("data/california_streamflow.xlsx", "streamflow")
names(stream_data)
```

OUTPUT

```
[1] "lake_sabrina" "pine_creek"   "rock_creek"   "runoff_volume"
```

CODE

```
s.mod_full <- lm(runoff_volume ~ lake_sabrina + pine_creek + rock_creek, data=stream_data)
s.mod_full <- lm(runoff_volume ~ ., data=stream_data) ## you can also use the . to indicate use all
variables
summary(s.mod_full)
```

OUTPUT

```
Call:
lm(formula = runoff_volume ~ ., data = stream_data)
```

Residuals:

Min	1Q	Median	3Q	Max
-0.09885	-0.03331	0.01025	0.03359	0.09495

Coefficients:

```

              Estimate Std. Error t value Pr(>|t|)
(Intercept)  3.25716    0.12360  26.352 < 2e-16 ***
lake_sabrina  0.05631    0.03756   1.499  0.14185
pine_creek   0.21085    0.06756   3.121  0.00339 **
rock_creek   0.43838    0.08798   4.983  1.32e-05 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.04861 on 39 degrees of freedom
Multiple R-squared:  0.8817,    Adjusted R-squared:  0.8726
F-statistic: 96.88 on 3 and 39 DF,  p-value: < 2.2e-16

```

Partial F-Tests

The above analysis tells us that both `pine_creek` & `rock_creek` are significant, according to the t-test, in the model and `lake_sabrina` is not? This involves performing Partial F-Tests as discussed in the lecture.

This can be done in **R** by using `anova()` on two model objects. To be able to compare the models and run the anova, you need to make objects of all the possible model combinations you want to compare.

```

CODE
s.mod_reduced <- lm(runoff_volume ~ rock_creek + pine_creek, data=stream_data)
anova(s.mod_reduced, s.mod_full)

```

The last row gives the results of the partial F-test.

Question 1

Should we remove `lake_sabrina` from the model?

Answer

Yes, we should remove `lake_sabrina` as the p-value is > 0.05 and opt for the simpler model.

Question 2

Is the p-value for the f-test the same as for the t-test?

Answer

Yes, P-values for the t-statistic and for the Partial F-statistic are related ($\text{Partial F} = t^2$)

Question 3

Write out the hypotheses you are testing.

Answer

$$H_0: \beta_{lake_sabrina} = 0$$

$$H_1: \beta_{lake_sabrina} \neq 0$$

Perform a Partial F-Test to work out if the removal of lake_sabrina and pine_creek improves upon the full model.

CODE

```
s.mod_reduced2 <- lm(runoff_volume ~ lake_sabrina + pine_creek, data=stream_data)
anova(s.mod_reduced2, s.mod_full)
```

OUTPUT

Analysis of Variance Table

Model 1: runoff_volume ~ lake_sabrina + pine_creek

Model 2: runoff_volume ~ lake_sabrina + pine_creek + rock_creek

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	40	0.150845				
2	39	0.092166	1	0.05868	24.83	1.321e-05 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Question 4

Which variable should be added to the model containing rock_creek?

Answer

lake_sabrina does not improve the model with only rock_creek ($\beta_{lake_sabrina} = 0$), so we can say that we should add pine_creek to the model containing rock_creek.

Remember: H0: No difference between the models, so choose the simplest H1: Full model is better

Question 5

Could things be even simpler? Perform a partial F-Test to see if a model containing rock_creek alone could be suitable.

CODE

```
s.mod_reduced3 <- lm(runoff_volume ~ rock_creek, data=stream_data)
anova(s.mod_reduced3, s.mod_full)
```

Answer

Fitting with only rock_creek does not improve model fit ($P < 0.05$) and so we can conclude that the better model is the one with pine_creek and rock_creek as predictors, with lake_sabrina removed.

Question 6

What is your optimal model?

Answer

The best model is: $runoff_{volume} = \beta_0 + \beta_1 rock_{creek} + \beta_2 pine_{creek} + error$

Review

- Simple linear regressions model the relationship between two variables
 - We can also make linear models with more than one predictor
- We can use histograms and correlation matrices to do some preliminary exploration of the data
- If any variables are skewed, we can transform them
- Looking at a correlation matrix to identify the best predictors (for both simple and multiple linear regression)
- Fit model using `lm()` function
- Check assumptions:
 - Collinearity (multiple linear regression only)
 - Linearity
 - Independence
 - Normality
 - Equal variance
- Use `summary()` to look at model output and interpret it
 - F-test : overall model significance
 - Coefficients table : individual predictors' significance
 - R^2 : How much variation in the data is explained by the model?

That's it for today! Great work fitting simple and multiple linear regression! Next week we jump into stepwise selection!

Attribution