

Supervised learning: Regression 1

Contents

Introduction	1
Regression in R	1
Plotting <code>lm()</code> in <code>ggplot</code>	4
Mean square error	5
Train-validation-test split	7
Programming exercise: cross-validation	9

Introduction

In this practical, you will learn how to perform regression analysis, how to plot with confidence and prediction intervals, how to calculate MSE, perform train-test splits, and write a function for cross validation.

Just like in the practical at the end of chapter 3 of the ISLR book, we will use the `Boston` dataset, which is in the `MASS` package that comes with R.

```
library(ISLR)
library(MASS)
library(tidyverse)
```

Regression in R

Regression is performed through the `lm()` function. It requires two arguments: a formula and data. A formula is a specific type of object that can be constructed like so:

```
some_formula <- outcome ~ predictor_1 + predictor_2
```

You can read it as “the outcome variable is a function of predictors 1 and 2”. As with other objects, you can check its class and even convert it to other classes, such as a character vector:

```
class(some_formula)
```

```
## [1] "formula"
```

```
as.character(some_formula)
```

```
## [1] "~" "outcome"
```

```
## [3] "predictor_1 + predictor_2"
```

You can estimate a linear model using `lm()` by specifying the outcome variable and the predictors in a formula and by inputting the dataset these variables should be taken from.

-
1. **Create a linear model object called `lm_ses` using the formula `medv ~ lstat` and the Boston dataset.**

You have now trained a regression model with `medv` (housing value) as the outcome/dependent variable and `lstat` (socio-economic status) as the predictor / independent variable.

Remember that a regression estimates β_0 (the intercept) and β_1 (the slope) in the following equation:

$$y = \beta_0 + \beta_1 \cdot x_1 + \epsilon$$

-
2. **Use the function `coef()` to extract the intercept and slope from the `lm_ses` object. Interpret the slope coefficient.**

-
3. **Use `summary()` to get a summary of the `lm_ses` object. What do you see? You can use the help file `?summary.lm`.**
-

We now have a model object `lm_ses` that represents the formula

$$\text{medv}_i = 34.55 - 0.95 * \text{lstat}_i + \epsilon_i$$

With this object, we can predict a new `medv` value by inputting its `lstat` value. The `predict()` method enables us to do this for the `lstat` values in the original dataset.

4. **Save the predicted y values to a variable called `y_pred`**

5. **Create a scatter plot with `y_pred` mapped to the x position and the true y value (`Boston$medv`) mapped to the y value. What do you see? What would this plot look like if the fit were perfect?**

We can also generate predictions from new data using the `newdat` argument in the `predict()` method. For that, we need to prepare a data frame with new values for the original predictors.

6. **Use the `seq()` function to generate a sequence of 1000 equally spaced values from 0 to 40. Store this vector in a data frame with (`data.frame()` or `tibble()`) as its column name `lstat`. Name the data frame `pred_dat`.**

7. **Use the newly created data frame as the `newdata` argument to a `predict()` call for `lm_ses`. Store it in a variable named `y_pred_new`.**

Plotting `lm()` in `ggplot`

A good way of understanding your model is by visualising it. We are going to walk through the construction of a plot with a fit line and prediction / confidence intervals from an `lm` object.

8. Create a scatter plot from the `Boston` dataset with `lstat` mapped to the x position and `medv` mapped to the y position. Store the plot in an object called `p_scatter`.

Now we're going to add a prediction line to this plot.

9. Add the vector `y_pred_new` to the `pred_dat` data frame with the name `medv`.

10. Add a `geom_line()` to `p_scatter`, with `pred_dat` as the data argument. What does this line represent?

11. The `interval` argument can be used to generate confidence or prediction intervals. Create a new object called `y_pred_95` using `predict()` (again with the `pred_dat` data) with the `interval` argument set to "confidence". What is in this object?

12. Create a data frame with 4 columns: `medv`, `lstat`, `lower`, and `upper`.

-
13. Add a `geom_ribbon()` to the plot with the data frame you just made. The ribbon geom requires three aesthetics: `x` (`lstat`, already mapped), `ymin` (lower), and `ymax` (upper). Add the ribbon below the `geom_line()` and the `geom_points()` of before to make sure those remain visible. Give it a nice colour and clean up the plot, too!

-
14. Explain in your own words what the ribbon represents.

-
15. Do the same thing, but now with the prediction interval instead of the confidence interval.

Mean square error

-
16. Write a function called `mse()` that takes in two vectors: true `y` values and predicted `y` values, and which outputs the mean square error.

Start like so:

```
mse <- function(y_true, y_pred) {  
  # your function here  
}
```

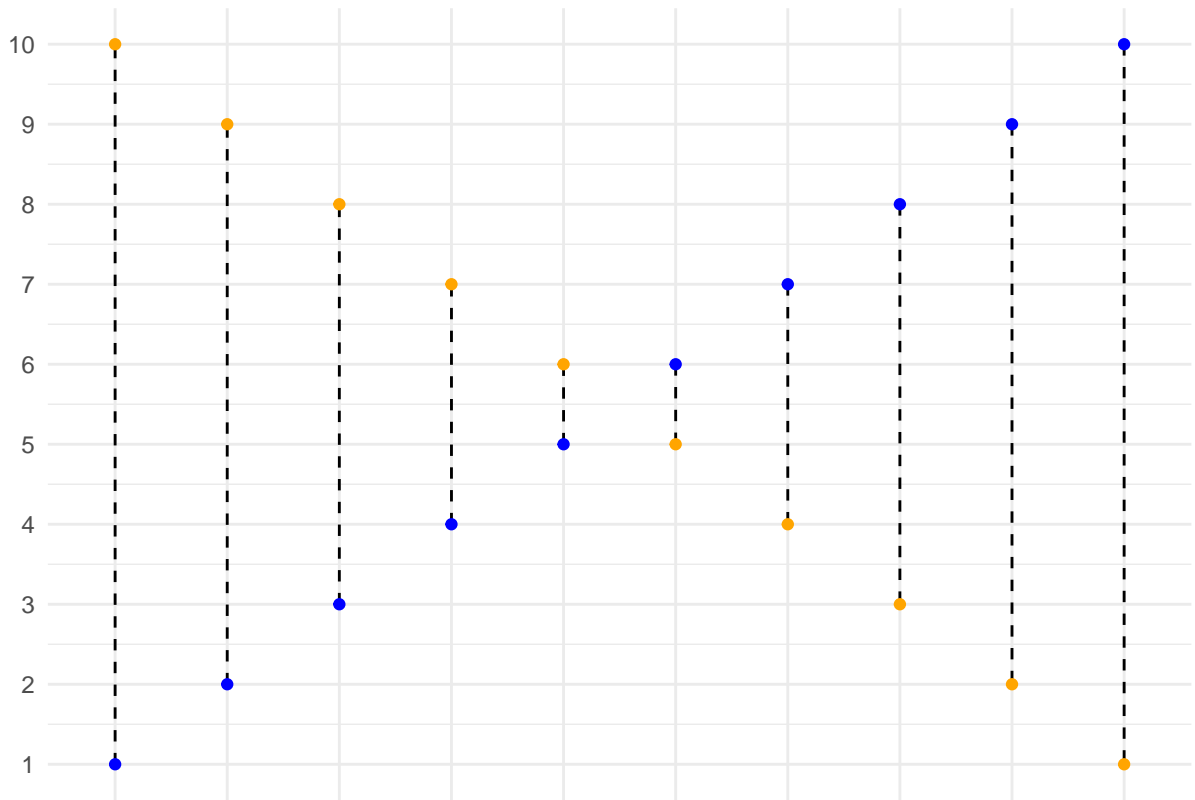
[Wikipedia](#) may help for the formula.

-
17. Make sure your `mse()` function works correctly by running the following code.
-

```
mse(1:10, 10:1)
```

```
## [1] 33
```

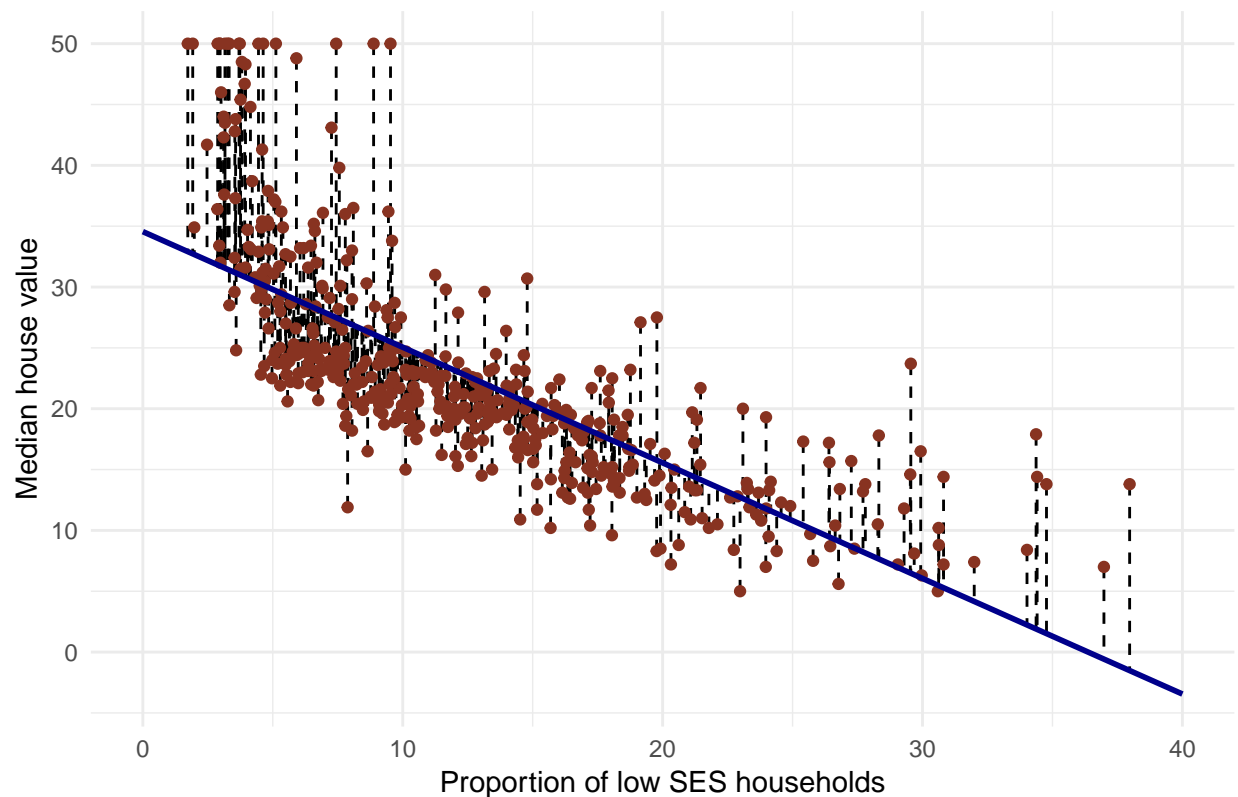
You have now calculated the mean squared length of the dashed lines below.



18. Calculate the mean square error of the `lm_ses` model. Use the `medv` column as `y_true` and use the `predict()` method to generate `y_pred`.

You have calculated the mean squared length of the dashed lines in the plot below.

Boston house prices: errors



Train-validation-test split

Now we will use the `sample()` function to randomly select observations from the `Boston` dataset to go into a training, test, and validation set. The training set will be used to fit our model, the validation set will be used to calculate the out-of-sample prediction error during model building, and the test set will be used to estimate the true out-of-sample MSE.

-
19. The `Boston` dataset has 506 observations. Use `c()` and `rep()` to create a vector with 253 times the word “train”, 152 times the word “validation”, and 101 times the word “test”. Call this vector `splits`.

-
20. Use the function `sample()` to randomly order this vector and add it to the `Boston` dataset using `mutate()`. Assign the newly created dataset to a variable called `boston_master`.

21. Now use `filter()` to create a training, validation, and test set from the `boston_master` data. Call these datasets `boston_train`, `boston_valid`, and `boston_test`.

We will set aside the `boston_test` dataset for now.

22. Train a linear regression model called `model_1` using the training dataset. Use the formula `medv ~ lstat` like in the first `lm()` exercise. Use `summary()` to check that this object is as you expect.

23. Calculate the MSE with this object. Save this value as `model_1_mse_train`.

24. Now calculate the MSE on the validation set and assign it to variable `model_1_mse_valid`. Hint: use the `newdata` argument in `predict()`.

This is the estimated out-of-sample mean squared error.

25. Create a second model `model_2` for the train data which includes `age` and `tax` as predictors. Calculate the train and validation MSE.

26. **Compare model 1 and model 2 in terms of their training and validation MSE. Which would you choose and why?**

27. **Calculate the test MSE for the model of your choice in the previous question. What does this number tell you?**

Programming exercise: cross-validation

This is an advanced exercise. Some components we have seen before in this and previous practicals, but some things will be completely new. Try to complete it by yourself, but don't worry if you get stuck. If you don't know about `for` loops in R, read up on those before you start the exercise.

Use help in this order:

- R help files
- Internet search & stack exchange
- Your peers
- The answer, which shows one solution

You may also just read the answer and try to understand what happens in each step.

28. **Create a function that performs k-fold cross-validation for linear models.**

Inputs:

- `formula`: a formula just as in the `lm()` function
- `dataset`: a data frame
- `k`: the number of folds for cross validation
- any other arguments you need necessary

Outputs:

- Mean square error averaged over folds
-

29. **Use your function to perform 9-fold cross validation with a linear model with as its formula $\text{medv} \sim \text{lstat} + \text{age} + \text{tax}$. Compare it to a model with as formula $\text{medv} \sim \text{lstat} + \text{I}(\text{lstat}^2) + \text{age} + \text{tax}$.**
