

MSiA 400 Lab Assignment 2

Oct 27, 2014

- Due: 11:59pm Nov 3, 2014
- This is an open book assignment.
- Please submit one report file that includes : short answer, related code and print for each problem if necessary.

Problem 1

Data set *bostonhousing.txt*, created by Harrison and Rubinfeld (1978), concerns housing values in suburbs of Boston. The attributes include

MEDV	Median value of owner-occupied homes in \$1000's
CRIM	per capita crime rate by town
ZN	proportion of residential land zoned for lots over 25,000 sq.ft.
INDUS	proportion of non-retail business acres per town
CHAS	Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)
NOX	nitric oxides concentration (parts per 10 million)
RM	average number of rooms per dwelling
AGE	proportion of owner-occupied units built prior to 1940
DIS	weighted distances to five Boston employment centres
RAD	index of accessibility to radial highways
TAX	full-value property-tax rate per \$10,000
PTRATIO	pupil-teacher ratio by town
B	$1000(Bk - 0.63)^2$ where Bk is the proportion of blacks by town
LSTAT	% lower status of the population,

in which MEDV is the response variable. The summary of the data set is below.

Name of the data set	bostonhousing
Number of observations	506
Number of attributes	14 (1 response variable and 13 explanatory variables)

problem 1(a)

Build regression model `reg` and display `summary()` of the model. Pick two explanatory variables that are least likely to be in the best model, and support your suggestion in one sentence.

problem 1(b)

Build regression model `reg.picked` by excluding the two explanatory variables selected in problem 1(a). Display `summary()` of the model.

problem 1(c)

For a regression model, the mean squared error (MSE) is defined as $\frac{SSE}{n-1-p}$, in which p is the number of explanatory variables used in the model. The mean absolute error (MAE) is similarly defined: $\frac{SAE}{n-1-p}$. Display *MSE* and *MAE* for regression models `reg` and `reg.picked` from the previous problems. Based on *MSE* and *MAE*, pick one model you prefer.

problem 1(d)

Run `step()` using regression model `reg` in problem 1(a). Compare the model with `reg.picked` in problem 1(b).

Problem 2

Import `labdata.txt`. The summary of the data set is below.

Name of the data set	labdata
Number of observations	400
Number of attributes	9 (1 response variable and 8 explanatory variables)

Column `y` is the response variable and remaining attributes `x1,x2,...` are the explanatory variables.

problem 2(a)

Build regression model `reg` and display `summary()` of the model

problem 2(b)

For each explanatory variable, plot it against the response variable. Based on the scatter plots, pick one variable that is most likely to be used in a piecewise regression model. Attach one plot associated with the variable you pick.

problem 2(c)

Calculate the mean of the variable you pick in problem 2(b) and build piecewise regression model `reg.piece` using the mean. Is model `reg.piece` better than model `reg` in problem 2(a)? Support your argument in one sentence.

References

Harrison, D. and Rubinfeld, D.L. (1978) Hedonic prices and the demand for clean air. *Economics & Management*, **5**, 81–102.