

Regularization and logistic regression

Question 1

Suppose we fit “Lasso Regression” to a data set, which has 100 features (X_1, X_2, \dots, X_{100}). Now, we rescale one of these feature by multiplying with 10 (say that feature is X_1), and then refit Lasso regression with the same regularization parameter.

Now, which of the following option will be correct?

- A. It is more likely for X_1 to be excluded from the model
- B. It is more likely for X_1 to be included in the model
- C. Can't say
- D. None of these

Question 2

Suppose you have fitted a multiple regression model on a dataset. Now, you are using Ridge regression with tuning parameter λ to reduce its complexity. Choose the options below which describes relationship of bias and variance with λ .

- A. In case of very large λ ; bias is low, variance is low
- B. In case of very large λ ; bias is low, variance is high
- C. In case of very large λ ; bias is high, variance is low
- D. In case of very large λ ; bias is high, variance is high

Question 3

Write a function to realize gradient descent in R. Understand how learning rate affects convergence.

Question 4

A five year follow-up study on 600 disease free subjects was carried out to assess the effect of whether having exposure E or not (of smoking for example) on the development (or not) of a certain disease. The variables AGE (continuous) and obesity status (boolean), which were determined at the start of the follow-up and were to be considered as control variables in analyzing the data.

- (1) State the logit form of a logistic regression model that assesses the effect of the 0/1 exposure variable E controlling for the confounding effects of AGE and OBS and the interaction effects of AGE with E and OBS with E .
- (2) Given above model you have, give a formula for the odds ratio for the exposure-disease relationship that controls for the confounding and interactive effects of AGE and OBS .
- (3) Now use the formula from above to write an expression for the estimated odds ratio for the exposure-disease relationship when $AGE=40$ and $OBS=1$.

Question 5

Build the best logistic regression model to predict loan will be default (delay) or not. Add regularization to control for multicollinearity.