

**7.3 Predicting Housing Median Prices.** The file *BostonHousing.csv* contains information on over 500 census tracts in Boston, where for each tract multiple variables are recorded. The last column (CAT.MEDV) was derived from MEDV, such that it obtains the value 1 if  $\text{MEDV} > 30$  and 0 otherwise. Consider the goal of predicting the median value (MEDV) of a tract, given the information in the first 12 columns.

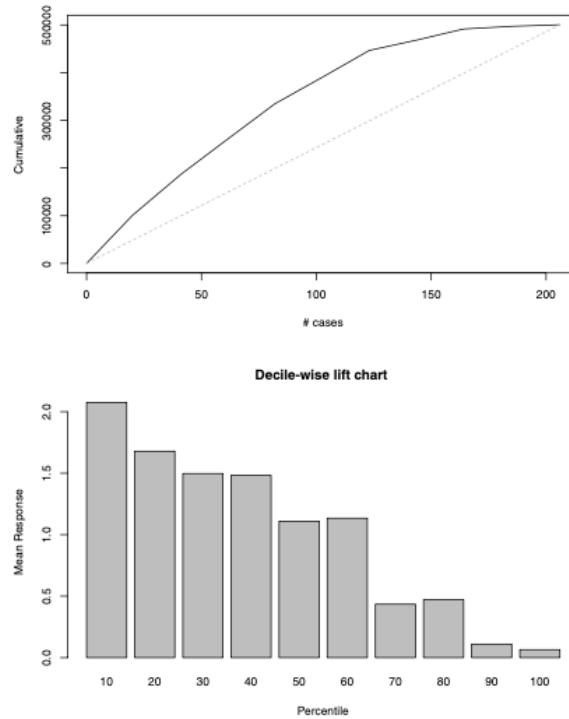
Partition the data into training (60%) and validation (40%) sets.

- a. Perform a  $k$ -NN prediction with all 12 predictors (ignore the CAT.MEDV column), trying values of  $k$  from 1 to 5. Make sure to normalize the data, and choose function `knn()` from package `class` rather than package `FNN`. To make sure R is using the `class` package (when both packages are loaded), use `class::knn()`. What is the best  $k$ ? What does it mean?
- b. Predict the MEDV for a tract with the following information, using the best  $k$ :

CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	LSTAT
0.2	0	7	0	0.538	6	62	4.7	4	307	21	10

- c. If we used the above  $k$ -NN algorithm to score the training data, what would be the error of the training set?
- d. Why is the validation data error overly optimistic compared to the error rate when applying this  $k$ -NN predictor to new data?
- e. If the purpose is to predict MEDV for several thousands of new tracts, what would be the disadvantage of using  $k$ -NN prediction? List the operations that the algorithm goes through in order to produce each prediction.

- 5.6** A firm that sells software services has been piloting a new product and has records of 500 customers who have either bought the services or decided not to. The target value is the estimated profit from each sale (excluding sales costs). The global mean is \$2128. However, the cost of the sales effort is not cheap—the company figures it comes to \$2500 for each of the 500 customers (whether they buy or not). The firm developed a predictive model in hopes of being able to identify the top spenders in the future. The lift and decile charts for the validation set are shown in Figure 5.13.



**FIGURE 5.13** LIFT AND DECILE-WISE LIFT CHARTS FOR SOFTWARE SERVICES PRODUCT SALES

- If the company begins working with a new set of 1000 leads to sell the same services, similar to the 500 in the pilot study, without any use of predictive modeling to target sales efforts, what is the estimated profit?
- If the firm wants the average profit on each sale to at least double the sales effort cost, and applies an appropriate cutoff with this predictive model to a new set of 1000 leads, how far down the new list of 1000 should it proceed (how many deciles)?
- Still considering the new list of 1000 leads, if the company applies this predictive model with a lower cutoff of \$2500, how far should it proceed down the ranked leads, in terms of deciles?

**10.1 Financial Condition of Banks.** The file *Banks.csv* includes data on a sample of 20 banks. The “Financial Condition” column records the judgment of an expert on the financial condition of each bank. This outcome variable takes one of two possible values—*weak* or *strong*—according to the financial condition of the bank. The predictors are two ratios used in the financial analysis of banks:  $\text{TotLns\&Lses}/\text{Assets}$  is the ratio of total loans and leases to total assets and  $\text{TotExp}/\text{Assets}$  is the ratio of total expenses to total assets. The target is to use the two ratios for classifying the financial condition of a new bank.

Run a logistic regression model (on the entire dataset) that models the status of a bank as a function of the two financial measures provided. Specify the *success* class as *weak* (this is similar to creating a dummy that is 1 for financially weak banks and 0 otherwise), and use the default cutoff value of 0.5.

- a. Write the estimated equation that associates the financial condition of a bank with its two predictors in three formats:
  - i. The logit as a function of the predictors
  - ii. The odds as a function of the predictors
  - iii. The probability as a function of the predictors
- b. Consider a new bank whose total loans and leases/assets ratio = 0.6 and total expenses/assets ratio = 0.11. From your logistic regression model, estimate the following four quantities for this bank (use R to do all the intermediate calculations; show your final answers to four decimal places): the logit, the odds, the probability of being financially weak, and the classification of the bank (use cutoff = 0.5).
- c. The cutoff value of 0.5 is used in conjunction with the probability of being financially weak. Compute the threshold that should be used if we want to make a classification based on the odds of being financially weak, and the threshold for the corresponding logit.
- d. Interpret the estimated coefficient for the total loans & leases to total assets ratio ( $\text{TotLns\&Lses}/\text{Assets}$ ) in terms of the odds of being financially weak.
- e. When a bank that is in poor financial condition is misclassified as financially strong, the misclassification cost is much higher than when a financially strong bank is misclassified as weak. To minimize the expected cost of misclassification, should the cutoff value for classification (which is currently at 0.5) be increased or decreased?