

Python Companion to ISLR

Naresh Gurbuxani

May 7, 2019

Contents

1	Introduction	1
2	Statistical Learning	4

1 Introduction

Figure 1 shows graphs of Wage versus three variables.

Figure 2 shows boxplots of previous days' percentage changes in S&P 500 grouped according to today's change Up or Down.

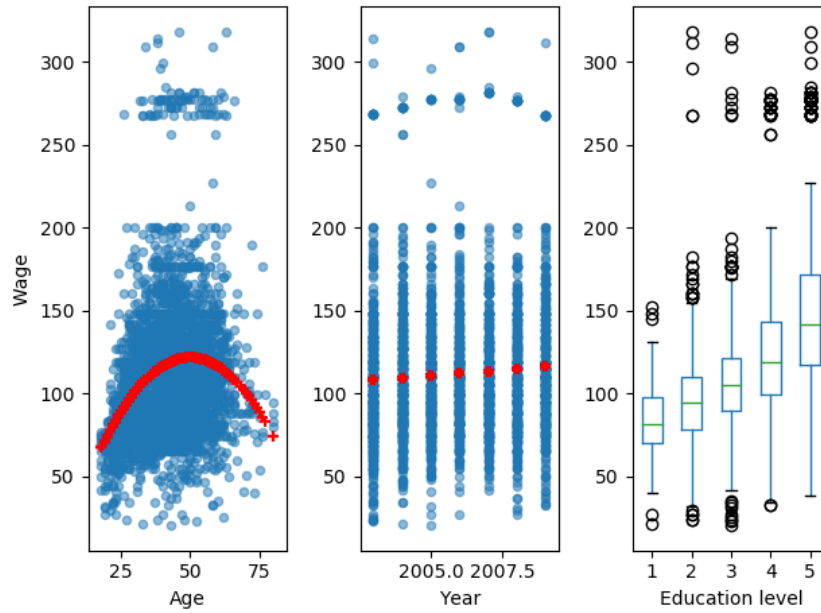


Figure 1: **Wage** data, which contains income survey information for males from the central Atlantic region of the United States. Left: **wage** as a function of **age**. On average, **wage** increases with **age** until about 60 years of age, at which point it begins to decline. Center: **wage** as a function of **year**. There is a slow but steady increase of approximately \$10,000 in the average **wage** between 2003 and 2009. Right: Boxplots displaying **wage** as a function of **education**, with 1 indicating the lowest level (no highschool diploma) and 5 the highest level (an advanced graduate degree). On average, **wage** increases with the level of **education**.

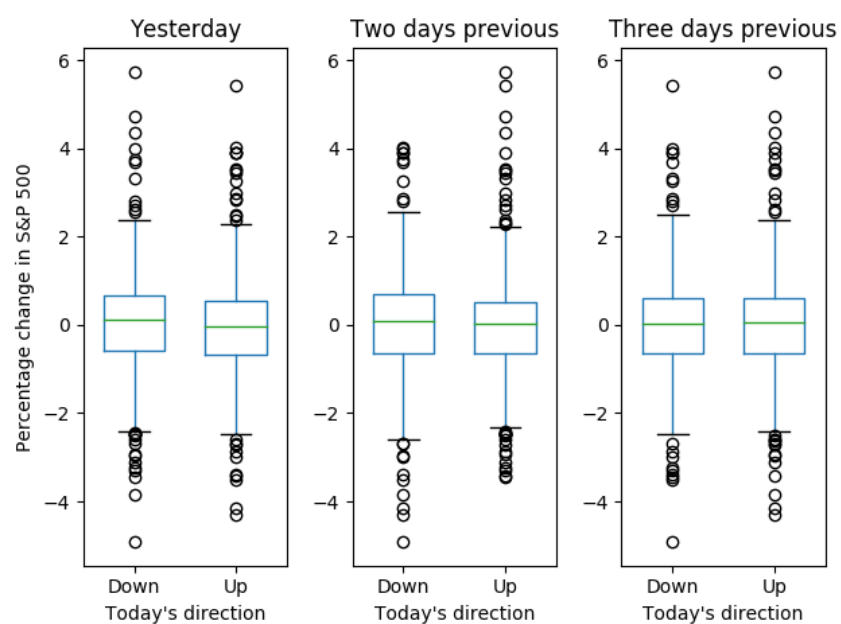


Figure 2: Left: Boxplots of the previous day's percentage change in the S&P 500 index for the days for which the market increased or decreased, obtained from the **Smarket** data. Center and Right: Same as left panel, but the percentage changes for two and three days previous are shown.

2 Statistical Learning

Figure 3 shows scatter plots of `sales` versus `TV`, `radio`, and `newspaper` advertising. In each panel, the figure also includes an OLS regression line.

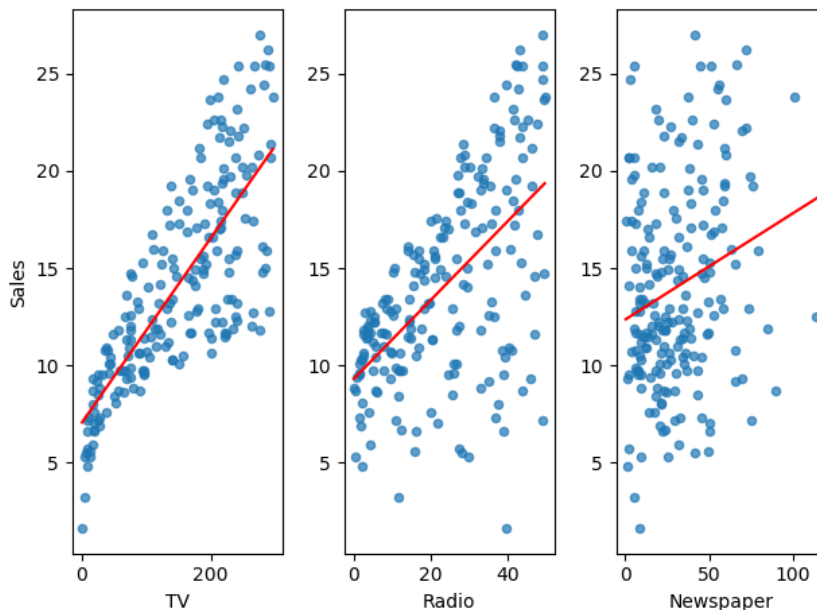


Figure 3: The `Advertising` data set. The plot displays `sales`, in thousands of units, as a function of `TV`, `radio`, and `newspaper` budgets, in thousands of dollars, for 200 different markets. In each plot we show the simple least squares fit of `sales` to that variable. In other words, each red line represents a simple model that can be used to predict `sales` using `TV`, `radio`, and `newspaper`, respectively.

Figure 4 is a plot of `Income` versus `Years of Education` from the `Income` data set. In the left panel, the “true” function (given by blue line) is actually my guess.

Figure 5 is a plot of `Income` versus `Years of Education` and `Seniority` from the `Income` data set. Since the book does not provide the true values of `Income`, “true” values shown in the plot are actually third order polynomial fit.

Figure 6 shows an example of the parametric approach applied to the `Income` data from previous figure.

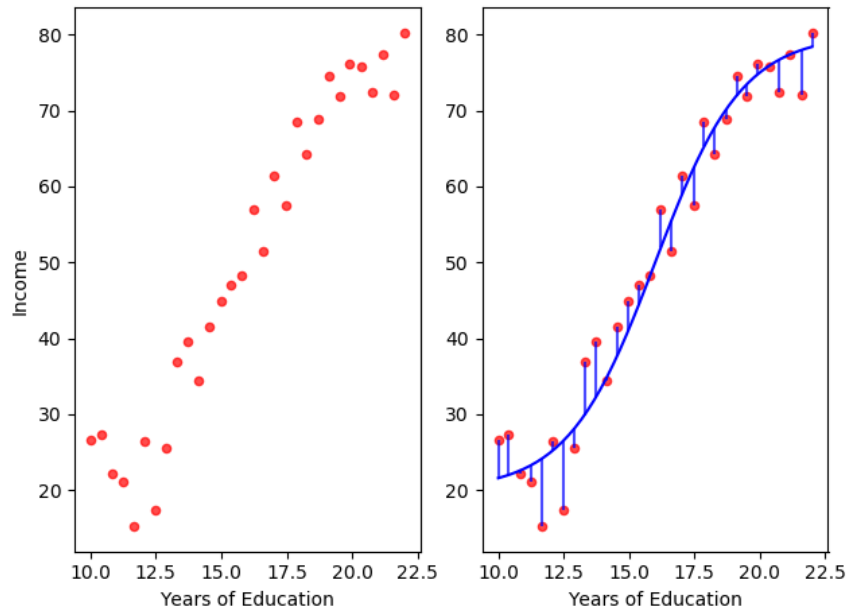


Figure 4: The `Income` data set. Left: The red dots are the observed values of `income` (in tens of thousands of dollars) and `years of education` for 30 individuals. Right: The blue curve represents the true underlying relationship between `income` and `years of education`, which is generally unknown (but is known in this case because the data are simulated). The vertical lines represent the error associated with each observation. Note that some of the errors are positive (when an observation lies above the blue curve) and some are negative (when an observation lies below the curve). Overall, these errors have approximately mean zero.

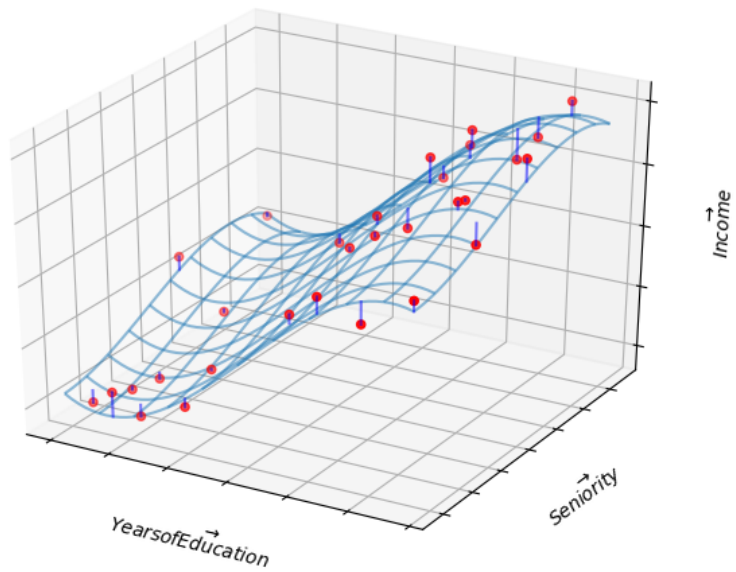


Figure 5: The plot displays `income` as a function of `years of education` and `seniority` in the `Income` data set. The blue surface represents the true underlying relationship between `income` and `years of education` and `seniority`, which is known since the data are simulated. The red dots indicate the observed values of these quantities for 30 individuals.

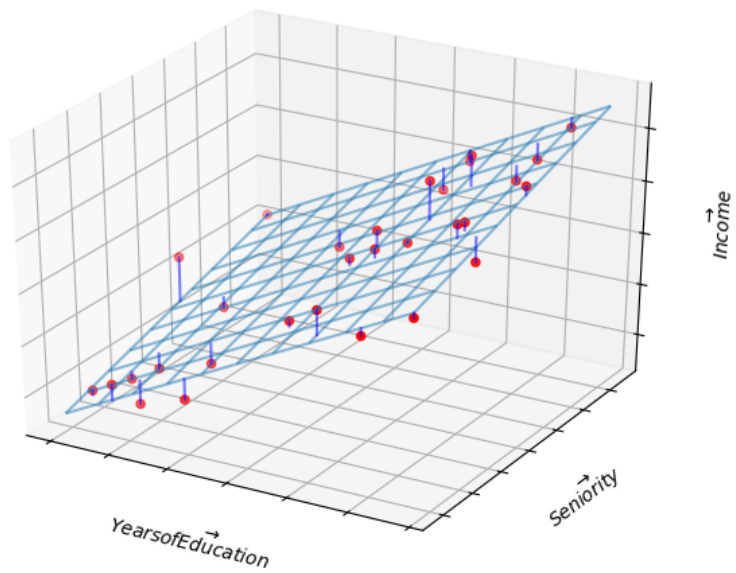


Figure 6: A linear model fit by least squares to the `Income` data from figure 5. The observations are shown in red, and the blue plane indicates the least squares fit to the data.

Figure 7 provides an illustration of the trade-off between flexibility and interpretability for some of the methods covered in this book.

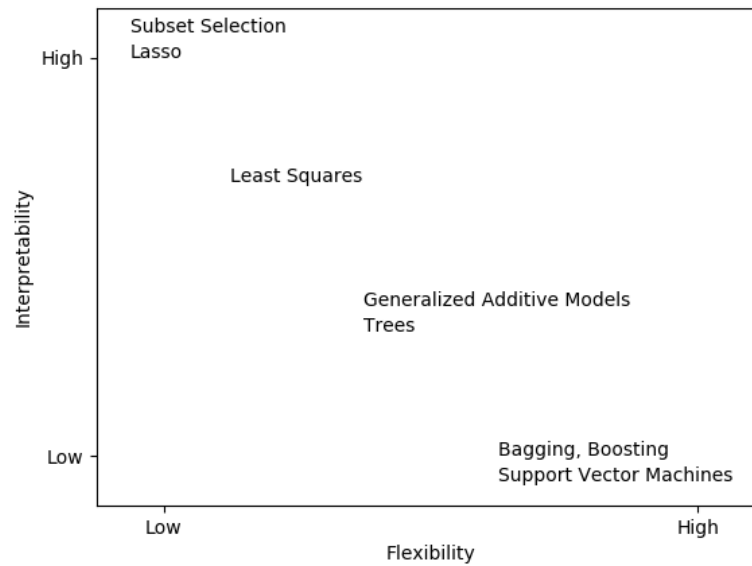


Figure 7: A representation of the tradeoff between flexibility and interpretability, using different statistical learning methods. In general, as the flexibility of a method increases, its interpretability decreases.

Figure 8 provides a simple illustration of the clustering problem.

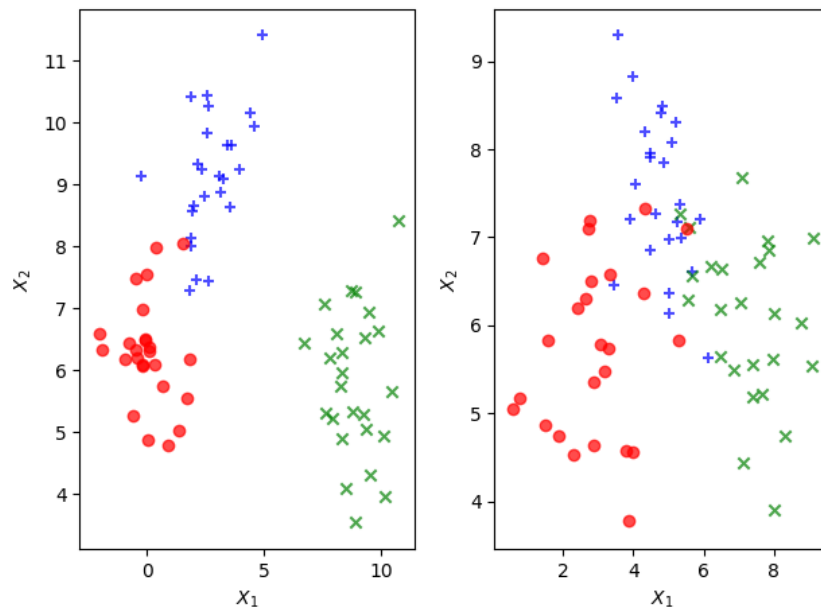


Figure 8: A clustering data set involving three groups. Each group is shown using a different colored symbol. Left: The three groups are well-separated. In this setting, a clustering approach should successfully identify the three groups. Right: There is some overlap among the groups. Now the clustering task is more challenging.