

Chapter 3 – Linear regression

Wenjing Liao

School of Mathematics, Georgia Institute of Technology

Math 4803

Fall 2019

Outline

- 1 3.1 Simple linear regression
- 2 3.2 Multiple linear regression
- 3 Review – linear system of equations
- 4 3.3 Considerations in the regression model
- 5 3.5 Comparison of linear regression with KNN regression

Simple linear regression

Variables: (X, Y) where $X, Y \in \mathbb{R}$

Linear relation:

$$Y \approx \beta_0 + \beta_1 X$$

$$\text{sales} \approx \beta_0 + \beta_1 \times \text{TV}$$

How to estimate β_0 and β_1 ?

Advertising data: sales and TV

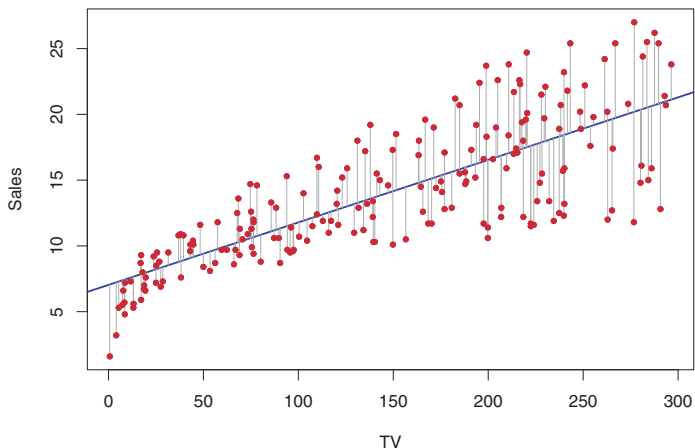


FIGURE 3.1. For the Advertising data, the least squares fit for the regression of sales onto TV is shown. The fit is found by minimizing the sum of squared errors. Each grey line segment represents an error, and the fit makes a compromise by averaging their squares. In this case a linear fit captures the essence of the relationship, although it is somewhat deficient in the left of the plot.

Estimating the coefficients

Data: $(x_i, y_i), i = 1, \dots, n$

Residual sum of squares: With coefficients $\hat{\beta}_0, \hat{\beta}_1$

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

$$e_i = y_i - \hat{y}_i$$

$$\text{RSS} = e_1^2 + \dots + e_n^2$$

Coefficient estimation: $(\hat{\beta}_0, \hat{\beta}_1)$ minimizes RSS

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad (3.4)$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x},$$

where $\bar{y} \equiv \frac{1}{n} \sum_{i=1}^n y_i$ and $\bar{x} \equiv \frac{1}{n} \sum_{i=1}^n x_i$ are the sample means. In other words, (3.4) defines the *least squares coefficient estimates* for simple linear regression.

Optimization

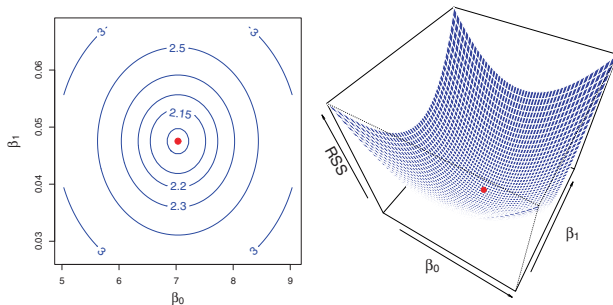


FIGURE 3.2. Contour and three-dimensional plots of the RSS on the Advertising data, using sales as the response and TV as the predictor. The red dots correspond to the least squares estimates $\hat{\beta}_0$ and $\hat{\beta}_1$, given by (3.4).

Outline

- 1 3.1 Simple linear regression
- 2 3.2 Multiple linear regression**
- 3 Review – linear system of equations
- 4 3.3 Considerations in the regression model
- 5 3.5 Comparison of linear regression with KNN regression

Simple linear regression

Variables: (X, Y) where $X \in \mathbb{R}^p$ and $Y \in \mathbb{R}$

Linear relation:

$$Y \approx \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$

$$\text{sales} = \beta_0 + \beta_1 \times \text{TV} + \beta_2 \times \text{radio} + \beta_3 \times \text{newspaper} + \epsilon.$$

How to estimate the coefficients?

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \dots + \hat{\beta}_p x_p.$$

RSS:

$$\begin{aligned} \text{RSS} &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\ &= \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \hat{\beta}_2 x_{i2} - \dots - \hat{\beta}_p x_{ip})^2. \end{aligned}$$

Example

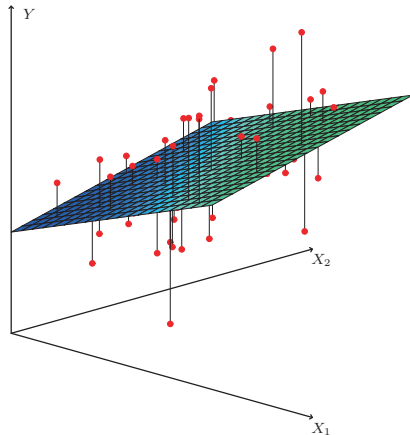


FIGURE 3.4. In a three-dimensional setting, with two predictors and one response, the least squares regression line becomes a plane. The plane is chosen to minimize the sum of the squared vertical distances between each observation (shown in red) and the plane.

Outline

- 1 3.1 Simple linear regression
- 2 3.2 Multiple linear regression
- 3 Review – linear system of equations**
- 4 3.3 Considerations in the regression model
- 5 3.5 Comparison of linear regression with KNN regression

Outline

- 1 3.1 Simple linear regression
- 2 3.2 Multiple linear regression
- 3 Review – linear system of equations
- 4 3.3 Considerations in the regression model**
- 5 3.5 Comparison of linear regression with KNN regression

Collinearity

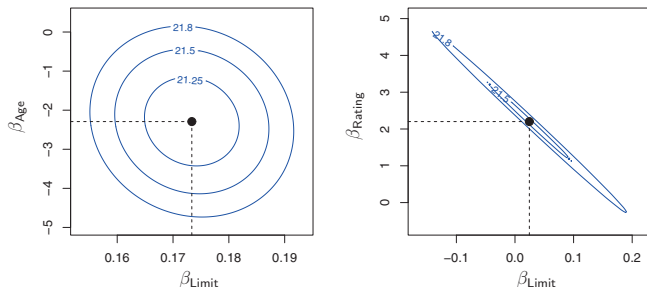
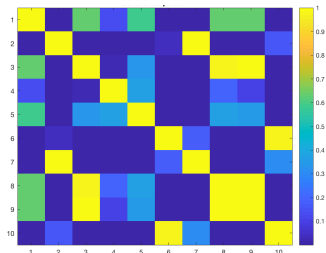


FIGURE 3.15. Contour plots for the RSS values as a function of the parameters β for various regressions involving the **Credit** data set. In each plot, the black dots represent the coefficient values corresponding to the minimum RSS. Left: A contour plot of RSS for the regression of **balance** onto **age** and **limit**. The minimum value is well defined. Right: A contour plot of RSS for the regression of **balance** onto **rating** and **limit**. Because of the collinearity, there are many pairs $(\beta_{\text{Limit}}, \beta_{\text{Rating}})$ with a similar value for RSS.

Correlation matrix

- Correlation matrix helps to detect the collinearity between two columns of A .

$$\text{Cor}(X, Y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}},$$



- Correlation matrix does not help to detect multicollinearity, for example, $\vec{v}_1 + 2\vec{v}_2 - \vec{v}_3 = 0$.

How to handle collinearity?

Reference: <https://en.wikipedia.org/wiki/Multicollinearity>

- Drop one of the variables
- Obtain more data, if possible
- Mean-center the predictor variables
- Ridge regression

$$\min_x \|Ax - b\|_2^2 + \|x\|_2^2$$

Qualitative variables

Credit data set

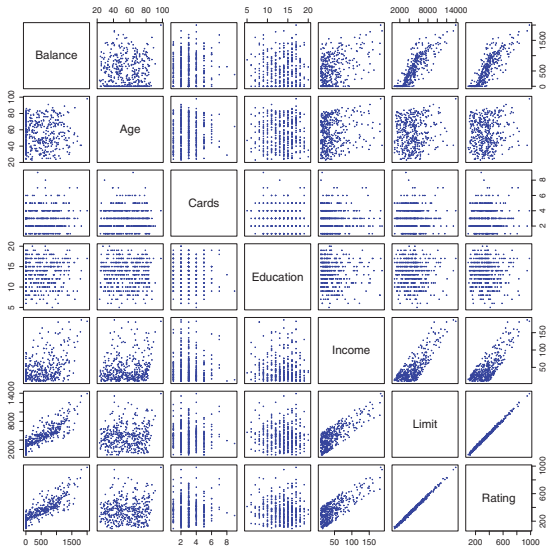


FIGURE 3.6. The Credit data set contains information about balance, age, cards, education, income, limit, and rating for a number of potential customers.

Gender variable

Quantify gender:

$$x_i = \begin{cases} 1 & \text{if } i\text{th person is female} \\ 0 & \text{if } i\text{th person is male,} \end{cases}$$

Credit balance versus gender:

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i = \begin{cases} \beta_0 + \beta_1 + \epsilon_i & \text{if } i\text{th person is female} \\ \beta_0 + \epsilon_i & \text{if } i\text{th person is male.} \end{cases}$$

Results: $\beta_0 = 509.80, \beta_1 = 17.73$

Average debt for males $\beta_0 = 509.80$

Average debt for females $\beta_0 + \beta_1 = 509.80 + 19.73 = 529.53$

An alternative way

Quantify gender:

$$x_i = \begin{cases} 1 & \text{if } i\text{th person is female} \\ -1 & \text{if } i\text{th person is male} \end{cases}$$

Credit balance versus gender:

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i = \begin{cases} \beta_0 + \beta_1 + \epsilon_i & \text{if } i\text{th person is female} \\ \beta_0 - \beta_1 + \epsilon_i & \text{if } i\text{th person is male.} \end{cases}$$

Results: $\beta_0 = 519.665$, $\beta_1 = 9.865$

Average debt for males $\beta_0 - \beta_1 = 509.80$

Average debt for females $\beta_0 + \beta_1 = 529.53$

Qualitative variables with more than two levels

$$x_{i1} = \begin{cases} 1 & \text{if } i\text{th person is Asian} \\ 0 & \text{if } i\text{th person is not Asian,} \end{cases}$$

$$x_{i2} = \begin{cases} 1 & \text{if } i\text{th person is Caucasian} \\ 0 & \text{if } i\text{th person is not Caucasian.} \end{cases}$$

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \epsilon_i = \begin{cases} \beta_0 + \beta_1 + \epsilon_i & \text{if } i\text{th person is Asian} \\ \beta_0 + \beta_2 + \epsilon_i & \text{if } i\text{th person is Caucasian} \\ \beta_0 + \epsilon_i & \text{if } i\text{th person is African American.} \end{cases}$$

Extensions of the linear model: incorporating interaction terms

Linear model:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon.$$

Incorporating product terms:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + \epsilon.$$

$$\begin{aligned} Y &= \beta_0 + (\beta_1 + \beta_3 X_2) X_1 + \beta_2 X_2 + \epsilon \\ &= \beta_0 + \tilde{\beta}_1 X_1 + \beta_2 X_2 + \epsilon \end{aligned}$$

Incorporating interaction terms

Linear model: $\text{balance} \approx \beta_0 + \beta_1 \times \text{income} + \beta_2 \times \text{student}$

$$\begin{aligned} \text{balance}_i &\approx \beta_0 + \beta_1 \times \text{income}_i + \begin{cases} \beta_2 & \text{if } i\text{th person is a student} \\ 0 & \text{if } i\text{th person is not a student} \end{cases} \\ &= \beta_1 \times \text{income}_i + \begin{cases} \beta_0 + \beta_2 & \text{if } i\text{th person is a student} \\ \beta_0 & \text{if } i\text{th person is not a student.} \end{cases} \end{aligned}$$

Incorporating product terms:

$\text{balance} \approx \beta_0 + \beta_1 \times \text{income} + \beta_2 \times \text{student} + \beta_3 \times \text{income} \times \text{student}$

$$\begin{aligned} \text{balance}_i &\approx \beta_0 + \beta_1 \times \text{income}_i + \begin{cases} \beta_2 + \beta_3 \times \text{income}_i & \text{if student} \\ 0 & \text{if not student} \end{cases} \\ &= \begin{cases} (\beta_0 + \beta_2) + (\beta_1 + \beta_3) \times \text{income}_i & \text{if student} \\ \beta_0 + \beta_1 \times \text{income}_i & \text{if not student} \end{cases} \end{aligned}$$

Incorporating interaction terms

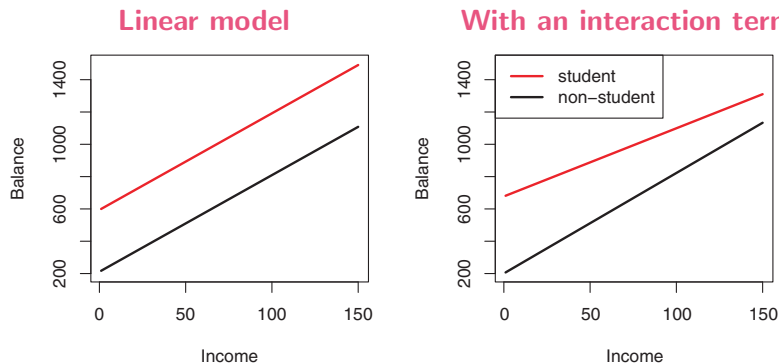


FIGURE 3.7. For the **Credit** data, the least squares lines are shown for prediction of **balance** from **income** for students and non-students. Left: The model (3.34) was fit. There is no interaction between **income** and **student**. Right: The model (3.35) was fit. There is an interaction term between **income** and **student**.

Nonlinear relations

Linear: $Y = \beta_0 + \beta_1 X$

Quadratic: $Y = \beta_0 + \beta_1 X + \beta_2 X^2$

Polynomials of degree p : $Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \dots + \beta_p X^p$

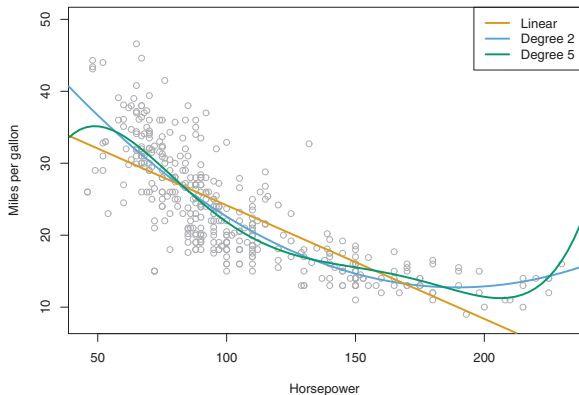


FIGURE 3.8. The `Auto` data set. For a number of cars, `mpg` and `horsepower` are shown. The linear regression fit is shown in orange. The linear regression fit for a model that includes `horsepower`² is shown as a blue curve. The linear regression fit for a model that includes all polynomials of `horsepower` up to fifth-degree is shown in green.

Outline

- 1 3.1 Simple linear regression
- 2 3.2 Multiple linear regression
- 3 Review – linear system of equations
- 4 3.3 Considerations in the regression model
- 5 3.5 Comparison of linear regression with KNN regression**

KNN regression

At X_0 , $\hat{Y}_0 = \frac{1}{K} \sum_{i \in \mathcal{N}_0} Y_i$ where \mathcal{N}_0 contains the K points in the training data that are closest to X_0

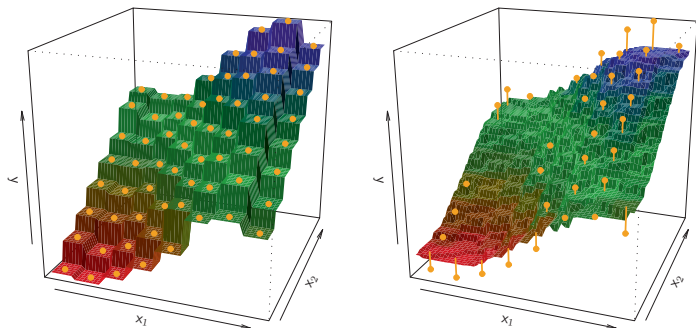


FIGURE 3.16. Plots of $\hat{f}(X)$ using KNN regression on a two-dimensional data set with 64 observations (orange dots). Left: $K = 1$ results in a rough step function fit. Right: $K = 9$ produces a much smoother fit.

KNN regression

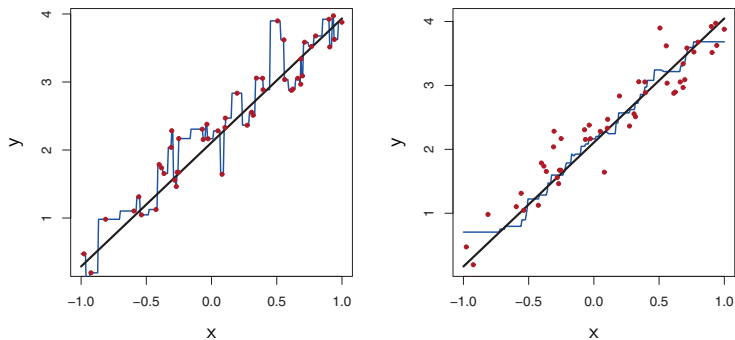


FIGURE 3.17. Plots of $\hat{f}(X)$ using KNN regression on a one-dimensional data set with 100 observations. The true relationship is given by the black solid line. Left: The blue curve corresponds to $K = 1$ and interpolates (i.e. passes directly through) the training data. Right: The blue curve corresponds to $K = 9$, and represents a smoother fit.

Linear regression

- Linear regression works well if the underlying function is indeed linear.

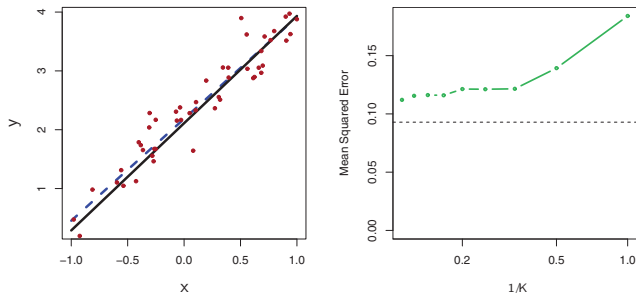


FIGURE 3.18. The same data set shown in Figure 3.17 is investigated further. Left: The blue dashed line is the least squares fit to the data. Since $f(X)$ is in fact linear (displayed as the black line), the least squares regression line provides a very good estimate of $f(X)$. Right: The dashed horizontal line represents the least squares test set MSE, while the green solid line corresponds to the MSE for KNN as a function of $1/K$ (on the log scale). Linear regression achieves a lower test MSE than does KNN regression, since $f(X)$ is in fact linear. For KNN regression, the best results occur with a very large value of K , corresponding to a small value of $1/K$.

Linear regression

- Linear regression may not work well if the function is nonlinear.

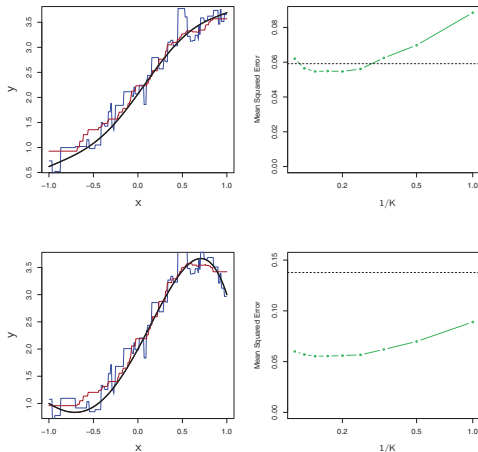


FIGURE 3.19. Top Left: In a setting with a slightly non-linear relationship between X and Y (solid black line), the KNN fits with $K = 1$ (blue) and $K = 9$ (red) are displayed. Top Right: For the slightly non-linear data, the test set MSE for least squares regression (horizontal black) and KNN with various values of $1/K$ (green) are displayed. Bottom Left and Bottom Right: As in the top panel, but with a strongly non-linear relationship between X and Y .

Curse of dimensionality for KNN regression

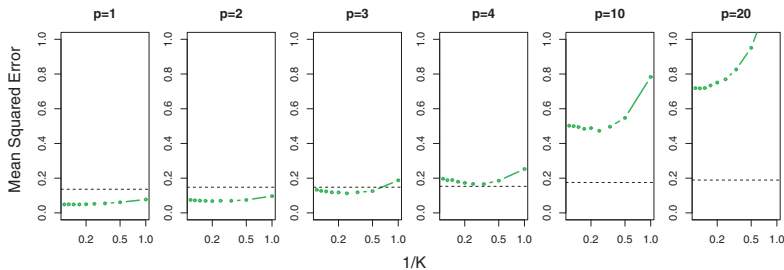


FIGURE 3.20. Test MSE for linear regression (black dashed lines) and KNN (green curves) as the number of variables p increases. The true function is non-linear in the first variable, as in the lower panel in Figure 3.19, and does not depend on the additional variables. The performance of linear regression deteriorates slowly in the presence of these additional noise variables, whereas KNN's performance degrades much more quickly as p increases.

Reference

Textbook: James, Gareth, Daniela Witten, Trevor Hastie and Robert Tibshirani, An introduction to statistical learning. Vol. 112, New York: Springer, 2013