

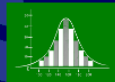
Introduction to Data Science

Winter Semester 2018/19

Oliver Ernst

TU Chemnitz, Fakultät für Mathematik, Professur Numerische Mathematik

Lecture Slides



Contents I

① What is Data Science?

② Learning Theory

2.1 What is Statistical Learning?

2.2 Assessing Model Accuracy

③ Linear Regression

3.1 Simple Linear Regression

3.2 Multiple Linear Regression

3.3 Other Considerations in the Regression Model

3.4 Revisiting the Marketing Data Questions

3.5 Linear Regression vs. K -Nearest Neighbors

④ Classification

4.1 Overview of Classification

4.2 Why Not Linear Regression?

4.3 Logistic Regression

4.4 Linear Discriminant Analysis

4.5 A Comparison of Classification Methods

⑤ Resampling Methods

Contents II

5.1 Cross Validation

5.2 The Bootstrap

6 Linear Model Selection and Regularization

6.1 Subset Selection

6.2 Shrinkage Methods

6.3 Dimension Reduction Methods

6.4 Considerations in High Dimensions

6.5 Miscellanea

7 Nonlinear Regression Models

7.1 Polynomial Regression

7.2 Step Functions

7.3 Regression Splines

7.4 Smoothing Splines

7.5 Generalized Additive Models

8 Tree-Based Methods

8.1 Decision Tree Fundamentals

8.2 Bagging, Random Forests and Boosting

Contents III

9 Support Vector Machines

- 9.1 Maximal Margin Classifier
- 9.2 Support Vector Classifiers
- 9.3 Support Vector Machines
- 9.4 SVMs with More than Two Classes
- 9.5 Relationship to Logistic Regression

10 Unsupervised Learning

- 10.1 Principal Components Analysis
- 10.2 Clustering Methods

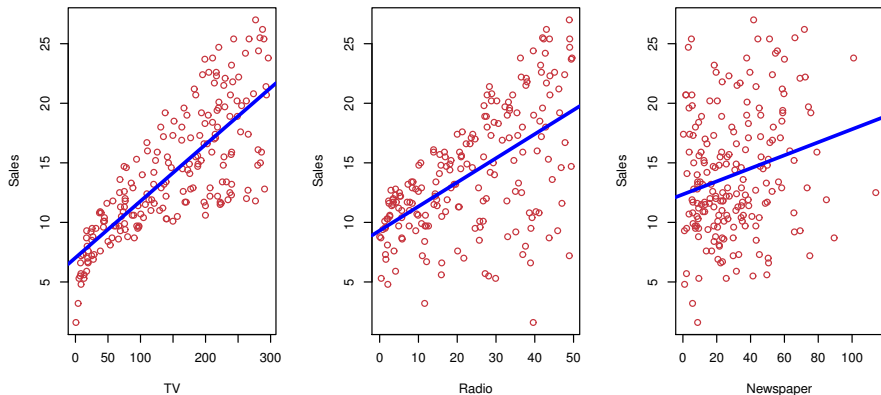
③ Linear Regression

- 3.1 Simple Linear Regression
- 3.2 Multiple Linear Regression
- 3.3 Other Considerations in the Regression Model
- 3.4 Revisiting the Marketing Data Questions
- 3.5 Linear Regression vs. K -Nearest Neighbors

Linear Regression

Advertising again

Recall advertising data set from Slide 28:



We will use the simple and well-established statistical learning technique known as **linear regression** to answer the following questions:

Linear Regression

Questions about advertising data set

- ❶ Is there a relationship between advertising budget and sales?
Otherwise, why bother?
- ❷ How strong is this relationship between advertising budget and sales?
Prediction possibly better than random guess?
- ❸ Which media contribute to sales?
Separate individual contributions
- ❹ How accurately can we estimate the effect of each medium on sales?
Euro by Euro?
- ❺ How accurately can we predict future sales?
Precise prediction for each medium?
- ❻ Is the relationship linear?
If yes, linear regression appropriate (possibly after transforming data)
- ❼ Is there synergy among the advertising media?
*Called **interaction effect** in statistics.*

③ Linear Regression

3.1 Simple Linear Regression

3.2 Multiple Linear Regression

3.3 Other Considerations in the Regression Model

3.4 Revisiting the Marketing Data Questions

3.5 Linear Regression vs. K -Nearest Neighbors

Simple Linear Regression

Definition, terminology, notation

Linear model for quantitative response Y of single predictor X :

$$Y \approx \beta_0 + \beta_1 X. \quad (3.1)$$

Statistician: "We are regressing Y onto X ."

E.g., with predictor **TV** advertising and response **sales**,

$$\mathbf{sales} \approx \beta_0 + \beta_1 \times \mathbf{TV}.$$

The values of **coefficients** or **parameters** β_0, β_1 obtained from fitting to the training data are denoted by $\hat{\beta}_0, \hat{\beta}_1$, leading to the prediction values

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x \quad (3.2)$$

when $X = x$, where the hat on \hat{y} denotes the predicted value of the response.

Simple Linear Regression

Estimating the coefficients

Determining **intercept** $\hat{\beta}_0$ and **slope** $\hat{\beta}_1$ in (3.1) amounts to choosing these parameters such that the **residuals** or **data misfits**

$$r_i := y_i - \hat{y}_i = y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i), \quad i = 1, \dots, n,$$

are minimized.

There are many options for defining smallness here, in **least squares estimation** this is measured by the **residual sum of squares (RSS)**

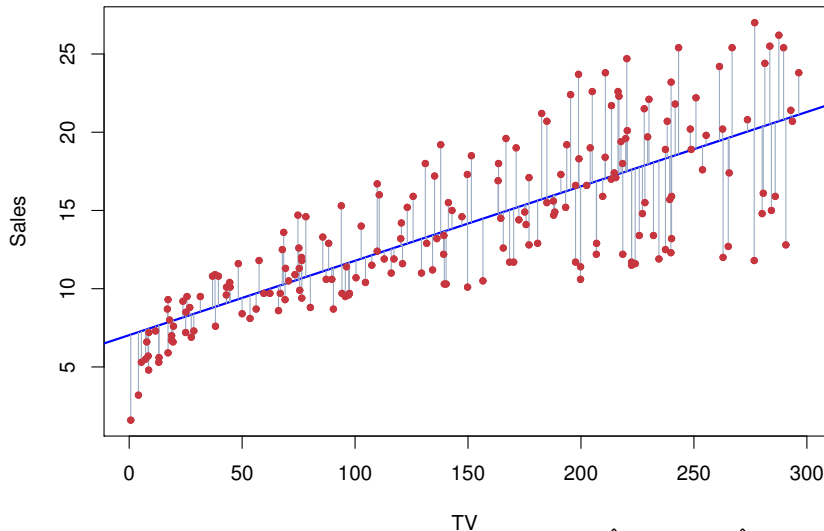
$$\text{RSS} := r_1^2 + \dots + r_n^2 = (y_1 - \hat{\beta}_0 - \hat{\beta}_1 x_1)^2 + \dots + (y_n - \hat{\beta}_0 - \hat{\beta}_1 x_n)^2. \quad (3.3)$$

An easy calculation reveals

$$\begin{aligned} \hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x}, & \bar{x} &:= \frac{1}{n} \sum_{i=1}^n x_i, \\ \hat{\beta}_1 &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}, & \bar{y} &:= \frac{1}{n} \sum_{i=1}^n y_i. \end{aligned} \quad (3.4)$$

Simple Linear Regression

Example: LS fit for advertising data

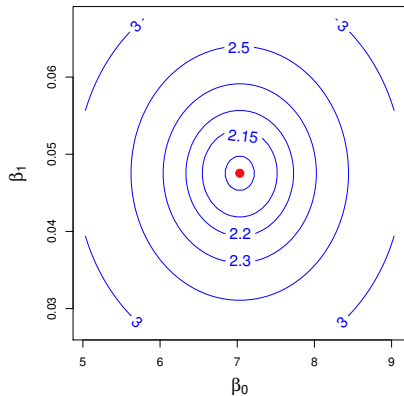


$$\hat{\beta}_0 = 7.03, \hat{\beta}_1 = 0.0475$$

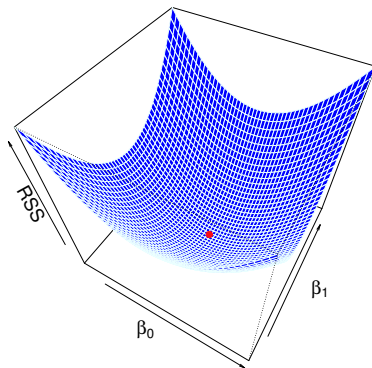
Simple Linear Regression

Example: LS fit for advertising data

LS fit of sales vs. TV budget: RSS as a function of (β_0, β_1)



Left: Level curves.



Right: Surface plot.

Simple Linear Regression

Assessing the accuracy of the coefficient estimates

Linear regression yields a linear model

$$Y = \beta_0 + \beta_1 X + \varepsilon \quad (3.5)$$

where β_0 : intercept

β_1 : slope

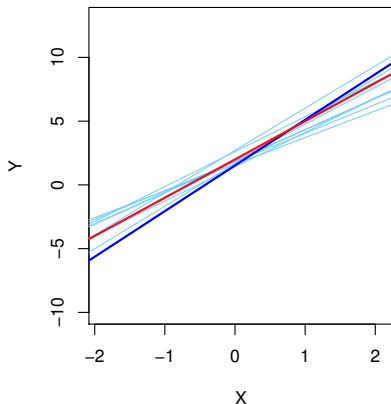
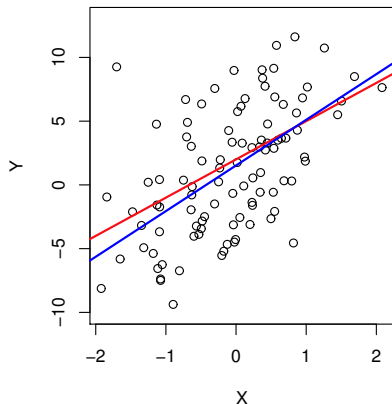
ε : model error, modeled as centered random variable,
independent of X .

Model (3.5) defines the **population regression line**, the best linear approximation to the true (generally unknown) relationship between X and Y .

The linear relation (3.2) containing the coefficients $\hat{\beta}_0, \hat{\beta}_1$ estimated from a given data set is called the **least squares line**.

Simple Linear Regression

Example: population regression line, least squares line



- Left: Simulated data set ($n = 100$) from model $f(X) = 2 + 3X$.
Red line: population regression line (true model).
Blue line: least squares line from data (black dots).
- Right: Additionally ten (light blue) least squares lines obtained from ten separate randomly generated data sets from same model; seen to average to the red line.

Simple Linear Regression

Analogy: estimation of mean

- Standard statistical approach: use information contained in a sample to estimate characteristics of a large (possibly infinite) population.
- Example: approximate **population mean** μ (expectation, expected value) of random variable Y from observations y_1, \dots, y_n by **sample mean**
 $\hat{\mu} := \bar{y} := \frac{1}{n} \sum_{i=1}^n y_i$.
- Just like $\hat{\mu} \approx \mu$ but, in general, $\hat{\mu} \neq \mu$, the coefficients $\hat{\beta}_0, \hat{\beta}_1$ defining the least squares line are estimates of the true values β_0, β_1 of the model.
- Sample mean $\hat{\mu}$ is an **unbiased estimator** of μ , i.e., it does not *systematically* over- or underestimate the true value μ .
Same holds for estimators $\hat{\beta}_0, \hat{\beta}_1$.
- How accurate is $\hat{\mu} \approx \mu$?
Standard error⁴ of $\hat{\mu}$, denoted $\text{SE}(\hat{\mu})$, satisfies

$$\text{Var } \hat{\mu} = \text{SE}(\hat{\mu})^2 = \frac{\sigma^2}{n}, \quad \text{where } \sigma^2 = \text{Var } Y. \quad (3.6)$$

⁴Standard deviation of the sample distribution, i.e., average amount $\hat{\mu}$ differs from μ .

Simple Linear Regression

Standard error of regression coefficients

For the regression coefficients (assuming uncorrelated observation errors)

$$\begin{aligned} \text{SE}(\hat{\beta}_0)^2 &= \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right], \\ \text{SE}(\hat{\beta}_1)^2 &= \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}, \end{aligned} \quad \sigma^2 = \mathbf{Var} \, \varepsilon. \quad (3.7)$$

- $\text{SE}(\hat{\beta}_1)$ smaller when x_i more spread out (provides more leverage to estimate slope).
- $\text{SE}(\hat{\beta}_0) = \text{SE}(\hat{\mu})$ if $\bar{x} = 0$. (Then $\hat{\beta}_0 = \bar{y}$.)
- σ generally unknown, can be estimated from the data by **residual standard error**

$$\text{RSE} := \sqrt{\frac{\text{RSS}}{n-2}}.$$

When RSE used in place of σ , should write $\widehat{\text{SE}}(\hat{\beta}_1)$.

Simple Linear Regression

Confidence intervals

- 95% **confidence interval**: range of values containing true unknown value of parameter with probability 95%.
- For linear regression: 95% CI for β_1 approximately

$$\hat{\beta}_1 \pm 2 \cdot \text{SE}(\hat{\beta}_1), \quad (3.8)$$

i.e., with probability 95%,

$$\beta_1 \in [\hat{\beta}_1 - 2 \cdot \text{SE}(\hat{\beta}_1), \hat{\beta}_1 + 2 \cdot \text{SE}(\hat{\beta}_1)]. \quad (3.9)$$

- Similarly, for β_0 , 95% CI approximately given by

$$\hat{\beta}_0 \pm 2 \cdot \text{SE}(\hat{\beta}_0). \quad (3.10)$$

- For advertising example: with 95% probability

$$\beta_0 \in [6.130, 7.935], \quad \beta_1 \in [0.042, 0.053].$$

Simple Linear Regression

Hypothesis tests

Use SE to test **null hypothesis**

$$H_0 : \text{no relationship between } X \text{ and } Y \quad (3.11)$$

and **alternative hypothesis**

$$H_a : \text{some relationship between } X \text{ and } Y \quad (3.12)$$

or, mathematically,

$$H_0 : \beta_1 = 0 \quad \text{vs.} \quad H_a : \beta_1 \neq 0.$$

- Reject H_0 if $\hat{\beta}_1$ sufficiently far from 0 relative to $\text{SE}(\hat{\beta}_1)$.
- **t-statistic**

$$t = \frac{\hat{\beta}_1 - 0}{\text{SE}(\hat{\beta}_1)} \quad (3.13)$$

measures distance of $\hat{\beta}_1$ from 0 in # standard deviations.

Simple Linear Regression

Hypothesis tests

- $\beta_1 = 0$ implies t follows t -distribution with $n - 2$ degrees of freedom.
- We compute probability of observing $|t|$ or larger under assumption $\beta_1 = 0$, its **p -value**.
- Small p -value: unlikely to observe substantial relation between X and Y due to purely random variation, unless the two actually are related.
- In this case we reject H_0 .
- Typical cutoffs for p -value: 1%, 5%; for $n = 30$ corresponds to t -statistic (3.13) values 2 and 2.75. respectively.

For **TV** sales data in advertising data set:

	Estimate	SE	t -statistic	p -value
β_0	7.0325	0.4578	15.36	< 0.0001
β_1	0.0475	0.0027	17.67	< 0.0001

Simple Linear Regression

Reminder: Student's t distribution

- Given X_1, \dots, X_n i.i.d. $\sim N(\mu, \sigma^2)$
- Sample mean:

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i.$$

- (Bessel corrected) sample variance:

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

- RV

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$$

distributed according to $N(0, 1)$.

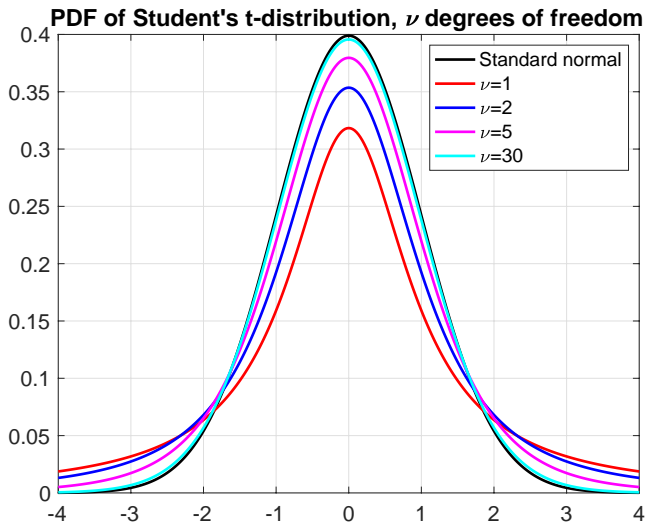
- RV

$$\frac{\bar{X} - \mu}{S/\sqrt{n}}$$

distributed according to Student's t -distribution with $n - 1$ DoF.

Simple Linear Regression

Student's t distribution



Simple Linear Regression

Assessing model accuracy

- **Residual standard error**: estimate of standard deviation of ε (model error)

$$\text{RSE} = \sqrt{\frac{\text{RSS}}{n-2}} = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2}. \quad (3.14)$$

- For **TV** data $\text{RSS} = 3.26$, i.e., deviation of sales from true regression line on average by 3,260 units (even if exact β_0, β_1 known).
Corresponds to $3,260/14,000 = 23\%$ error relative to mean value of all sales.
- RSE measures *lack of model fit*.

Simple Linear Regression

Assessing model accuracy

- **R^2 statistic**: alternative measure of fit: proportion of variance explained.
- $\in [0, 1]$, independent of scale of Y .
- Defined in terms of **total sum of squares (TSS)** as

$$R^2 = \frac{\text{TSS} - \text{RSS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}}, \quad \text{TSS} = \sum_{i=1}^n (y_i - \bar{y})^2. \quad (3.15)$$

- TSS : total variance in response Y ,
RSS : amount of variability left unexplained after regression,
TSS – RSS : response variability explained by regression model,
 R^2 : proportion of variability in Y explained using X .
- $R^2 \approx 0$: linear model wrong, high model error variance.
- For **TV** data $R^2 = 0.61$: 2/3 of **sales** variability explained by (linear regression on) TV budget.
- $R^2 \in [0, 1]$, but sufficient value problem dependent.

Simple Linear Regression

Correlation

- Measure of linear relationship between X and Y : **(sample) correlation**:

$$\text{Cor}(X, Y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}. \quad (3.16)$$

- In simple linear regression: $\text{Cor}(X, Y)^2 = R^2$.
- Correlation expresses association between single pair of variables; R^2 between larger number of variables in multivariate linear regression.

③ Linear Regression

3.1 Simple Linear Regression

3.2 Multiple Linear Regression

3.3 Other Considerations in the Regression Model

3.4 Revisiting the Marketing Data Questions

3.5 Linear Regression vs. K -Nearest Neighbors

Multiple Linear Regression

Justification

- $p > 1$ predictor variables
(as in **advertising** data set: **TV**, **newspaper**, **radio**)
- Easiest option: simple linear regression for each

For **radio** sales data in advertising data set:

	Estimate	SE	<i>t</i> -statistic	<i>p</i> -value
β_0	9.312	0.563	16.54	< 0.0001
β_1	0.203	0.020	9.92	< 0.0001

For **newspaper** sales data in advertising data set:

	Estimate	SE	<i>t</i> -statistic	<i>p</i> -value
β_0	12.351	0.621	19.88	< 0.0001
β_1	0.055	0.017	3.30	< 0.00115

Multiple Linear Regression

Justification

- How to predict total sales given 3 budgets?
- Each separate regression equation ignores the other 2 media.
- For correlated media budgets this can lead to misleading estimates of individual media effects.

Multiple linear regression model for p predictor variables:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p + \varepsilon \quad (3.17)$$

β_j : average effect on Y of 1-unit increase in X_j holding other predictors fixed.

In advertising example:

$$\text{sales} = \beta_0 + \beta_1 \times \text{TV} + \beta_2 \times \text{radio} + \beta_3 \times \text{newspaper} \quad (3.18)$$

Multiple Linear Regression

Estimating the coefficients

- Given estimates $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$, obtain prediction formula

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p. \quad (3.19)$$

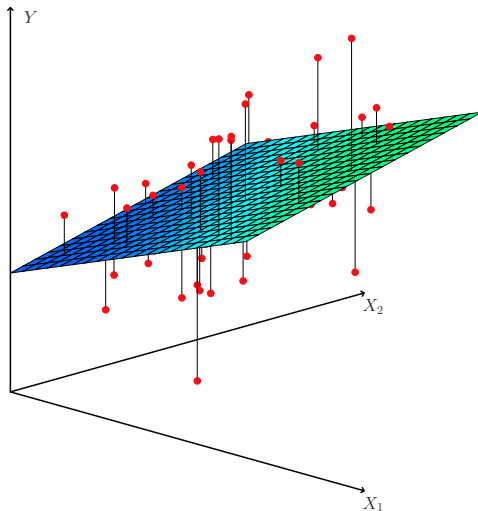
- Same fitting approach: choose $\{\hat{\beta}_j\}_{j=0}^p$ to minimize

$$\text{RSS} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i,1} - \dots - \hat{\beta}_p x_{i,p})^2, \quad (3.20)$$

yielding the **multiple least squares regression coefficients**

Multiple Linear Regression

Example: multiple linear regression, 2 predictors, 1 response



Multiple Linear Regression

Numerical methods for least squares fitting

- Determining the coefficients $\{\hat{\beta}_j\}_{j=0}^p$ to minimize the RSS in (3.20) is equivalent to minimizing $\|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|_2^2$, where we have introduced the notation

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} 1 & x_{1,1} & \dots & x_{1,p} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n,1} & \dots & x_{n,p} \end{bmatrix}, \quad \hat{\boldsymbol{\beta}} = \begin{bmatrix} \hat{\beta}_0 \\ \vdots \\ \hat{\beta}_p \end{bmatrix}$$

for the vector $\mathbf{y} \in \mathbb{R}^n$ of response observations, the matrix $\mathbf{X} \in \mathbb{R}^{n \times (p+1)}$ of predictor observations and vector $\hat{\boldsymbol{\beta}} \in \mathbb{R}^{p+1}$ of coefficient estimates.

- The problem of finding a vector $\mathbf{x} \in \mathbb{R}^n$ such that $\mathbf{b} \approx \mathbf{A}\mathbf{x}$ for given $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$ is called a **linear regression problem**.
- One (of many) possible approaches for achieving this is choosing \mathbf{x} to minimize $\|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2$, which is a **linear least squares problem**.

Multiple Linear Regression

Numerical methods for least squares fitting

- A somewhat more general fitting approach using a model

$$y \approx \beta_0 + \beta_1 f_1(\mathbf{x}) + \cdots + \beta_p f_p(\mathbf{x})$$

with fixed **regression functions** $\{f_j\}_{j=1}^p$ also leads to a linear regression problem, where now $[\mathbf{X}]_{i,j} = f_j(x_i)$.

- A linear least squares problem $\|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2 \rightarrow \min$ with $m \geq n$ has a unique solution if the columns of \mathbf{A} are linearly independent, i.e., when \mathbf{A} has full rank, given by $\mathbf{x} = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b}$.
In this case the solution can be computed using a Cholesky decomposition.
- In the (nearly) rank-deficient case, more sophisticated techniques of numerical linear algebra like the QR decomposition or the SVD are required to obtain a (stable) solution.
- When \mathbf{A} is large and sparse or structured, iterative methods such as CGLS or LSQR can be employed which require only matrix-vector products in place of manipulations of matrix entries.

Multiple Linear Regression

Advertising data

	Estimate	SE	<i>t</i> -statistic	<i>p</i> -value
β_0	2.939	0.3119	9.42	< 0.0001
β_1 (TV)	0.046	0.0014	32.81	< 0.0001
β_2 (radio)	0.189	0.0086	21.89	< 0.0001
β_3 (newspaper)	-0.001	0.0059	-0.18	0.8599

- Newspaper slope differs from simple regression.
Small estimate, *p*-value no longer significant.
- Now no relation between sales and newspaper budget. Contradiction?

Multiple Linear Regression

Advertising data

Correlation matrix:

	TV	radio	newspaper	sales
TV	1.0000	0.0548	0.0567	0.7822
radio		1.0000	0.3541	0.5762
newspaper			1.0000	0.2283
sales				1.0000

- Correlation between newspaper and radio: ≈ 0.35 :
Tend to spend more on radio ads where more is spent on newspaper ads.
- If correct, i.e., $\beta_{\text{newspaper}} \approx 0$, $\beta_{\text{radio}} > 0$, radio increased sales, and where radio budget high, newspaper budget tends to also be high.
- Simple linear regression: indicates newspaper associated with higher sales.
Multiple regression reveals no such affect.
- Newspaper receives credit for radio's affect on sales.
Sales due to newspaper advertising is a **surrogate** for sales due to radio advertising.

Multiple Linear Regression

Absurd example, same effect

- Counterintuitive but not uncommon. Consider following (absurd) example.
- Data on **shark attacks** versus **ice cream sales** at beach community would show similar positive relationship as **newspaper** and **radio** ads.
- Should one ban ice cream sales to reduce risk of shark attacks?
- Answer: High temperatures cause both (more people at beach for shark encounters, more ice cream customers).
- Multiple regression reveals ice cream sales not a predictor for shark attacks after adjusting for temperature.

Multiple Linear Regression

Questions to consider

- 1 Is at least one of the predictors X_1, X_2, \dots, X_p useful in predicting the response?
- 2 Do all predictors help to explain Y , or is only a subset of the predictors useful?
- 3 How well does the model fit the data?
- 4 Given a set of predictor values, what response value should we predict, and how accurate is our prediction?

Multiple Linear Regression

(1) Is there a relationship between response and predictors?

- As for simple regression, perform statistical hypothesis test: null hypothesis

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_p = 0$$

versus alternative

$$H_a : \text{at least one } \beta_j (j = 1, \dots, p) \text{ is nonzero.}$$

- Such a test can be based on the **F-statistic**

$$F = \frac{(\text{TSS} - \text{RSS})/p}{\text{RSS}/(n - p - 1)} \quad (3.21)$$

where, as before,

$$\text{TSS} = \sum_{i=1}^n (y_i - \bar{y})^2, \quad \text{RSS} = \sum_{i=1}^n (y_i - \hat{y}_i)^2.$$

Multiple Linear Regression

(1) Is there a relationship between response and predictors?

$$F = \frac{(\text{TSS} - \text{RSS})/p}{\text{RSS}/(n - p - 1)}$$

- Under linear model assumption, can show

$$\mathbf{E} \left[\frac{\text{RSS}}{n - p - 1} \right] = \sigma^2.$$

- If also H_0 is true, can show

$$\mathbf{E} \left[\frac{\text{TSS} - \text{RSS}}{p} \right] = \sigma^2.$$

- Hence $F \approx 1$ if no relationship between response and predictors. Alternatively, if H_a true, $\mathbf{E}[(\text{TSS} - \text{RSS})/p] > \sigma^2$, hence $F > 1$.

Multiple Linear Regression

(1) Is there a relationship between response and predictors?

Statistics for multiple regression of **sales** onto **radio**, **TV** and **newspaper** in the advertising data set:

Quantity	Value
RSE	1.69
R^2	0.897
F	570

- $F \gg 1$ strong evidence against H_0 .
- Proper threshold value for F depends on n, p .
Larger F needed to reject H_0 for small n .
- H_0 true, ε_i Gaussian, then F follows **F-distribution**; calculate p -value using statistical software.
- Here, p -value ≈ 0 for $F = 590$ in this example, hence we safely reject H_0 .

Multiple Linear Regression

(1) Is there a relationship between response and predictors?

- To test whether *subset* of last $q < p$ coefficients relevant, use null hypothesis

$$H_0 : \beta_{p-q+1} = \beta_{p-q+2} = \cdots = \beta_p = 0.$$

- Fit model using all variables *except* last q , obtaining residual sum of squares RSS_0 .
- Appropriate F -statistic now

$$F = \frac{(RSS_0 - RSS)/q}{RSS/(n - p - 1)}$$

- For multiple regression, t -statistic and p values for each variable indicate whether each predictor related to response after adjusting for the remaining variables.

Equivalent to F -test omitting single variable ($q = 1$).

Reports partial effect of adding each variable.

Multiple Linear Regression

(1) Is there a relationship between response and predictors?

What does F statistic tell us that individual p -values don't?

- Does single small p -value indicate at least one variable relevant? No.
- Example: $p = 100$, $H_0 : \beta_1 = \dots = \beta_p = 0$ true.
Then by chance, 5% of p -values below 0.05.
Almost guaranteed that $p < 0.05$ for at least one variable by chance.
- Thus, for large p , looking only at p -values of individual t -statistics tends to discover spurious relationships.
- For F -statistic, if H_0 true, only 5% chance of p -value < 0.05 independently of n, p .

Note: F -statistic approach works for $p < n$.

Multiple Linear Regression

(2) Deciding on important variables

- Typically, not all predictors related to response (**variable selection** problem).
- One approach: try all possible models, select best one. Criteria?
Mallow's C_p , **Akaike information criterion (AIC)**,
Bayesian information criterion (BIC) (later)
- For p large, trying 2^p models with subsets of variables impractical.
- *Forward selection*: Start with null model (only β_0), fit p simple regressions, add variable leading to lowest RSS, then add variable leading to two-variable model with lowest RSS, continue until stopping criterion met.
- *Backward selection*: Start with full model, remove variable with largest p -value, fit new $(p - 1)$ -variable model, keep removing least significant variable, until stopping criterion met.
- *Mixed selection*: Start with null model, adding variables with best fit one-by-one, remove variables whenever its p -value rises above threshold, until model contains only variables with low p -values and excludes those with high p -value.

Multiple Linear Regression

(3) Model fit

RSE, R^2 computed and interpreted as in simple linear regression.

- $R^2 = \text{Cor}(X, Y)^2$ for simple linear regression.
- $R^2 = \text{Cor}(\hat{Y}, Y)^2$ for multiple linear regression, maximized by fitted model.
- $R^2 \approx 1$: model explains large portion of response variance.
- *Advertising example:*

$$\{\mathbf{TV}, \mathbf{radio}, \mathbf{newspaper}\} \quad R^2 = 0.8972$$

$$\{\mathbf{TV}, \mathbf{radio}\} \quad R^2 = 0.89719$$

Small *increase* on including **newspaper** (even though **newspaper** not significant)

- Note: R^2 always increases when variables are added.
- Tiny increase in R^2 on including **newspaper** more evidence this variable can be dropped.
- Including redundant variables promotes overfitting.

Multiple Linear Regression

(3) Model fit

- *Advertising example:*

$$\{\mathbf{TV}\} \quad R^2 = 0.61$$

$$\{\mathbf{TV}, \mathbf{radio}\} \quad R^2 = 0.89719$$

Substantial improvement on adding **radio**.

(Could also look at p -value of **radio**'s coefficient in last model.)

- *Advertising example:*

$$\{\mathbf{TV}, \mathbf{radio}, \mathbf{newspaper}\} \quad \text{RSE} = 1.686$$

$$\{\mathbf{TV}, \mathbf{radio}\} \quad \text{RSE} = 1.681$$

$$\{\mathbf{TV}\} \quad \text{RSE} = 3.26$$

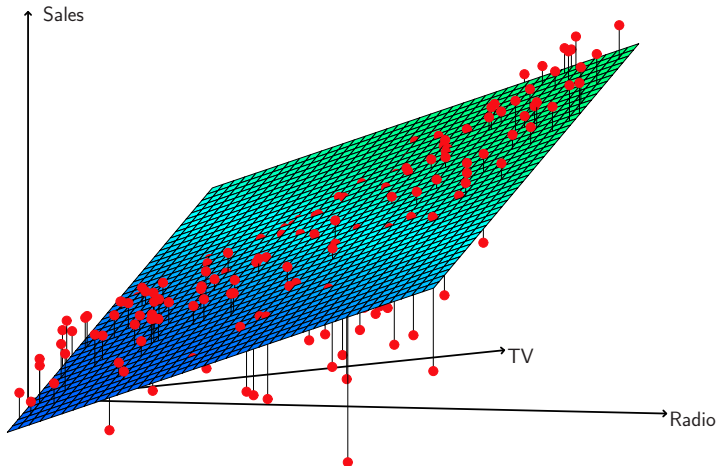
- Note: for multiple linear regression RSE defined as

$$\text{RSE} = \sqrt{\frac{\text{RSS}}{n - p - 1}}.$$

Multiple Linear Regression

(3) Model fit

$\{\text{TV}, \text{radio}\}$



Multiple Linear Regression

(3) Model fit

Previous figure:

- Some observations above, some below least squares regression plane.
- Linear model overestimates sales where most of budget spent either exclusively on **TV** or **radio**.
- Underestimation where budget split between two media.
- Such *nonlinear pattern* not reflected by linear model; suggests *synergy* effect between these two media.

Multiple Linear Regression

(4) Predictions

We note three sources of prediction uncertainty:

- 1 Reducible error: $\hat{Y} \approx f(X)$ since $\hat{\beta}_j \approx \beta_j$.
Can construct confidence intervals to ascertain closeness \hat{Y} to $f(X)$.
- 2 *Model bias*: linear model can only yield best *linear* approximation.
- 3 Irreducible error: $Y = f(X) + \varepsilon$.
Assess prediction error with **prediction intervals**: incorporate both reducible and irreducible errors.

Example: Prediction using $\{\mathbf{TV}, \mathbf{radio}\}$ model.

$X_{\text{TV}} = 100\,000$ \$, $X_{\text{radio}} = 20\,000$ \$.

Confidence interval on **sales**: 95% confidence interval : [10.985, 11.528].

Prediction interval on **sales**: 95% prediction interval : [7.930, 14.580].

Increased uncertainty about sales for given city in contrast with average sales over many locations.

③ Linear Regression

3.1 Simple Linear Regression

3.2 Multiple Linear Regression

3.3 Other Considerations in the Regression Model

3.4 Revisiting the Marketing Data Questions

3.5 Linear Regression vs. K -Nearest Neighbors

Other Considerations in the Regression Model

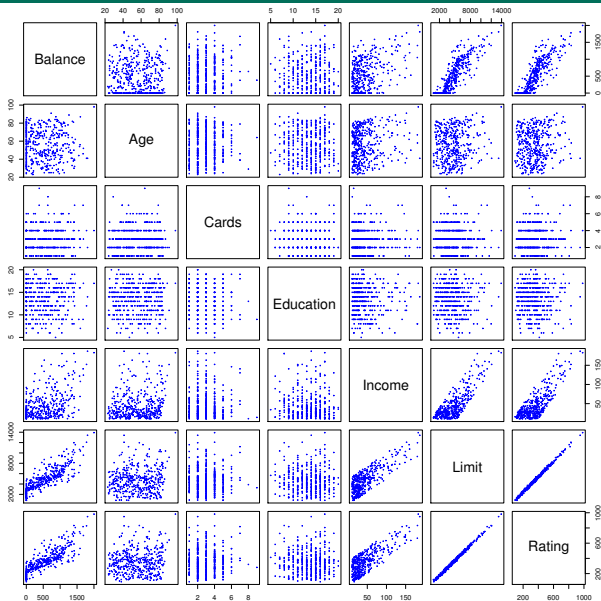
Qualitative predictors

Credit data set:

- Quantitative predictors:
 - **balance**: average credit card debt for a number of individuals
 - **age**
 - **cards** (# credit cards)
 - **education** (years of education)
 - **income** (in thousands of dollars)
 - **limit** (credit limit)
 - **rating** (credit rating)
- Qualitative predictors:
 - **gender**
 - **student** (student status)
 - **status** (marital status)
 - **ethnicity** (Caucasian, African American or Asian)

Other Considerations in the Regression Model

Qualitative predictors



Other Considerations in the Regression Model

Two-valued predictors

- Goal: investigate differences in credit card balance between males/females.
- **Gender** (qualitative variable, factor) represented with **indicator** (dummy variable)

$$x_i = \begin{cases} 1 & \text{if } i\text{-th person female,} \\ 0 & \text{if } i\text{-th person male.} \end{cases} \quad (3.22)$$

- Using x_i in regression equation results in model

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i = \begin{cases} \beta_0 + \beta_1 + \varepsilon_i & \text{if } i\text{-th person female,} \\ \beta_0 + \varepsilon_i & \text{if } i\text{-th person male.} \end{cases} \quad (3.23)$$

- Interpretation

β_0 : average credit card balance among males,

$\beta_0 + \beta_1$: average credit card balance among females,

β_1 : average difference in credit card balance male/female.

Other Considerations in the Regression Model

Two-valued predictors

	Coefficient	Standard error	<i>t</i> -statistic	<i>p</i> -value
β_0	509.80	33.13	15.389	< 0.0001
β_1	19.73	46.05	0.429	0.6690

- Average credit card debt males: \$509.80.
- Average additional credit card debt females: \$19.73.
- Total average female credit card debt: \$529.53.
- High *p* value for dummy variable. Conclusion?
Gender not a statistically significant factor for credit card debt.
- Switching male/female coding yields estimates

$$\hat{\beta}_0 = \$529.53, \quad \hat{\beta}_1 = \$ - 19.73, \quad \hat{\beta}_0 + \hat{\beta}_1 = \$509.80.$$

Other Considerations in the Regression Model

Two-valued predictors

Another alternative coding of two-valued **gender** predictor:

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

Results in model

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

with interpretation

β_0 : average credit card balance (ignoring gender),

β_1 : amount females are above/males below this average,

giving estimates

$\hat{\beta}_0 = \$519.665$ (half way between male and female averages)

$\hat{\beta}_1 = \$ 9.865$ (half of \$19.63, average male/female difference).

Other Considerations in the Regression Model

Multi-valued qualitative predictors

To encode **ethnicity** $\in \{\text{Caucasian, African American, Asian}\}$, use multiple dummy variables ($\# \text{ values} - 1$)

$$x_{i,1} = \begin{cases} 1 & \text{if } i\text{-th person Asian,} \\ 0 & \text{if } i\text{-th person not Asian,} \end{cases} \quad (3.24)$$

$$x_{i,2} = \begin{cases} 1 & \text{if } i\text{-th person Caucasian,} \\ 0 & \text{if } i\text{-th person not Caucasian,} \end{cases} \quad (3.25)$$

resulting in model

$$y_i = \beta_0 + \beta_1 x_{i,1} + \beta_2 x_{i,2} + \varepsilon_i = \begin{cases} \beta_0 + \beta_1 + \varepsilon_i & \text{if } i\text{-th person Asian} \\ \beta_0 + \beta_2 + \varepsilon_i & \text{if } i\text{-th person Caucasian} \\ \beta_0 + \varepsilon_i & \text{if } i\text{-th person African American} \end{cases} \quad (3.26)$$

Interpretation: β_0 : average credit card balance for African Americans (**baseline**),
 β_1 : difference between Asian and African Americans,
 β_2 : difference between Caucasian and African Americans

Other Considerations in the Regression Model

Multi-valued qualitative predictors

	Coefficient	Standard error	<i>t</i> -statistic	<i>p</i> -value
β_0	531.00	46.32	11.464	< 0.0001
β_1 (Asian)	-18.69	65.02	-0.287	0.7740
β_2 (Caucasian)	-12.50	56.68	-0.221	0.8260

- Estimated balance for African Americans (baseline): \$531.00.
- Asians estimated to have \$18.69 less debt than African Americans.
- Caucasians estimated to have \$12.50 less debt than African Americans.
- β_1, β_2 have high *p*-values, indicating no statistical significance for ethnicity as factor in credit card balance.
- Coefficients and *p*-values depend on coding, result does not.
F-test to reject $H_0 : \beta_1 = \beta_2 = 0$ has *p*-value 0.96 (cannot reject).
- Dummy variable approach works for combining qualitative and quantitative predictors.
(Other coding schemes for qualitative variables possible.)

Other Considerations in the Regression Model

Extending the linear model

- Restrictive assumptions in linear model: **linearity**, **additivity**.
- Additivity: effect on Y of changing X_j independent of remaining variables.
- Linearity: rate of change in Y with respect to X_j constant in X_j .
- Recall advertising data set: indication that higher **radio** budget made effect of **TV** spending stronger (interaction effect, synergy).
- Add **interaction term** to two-predictor model:

$$\begin{aligned} Y &= \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + \varepsilon \\ &= \beta_0 + (\beta_1 + \beta_3 X_2) X_1 + \beta_2 X_2 + \varepsilon \\ &= \beta_0 + \tilde{\beta}_1 X_1 + \beta_2 X_2 + \varepsilon, \end{aligned} \quad \tilde{\beta}_1 := \beta_1 + \beta_3 X_2.$$

$\tilde{\beta}_1$ changes with X_2 , hence effect of X_1 on Y changes with X_2 .

Other Considerations in the Regression Model

Extending the linear model: factory example

Example: factory productivity.

- Predict # produced **units** based on # production **lines** and # **workers**.
- Expected: increase in # production **lines** will depend on # **workers**.
- In linear model of **units**, include interaction term between **lines** and **workers**:

$$\begin{aligned}\text{units} &\approx 1.2 + 3.4 \times \text{lines} + 0.22 \times \text{workers} + 1.4 \times (\text{lines} \times \text{workers}) \\ &= 1.2 + (3.4 + 1.4 \times \text{workers}) \times \text{lines} + 0.22 \times \text{workers}\end{aligned}$$

- Adding additional line will increase # produced units by $3.4 + 1.4 \times \text{workers}$. The more **workers**, the stronger the effect of adding a **line**.

Other Considerations in the Regression Model

Extending the linear model: advertising example

Linear model for **sales** predicted by interacting **TV**, **radio** terms:

$$\begin{aligned}\text{sales} &= \beta_0 + \beta_1 \times \text{TV} + \beta_2 \times \text{radio} + \beta_3 \times (\text{radio} \times \text{TV}) + \varepsilon \\ &= \beta_0 + (\beta_1 + \beta_3 \times \text{radio}) \times \text{TV} + \beta_2 \times \text{radio} + \varepsilon\end{aligned}\tag{3.27}$$

Interpretation of β_3 : increase in effectiveness of TV advertising for one-unit increase in radio advertising.

	Coefficient	Standard error	<i>t</i> -statistic	<i>p</i> -value
β_0	6.7502	0.248	27.23	< 0.0001
β_1	0.0191	0.002	12.70	< 0.0001
β_2	0.0289	0.009	3.24	0.0014
β_3	0.0011	0.000	20.73	< 0.0001

- Model with interaction term superior to that including only **main effects**.
- Low *p*-value of interaction term strong evidence for rejecting $H_0 : \beta_3 = 0$.

Other Considerations in the Regression Model

Extending the linear model: advertising example

- Model (3.27) has $R^2 = 96.8\%$
(vs. $R^2 = 89.7\%$ for model without interaction term).
- Interpretation: of the variability remaining after fitting the model without interaction term,

$$\frac{96.8\% - 89.7\%}{100\% - 89.7\%} = 69\%$$

is explained by model (3.27) which includes the interaction term.

- \$1000 increase in TV budget associated with sales increase of $(\hat{\beta}_1 + \hat{\beta}_3 \times \mathbf{radio}) \times 1000 = 19 + 1.1 \times \mathbf{radio}$ units.
\$1000 increase in radio budget associated with sales increase of $(\hat{\beta}_2 + \hat{\beta}_3 \times \mathbf{TV}) \times 1000 = 29 + 1.1 \times \mathbf{TV}$ units.
- Hierarchical principle:** for every interaction term, include all associated main effects, even if the p values of their coefficients not significant.
Rationale: If $X_1 X_2$ related to response, vanishing coefficients for X_1 , X_2 unimportant. $X_1 X_2$ typically correlated with X_1 , X_2 ; leaving these out alters meaning of interaction.

Other Considerations in the Regression Model

Extending the linear model: credit example

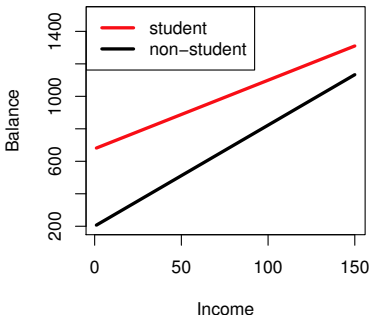
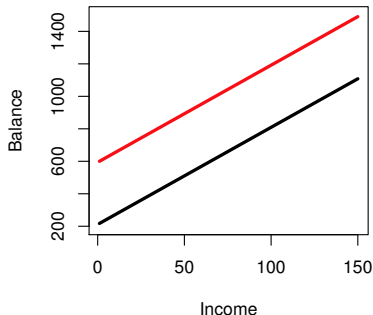
Credit data set: predict **balance** using **income** (quantitative) and **student** (qualitative). Without interaction term:

$$\begin{aligned}\mathbf{balance}_i &\approx \beta_0 + \beta_1 \times \mathbf{income}_i + \begin{cases} \beta_2 & \text{if } i\text{-th person student} \\ 0 & \text{otherwise} \end{cases} \\ &= \beta_1 \times \mathbf{income}_i + \begin{cases} \beta_0 + \beta_2 & \text{if } i\text{-th person student} \\ \beta_0 & \text{otherwise.} \end{cases}\end{aligned}\tag{3.28}$$

- Results in fitting two parallel lines to data (one each for students and non-students).
- Parallel implies: average affect on balance of one-unit increase in **income** independent of **Student** status.
- Reflects model shortcoming: change in **income** may have very different effect on credit card balance for students and non-students.

Other Considerations in the Regression Model

Extending the linear model: credit example



With interaction term: multiply **income** with dummy variable for **student**

$$\begin{aligned}\text{balance}_i &\approx \beta_0 + \beta_1 \times \text{income}_i + \begin{cases} \beta_2 + \beta_3 \times \text{income}_i & \text{if } i\text{-th person student} \\ 0 & \text{otherwise} \end{cases} \\ &= \begin{cases} (\beta_0 + \beta_2) + (\beta_1 + \beta_3) \times \text{income}_i & \text{if } i\text{-th person student} \\ \beta_0 + \beta_1 \times \text{income}_i & \text{otherwise.} \end{cases}\end{aligned}$$

(3.29)

Other Considerations in the Regression Model

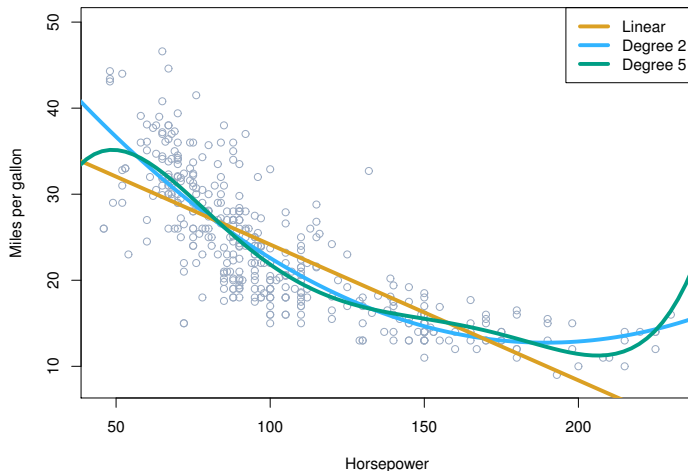
Extending the linear model: credit example

- Now the two lines have different intercepts and different slopes.
- Slope for students lower, indicates increases in income associated with smaller increase in credit card balance than for non-students.

Other Considerations in the Regression Model

Extending the linear model: nonlinear relationships

Polynomial regression vs. linear regression:



Auto data set showing **mpg** (miles per gallon) versus **horsepower** for different cars.

Other Considerations in the Regression Model

Extending the linear model: nonlinear relationships

Since the data seem to suggest curved relationship, add quadratic term:

$$\text{mpg} = \beta_0 + \beta_1 \times \text{horsepower} + \beta_2 \times \text{horsepower}^2 + \varepsilon. \quad (3.30)$$

	Coefficient	Standard error	<i>t</i> -statistic	<i>p</i> -value
β_0	56.9001	1.8004	31.6	< 0.0001
β_1	-0.4662	0.0311	-15.0	< 0.0001
β_2	0.0012	0.0001	10.1	< 0.0001

- Linear fit has $R^2 = 0.606$, quadratic fit has $R^2 = 0.688$.
- *p*-value for quadratic term highly significant.
- Degree 5 fit more oscillatory, doesn't appear to explain data any better than quadratic.

Other Considerations in the Regression Model

Potential problems

Most common problems when fitting a linear regression model to a data set:
(identification and solution as much an art as a science)

- 1 Nonlinear dependence of response on predictors
- 2 Correlated error terms
- 3 Non-constant variance of error terms
- 4 Outliers
- 5 High-leverage points
- 6 Collinearity

Other Considerations in the Regression Model

Potential problems: (1) Nonlinear dependence

Inference and prediction from linear regression model suspect when true model nonlinear.

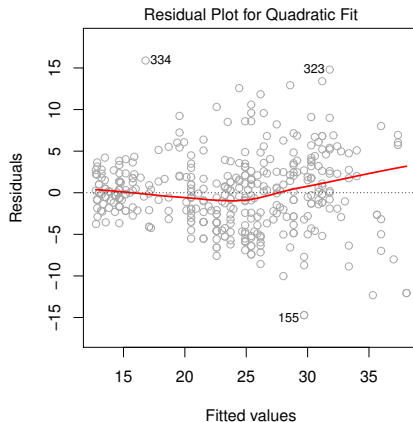
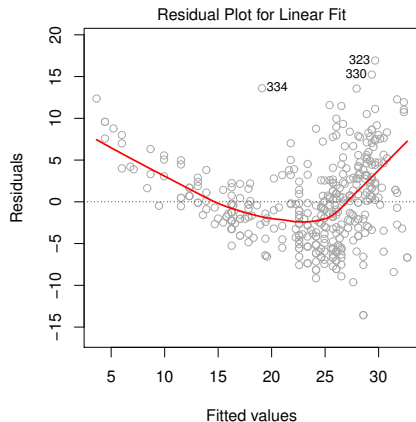
- Identifying nonlinearity aided by **residual plots**

$$e_i = y_i - \hat{y}_i \quad \text{against predictors } x_i.$$

- For multiple regression models, plot residuals against predicted (fitted) values \hat{y}_i .
- Ideal picture: no discernible pattern.
- Pattern indicates possible problem with model.
- When nonlinearity is suggested, introduce nonlinear functions of predictors as **regression functions** into the regression model.

Other Considerations in the Regression Model

Potential problems: (1) Nonlinear dependence



Residuals versus predicted values for **Auto** data set.

Red line is smooth fit to residuals to aid in identifying trends.

Left: linear regression of **mpg** on **horsepower** (strong pattern).

Right: linear regression of **mpg** on **horsepower** and **horsepower**² (little pattern).

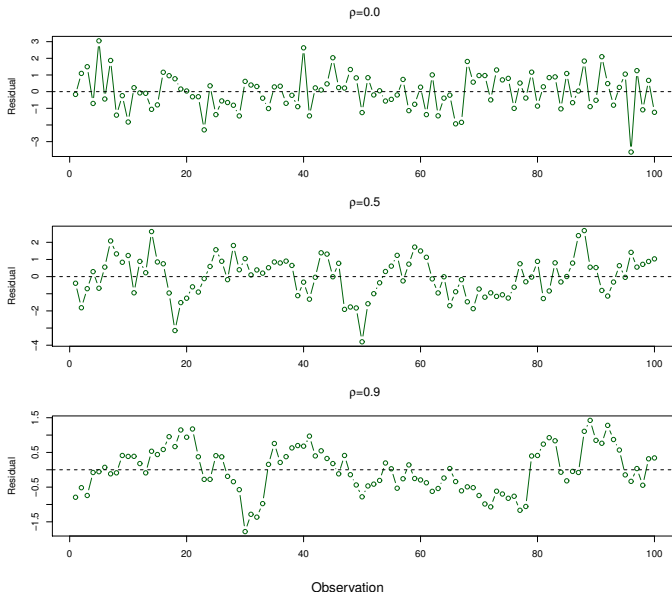
Other Considerations in the Regression Model

Potential problems: (2) Correlated error terms

- Linear regression assumes uncorrelated errors ε_i .
- Computation of SE for coefficient estimates, fitted values, based on this assumption. Otherwise estimated SE tend to underestimate true SE, confidence and prediction intervals too optimistic (narrow), p -values lower than they should be.
- Extreme example: double data (observations, error terms identical in pairs). SE calculations use sample size $2n$ in place of n , hence CI narrower by factor of $\sqrt{2}$.
- Detection for **time series**: plot residuals as function of time. No correlations implies no visible pattern; correlations lead to **tracking** of residuals.
- Example (next slide): time series with error correlation $\rho = 0, 0.5, 0.9$
- Example: study of persons' heights predicted from their weights. Uncorrelatedness assumption violated if, e.g., individuals related, same diet or environmental factors.

Other Considerations in the Regression Model

Potential problems: (2) Correlated error terms



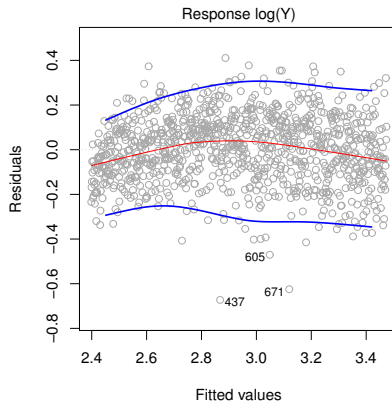
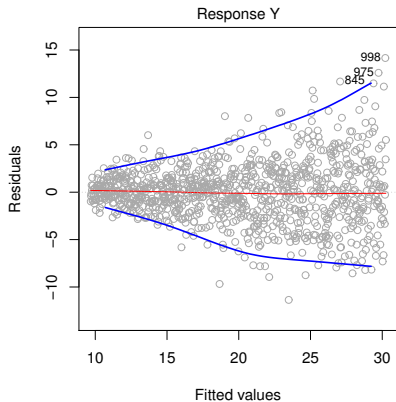
Other Considerations in the Regression Model

Potential problems: (3) Non-constant variance of error terms

- SE, CI, hypothesis tests associated with linear model rely on assumption **$\text{Var } \varepsilon_i = \sigma^2 (\forall i)$** .
- Non-constant error variance (**heteroscedasticity**), e.g. increase with response value, leads to *funnel-shaped* residual plot.
- Possible solution: transform response Y using concave function such as $\log Y$ or \sqrt{Y} , leads to damping of larger responses, reducing heteroscedasticity.
- When variation of response variance known, e.g., i -th response average of n_i observations which are uncorrelated with variance σ^2 , then average has variance $\sigma_i^2 = \sigma^2/n_i$. Remedy: **weighted least squares** with weights proportional to inverse variances.

Other Considerations in the Regression Model

Potential problems: (3) Non-constant variance of error terms

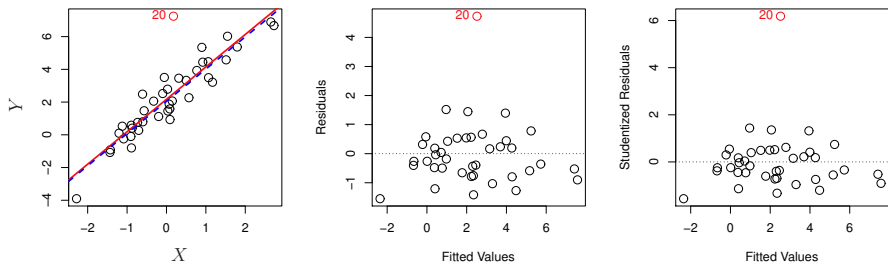


Residual plots. Red: smooth fit of residuals. Blue: track outer quantiles of residuals.
Left: funnel shape indicating heteroscedasticity.
Right: After log-transforming response, heteroscedasticity removed.

Other Considerations in the Regression Model

Potential problems: (4) Outliers

- **Outlier**: point where y_i far from value predicted by model.
- Possible causes: observation errors.



Left: red solid line: least squares line with outlier, blue: without.

Center: Residual plot identifies outlier.

Right: Outlier seen to have studentized residual (divide e_i by its estimated standard error) of 6 (between -3 and 3 expected).

R^2 declines from 0.892 to 0.805 on including outlier.

Other Considerations in the Regression Model

Potential problems: (5) High-leverage points

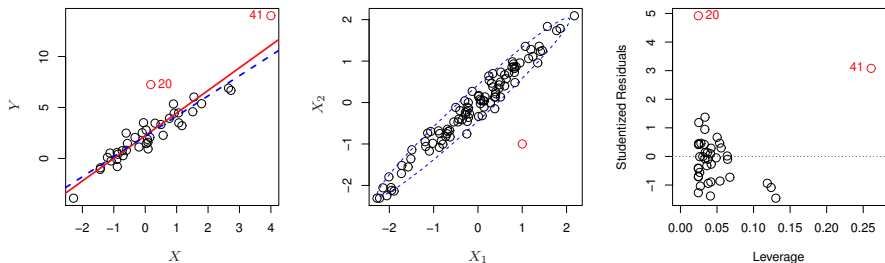
- Outliers: observations where y_i is unusual given x_i .
- Observations with **high leverage** have unusual value for x_i .
- If least squares line strongly affected by certain points, problems with these may invalidate entire fit, hence important to identify such observations.
- Simple linear regression: extremal x-values; multiple linear regression: in range of all other observation coordinates, but unusual (difficult to detect for more than two predictors).
- Large value of **leverage statistic** indicates high leverage.
For simple linear regression:

$$h_i = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{i'=1}^n (x_{i'} - \bar{x})^2} \in \left(\frac{1}{n}, 1 \right). \quad (3.31)$$

Average value always $\frac{p+1}{n}$, deviation from average indicates high leverage.

Other Considerations in the Regression Model

Potential problems: (5) High-leverage points



Left: Same data as previous figure, with added observation 41 (red) of high leverage. Red solid line is least squares fit with, blue dashed without observation 41.

Center: two predictor variables, most observations within blue dashed ellipse, red observation distinctly outside.

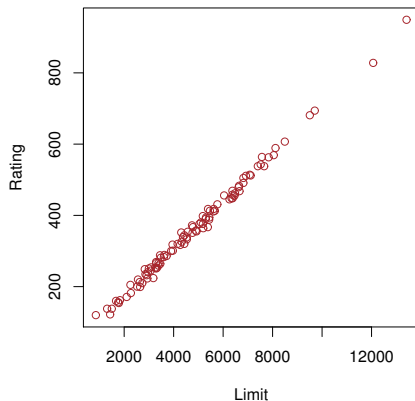
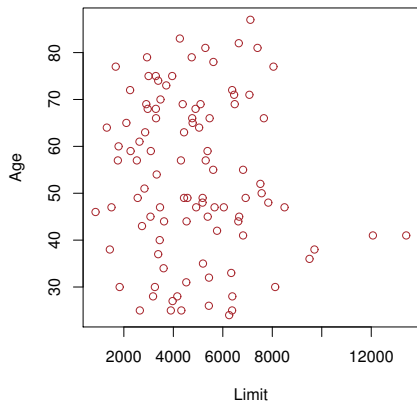
Right: same data as in left panel, studentized residuals vs. leverage statistic. Observation 41 has high leverage and high residual, i.e., outlier *and* high-leverage point.

Outlier observation 20 has low leverage.

Other Considerations in the Regression Model

Potential problems: (6) Collinearity

Collinearity: two or more predictor variables closely related.

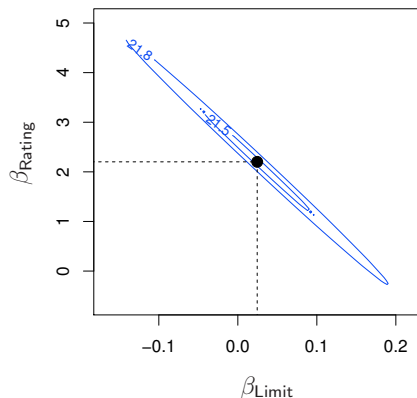
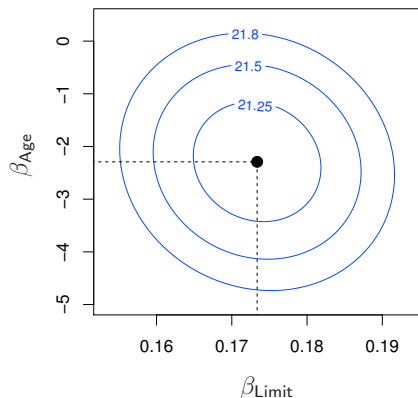


From **Credit** data set. Left: **limit** vs. **age**. Right: **limit** vs. **rating** (strongly collinear).

Other Considerations in the Regression Model

Potential problems: (6) Collinearity

Difficult to separate individual effects of collinear variables on response.



Contour plot of RSS associated with different coefficient estimates for **Credit** data set. Axes scaled to include 4 SE on either side of optimum.

Left: for regression of **balance** on **limit** and **age**.

Right: for regression of **balance** on **limit** and **rating**.

Other Considerations in the Regression Model

Potential problems: (6) Collinearity

- Collinearity increases SE, hence reduces t -statistic, and we will more likely fail to reject $H_0 : \beta_j = 0$. This reduces the **power** of the hypothesis test, i.e., the probability of correctly detecting a nonzero coefficient.

	Coefficient	Standard error	t -statistic	p -value
Model 1				
β_0	-173.411	43.828	-3.957	< 0.0001
β_1 (age)	-2.292	0.672	-3.407	0.0007
β_2 (limit)	0.173	0.005	34.496	< 0.0001
Model 2				
β_0	-377.537	45.254	-8.343	< 0.0001
β_1 (rating)	2.202	0.952	2.312	0.0213
β_2 (limit)	0.025	0.064	0.384	0.7012

- Model 1: **age**, **limit** both highly significant.
Model 2: collinearity between **rating** and **limit** increases SE for **limit** coefficient by factor 12, p -value increases to 0.701. Collinearity masks importance of **limit** variable.

Other Considerations in the Regression Model

Potential problems: (6) Collinearity

- Important to detect collinearity when fitting a model.
- Correlation matrix may give indication.
- **Multicollinearity**: collinearity between 3 or more variables which each have low pairwise correlation.
- **Variance inflation factor (VIF)**: ratio of variance of $\hat{\beta}_j$ when fitting the full model and variance of $\hat{\beta}_j$ when fitted separately.
- $VIF \geq 1$, minimum at complete absence of collinearity.
Problematic if VIF exceeds 5 or 10.

$$VIF(\hat{\beta}_j) = \frac{1}{1 - R_{X_j|X_{-j}}^2}$$

$R_{X_j|X_{-j}}^2$: R^2 from regression of X_j onto all other predictors.

- In **Credit** data example: predictors have VIF values of 1.01, 160.67, 160.59.
- Remedies: drop problematic variables, combine collinear variables into single predictor.

③ Linear Regression

- 3.1 Simple Linear Regression
- 3.2 Multiple Linear Regression
- 3.3 Other Considerations in the Regression Model
- 3.4 Revisiting the Marketing Data Questions
- 3.5 Linear Regression vs. K -Nearest Neighbors

Revisiting the Marketing Data Questions

Recall the seven questions relating to the **Advertising** data set we set out to answer on Slide 71:

- ➊ Is there a relationship between advertising budget and sales?
- ➋ How strong is this relationship between advertising budget and sales?
- ➌ Which media contribute to sales?
- ➍ How accurately can we estimate the effect of each medium on sales?
- ➎ How accurately can we predict future sales?
- ➏ Is the relationship linear?
- ➐ Is there synergy among the advertising media?

We revisit each in turn.

Revisiting the Marketing Data Questions

- ① Is there a relationship between advertising budget and sales?
 - Fit multiple regression model of **sales** onto **TV**, **radio** and **newspaper**.
 - Test hypothesis $H_0 : \beta_{\text{TV}} = \beta_{\text{radio}} = \beta_{\text{newspaper}} = 0$.
 - Rejection/non-rejection based on F -statistic (Slide 100).
 - For advertising data: low p -value of F -statistic (table on Slide 102) strong evidence for rejecting H_0 .

Revisiting the Marketing Data Questions

- ② How strong is this relationship between advertising budget and sales?
 - Measure of model error: RSE (see Slide 80), estimates standard deviation of response from (true) population regression line.
 - Advertising data:
For multiple regression model of **sales** on **TV** and **radio**, $RSE = 1,681$ units (Slide 107).
Relative to response sample mean of 14,022 units, this is an error of 12%.
 - Measure of model error: R^2 (Slide 87), measures proportion of response variability explained by model.
 - Advertising data:
For multiple regression model of **sales** on **TV**, **radio** and **newspaper**, $R^2 = 0.897$, i.e., $\approx 90\%$ of **sales** variability explained by multiple linear regression model (Slide 102).

Revisiting the Marketing Data Questions

③ Which media contribute to sales?

- p -values of t -statistic in multiple regression model of **sales** on **TV**, **radio** and **newspaper**: small for **TV** and **radio**, large for **newspaper**.
- Suggest only **TV** and **radio** budgets related to **sales**.

Revisiting the Marketing Data Questions

- ④ How accurately can we estimate the effect of each medium on sales?
 - Confidence intervals for β_j constructed from SE of $\hat{\beta}_j$.
 - Advertising data: 95%-confidence intervals for multiple regression coefficients are

TV	(0.043, 0.049)
radio	(0.172, 0.206)
newspaper	(-0.013, 0.011)

- Wide SE due to collinearity? (Slide 138).
VIF scores for **TV**, **radio** and **newspaper** are 1.005, 1.145, 1.145, so not likely.
- Separate simple regressions of **sales** on **TV**, **radio** and **newspaper** show strong association of **TV** and **radio** with **sales**, mild association of **newspaper** with **sales**, when remaining two predictors ignored.

Revisiting the Marketing Data Questions

5 How accurately can we predict future sales?

- Can use (3.19) for prediction.
- Prediction intervals assess accuracy of predicting individual responses
 $Y = f(X) + \varepsilon$.
- Confidence intervals assess accuracy of predicting average responses
 $Y = f(X)$.
- Former always wider due to accounting for additional variability due to irreducible error ε .

Revisiting the Marketing Data Questions

⑥ Is the relationship linear?

- Identify nonlinearity using residual plots of linear model (Slide 129).
- Advertising data:
Nonlinear effects visible in figure on Slide 108.
- Discussed regression functions which are nonlinear in the predictor variables.

Revisiting the Marketing Data Questions

- 7 Is there synergy among the advertising media?
 - Non-additive relationships modeled by interaction term in model (Slide 119).
 - Presence of interaction (synergy) confirmed by small p -value of interaction term.
 - Advertising data:
Including interaction term increased R^2 from $\approx 90\%$ to $\approx 97\%$.

③ Linear Regression

- 3.1 Simple Linear Regression
- 3.2 Multiple Linear Regression
- 3.3 Other Considerations in the Regression Model
- 3.4 Revisiting the Marketing Data Questions
- 3.5 Linear Regression vs. K -Nearest Neighbors

Linear Regression vs. K -Nearest Neighbors

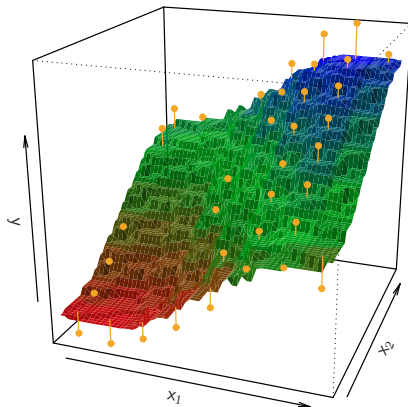
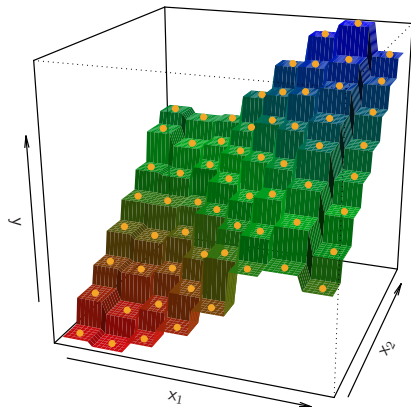
Non-parametric approach

- Linear regression is a parametric method.
- Non-parametric methods make no strong a priori assumptions on functional form of model $Y \approx f(X)$, more flexibility in adapting to data.
- Here: **K -nearest neighbors (KNN) regression** (Cf. KNN classifier in Chapter 2).
- Given prediction point x_0 , first determine the set \mathcal{N}_0 consisting of the K ($K \in \mathbb{N}$) training observations closest to x_0 .
- Predict \hat{y}_0 to be average training response in \mathcal{N}_0 , i.e.,

$$\hat{f}(x_0) = \frac{1}{K} \sum_{x_i \in \mathcal{N}_0} y_i.$$

Linear Regression vs. K -Nearest Neighbors

Non-parametric approach



Two KNN fits on a data set with 64 observations using $p = 2$ predictors.

Left: $K = 1$. Interpolation, rough step-like function.

Right: $K = 9$. Not interpolatory, smoother.

Linear Regression vs. K -Nearest Neighbors

Tuning K

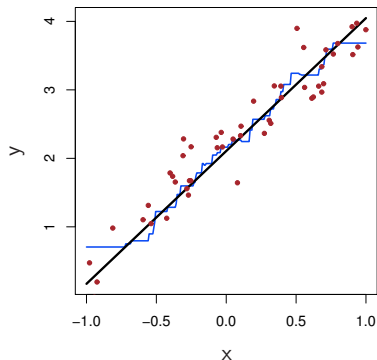
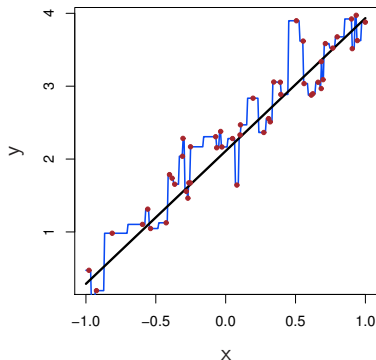
- Flexibility of model controlled by K : less flexible. smoother fit, for large K .
- Bias-variance tradeoff.
- Flexible model: low bias, high variance
(prediction depends on only one nearby observation).
Unflexible model: high bias, low variance (changing one observation has smaller effect, averaging introduces bias).
- Optimal value of K ? (later)

Linear Regression vs. K -Nearest Neighbors

Parametric vs. non-parametric

Q: In what setting will a parametric approach outperform a non-parametric approach?

A: Depends on how closely assumed form of f matches true form.

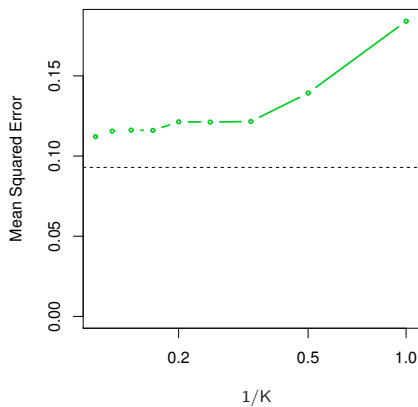
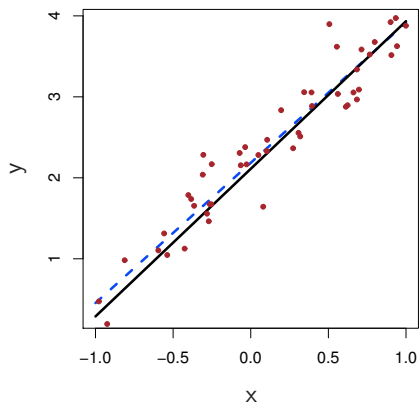


1D data, 100 observations (red), linear true model (black), KNN regression (blue).

Left: $K = 1$, right: $K = 9$.

Linear Regression vs. K -Nearest Neighbors

Parametric vs. non-parametric

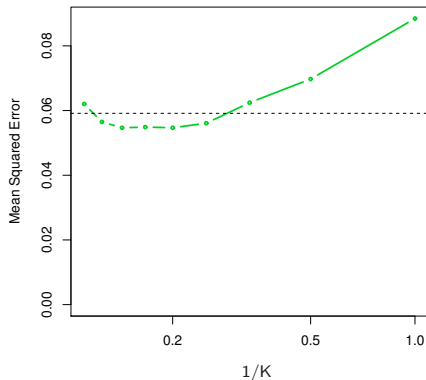
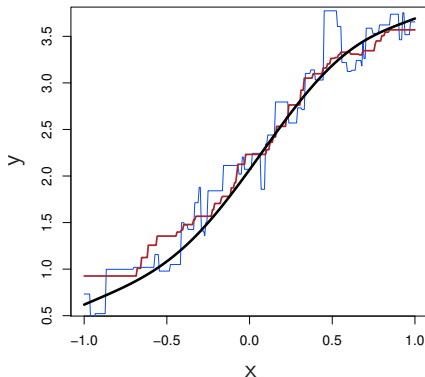


Left: same data, linear regression fit.

Right: test set MSE for linear regression (dotted line) and KNN for different values of K (plotted against $1/K$).

Linear Regression vs. K -Nearest Neighbors

Parametric vs. non-parametric

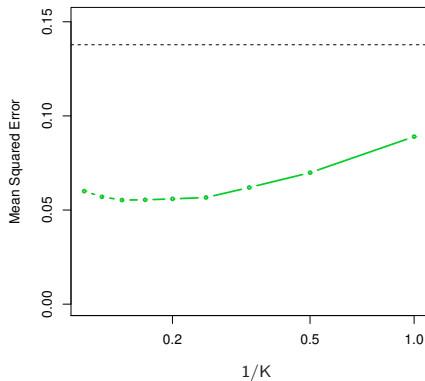
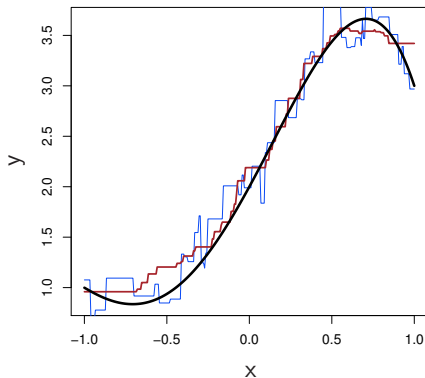


Left: slightly nonlinear data, true model (black), KNN regression with $K = 1$ (blue) and $K = 9$ (red).

Right: test set MSE for linear regression (dotted line) and KNN (against $1/K$). KNN wins for $K \geq 4$.

Linear Regression vs. K -Nearest Neighbors

Parametric vs. non-parametric

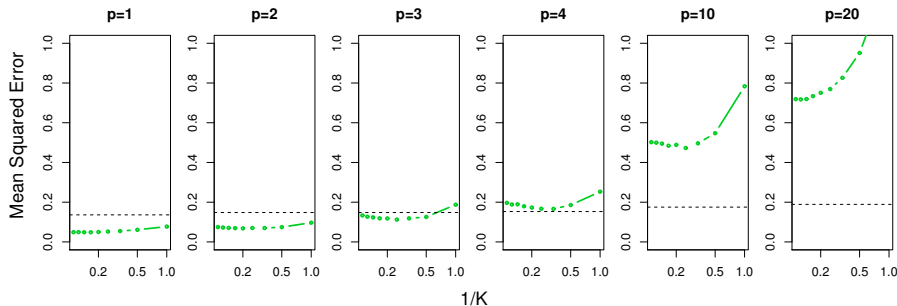


Left: strongly nonlinear data, true model (black), KNN regression with $K = 1$ (blue) and $K = 9$ (red).

Right: test set MSE for linear regression (dotted line) and KNN (against $1/K$). KNN wins for all K displayed.

Linear Regression vs. K -Nearest Neighbors

Parametric vs. non-parametric



Strongly nonlinear case, added noise predictors not associated with response. Linear regression MSE deteriorates only slightly as p rises, KNN regression MSE much more sensitive.

- For $p = 1$ KNN seems at most slightly worse than linear regression. For $p > 1$ this is no longer true.
- **Curse of dimensionality**: for $p = 20$, many of the 100 observations have no nearby observations.

Linear Regression vs. K -Nearest Neighbors

Parametric vs. non-parametric

- General rule: parametric methods tend to outperform non-parametric methods when there is a small number of observations per predictor.
- Even for small p , parametric methods offer the added advantage of better interpretability.