# ISLR Notes

TBD

2021

# Contents

# About

Notes and solutions for the exercises in the book: *An Introduction to Statistical Learning with Applications in R (1st edition)* by Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani (website: https://www.statlearning.com/)

**License**

This work, as a whole, is licensed under a Attribution-NonCommercial-ShareAlike 4.0 International License

# Chapter 1

# Introduction

## 1.1 An Overview of Statistical Learning

"Statistical learning refers to a vast set of tools for understanding data."

- Supervised: Using statistical models to **predict** or **estimate outputs** based on **inputs**.
- Unsupervised: Finding relationships between variables and structure in the data
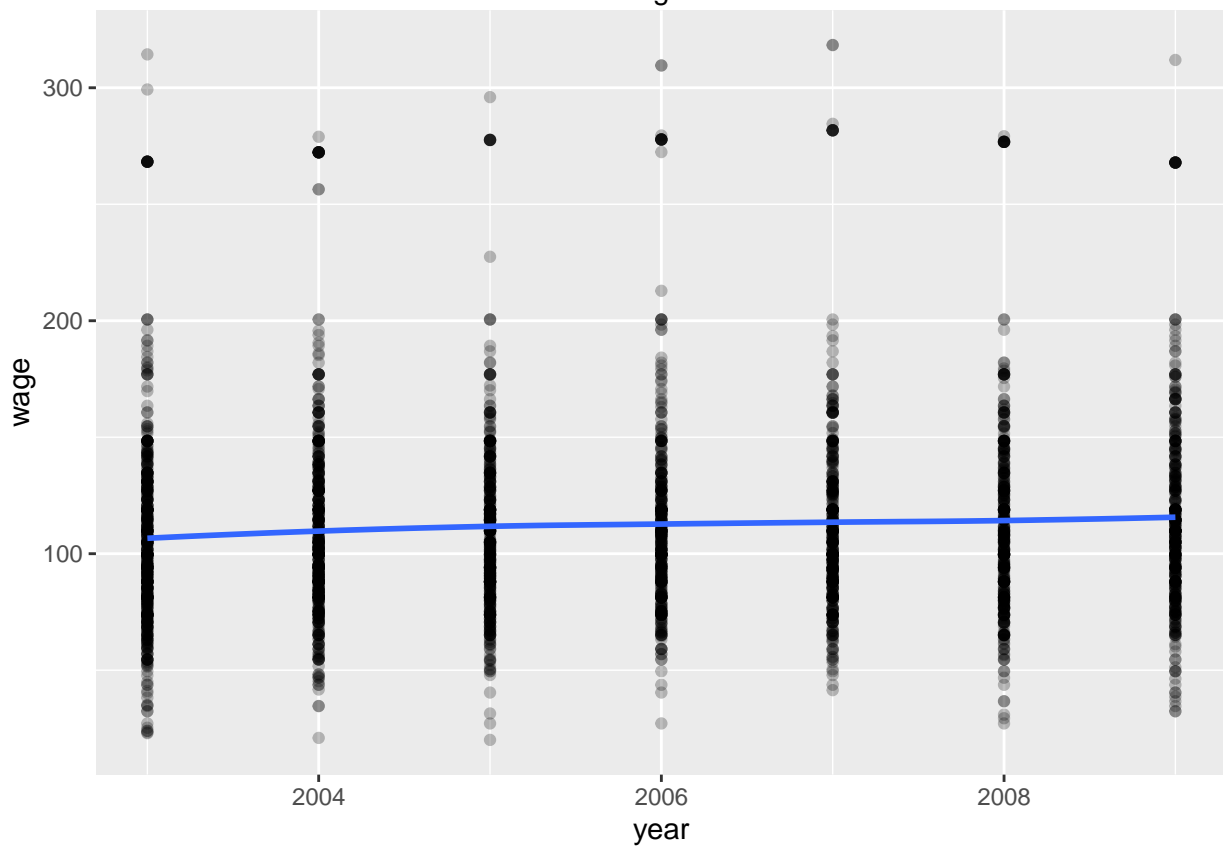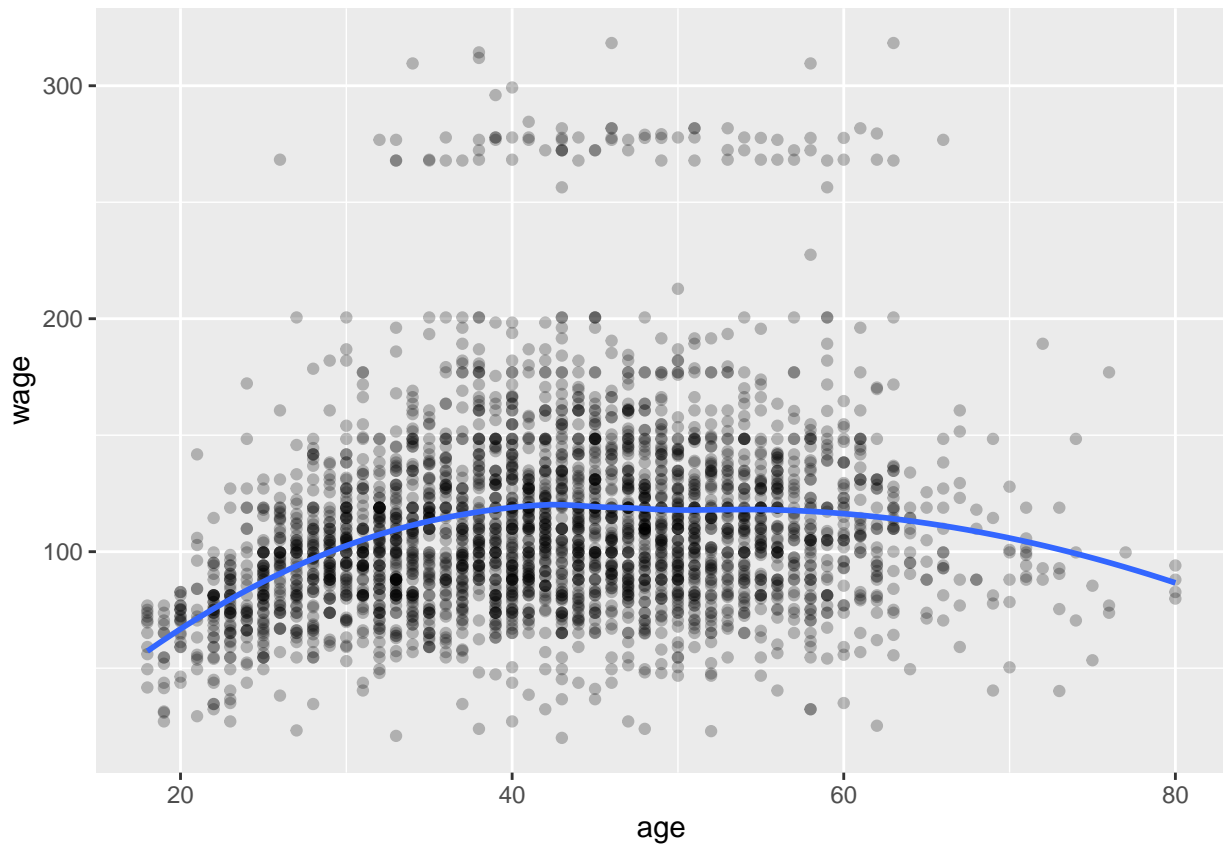
## 1.2 Data sets

Example data used in the book

- Wages
- Stock Market Data
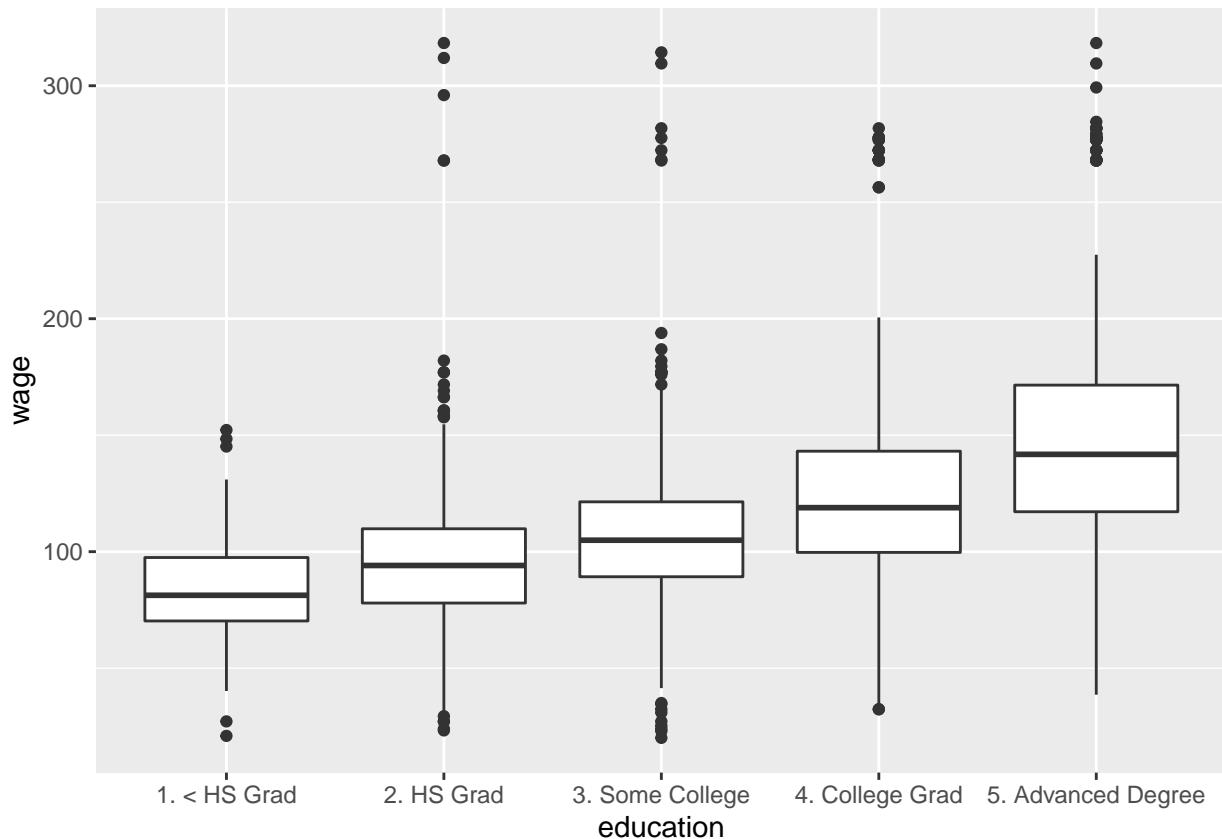- Gene Expression Data

### 1.2.1 Wages

Used for regression problem examples such as predicting wage based on age and education

```
glimpse(Wage)
```

```
## Rows: 3,000
## Columns: 11
## $ year       <int> 2006, 2004, 2003, 2003, 2005, 2008, 2009, 2008, 2006, 20...
## $ age        <int> 18, 24, 45, 43, 50, 54, 44, 30, 41, 52, 45, 34, 35, 39, ...
## $ maritl     <fct> 1. Never Married, 1. Never Married, 2. Married, 2. Marri...
## $ race       <fct> 1. White, 1. White, 1. White, 3. Asian, 1. White, 1. Whi...
## $ education  <fct> 1. < HS Grad, 4. College Grad, 3. Some College, 4. Colle...
## $ region     <fct> 2. Middle Atlantic, 2. Middle Atlantic, 2. Middle Atlant...
## $ jobclass   <fct> 1. Industrial, 2. Information, 1. Industrial, 2. Informa...
## $ health     <fct> 1. <=Good, 2. >=Very Good, 1. <=Good, 2. >=Very Good, 1....
## $ health_ins <fct> 2. No, 2. No, 1. Yes, 1. Yes, 1. Yes, 1. Yes, 1. Yes, 1....
## $ logwage    <dbl> 4.318063, 4.255273, 4.875061, 5.041393, 4.318063, 4.8450...
## $ wage       <dbl> 75.04315, 70.47602, 130.98218, 154.68529, 75.04315, 127....
```
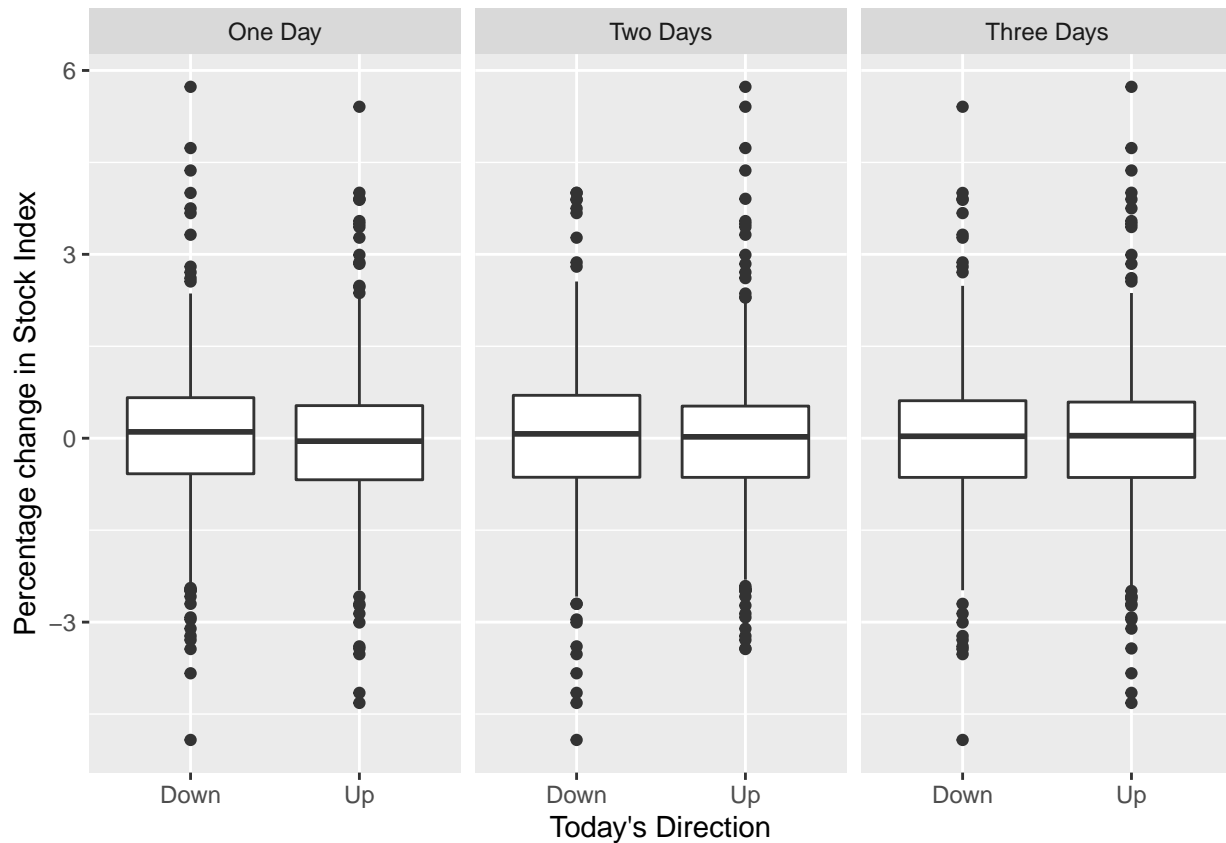
## 1.2.2 Stock Market Data

Used for classification problem examples with categorical or qualitative output, such as predicting whether a stock index will either increase or decrease on any given day.

Daily percentage change of S&P 500 stock index and 5 prior days

```
glimpse(Smarket)
```

```
## Rows: 1,250
## Columns: 9
## $ Year      <dbl> 2001, 2001, 2001, 2001, 2001, 2001, 2001, 2001, 2001, 200...
## $ Lag1      <dbl> 0.381, 0.959, 1.032, -0.623, 0.614, 0.213, 1.392, -0.403,...
## $ Lag2      <dbl> -0.192, 0.381, 0.959, 1.032, -0.623, 0.614, 0.213, 1.392,...
## $ Lag3      <dbl> -2.624, -0.192, 0.381, 0.959, 1.032, -0.623, 0.614, 0.213...
## $ Lag4      <dbl> -1.055, -2.624, -0.192, 0.381, 0.959, 1.032, -0.623, 0.61...
## $ Lag5      <dbl> 5.010, -1.055, -2.624, -0.192, 0.381, 0.959, 1.032, -0.62...
## $ Volume    <dbl> 1.1913, 1.2965, 1.4112, 1.2760, 1.2057, 1.3491, 1.4450, 1...
## $ Today     <dbl> 0.959, 1.032, -0.623, 0.614, 0.213, 1.392, -0.403, 0.027,...
## $ Direction <fct> Up, Up, Down, Up, Up, Up, Down, Up, Up, Up, Down, Down, U...
```

### 1.2.3   Gene Expression Data

Used for examples of clustering problems such as identifying related groups of cancer cells based on observed characteristics.

```
str(NCI60)
```

```
## List of 2
##  $ data: num [1:64, 1:6830] 0.3 0.68 0.94 0.28 0.485 ...
##   ..- attr(*, "dimnames")=List of 2
##   .. ..$ : chr [1:64] "V1" "V2" "V3" "V4" ...
##   .. ..$ : chr [1:6830] "1" "2" "3" "4" ...
##  $ labs: chr [1:64] "CNS" "CNS" "CNS" "RENAL" ...
```

## 1.3   History

A brief timeline for the development of statistical learning

- 1800's *Linear Regression* (*Method of Least Squares*)
- 1936 *Linear Discriminant Analysis* developed to predict qualitative values
- 1940s *Logistic Regression* developed to predict qualitative values
- 1970s *Generalized Linear Models* including both logistic and linear regression
- 1980s *Classification and Regression Trees*
- 1986 *Generalized Additive Models*
- Present day (2001) *Machine Learning*

## 1.4   Other Considerations

"How Eugenics Shaped Statistics: Exposing the damned lies of three science pioneers.

## 1.5   Matrix Notation

Conventions used in the book

- $n$ number of observations in a sample

- $p$ number of variables

- **X** an $n \times p$ matrix

  - where $x_{ij}$ represents the element in the $i$th row and the $j$th column.
  - $x_i$ represents a single observation (row) as a vector with length $p$. Note that vectors are written vertically by convention in math notation.
  - $\mathbf{x}_j$ represents a single variable (column) as a vector with length $n$. Note that the bold face font is used to distinguish columns ($\mathbf{x}_3$) from rows ($x_3$).

- The $^T$ superscript operator denotes the transpose of a matrix or vector, where row and column indices are reversed such that the resulting matrix or vector will have $p$ rows and/or $n$ columns.

Examples

- A matrix of elements

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix}$$

- A row vector

$$x_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{pmatrix}$$

- A column vector

$$\mathbf{x}_j = \begin{pmatrix} x_{1j} \\ x_{2j} \\ \vdots \\ x_{nj} \end{pmatrix}$$

- A matrix represented as a collection of column vectors

$$\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_j)$$

- A transposed matrix. Rows become columns and columns become rows

$$\mathbf{X}^T = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{p1} & x_{p2} & \dots & x_{pn} \end{pmatrix}$$

- A transposed row vector. Again, vector elements are listed vertically by default, so this presentation shows the new orientation.

$$x_i^T = (x_{i1}, x_{i2}, \ldots, x_{ip})$$

- A matrix represented as a collection of row vectors

$$\mathbf{X} = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{pmatrix}$$

# Chapter 2

# Statistical Learning

```
## [1] 1 3 2 5
```

## 2.13   2.3.2 Graphics

## 2.14   2.3.3 Indexing Data

## 2.15   2.3.4 Loading Data

## 2.16   2.3.5 Additional Graphical and Numerical Summaries

## 2.17   2.4 Exercises

## 2.18   Conceptual

1.
2.
3.
4.
5.
6.
7.

## 2.19   Applied

8.
9.
10.

# Chapter 3

# Linear Regression

# Chapter 4

# Classification

# Chapter 5

# Resampling Methods

# Chapter 6

# Model Selection and Regularization

# Chapter 7

# Moving Beyond Linearity

# Chapter 8

# Tree Based Methods

# Chapter 9

# Support Vector Machines

# Chapter 10

# Unsupervised Learning