

Data driven sampling (Lecture 6)

Regression

Hongwei YUAN

University of Macau

Model of regression

We first assume that the regression problem is one dimensional. We consider

$$Y = f(X) + \epsilon$$

where

- Y is called the output, X is the input, ϵ is noise term due to measurement errors.
- f is unknown function to be learnt.
- We aim to learn f by observed data $(x_1, y_1), \dots, (x_n, y_n)$.

Model of regression

We write

$$Y = f(X) + \epsilon$$

where

- f is unknown function to be learnt.
- In Statistics, one often assumes $f(x) = g(x, \theta)$ with g being known but θ unknown, e.g. $g(x, \theta) = \theta x$.
- When $f(x) = g(x, \theta)$, the problem turns to be finding θ by the observed data.
- The noise ϵ is often assumed to have a normal distribution $N(0, \sigma^2)$, usually the variance σ^2 is known.

Model of regression

We consider the regression:

$$Y = g(X, \theta) + \epsilon$$

where we have observed data $(x_1, y_1), \dots, (x_n, y_n)$, and ϵ has the distribution $N(0, \sigma^2)$.

- for an observed data (x, y) , the distribution is

$$p(x, y, \theta) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y-g(x,\theta))^2}{2\sigma^2}} p(x),$$

where $p(x)$ is the distribution that the data x satisfies.

- the likelihood function is

$$L_n(\theta) = \prod_{i=1}^n p(x_i, y_i, \theta) = \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^n \prod_{i=1}^n e^{-\frac{(y_i-g(x_i,\theta))^2}{2\sigma^2}} p(x_i).$$

Model of regression

- The log likelihood function is

$$l_n(\theta) = -\frac{1}{2} \sum_{i=1}^n (y_i - g(x_i, \theta))^2 + C + \sum_{i=1}^n \log p(x_i),$$

where $C = -n \log(\sqrt{2\pi}\sigma)$.

- Because $C + \sum_{i=1}^n \log p(x_i)$ does not depend on the θ , so

$$\begin{aligned} \arg \max_{\theta} l_n(\theta) &= \arg \max_{\theta} \left[-\frac{1}{2} \sum_{i=1}^n (y_i - g(x_i, \theta))^2 \right] \\ &= \arg \min_{\theta} \left[\frac{1}{2} \sum_{i=1}^n (y_i - g(x_i, \theta))^2 \right] \end{aligned}$$

Hence, the estimator of θ is

$$\hat{\theta}_n = \operatorname{argmin}_{\theta} \left[\frac{1}{2} \sum_{i=1}^n (y_i - g(x_i, \theta))^2 \right].$$

$$Y = g(X, \theta) + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$$

$$f(x, y; \theta) = f(x) \cdot p(y|x; \theta) = f(x) \cdot \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y - g(x, \theta))^2}{2\sigma^2}}$$

$$Y = g(x, \theta) + \varepsilon \sim N(g(x, \theta), \sigma^2)$$

$$L_n(\theta) = \prod_{i=1}^n p(x_i, y_i; \theta), \quad l_n(\theta) = \log L_n(\theta)$$

$$\hat{\theta}_n := \arg \max_{\theta} L_n(\theta) = \arg \max_{\theta} l_n(\theta)$$

$$l'_n(\theta) = 0 \Rightarrow \hat{\theta}_n = ?$$

$$\text{Example: } g(x, \theta) = \theta_1 x + \theta_0$$

$$= \arg \max_{\theta} -\frac{1}{2} \sum_{i=1}^n (y_i - g(x_i, \theta))^2$$

$$\begin{aligned} \bar{l}_n(\theta_1, \theta_0) &= -\frac{1}{2} \sum_{i=1}^n (y_i - \theta_1 x_i - \theta_0)^2 \\ \partial \bar{l}_n / \partial \theta_1 &= \sum_{i=1}^n x_i y_i - \theta_1 x_i^2 - \theta_0 x_i = n \cdot \bar{xy} - n \cdot \bar{x}^2 \cdot \theta_1 - n \cdot \bar{x} \cdot \theta_0 \\ \partial \bar{l}_n / \partial \theta_0 &= \sum_{i=1}^n (y_i - \theta_1 x_i - \theta_0) \\ &= n \cdot \bar{y} - \theta_1 \cdot n \cdot \bar{x} - n \cdot \theta_0 \end{aligned}$$

$$\Rightarrow \bar{xy} - \bar{x}^2 \cdot \theta_1 - \bar{x} \cdot \theta_0 = 0 \dots \textcircled{1}$$

$$\bar{y} - \bar{x} \cdot \theta_1 - \theta_0 = 0 \dots \textcircled{2}$$

$$\textcircled{1} - \textcircled{2} \cdot \bar{x} \Rightarrow \bar{xy} - \bar{x}^2 \cdot \theta_1 - \bar{y} \cdot \bar{x} + \bar{x}^2 \cdot \theta_1 = 0$$

$$\Rightarrow \theta_1 = \boxed{\frac{\bar{xy} - \bar{x} \cdot \bar{y}}{\bar{x}^2 - \bar{x}^2}}$$

$$\theta_0 = \bar{y} - \bar{x} \cdot \theta_1 = \boxed{\bar{y} - \bar{x}}$$

One dimensional linear regression

- Suppose $g(x, \theta) = \theta_1 x + \theta_0$ with $\theta_1 \in \mathbb{R}$ and $\theta_0 \in \mathbb{R}$ being unknown parameters to be estimated, therefore,

$$\hat{\theta}_n = \arg \min_{\theta=(\theta_0, \theta_1) \in \mathbb{R}^2} \left[\frac{1}{2} \sum_{i=1}^n (y_i - \theta_1 x_i - \theta_0)^2 \right]. \quad (1)$$

- (Class Exercise) Show that $\hat{\theta}_n = (\hat{\theta}_{n,1}, \hat{\theta}_{n,0})$ with the form:

$$\hat{\theta}_{n,1} = \frac{\bar{xy} - \bar{x}\bar{y}}{\bar{x^2} - (\bar{x})^2},$$

$$\hat{\theta}_{n,0} = \bar{y} - \left(\frac{\bar{xy} - \bar{x}\bar{y}}{\bar{x^2} - (\bar{x})^2} \right) \bar{x},$$

where $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$, $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$, $\bar{xy} = \frac{1}{n} \sum_{i=1}^n x_i y_i$, $\bar{x^2} = \frac{1}{n} \sum_{i=1}^n x_i^2$.

One dimensional polynomial regression

Suppose here is K dimensional

Suppose $g(x, \theta) = \theta_K x^K + \dots + \theta_1 x + \theta_0$ with $\theta_K \in \mathbb{R}, \dots, \theta_1 \in \mathbb{R}$ and $\theta_0 \in \mathbb{R}$ being unknown parameters to be estimated, therefore,

$$\hat{\theta} = \arg \min_{\substack{\theta = (\theta_0, \dots, \theta_K) \in \mathbb{R}^{K+1} \\ \theta \in \mathbb{R}^{K+1}}} \left[\frac{1}{2} \sum_{i=1}^n (y_i - \theta_K x_i^K - \dots - \theta_1 x_i - \theta_0)^2 \right]. \quad (2)$$

Denote $f_n(\theta_0, \dots, \theta_K) = \frac{1}{2} \sum_{i=1}^n (y_i - \theta_K x_i^K - \dots - \theta_1 x_i - \theta_0)^2$, to find the minimizer of f_n , we differentiate it with respect to $\theta_0, \dots, \theta_K$ and let these partial derivatives be 0, i.e.

$$\frac{\partial f_n(\theta)}{\partial \theta_k} = 0, \quad k = 0, \dots, K,$$

which read as

$$\sum_{i=1}^n (y_i - \theta_K x_i^K - \dots - \theta_1 x_i - \theta_0) x_i^k = 0, \quad k = 0, \dots, K. \quad (3)$$

One dimensional polynomial regression

Write

$$\mathbf{X} = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^K \\ 1 & x_2 & x_2^2 & \dots & x_2^K \\ \dots & \dots & \dots & \dots & \dots \\ 1 & x_n & x_n^2 & \dots & x_n^K \end{bmatrix}$$

Assignment 1: Verify that (3) can be rewritten as a vector equation as

$$\mathbf{X}^\top \mathbf{y} = \mathbf{X}^\top \mathbf{X} \theta.$$

where $\mathbf{y} = [y_1, \dots, y_n]^\top$ and $\theta = [\theta_0, \dots, \theta_K]^\top$.

One dimensional polynomial regression

$$\text{s.t. } X^T X = P^T D P \quad \lambda_1, \dots, \lambda_{k+1} \text{ are eigenvalues of } X^T X$$

$$\exists \text{ invertible matrix } P \text{ and } D = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_{k+1} \end{bmatrix}$$

$X^T X$ is real symmetric $\Leftrightarrow \lambda_i \neq 0$, for all i

If the matrix $X^T X \in \mathbb{R}^{(K+1) \times (K+1)}$ is invertible, then

if $X^T X \in \mathbb{R}^{(K+1) \times (K+1)}$ is not invertible

Choose $\lambda \neq \lambda_i$

$$X^T y - \lambda \theta = (X^T X - \lambda I)\theta = P^{-1}(\hat{\theta} - \lambda I)P\theta$$

$$\Rightarrow \theta_{j+1} = (X^T X - \lambda I)^{-1}(X^T y - \lambda \theta_j)$$

do iteration

$$\theta_0, \theta_1, \theta_2, \dots, \theta_j, \dots$$

$$\begin{pmatrix} \lambda - \lambda & 0 \\ 0 & \lambda - \lambda \end{pmatrix}$$

$$\rightarrow \theta^*, X^T y = X^T X \theta^*$$

$$\sum_{j=1}^n h_{kj} h_{mj} x_j^2 = \frac{h_{kj}}{|h_{kj}|} \cdot \frac{h_{mj}}{|h_{mj}|} \cdot X_j^2 \frac{1}{\bar{x}_n} \cdot \frac{1}{\bar{x}_m}$$
$$= \frac{x_j^2}{n} \cdot \frac{h_{kj}}{|h_{kj}|} \cdot \frac{h_{mj}}{|h_{mj}|}$$

Multi-d LR.

d is the dimension of the unknown θ .
 n is number of samples. $d > n$

$$y_i = \theta_0 + \theta_1 x_{i1} + \dots + \theta_d x_{id} + \varepsilon_i, \quad i=1, \dots, n.$$

$$\min_{\theta} \|y - X\theta\|^2 \Rightarrow X^T X \hat{\theta} = X^T y \quad (\star)$$

If $X^T X$ is invertible, $\hat{\theta} = (X^T X)^{-1} X^T y$

If $X^T X$ is not invertible, try

$$X \in \mathbb{R}^{n \times (d+1)}$$

$$\theta_{j+1} = (X^T X + \lambda I)^{-1} (X^T y + \lambda \theta_j) \quad \text{error compensation to ridge regression}$$

If $\lim_{j \rightarrow \infty} \theta_j = \hat{\theta}$ exists, then $\hat{\theta}$ solves (\star) .

$$\theta_{j+1} - \theta_j = (X^T X + \lambda I)^{-1} \lambda (\theta_j - \theta_{j-1})$$

$$= P(D + \lambda I)^{-1} P^{-1} \lambda (\theta_j - \theta_{j-1})$$

$$\Rightarrow P^{-1}(\theta_{j+1} - \theta_j) = (D + \lambda I)^{-1} \lambda P^{-1}(\theta_j - \theta_{j-1})$$

$$\exists \delta \in (0, 1) \quad \|P^{-1}(\theta_{j+1} - \theta_j)\| \leq \delta \|P^{-1}(\theta_j - \theta_{j-1})\| ?$$

$$(D + \lambda I)^{-1} \lambda = \begin{pmatrix} \frac{\lambda}{\lambda + \lambda_1} & & & \\ & \ddots & & 0 \\ & & \ddots & \\ 0 & & & \frac{\lambda}{\lambda + \lambda_n} \end{pmatrix}. \quad \left| \frac{\lambda}{\lambda + \lambda_i} \right| \leq \delta, \text{ for all } i$$

this fails since
there are $\lambda_i = 0$.

$$X^T X = P D P^T, \quad P^T = P^{-1}$$

$$\theta_0, \theta_1, \dots, \theta_d, \dots$$

If $\exists \delta \in (0, 1)$ s.t.

$$\|\theta_{j+1} - \theta_j\| \leq \delta \|\theta_j - \theta_{j-1}\|, \text{ for all } j$$

then $\lim_{j \rightarrow \infty} \theta_j$ exists.

$$D = \begin{pmatrix} \lambda_1 & & & \\ 0 & \ddots & & 0 \\ & & \ddots & \\ 0 & & & \lambda_{d+1} \end{pmatrix}$$

Multi-dimensional linear regression

Many problems in Statistics can be recast as a multi-dimensional linear regression as the following: given a sequence of observed data

$$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n),$$

where

- $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$ for all i ,
- one would like to learn from these given data the changing tendency of y_i according to \mathbf{x}_i ,
- For instance, $\mathbf{x}_i = [x_{i1}, x_{i2}]$ is a vector of height and weight of i -th person and y_i is his blood pressure.

Linear regression: an optimization problem

- Linear regression uses a linear function to model the relation between \mathbf{x}_i and y_i , i.e., we aim to find some $\theta = (\theta_0, \theta_1, \dots, \theta_d)' \in \mathbb{R}^{d+1}$ such that

$$y_i = \theta_0 + x_{i1}\theta_1 + \dots + x_{id}\theta_d + \epsilon_i, \quad i = 1, \dots, n,$$

where ϵ_i is the noise (measurement error), it can be written as

$$\mathbf{y} = \mathbf{X}\theta + \epsilon \tag{4}$$

where $\epsilon = [\epsilon_1, \dots, \epsilon_n]^\top$, $\mathbf{X} = [\mathbf{1}, [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top] \in \mathbb{R}^{n \times (d+1)}$ with $\mathbf{1} = [1, \dots, 1]^\top \in \mathbb{R}^n$.

n data inputs
 $\overrightarrow{\mathbf{x}_i \in \mathbb{R}^d}$
for $i=1, 2, \dots, n$

- The least square (LE) regression is to find the 'best' θ for the equation (4).

Least square regression

- The linear equation (4) can be written as

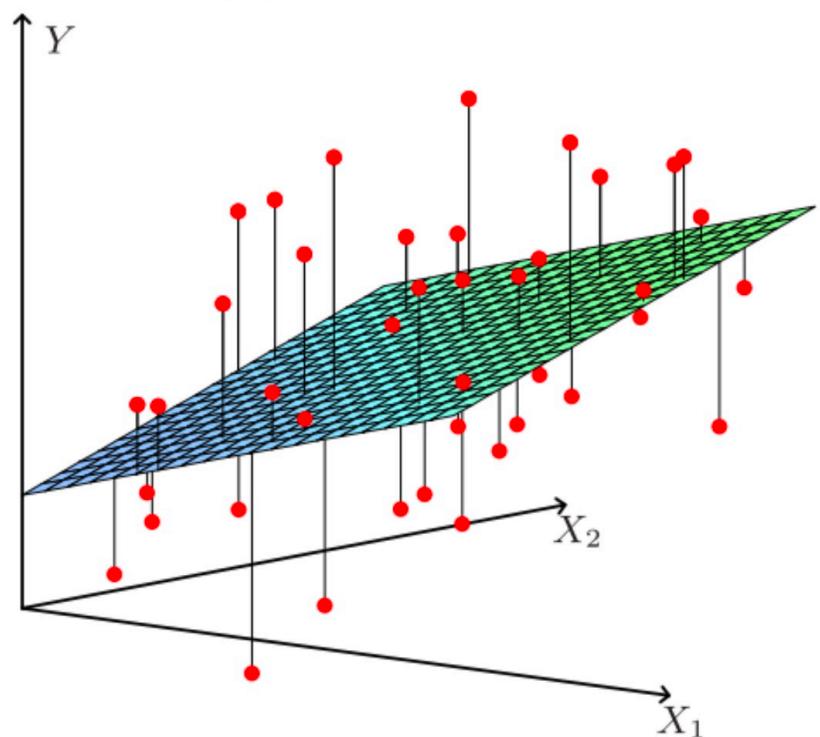
$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1d} \\ 1 & x_{21} & x_{22} & \dots & x_{2d} \\ \dots & \dots & \dots & \dots & \dots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nd} \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

- Inspired by (2), LS regression aims to minimize the following question

$$\min_{\theta} \|\mathbf{y} - \mathbf{X}\theta\|^2 = \min_{\theta} \sum_{i=1}^n (y_i - \theta_0 - x_{i1}\theta_1 - \dots - x_{id}\theta_d)^2$$

where $\|\cdot\|$ is the Euclidean distance in \mathbb{R}^n , i.e., for all $x \in \mathbb{R}^n$, $\|x\| = \sqrt{x_1^2 + \dots + x_n^2}$.

$$f(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$



Least square regression

- How to solve this problem? By differentiate $\|\mathbf{y} - \mathbf{X}\theta\|^2$ with respect to θ , i.e. for $j = 0, 1, \dots, d$,

$$\frac{\partial}{\partial \theta_j} \|\mathbf{y} - \mathbf{X}\theta\|^2 = 0.$$

$$\tilde{L}_n(\theta)$$

$$\hat{\theta}_n := \arg \min_{\theta} \tilde{L}_n(\theta)$$

$$\frac{\partial}{\partial \theta_j} \tilde{L}_n(\theta) = 0, j=0, 1, \dots, k.$$

which leads to

$$\mathbf{X}^\top \mathbf{y} = \mathbf{X}^\top \mathbf{X} \theta.$$

- As long as $\mathbf{X}^\top \mathbf{X}$ invertible, we get

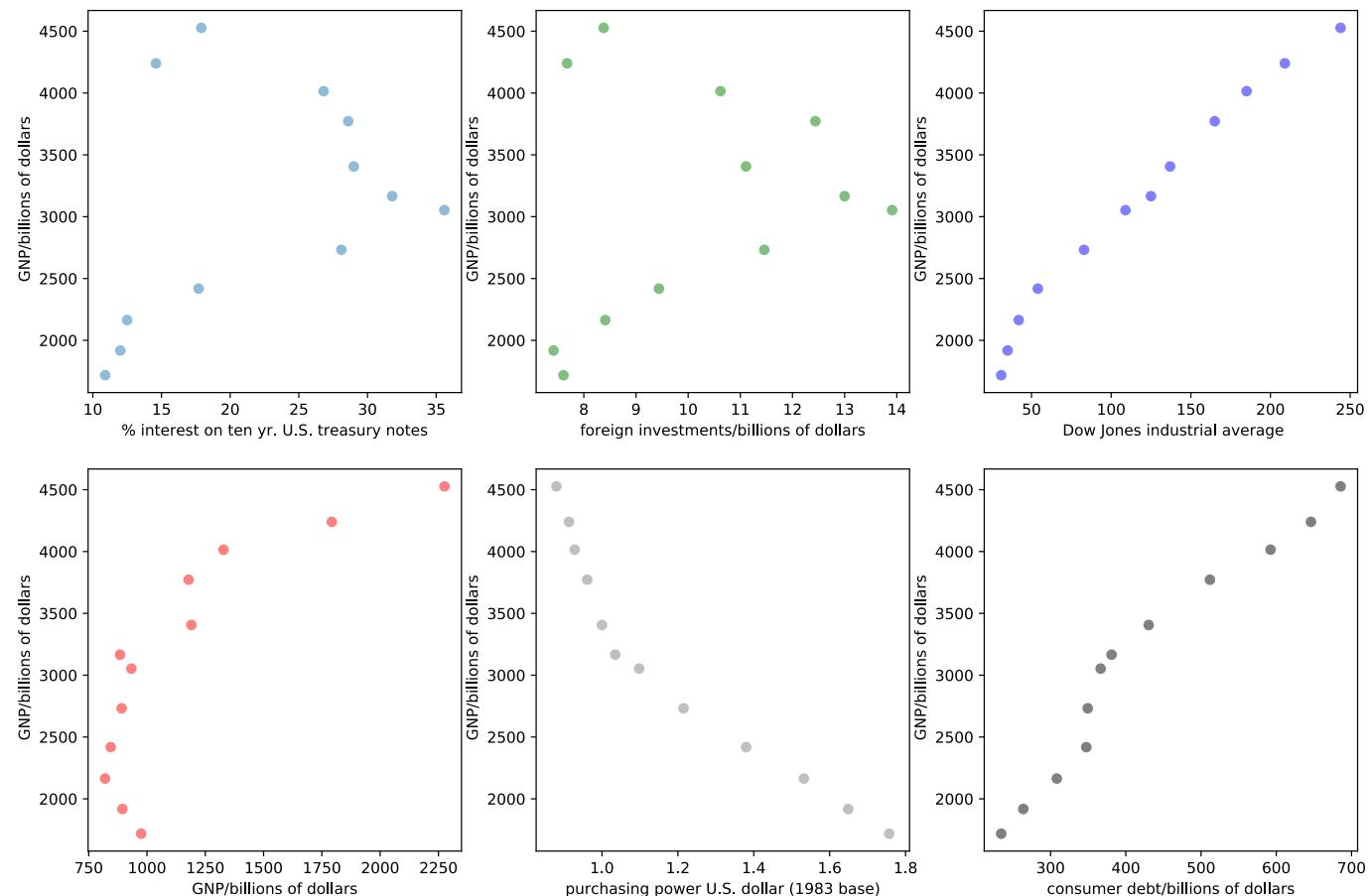
$$\hat{\theta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.$$

Example 1: US Gross National Product (GNP)

U.S. Economy Case Study U.S.. economic data 1976 to 1987

- $Y = \text{GNP/billions of dollars}$
- $X_1 = \text{dollars/barrel crude oil}$
- $X_2 = \% \text{ interest on ten yr. U.S. treasury notes}$
- $X_3 = \text{foreign investments/billions of dollars}$
- $X_4 = \text{Dow Jones industrial average}$
- $X_5 = \text{purchasing power U.S. dollar (1983 base)}$
- $X_6 = \text{consumer debt/billions of dollars}$

Example 1: US Gross National Product (GNP)



$$Y = 4.9 * X_1 + 1.1 * X_2 + 5.1 * X_3 - 1.1 * X_4 - 8.7 * X_5 + 2.1 * X_6 + 2531.6$$

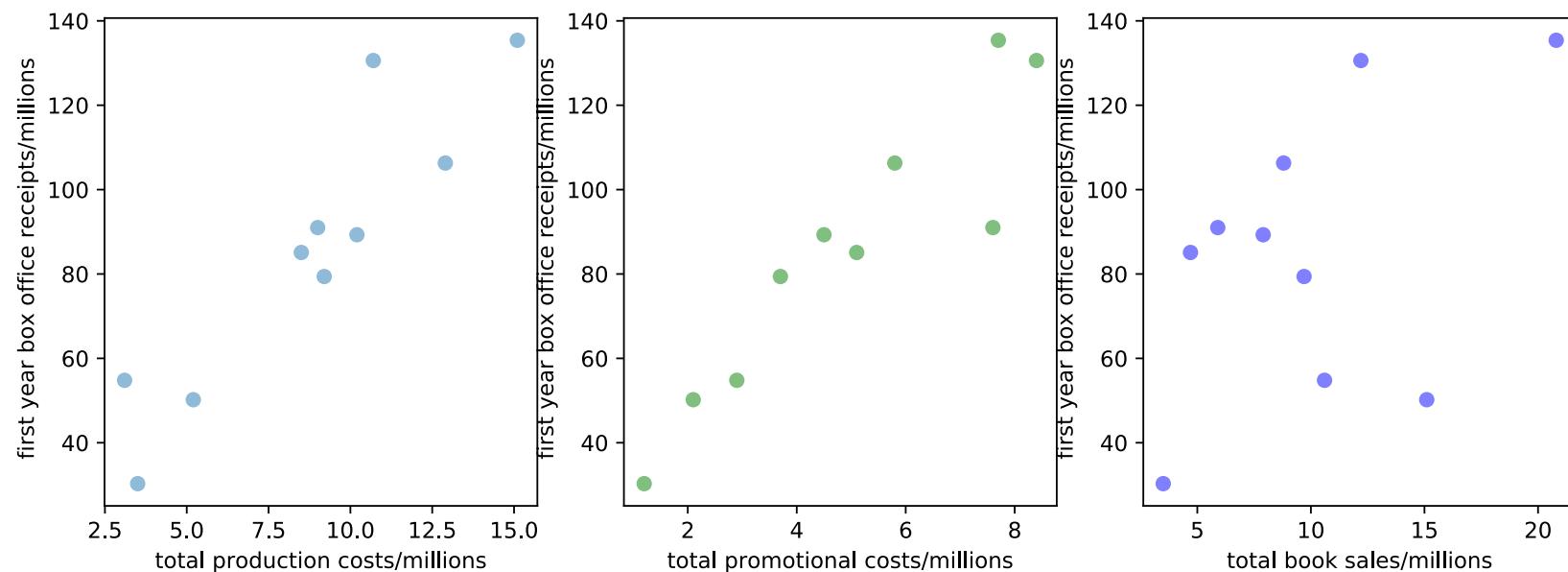
Example 1: US Gross National Product (GNP)

$$\begin{aligned} \text{GNP} = & 4.9 * \text{'CRUDE'} + 1.1 * \text{'INTEREST'} + 5.1 * \text{'FOREIGN'} - 1.1 * \text{'DJIA'} \\ & - 8.7 * \text{'PURCHASE'} + 2.1 * \text{'CONSUMER'} + 2531.6 \end{aligned}$$

Example 2: Hollywood Movies Box Office Receipts

- Y = first year box office receipts/millions,
- X_1 = total production costs/millions,
- X_2 = total promotional costs/millions,
- X_3 = total book sales/millions

Example 2: Hollywood Movies Box Office Receipts

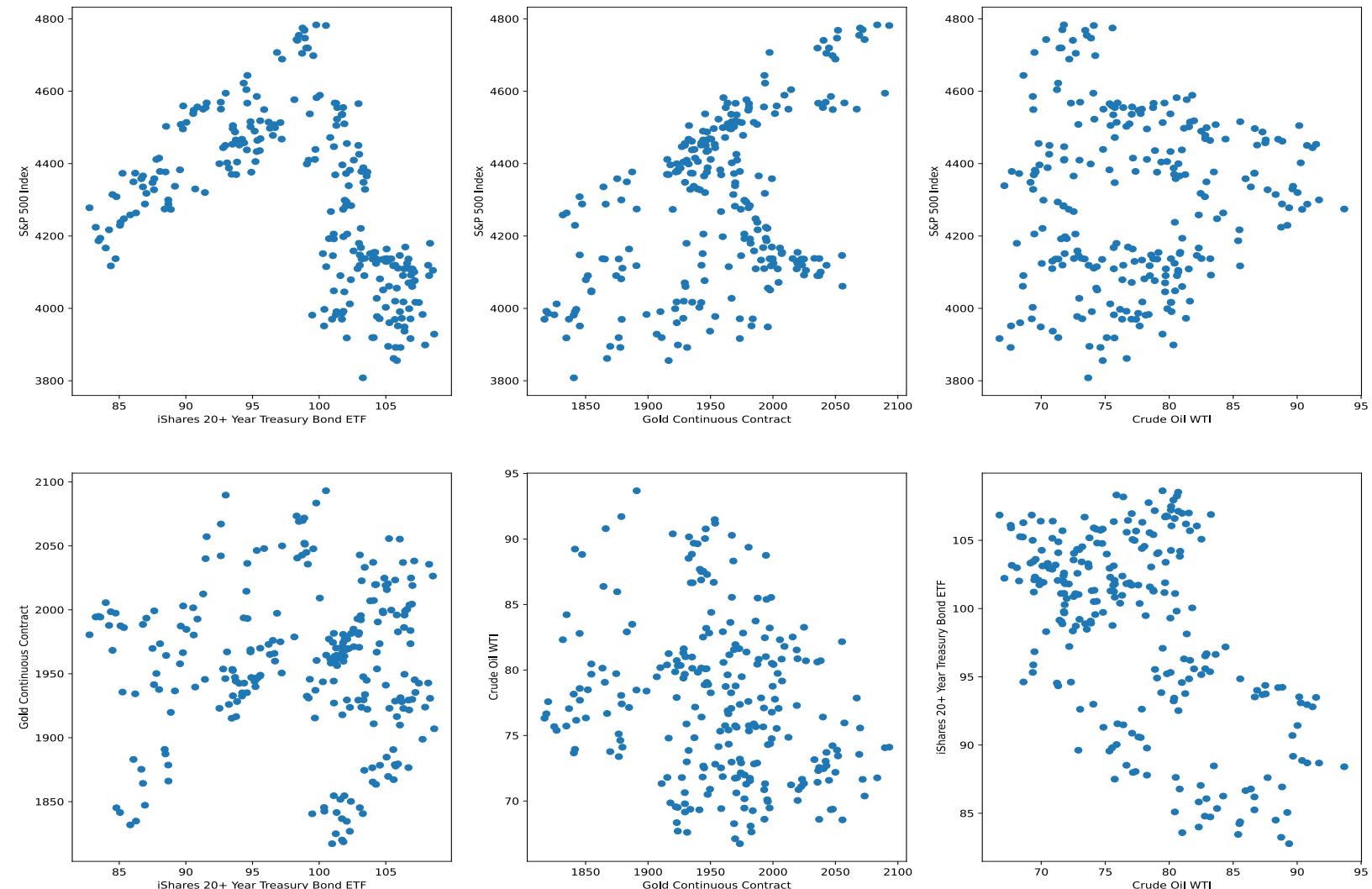


$$Y = 3.7 * X_1 + 7.6 * X_2 + 0.8 * X_3 + 7.7.$$

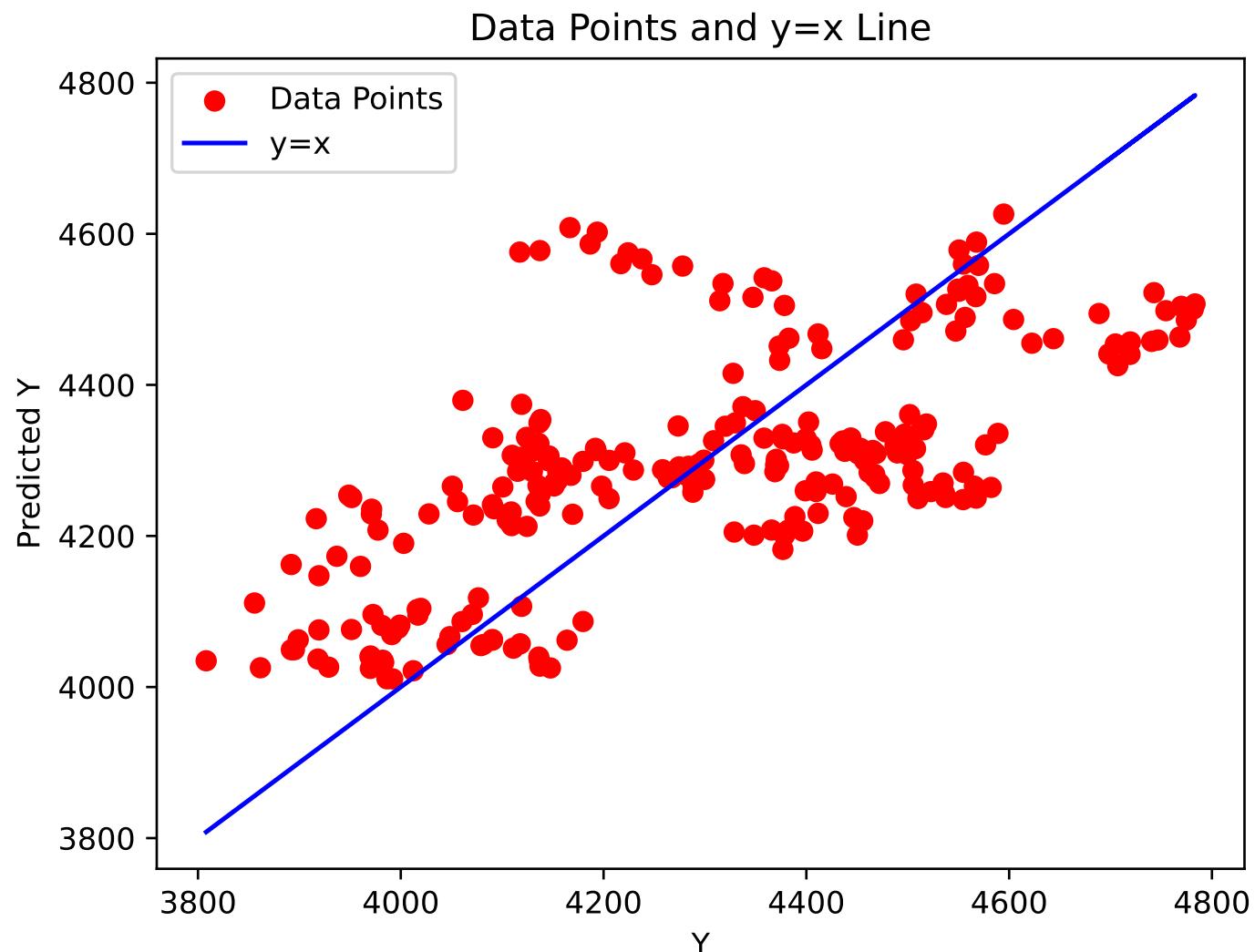
Example 3: S&P 500 Index

- $Y = \text{S\&P 500 Index}$,
- $X_1 = \text{iShares 20+ Year Treasury Bond ETF}$,
- $X_2 = \text{Gold Continuous Contract}$,
- $X_3 = \text{Crude Oil WTI}$

Example 3: S&P 500 Index



Example 3: S&P 500 Index

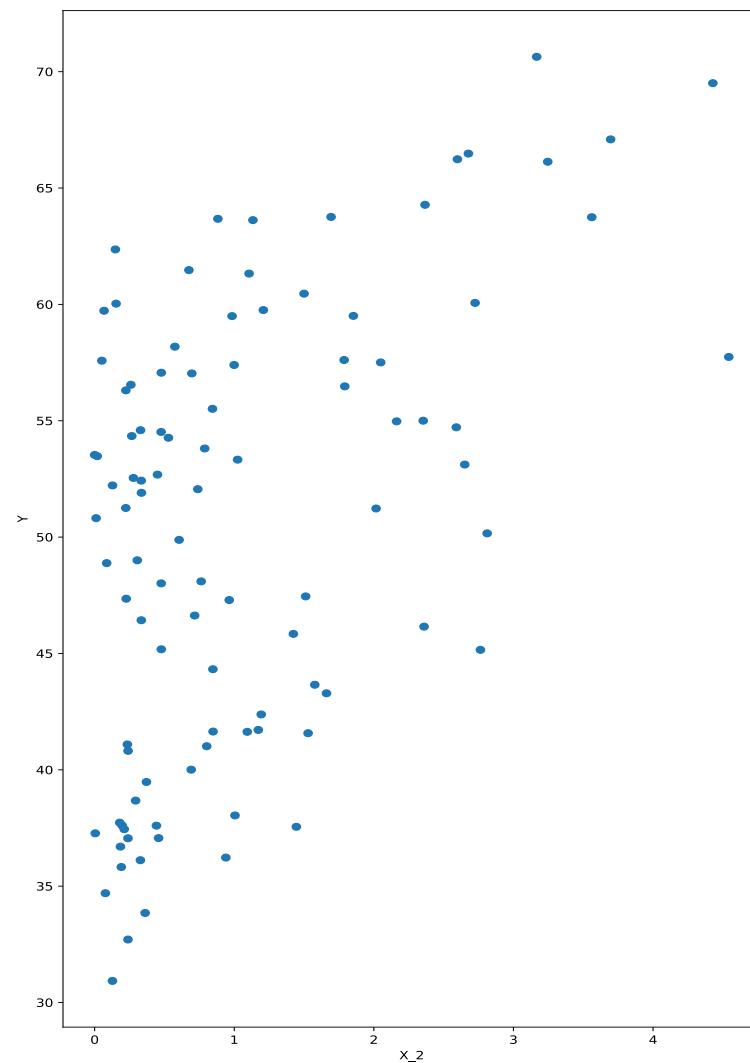
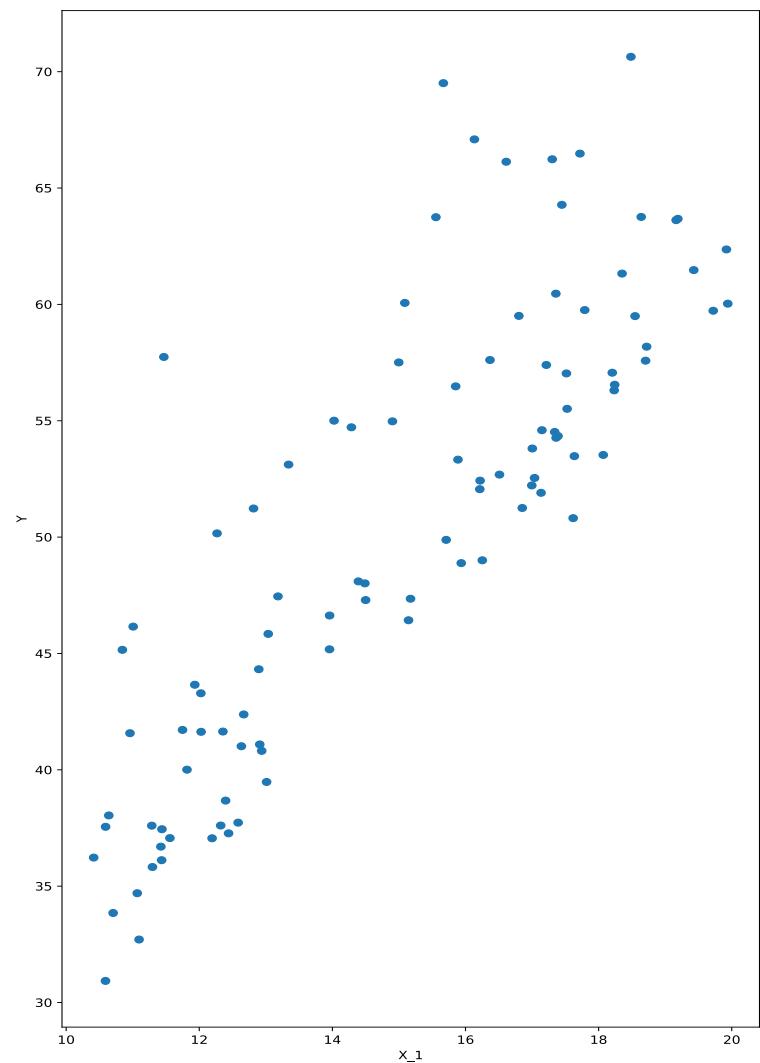


$$Y = -17.47 * X_1 + 1.66 * X_2 - 4.33 * X_3 + 33104.34.$$

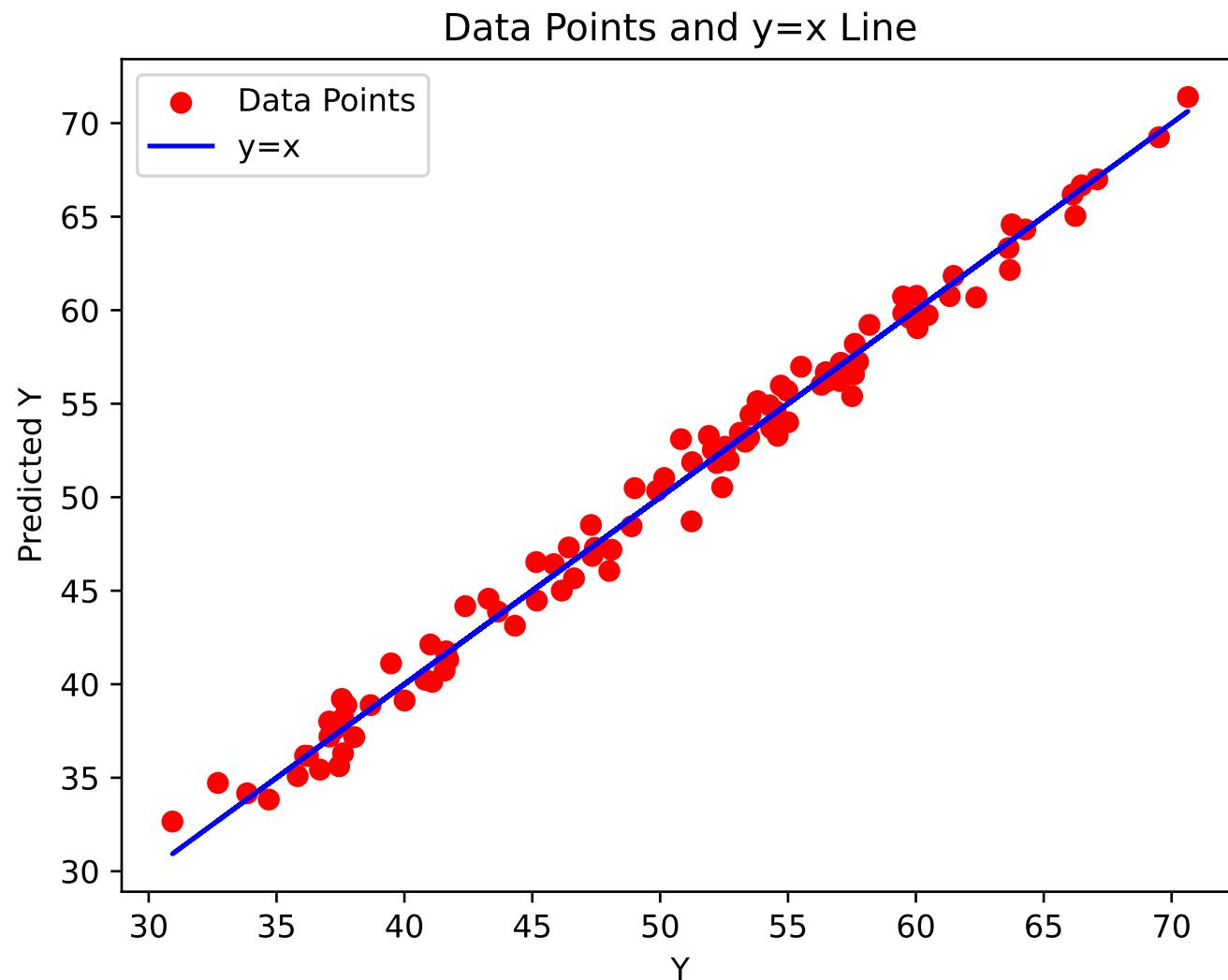
Example 4: Simulation

- $Y = 3*X1 + 5*X2 + \epsilon,$
- $X1 = \text{Uniform}(10, 20),$
- $X2 = \text{Exp}(1),$
- $\epsilon = N(0, 1)$

Example 4: Simulation



Example 4: S&P 500 Index



$$Y = 2.99 * X_1 + 4.97 * X_2 + 0.32.$$

Python codes for example 1

```
# U.S.. economic data 1976 to 1987
# X1 = dollars/barrel crude oil
# X2 = % interest on ten yr. U.S. treasury notes
# X3 = foreign investments/billions of dollars
# X4 = Dow Jones industrial average
# X5 = GNP/billions of dollars
# X6 = purchasing power U.S. dollar (1983 base)
# X7 = consumer debt/billions of dollars
data = pd.read_csv("U.S._Economy.csv")
X = data[['CRUDE','INTEREST','FOREIGN','DJIA','PURCHASE','CONSUMER']]
y = data['GNP']
fig = plt.figure(1,figsize=[15.,10.])
ax = plt.subplot(231)
plt.scatter(X['CRUDE'],y,alpha=0.5)
plt.xlabel('% interest on ten yr. U.S. treasury notes')
plt.ylabel(' GNP/billions of dollars')
plt.subplot(232)
plt.scatter(X['INTEREST'],y,c='green',alpha=0.5)
plt.xlabel('foreign investments/billions of dollars')
plt.ylabel(' GNP/billions of dollars')
plt.subplot(233)
plt.scatter(X['FOREIGN'],y,c='blue',alpha=0.5)
plt.xlabel(' Dow Jones industrial average')
plt.ylabel(' GNP/billions of dollars')
plt.subplot(234)
plt.scatter(X['DJIA'],y,c='red',alpha=0.5)
plt.xlabel('GNP/billions of dollars')
plt.ylabel(' GNP/billions of dollars')
plt.subplot(235)
plt.scatter(X['PURCHASE'],y,c='grey',alpha=0.5)
plt.xlabel('purchasing power U.S. dollar (1983 base)')
plt.ylabel(' GNP/billions of dollars')
plt.subplot(236)
plt.scatter(X['CONSUMER'],y,c='black',alpha=0.5)
plt.xlabel('consumer debt/billions of dollars')
plt.ylabel(' GNP/billions of dollars')
plt.savefig("U.S._Economy scatter plots.pdf")
plt.show()
regr10 = LinearRegression()
regr10.fit(X,y)
predicts = regr10.predict(X) #predicted values
coef10= regr10.coef_
intercept10 = regr10.intercept_
print(coef10, intercept10)
```

Python codes for example 2

```
import matplotlib.pyplot as plt
import pandas as pd
import numpy as np
from sklearn.linear_model import LinearRegression
# Hollywood Movies data,
# X1 = first year box office receipts/millions,
# X2 = total production costs/millions,
# X3 = total promotional costs/millions,
# X4 = total book sales/millions
data = pd.read_csv("Hollywood_Movies.csv")
X = data[['X2','X3','X4']]
y = data['X1']
#plot the scatters
fig = plt.figure(1,figsize=[12.,4.])
ax = plt.subplot(131)
plt.scatter(X['X2'],y,alpha=0.5)
plt.xlabel('total production costs/millions')
plt.ylabel(' first year box office receipts/millions')
plt.subplot(132)
plt.scatter(X['X3'],y,c='green',alpha=0.5)
plt.xlabel('total promotional costs/millions')
plt.ylabel(' first year box office receipts/millions')
plt.subplot(133)
plt.scatter(X['X4'],y,c='blue',alpha=0.5)
plt.xlabel('total book sales/millions')
plt.ylabel(' first year box office receipts/millions')
plt.savefig("Hollywood_Movies scatter plots.pdf")
plt.show()
# build regression model
regr4 = LinearRegression()
regr4.fit(X,y)
predicts = regr4.predict(X) #predicted values
coef4 = regr4.coef_
intercept4 = regr4.intercept_
print(coef4, intercept4)
```

High Dimensional Linear regression

- Recall the linear regression:

$$y_i = \theta_0 + x_{i1}\theta_1 + \dots + x_{id}\theta_d + \epsilon_i, \quad i = 1, \dots, n,$$

where ϵ_i is the noise (measurement error), it can be written as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta} + \boldsymbol{\epsilon}. \quad (5)$$

- As long as $\mathbf{X}^\top \mathbf{X}$ invertible, we get $\hat{\boldsymbol{\theta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$.
- In high dimensional regression problem, we often have the following case:

$$d > n,$$

i.e. the dimension of parameter is larger than the sample size. Then $\mathbf{X}^\top \mathbf{X}$ is **NOT** invertible.

Ridge Linear Regression

- We consider the linear regression:

$$\mathbf{y} = \mathbf{X}\theta + \epsilon, \quad (6)$$

with the restriction: for some $K > 0$, $\|\theta\| \leq K$.

- By Lagrangian duality, to estimate the parameter θ , we solve the following optimization problem:

$$\min_{\theta} [\|\mathbf{y} - \mathbf{X}\theta\|^2 + \lambda\|\theta\|^2], \quad (7)$$

where $\lambda > 0$ is some tuning parameter to be chosen.

- Assignment 4: For (7), we have

$$\hat{\theta} = (\lambda I_d + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}, \quad (8)$$

where I_d is the $d \times d$ identity matrix.

Ridge Regression

$$= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{y} + \lambda \mathbf{0})$$

$\lim_{j \rightarrow \infty} \theta_j = \hat{\theta}$ exists, then $\hat{\theta}$ solves (*).

$$\Rightarrow \mathbf{P}^T, \quad \mathbf{P}^T = \mathbf{P}^{-1}$$

$$\mathbf{P} = \begin{pmatrix} \theta_0 & \theta_1 & \dots & \theta_j & \dots \\ 0 & \ddots & 0 & \ddots & \vdots \\ \vdots & & \ddots & & 0 \\ 0 & \ddots & & \lambda_{d+1} & \end{pmatrix}$$

If $\exists \delta \in (0, 1)$, s.t.
 $\|\theta_{j+1} - \theta_j\| \leq \delta \cdot \|\theta_j - \theta_{j-1}\|$, for all j ,
then $\lim_{j \rightarrow \infty} \theta_j$ exists.

$$(\mathbf{D} + \lambda \mathbf{I})^T \lambda = \begin{pmatrix} \frac{\lambda}{\lambda + \lambda_1} & 0 & \dots & 0 \\ 0 & \frac{\lambda}{\lambda + \lambda_2} & \dots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \frac{\lambda}{\lambda + \lambda_{d+1}} \end{pmatrix}$$

$|\frac{\lambda}{\lambda + \lambda_i}| \leq \delta$, for all i .
this fails since there are $\lambda_i = 0$.

\Rightarrow Ridge LR.

$$\|\theta\|_2 = \sqrt{\theta_0^2 + \theta_1^2 + \dots + \theta_d^2}$$

$$\min_{\lambda} \|\mathbf{y} - \mathbf{X}(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}\|^2$$

$$\frac{\partial}{\partial \theta_j} J_R(\theta) = 0, \text{ for all } j$$

$$\Rightarrow \mathbf{X}^T (\mathbf{y} - \mathbf{X} \theta) = \lambda \theta$$

$$\Rightarrow (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \theta = \mathbf{X}^T \mathbf{y} = \hat{\theta} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

Advantages:

1. Handles Multicollinearity: Ridge Regression is effective in dealing with multicollinearity, which occurs when predictors are highly correlated. It reduces the impact of correlated predictors by shrinking their coefficients without eliminating them entirely.
2. Stable Solution: Ridge Regression provides a more stable solution compared to Lasso Regression, as it does not force coefficients to be exactly zero. This makes it more suitable when all predictors are potentially relevant.

$$\min_{\lambda} \|y - X \cdot (X^T X + \lambda I)^{-1} X^T y\|^2 \text{ for } \alpha \in (0, 1), \text{ suitable } \lambda.$$

$$\theta_{j+1} = (X^T X + \lambda I)^{-1} (X^T y + \alpha \lambda \theta_j) \Rightarrow \lim_{j \rightarrow \infty} \theta_j \text{ exists.}$$

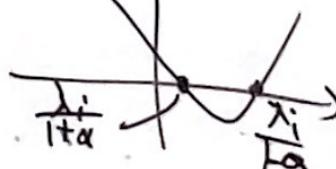
$$\|\bar{P}^{-1}(\theta_{j+1} - \theta_j)\| = \left\| \begin{pmatrix} \frac{\alpha \lambda}{\lambda + \lambda_1} & 0 & \dots & 0 \\ 0 & \ddots & \ddots & 0 \\ \vdots & \ddots & \ddots & \frac{\alpha \lambda}{\lambda + \lambda_{d+1}} \end{pmatrix} \bar{P}^{-1}(\theta_j - \theta_{j-1}) \right\| \\ \leq \delta \|\bar{P}(\theta_j - \theta_{j-1})\|$$

$\alpha = 0 \Rightarrow$ Ridge L.R.

$$\left| \frac{\alpha \lambda}{\lambda + \lambda_i} \right| < 1$$

$\alpha = 1 \Rightarrow$ Original L.R..

$$\underline{\alpha \in (0, 1)}, \theta = (X^T X + \lambda I)^{-1} (X^T y + \alpha \lambda \theta) \Leftrightarrow (1 - \alpha^2) \lambda^2 + 2\lambda_i \lambda + \lambda_i^2 > 0 \\ = (1 - \alpha^2) \cdot \left(\lambda + \frac{\lambda_i}{1 + \alpha} \right) \cdot \left(\lambda + \frac{\lambda_i}{1 - \alpha} \right) > 0$$



$$\lambda > \frac{\lambda_i}{1 - \alpha} \quad \text{or} \quad \lambda < \frac{\lambda_i}{1 + \alpha}$$

\Rightarrow Ridge L.R.

$$\|\theta\|_2 = \sqrt{\theta_0^2 + \theta_1^2 + \theta_2^2 + \dots + \theta_d^2}$$

Ridge Regression

Disadvantages:

$$\begin{aligned} y &= \chi_1 + 2\chi_2 + \chi_3 + \varepsilon = 2\chi_3 + 2 = 2\chi_1 + 4\chi_2 + \varepsilon \\ \chi_3 &= \chi_1 + 2\chi_2 \end{aligned}$$

1. Does Not Perform Feature Selection: Unlike Lasso Regression, Ridge Regression does not perform feature selection. It shrinks the coefficients of all predictors but does not force any of them to become exactly zero. Therefore, it does not provide a sparse model with selected predictors.
2. Less Interpretable: Due to the shrinkage effect, Ridge Regression may not provide a readily interpretable model, as the coefficients are not easily distinguishable between important and unimportant predictors.

Lasso Linear Regression

The least absolute shrinkage and selection operator (lasso) was put forward by Tibshirani (1996).

- We consider the linear regression:

$$\mathbf{y} = \mathbf{X}\theta + \epsilon, \quad (9)$$

with the restriction: **for some $K > 0$, $\sum_{i=0}^d |\theta_i| \leq K$.**

- By Lagrangian duality, to estimate the parameter θ , we solve the following optimization problem:

$$\min_{\theta} \left[\|\mathbf{y} - \mathbf{X}\theta\|^2 + \lambda \sum_{i=0}^d |\theta_i| \right],$$

where $\lambda > 0$ is some tuning parameter to be chosen.

Lasso Regression

Advantages:

1. Feature Selection: Lasso Regression performs feature selection by driving some coefficients to exactly zero. This makes it useful when dealing with high-dimensional datasets with many irrelevant features, as it helps in identifying the most important predictors.
2. Simplicity: Lasso Regression provides a simple and interpretable model by shrinking less important features to zero, effectively eliminating them from the model.

Lasso Regression

Disadvantages:

1. Unstable Solution: Lasso Regression can be sensitive to the presence of highly correlated predictors. In the presence of multicollinearity, it tends to arbitrarily select one variable over another, leading to an unstable solution.
2. Biased Estimates: Lasso Regression tends to introduce bias in coefficient estimates, particularly when dealing with small sample sizes or when there is a high degree of multicollinearity among predictors.

Lasso Regression v.s. Ridge Regression

In summary, Lasso Regression offers feature selection and simplicity at the cost of potential instability and bias. On the other hand, Ridge Regression handles multicollinearity and provides a stable solution but lacks feature selection capability and may be less interpretable. The choice between the two techniques depends on the specific requirements of the problem and the nature of the dataset.

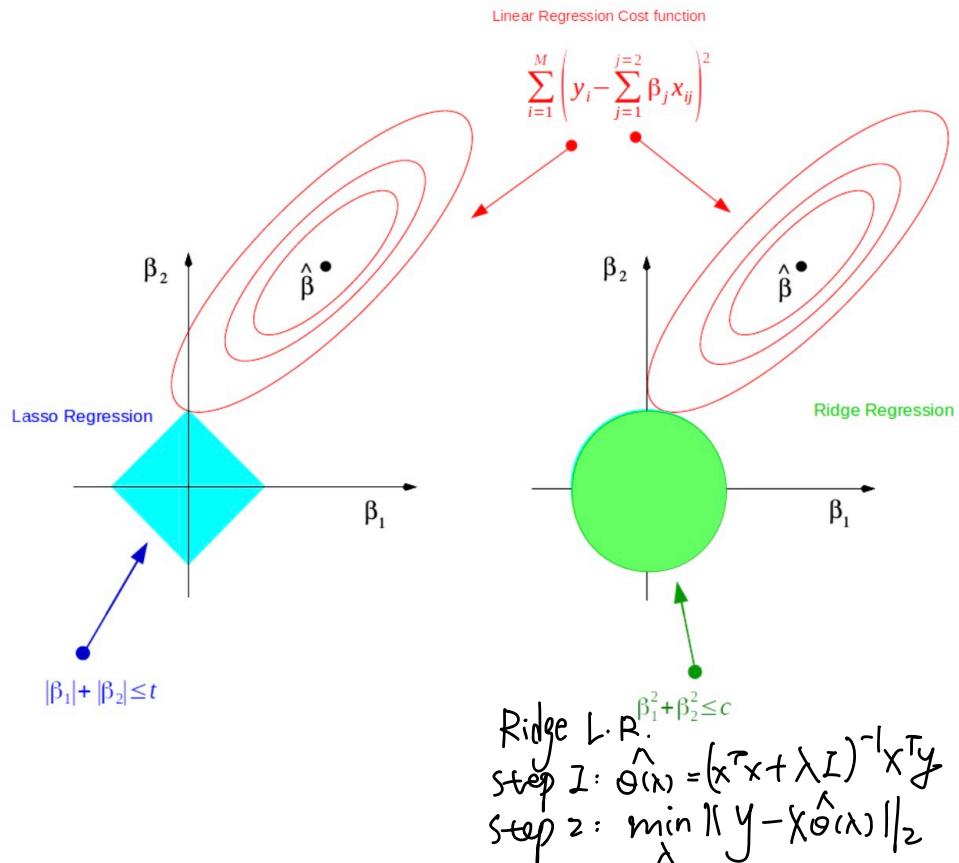
Lasso Regression v.s. Ridge Regression

$$\text{Ridge: } \min_{\theta} \frac{\|y - X\theta\|_2^2 + \lambda \|\theta\|_2^2}{= \text{MSF}}$$

$$\text{Lasso: } \min_{\theta} \|y - X\theta\|_2^2 + \lambda\|\theta\|_1$$

Geometric interpretation of Lasso and ridge regressions:

Dimension Reduction of Feature Space with LASSO



one case: linear dependence

$$\begin{cases} y = x_1 + x_2 \\ x_1 = 2x_2 \end{cases} \Rightarrow x_2 = 1.5x_1$$

$$y = \theta_1 x_1 + \theta_2 x_2$$

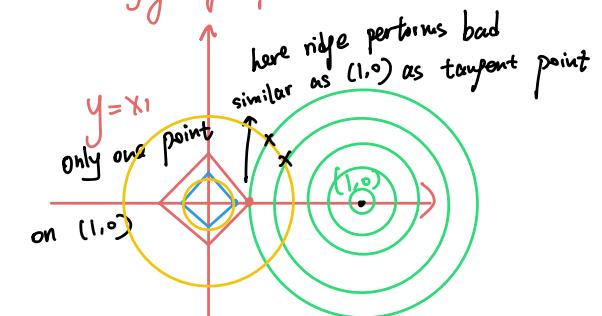
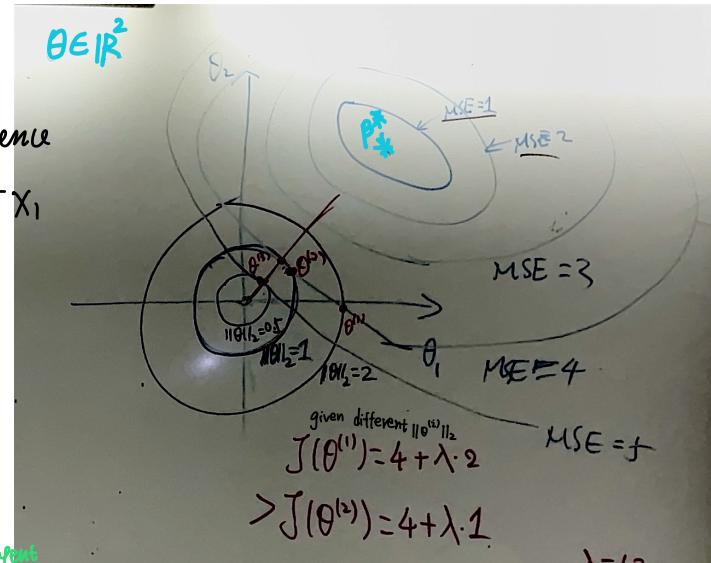
if Lasso L.R. much likely the point tangent
is on the 4 points, thus some $\beta_i = 0$, and no
change to the 4 points randomly. So Ridge L.R. is T

$$\hat{\theta}^{(1)} = (X^T X + \lambda^{(1)} I)^{-1} X^T y$$

$$\theta^{(2)} = (X^T X + \lambda^{(2)} I)^{-1} X^T y$$

if $\lambda^{(2)} > \lambda^{(1)}$

$$\| \theta^{(2)} \|_2 < \| \theta^{(1)} \|_2$$



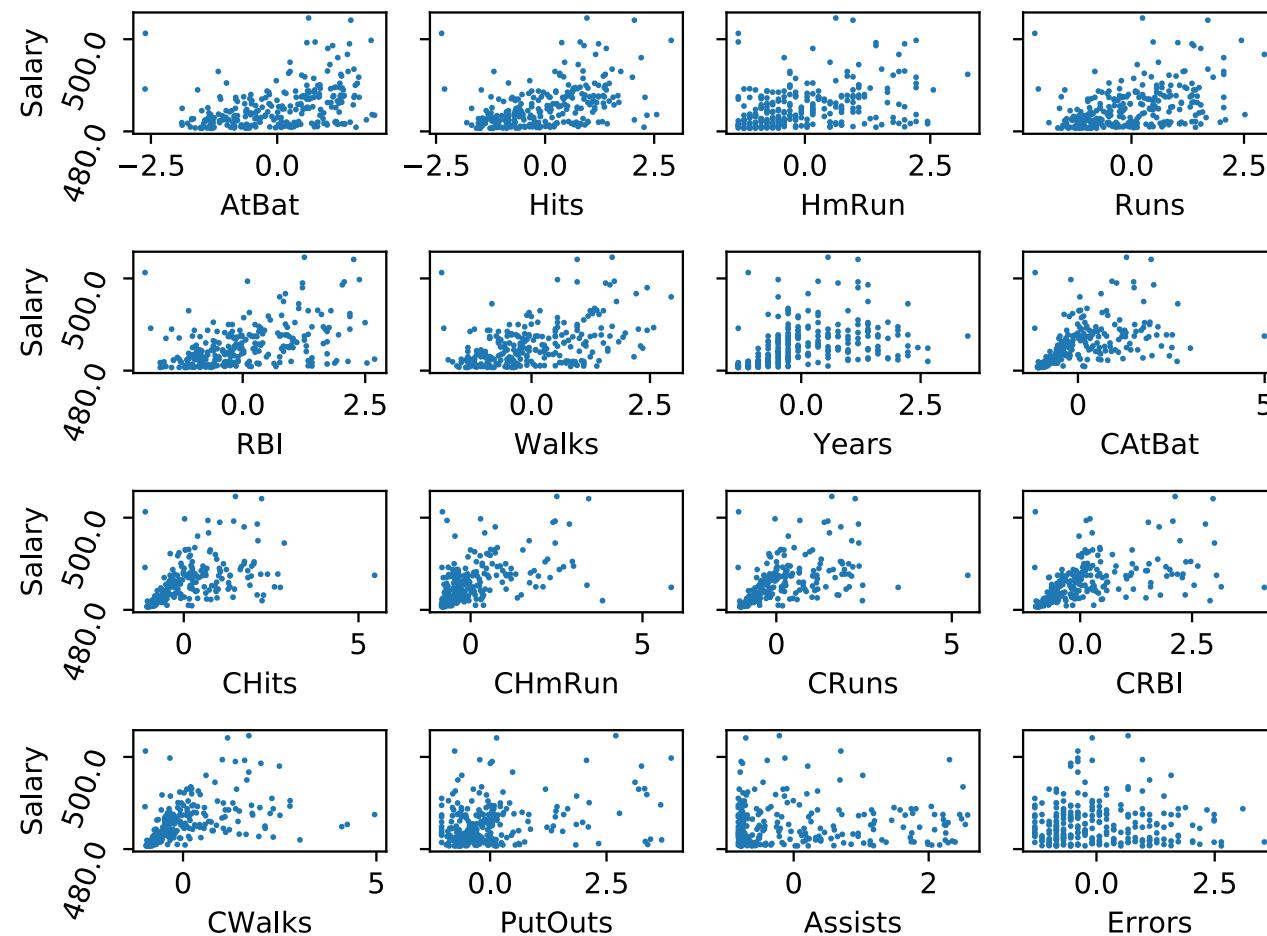
Example 3: Major League Baseball Data from the 1986 and 1987 seasons

- $Y =$ Salary 1987 annual salary on opening day in thousands of dollars,
- $X_1 =$ Number of times at bat in 1986 (AtBat),
- $X_2 =$ Number of hits in 1986 (Hits),
- $X_3 =$ Number of home runs in 1986 (HmRun),
- $X_4 =$ Number of runs in 1986 (Runs),
- $X_5 =$ Number of runs batted in 1986 (RBI),
- $X_6 =$ Number of walks in 1986 (Walks),
- $X_7 =$ Number of years in the major leagues (Years),

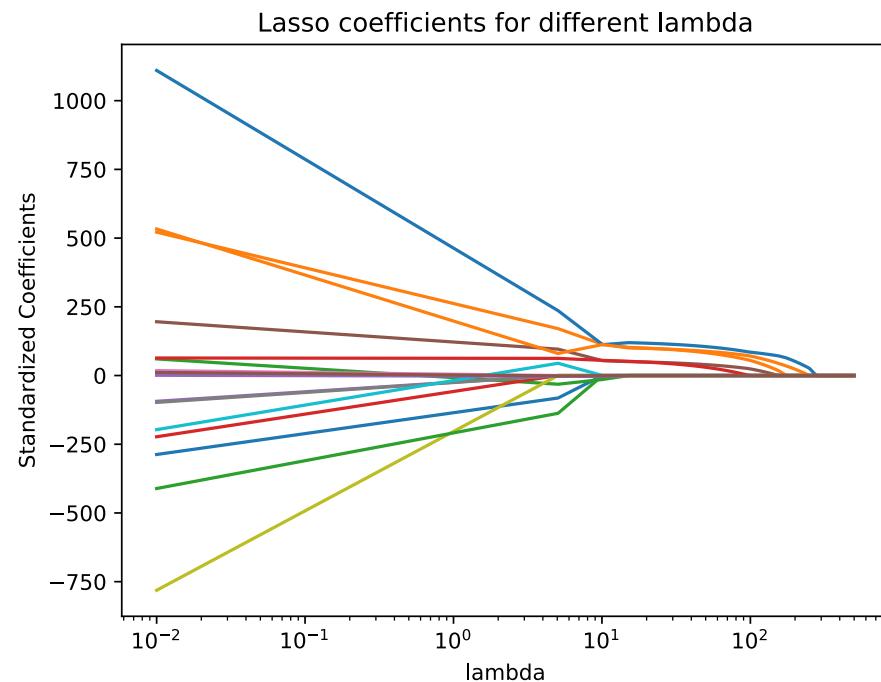
Example 3: Major League Baseball Data from the 1986 and 1987 seasons

- X_8 = Number of times at bat during his career (CAtBat),
- X_9 = Number of hits during his career (CHits),
- X_{10} = Number of home runs during his career (CHmRun),
- X_{11} = Number of runs during his career (CRuns),
- X_{12} = Number of runs batted in during his career (CRBI),
- X_{13} = Number of walks during his career (CWalks),
- X_{14} = Number of put outs in 1986 (PutOuts),
- X_{15} = Number of assists in 1986 (Assists),
- X_{16} = Number of errors in 1986 (Errors).

Example 3: Major League Baseball Data from the 1986 and 1987 seasons



Example 3: Major League Baseball Data from the 1986 and 1987 seasons



- Tuning parameters λ need to be carefully chosen.
- By 5-fold cross-validation, the optimal tuning parameter λ is 24.15.
- $Y = 97.48 * X_2 + 47.78 * X_6 + 115.97 * X_{11} + 97.28 * X_{12} + 46.16 * X_{14}$.

Python codes for example 3

```
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import Lasso
from sklearn.metrics import mean_squared_error
import numpy as np
import matplotlib.pyplot as plt
from sklearn.linear_model import LassoCV
df = pd.read_csv('Hitters.csv')
# data processing and delete the missing values
df = df.dropna()
# Convert the categorical variables to dummy variables
y = df['Salary']
X = df.drop(['Salary', 'League', 'Division', 'NewLeague'], axis=1).astype('float64')
# Standardizing X
list_numerical = X.columns
scaler1 = StandardScaler().fit(X[list_numerical])
X = scaler1.transform(X)
# Split the data into training data and test data
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=10)
# Training the Lasso regression
reg = Lasso(alpha=1)
reg.fit(X_train, y_train)
beta = reg.coef_
print('Coefficients', beta)
# Test data
pred = reg.predict(X_test)
mse_test = mean_squared_error(y_test, pred)
```

Python codes for example 3

```
# Plot of Lasso coefficients for different alpha
alphas = np.linspace(0.01,500,100)
lasso = Lasso(max_iter=10000)
coefs = []
for a in alphas:
    lasso.set_params(alpha=a)
    lasso.fit(X_train, y_train)
    coefs.append(lasso.coef_)
ax = plt.gca()
ax.plot(alphas, coefs)
ax.set_xscale('log')
plt.axis('tight')
plt.xlabel('alpha')
plt.ylabel('Standardized Coefficients')
plt.title('Lasso coefficients for different alpha')
plt.savefig('Lasso_coefficients.pdf')
plt.show()
# Lasso with 5-fold cross-validation
model = LassoCV(cv=5, random_state=0, max_iter=10000)
model.fit(X_train, y_train)
print('Optimal alpha', model.alpha_)
#Best model
lasso_best = Lasso(alpha=model.alpha_)
lasso_best.fit(X_train, y_train)
coefs = lasso_best.coef_
print(coefs)
best_mse_test = mean_squared_error(y_test, lasso_best.predict(X_test))
print('Best MSE test set', round(best_mse_test, 2))
```