

# Linear regression

Xueyi He, Tingyin Wang

School of Mathematics  
Sun Yat-sen University

August 19, 2020

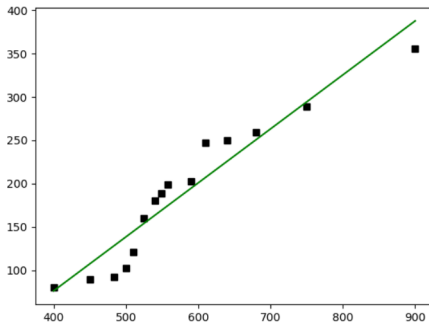
# Table of Contents

- 1 Introduction
- 2 Model specification
- 3 Some properties
- 4 Robust linear regression
- 5 Subset Selection Methods
- 6 Ridge Regression and PCA
- 7 Appendix

# Table of Contents

- 1 Introduction
- 2 Model specification
- 3 Some properties
- 4 Robust linear regression
- 5 Subset Selection Methods
- 6 Ridge Regression and PCA
- 7 Appendix

# Introduction



Given  $(x_i, y_i)$ ,  $i = 1, \dots, n$ , suppose

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i,$$

where  $E(\epsilon_i) = 0$ ,  $\text{var}(\epsilon_i) = \sigma^2$ . Then,

$$\hat{\beta}_1 = b_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}, \hat{\beta}_0 = b_0 = \bar{y} - b_1 \bar{x}.$$

# Introduction

Linear regression is the “work horse” of statistics and (supervised) machine learning. When augmented with kernels or other forms of basis function expansion, it can model also non-linear relationships. And when the Gaussian output is replaced with a Bernoulli or multinoulli distribution, it can be used for classification.

# Table of Contents

- 1 Introduction
- 2 Model specification**
- 3 Some properties
- 4 Robust linear regression
- 5 Subset Selection Methods
- 6 Ridge Regression and PCA
- 7 Appendix

# Model specification

- Generative Model:

This asserts that **the response is a linear function of the inputs**. This can be written as follows:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}.$$

where  $\boldsymbol{\epsilon}$  is the residual error between our linear predictions and the true response.

If we assume that  $\boldsymbol{\epsilon}$  has a normal distribution,

$$\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}),$$

we have:

$$\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}).$$

The log likelihood can be written:

$$\ell(\boldsymbol{\theta}) = \log\{(2\pi\sigma^2)^{-n/2}\} - \frac{1}{2\sigma^2}(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}).$$

# Model specification

- Generative Model:

$$\text{minimize : } Q = (Y - X\beta)^T(Y - X\beta).$$

$$\begin{aligned} Q &= (Y - X\beta)^T(Y - X\beta) \\ &= Y^T Y - Y^T X\beta - \beta^T X^T Y + \beta^T X^T X\beta \\ &= Y^T Y - 2Y^T X\beta + \beta^T X^T X\beta, \end{aligned}$$

$$\frac{\partial Q}{\partial \beta} = -2X^T Y + 2X^T X\beta,$$

$$\frac{\partial Q}{\partial \beta} = 0 \Rightarrow X^T X\hat{\beta} = X^T Y,$$

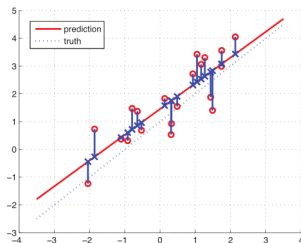
$$\Rightarrow \hat{\beta} = (X^T X)^{-1} X^T Y. \quad (1)$$



# Model specification

The corresponding solution  $\hat{\beta}$  to this linear system of equations is called the ordinary least squares or OLS solution, which is given by

$$\hat{\beta}_{OLS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}.$$



- Discriminative Model:

We seek a function  $f(X)$  for predicting  $Y$  given values of the input  $X$ .

This theory requires a loss function  $L(Y, f(X))$  for penalizing errors in prediction, and by far the most common and convenient is squared error loss:  $L(Y, f(X)) = (Y - f(X))^2$ . This leads us to a criterion for choosing  $f$

$$\begin{aligned}\text{EPE}(f) &= \mathbb{E}(Y - f(X))^2 \\ &= \int [y - f(x)]^2 \Pr(dx, dy),\end{aligned}$$

the expected (squared) prediction error .

# Model specification

- Discriminative Model:

By conditioning on  $X$ , we can write EPE as

$$\text{EPE}(f) = \mathbb{E}_X \mathbb{E}_{Y|X} ([Y - f(X)]^2 | X),$$

and we see that it suffices to minimize EPE pointwise:

$$f(x) = \operatorname{argmin}_c \mathbb{E}_{Y|X} ([Y - c]^2 | X = x).$$

The solution is

$$f(x) = \mathbb{E}(Y | X = x),$$

the conditional expectation, also known as the regression function.

Thus the best prediction of  $Y$  at any point  $X = x$  is the conditional mean, when best is measured by average squared error.

# Model specification

- Discriminative Model:

We settle for

$$\hat{f}(x) = \text{Ave}(y_i \mid x_i \in N_k(x)),$$

where "Ave" denotes average, and  $N_k(x)$  is the neighborhood containing the  $k$  points in domain closest to  $x$ .

Two approximations are happening here:

- (1) expectation is approximated by averaging over sample data;
- (2) conditioning at a point is relaxed to conditioning on some region "close" to the target point.

If the linear or some more structured model is appropriate, then we can usually get a more stable estimate than  $k$ -nearest neighbors.

- Discriminative Model:

Linear regression assumes that the regression function  $f(X)$  is approximately linear in its arguments:

$$f(X) \approx X\beta.$$

Then, we need to minimize  $R(\beta) = E\{(Y - X\beta)^T(Y - X\beta)\}$ . The solution is

$$\hat{\beta} = \left(E\{X^T X\}\right)^{-1} E\{X^T Y\}. \quad (2)$$

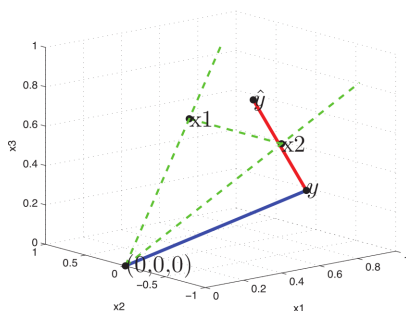
The least squares solution (1) amounts to replacing the expectation in (2) by averages over the training data.

# Geometric interpretation

The corresponding solution  $\hat{\beta}$  to this linear system of equations is called the ordinary least squares or OLS solution, which is given by

$$\underset{\hat{Y} \in \text{span}(\{\tilde{x}_1, \dots, \tilde{x}_p\})}{\operatorname{argmin}} \quad \|Y - \hat{Y}\|_2$$

$$\hat{Y} = X(X^T X)^{-1} X^T Y$$



# Regression by Successive Orthogonalization

Define

$$\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^N x_i y_i = \mathbf{x}^T \mathbf{y}.$$

Then in the univariate regression analysis,

$$\hat{\beta} = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle},$$
$$\mathbf{e} = \mathbf{y} - \mathbf{x} \hat{\beta}.$$

Suppose that the inputs  $x_1, x_2, \dots, x_p$  (the columns of the data matrix  $\mathbf{X}$ ) are orthogonal; that is  $\langle x_j, x_k \rangle = 0$  for all  $j \neq k$ . Then it is easy to check that the multiple least squares estimates  $\beta_j$  are equal to  $\langle x_j, \mathbf{y} \rangle / \langle x_j, x_j \rangle$  (the univariate estimates).

In other words, when the inputs are orthogonal, they have no effect on each other's parameter estimates in the model.

# Regression by Successive Orthogonalization

The orthogonalization does not change the subspace spanned by  $\{\mathbf{x}_j\}, j = 1, \dots, p$ , it simply produces an orthogonal basis for representing it.

---

**Algorithm 3.1** *Regression by Successive Orthogonalization.*

---

1. Initialize  $\mathbf{z}_0 = \mathbf{x}_0 = \mathbf{1}$ .

2. For  $j = 1, 2, \dots, p$

Regress  $\mathbf{x}_j$  on  $\mathbf{z}_0, \mathbf{z}_1, \dots, \mathbf{z}_{j-1}$  to produce coefficients  $\hat{\gamma}_{\ell j} = \langle \mathbf{z}_\ell, \mathbf{x}_j \rangle / \langle \mathbf{z}_\ell, \mathbf{z}_\ell \rangle$ ,  $\ell = 0, \dots, j-1$  and residual vector  $\mathbf{z}_j = \mathbf{x}_j - \sum_{k=0}^{j-1} \hat{\gamma}_{kj} \mathbf{z}_k$ .

3. Regress  $\mathbf{y}$  on the residual  $\mathbf{z}_p$  to give the estimate  $\hat{\beta}_p$ .

---

The result of this algorithm is

$$\hat{\beta}_p = \frac{\langle \mathbf{z}_p, \mathbf{y} \rangle}{\langle \mathbf{z}_p, \mathbf{z}_p \rangle}. \quad (3.28)$$



# Regression by Successive Orthogonalization

Univariate linear regression:  $y_i = a + bx_i + \varepsilon_i$ .

Let  $x_0 = [1, 1, \dots, 1]^T$ ,  $x_1$  is the explanatory variable.

Suppose  $x_0$  and  $x_1$  are not orthogonal, then for a simple univariate regression, the solution of  $a$  and  $b$ :

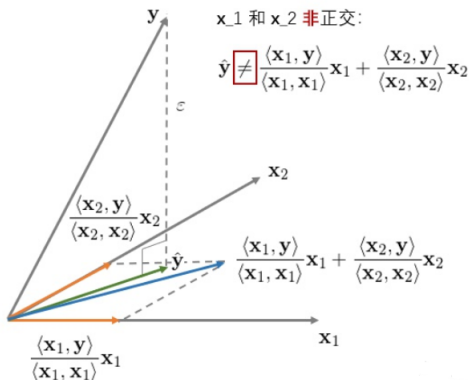
$$b = \frac{\sum x_i y_i - \frac{1}{n} \sum x_i \sum y_i}{\sum x_i^2 - \frac{1}{n} (\sum x_i)^2}, a = \frac{1}{n} \sum y_i - b \left( \frac{1}{n} \sum x_i \right).$$

# Regression by Successive Orthogonalization

Regression by Successive Orthogonalization: let  $z_0 = x_0$ , regress  $x_1$  on  $z_0$  to get the residual orthogonal with  $z_0$ , denoted as  $z_1$ .

$$\begin{aligned} z_1 &= x_1 - \frac{\langle z_0, x_1 \rangle}{\langle z_0, z_0 \rangle} z_0 = x_1 - \bar{x} \mathbf{1}, \\ b &= \frac{\langle z_1, y \rangle}{\langle z_1, z_1 \rangle} = \frac{\langle x_1 - \bar{x} \mathbf{1}, y \rangle}{\langle x_1 - \bar{x} \mathbf{1}, x_1 - \bar{x} \mathbf{1} \rangle} \\ &= \frac{\sum x_i y_i - \bar{x} \sum y_i}{\sum (x_i - \bar{x})^2} = \frac{\sum x_i y_i - \frac{1}{n} \sum x_i \sum y_i}{\sum (x_i^2 - 2x_i \bar{x} + \bar{x}^2)} \\ &= \frac{\sum x_i y_i - \frac{1}{n} \sum x_i \sum y_i}{\sum x_i^2 - 2(\sum x_i) \bar{x} + n\bar{x}^2} = \frac{\sum x_i y_i - \frac{1}{n} \sum x_i \sum y_i}{\sum x_i^2 - 2n\bar{x}^2 + n\bar{x}^2} \\ &= \frac{\sum x_i y_i - \frac{1}{n} \sum x_i \sum y_i}{\sum x_i^2 - n\bar{x}^2} = \frac{\sum x_i y_i - \frac{1}{n} \sum x_i \sum y_i}{\sum x_i^2 - \frac{1}{n} (\sum x_i)^2}. \end{aligned}$$

# Regression by Successive Orthogonalization



The multiple regression coefficient  $\hat{\beta}_j$  represents the additional contribution of  $x_j$  on  $y$ , after  $x_j$  has been adjusted for  $x_0, x_1, \dots, x_{j-1}, x_{j+1}, \dots, x_p$ .

# Multicollinearity

The estimator of  $\beta_p$ :

$$\hat{\beta}_p = \frac{\langle \mathbf{z}_p, \mathbf{y} \rangle}{\langle \mathbf{z}_p, \mathbf{z}_p \rangle}. \quad (3)$$

If  $\mathbf{x}_p$  is highly correlated with some of the other  $\mathbf{x}_k$  's, the residual vector  $\mathbf{z}_p$  will be close to zero, and from (3) the coefficient  $\hat{\beta}_p$  will be very unstable. From (3) we also obtain an alternate formula for the variance estimates  $\text{Var}(\hat{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2$ ,

$$\text{Var}(\hat{\beta}_p) = \frac{\sigma^2}{\langle \mathbf{z}_p, \mathbf{z}_p \rangle} = \frac{\sigma^2}{\|\mathbf{z}_p\|^2}.$$

In other words, the precision with which we can estimate  $\hat{\beta}_p$  depends on the length of the residual vector  $\mathbf{z}_p$ ; this represents how much of  $\mathbf{x}_p$  is unexplained by the other  $\mathbf{x}_k$  's.

# QR decomposition

We can represent step 2 of Algorithm 3.1 in matrix form:

$$X = Z\Gamma,$$

where  $Z$  has as columns the  $z_j$  (in order), and  $\Gamma$  is the upper triangular matrix with entries  $\hat{\gamma}_{kj}$ . Introducing the diagonal matrix  $D$  with  $j$  th diagonal entry  $D_{jj} = \|z_j\|$ , we get

$$\begin{aligned} X &= ZD^{-1}D\Gamma \\ &= QR, \end{aligned}$$

the so-called *QR* decomposition of  $X$ . Here  $Q$  is an  $N \times p$  orthogonal matrix,  $Q^T Q = I$ , and  $R$  is a  $p \times p$  upper triangular matrix. The QR decomposition represents a convenient orthogonal basis for the column space of  $X$ .

# QR decomposition

It is easy to see, for example, that the least squares solution is given by

$$\begin{aligned}QR\beta &= Y, \\R\beta &= Q^T Y, \\ \hat{\beta} &= R^{-1}Q^T Y, \\ \hat{Y} &= QQ^T Y.\end{aligned}$$

Equation is easy to solve because  $R$  is upper triangular.

Linear regression can be made to model non-linear relationships by replacing  $\mathbf{X}$  with some non-linear function of the inputs,  $\Phi(\mathbf{X})$ . That is,

$$p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{X}\hat{\boldsymbol{\beta}}, \sigma^2 \mathbf{I}),$$

$$p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta}) = \mathcal{N}(\hat{\boldsymbol{\beta}}\Phi(\mathbf{X}), \sigma^2 \mathbf{I}).$$

# Table of Contents

- 1 Introduction
- 2 Model specification
- 3 Some properties**
- 4 Robust linear regression
- 5 Subset Selection Methods
- 6 Ridge Regression and PCA
- 7 Appendix



# Some properties

Hat matrix:  $H = X(X^T X)^{-1} X^T$ , symmetric and idempotent.

Distribution of estimator:  $\hat{\beta} \sim N\left(\beta, \sigma^2 (X^T X)^{-1}\right)$ .

Distribution of fitted value:  $\hat{Y} = HY \sim N(X\beta, \sigma^2 H)$ .

Distribution of residuals:  $e \sim N(0, (I - H)\sigma^2)$ .

Analysis of variance:

$$SST = Y^T \left[ I - \left( \frac{1}{n} \right) J \right] Y \sim \sigma^2 \chi^2(n - 1, \delta_1),$$

$$SSE = Y^T [I - H] Y \sim \sigma^2 \chi^2(n - p, 0),$$

$$SSR = Y^T \left[ H - \left( \frac{1}{n} \right) J \right] Y \sim \sigma^2 \chi^2(p - 1, \delta_2).$$

# Some properties

Hat matrix:  $H = X(X^T X)^{-1} X^T$ .

Property of hat matrix  $H$  :

- $HY = \hat{Y}$ ,  $HX = X$ ,  $H\hat{Y} = \hat{Y}$ ,  $He = 0$   
 $HX = X(X^T X)^{-1} X^T X = X$ ,  $H\hat{Y} = HX\hat{\beta} = X\hat{\beta} = \hat{Y}$ .

- symmetric

$$H^T = \left( X(X^T X)^{-1} X^T \right)^T = X(XX)^{-1} X^T = H.$$

- idempotent

$$HH = X(X^T X)^{-1} X^T X (X^T X)^{-1} X^T = X(X^T X)^{-1} X^T = H.$$

$\text{rank}(H) = \text{trace}(H)$ .

# Some properties

Distribution of estimator:

$$\hat{\beta} = (X^T X)^{-1} X^T Y.$$

$$Y \sim N(X\beta, \sigma^2 I).$$

$$\begin{aligned} E\{\hat{\beta}\} &= (X^T X)^{-1} X^T E\{Y\} \\ &= (X^T X)^{-1} X^T X \beta = \beta, \end{aligned}$$

$$\begin{aligned} \text{var}\{\hat{\beta}\} &= (X^T X)^{-1} X^T \text{var}\{Y\} \left[ (X^T X)^{-1} X^T \right] \\ &= (X^T X)^{-1} X^T \sigma^2 X (X^T X)^{-1} = \sigma^2 (X^T X)^{-1}, \end{aligned}$$

$$\hat{\beta} \sim N\left(\beta, \sigma^2 (X^T X)^{-1}\right).$$

# Some properties

Distribution of fitted value:

$$\hat{Y} = X\hat{\beta} = X(X^T X)^{-1} X^T Y = HY.$$

$$E\{\hat{Y}\} = E\{HY\} = HE\{Y\} = HX\beta = X\beta,$$

$$\text{var}\{\hat{Y}\} = \text{var}\{HY\} = H \text{var}\{Y\} H^T = \sigma^2 H,$$

$$\hat{Y} = HY \sim N(X\beta, \sigma^2 H).$$

# Some properties

Distribution of residuals:

$$e = Y - \hat{Y} = Y - HY = (I - H)Y.$$

$$E\{e\} = E\{(I - H)Y\} = (I - H)E\{Y\} = (I - H)X\beta = X\beta - X\beta = 0,$$

$$\text{var}\{e\} = (I - H)\sigma^2 I (I - H)^T = \sigma^2(I - H).$$

$$e \sim N(0, (I - H)\sigma^2).$$

The matrix  $I - H$  is also symmetric and idempotent.

$$[I - H][I - H] = I - H - H + H^2 = I - H.$$

# Some properties

Analysis of variance:

$$SST = \sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n Y_i^2 - \frac{1}{n} \left( \sum_{i=1}^n Y_i \right)^2.$$

$$SST = Y^T Y - \left( \frac{1}{n} \right) Y^T J Y = Y^T \left[ I - \left( \frac{1}{n} \right) J \right] Y, \mathbf{J} = \mathbf{1}\mathbf{1}^T.$$

$$SSE = e^T e = [(I - H)Y]^T [(I - H)Y] = Y^T (I - H)^2 Y = Y^T (I - H) Y.$$

$$SSR = SST - SSE$$

$$= Y^T \left[ I - \left( \frac{1}{n} \right) J \right] Y - Y^T [I - H] Y = Y^T \left[ H - \left( \frac{1}{n} \right) J \right] Y.$$

# Some properties

Analysis of variance:

$$SST = Y^T \left[ I - \left( \frac{1}{n} \right) J \right] Y, \quad \text{rank} \left[ I - \left( \frac{1}{n} \right) J \right] = n - 1.$$

$$SSE = Y^T [I - H] Y, \quad \text{rank}[I - H] = n - p.$$

$$SSR = Y^T \left[ H - \left( \frac{1}{n} \right) J \right] Y, \quad \text{rank} \left[ H - \left( \frac{1}{n} \right) J \right] = p - 1.$$

$H, J/n, I - J/n, I - H, H - J/n$  are idempotent and symmetric

$$\text{rank}[H] = \text{tr} \left[ X (X^T X)^{-1} X^T \right] = \text{tr} \left[ (X^T X)^{-1} X^T X \right] = \text{tr} [I_p] = p.$$

$$\text{rank} \left[ I - \left( \frac{1}{n} \right) J \right] = n - 1, \text{rank}[I - H] = n - p, \text{rank} \left[ H - \left( \frac{1}{n} \right) J \right] = p - 1.$$

# Some properties

Analysis of variance:

$$SST = Y^T \left[ I - \left( \frac{1}{n} \right) J \right] Y \sim \sigma^2 \chi^2(n-1, \delta_1),$$

$$SSE = Y^T [I - H] Y \sim \sigma^2 \chi^2(n-p, 0),$$

$$SSR = Y^T \left[ H - \left( \frac{1}{n} \right) J \right] Y \sim \sigma^2 \chi^2(p-1, \delta_2).$$

Where

$$\delta_1 = \frac{1}{\sigma^2} (X\beta)^T \left[ I - \left( \frac{1}{n} \right) J \right] X\beta,$$

$$\delta_2 = \frac{1}{\sigma^2} (X\beta)^T \left[ H - \left( \frac{1}{n} \right) J \right] X\beta.$$

$SSR \perp SSE$  (See Cochran's Theorem in Appendix).



# Test whether $\beta_k = 0$

$$\hat{\beta} \sim N\left(\beta, \sigma^2 (X^T X)^{-1}\right).$$

$$\hat{\beta}_j \sim N(\beta_j, c_{jj}\sigma^2),$$

$$\frac{\hat{\beta}_j - \beta_j}{\sqrt{c_{jj}}\sigma} \sim N(0, 1),$$

$$\frac{SSE}{\sigma^2} \sim \chi^2(n - p),$$

$$\hat{\sigma} = \sqrt{\frac{SSE}{n - p}},$$

$$t_j = \frac{\hat{\beta}_j}{\sqrt{c_{jj}}\hat{\sigma}} \sim t(n - p).$$

$c_{jj}$  is  $j$ th diagonal element of  $(X^T X)^{-1}$ .

# Test whether $\beta_k = 0$ or several $\beta_k = 0$

$H_0 : \beta_q = \dots = \beta_{p-1} = 0, H_A : \text{At least one of } \beta_q \dots \beta_{p-1} \neq 0$

Full Model:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_{p-1} X_{i,p-1} + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2)$$

Reduced Model:

$$Y_i = \beta_0 + \beta_1 X_{i1} \dots + \beta_{q-1} X_{i,q-1} + \varepsilon_i \quad (q < p).$$

Full Model:  $SSE(F) = SSE(X_1, X_2, \dots, X_{p-1}), \quad df_F = n - p.$

Reduced Model:  $SSE(R) = SSE(X_1, X_2, \dots, X_{q-1}), \quad df_R = n - q.$

General Linear Test:

$$F^* = \frac{\left[ \frac{SSE(R) - SSE(F)}{df_R - df_F} \right]}{\left[ \frac{SSE(F)}{df_F} \right]} \stackrel{H_0}{\sim} F(p - q, n - p).$$

# R square and adjusted R square

The coefficient of multiple determination  $R^2$  is defined as:

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST} (0 \leq R^2 \leq 1).$$

$R^2$  always increases when there are more variables.

Adjusted  $R^2$  :

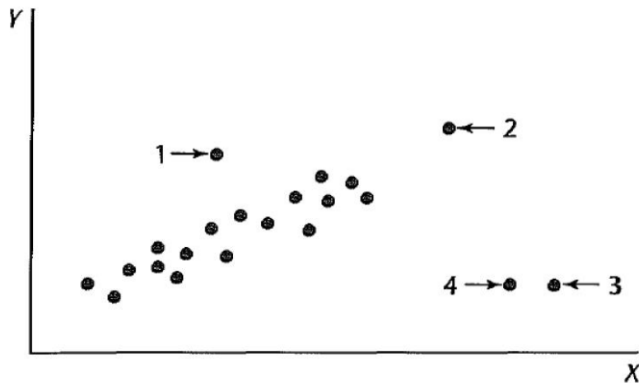
$$R_a^2 = 1 - \frac{[SSE/(n-p)]}{[SST/(n-1)]} = 1 - \frac{MSE}{MST} = 1 - \left( \frac{n-1}{n-p} \right) \frac{SSE}{SST}.$$

Adjusted  $R^2$  may decrease when  $p$  is large.

# Table of Contents

- 1 Introduction
- 2 Model specification
- 3 Some properties
- 4 Robust linear regression**
- 5 Subset Selection Methods
- 6 Ridge Regression and PCA
- 7 Appendix

# Outliers



# Outlying Y Observations

Residuals:

$$e \sim N(0, (I - H)\sigma^2)$$

Studentized residual:

$$\frac{e_i}{s\{e_i\}} = \frac{e_i}{\sqrt{MSE(1 - h_{ii})}}.$$

$$E\{e_i\} = 0; \quad \text{Let } h_{ij} = (i, j)^{th} \text{ element of } H = X(X^T X)^{-1}X^T,$$

$$\sigma^2\{e_i\} = \sigma^2(1 - h_{ii}), \quad \sigma\{e_i, e_j\} = -h_{ij}\sigma^2, \quad \forall i \neq j$$

$$s^2\{e_i\} = MSE(1 - h_{ii}), \quad s\{e_i, e_j\} = -h_{ij}MSE, \quad \forall i \neq j$$

The rough criterion is that the studentized residual is greater than 2 or less than -2.

# Outlying X-Cases

$$H = X (X^T X)^{-1} X^T, \quad h_{ij} = x_i^T (X^T X)^{-1} x_j, \\ 0 \leq h_{ii} \leq 1, \quad \sum_{i=1}^n h_{ii} = \text{trace}(H) = p.$$

- $h_{ii}$  known as the leverage of  $i$ th case. It is a measure of distance between the  $X_i$  value and the mean of the  $X$  values.
- Cases with  $X$ -levels close to the "center" of the sampled  $X$ -levels will have small leverages. Cases with "extreme" levels have large leverages
- Large leverage values:  $h_{ii} > 2p/n$ .

$$\hat{Y} = HY \Rightarrow \hat{Y}_i = \sum_{j=1}^n h_{ij} Y_j = h_{i1} Y_1 + h_{i2} Y_2 + \dots + h_{in} Y_n.$$

- Thus  $h_{ii}$  is a measure of how much  $Y_i$  is contributing to the prediction  $\hat{Y}_i$ .
- Cases with large leverages have the potential to "pull" the regression equation toward their observed  $Y$ -values.

# Identifying Influential Cases

Influence on all fitted values (Cook's Distance):

$$\begin{aligned} D_i &= \frac{\sum_{j=1}^n \left( \hat{Y}_j - \hat{Y}_{j(i)} \right)^2}{pMSE} = \frac{\left( \hat{Y} - \hat{Y}_{(i)} \right)^T \left( \hat{Y} - \hat{Y}_{(i)} \right)}{pMSE} \\ &= \frac{e_i^2}{pMSE} \left[ \frac{h_{ii}}{(1 - h_{ii})^2} \right] = \frac{h_{ii}}{(1 - h_{ii})} \frac{\tilde{e}_i^2}{p}. \text{(See Appendix)} \end{aligned}$$

where  $\hat{Y}_{j(i)}$  is the fitted value when regression is fit on the other  $n-1$  cases,  
 $\tilde{e}_i = \frac{e_i}{\sqrt{MSE(1-h_{ii})}}$  is studentized residual.

Problem cases are  $> F(0.50; p, n - p)$ .



# Robust linear regression

One way to achieve robustness to outliers is to replace the Gaussian distribution for the response variable with a distribution that has heavy tails. Such a distribution will assign higher likelihood to outliers, without having to perturb the straight line to “explain” them.

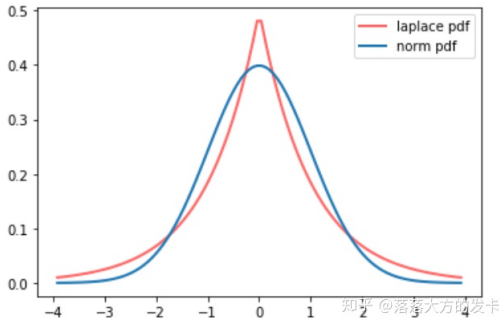
Gaussian distribution:

$$p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\beta}) = \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}).$$

Laplace distribution:

$$\begin{aligned} p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\beta}, b) &= \text{Lap}(\mathbf{Y}|\boldsymbol{\beta}^T \mathbf{X}, b) \\ &\propto \exp\left(-\frac{1}{b}|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}|\right). \end{aligned}$$

# Robust linear regression



The robustness arises from the use of  $|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}|$  instead of  $(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})$ . For simplicity, we will assume  $b$  is fixed. Let  $r_i \triangleq y_i - \boldsymbol{\beta}^T \mathbf{x}_i$  be the  $i$ 'th residual. The NLL(negative log likelihood) has the form

$$\ell(\boldsymbol{\beta}) = \sum_i |r_i(\boldsymbol{\beta})|.$$

# Robust linear regression

This is a non-linear objective function.

But we can convert the NLL to a linear objective, subject to linear constraints, using the following split variable trick.

First we define

$$r_i^+ = \frac{1}{2}(|r_i| + r_i) \geq 0, r_i^- = \frac{1}{2}(|r_i| - r_i) \geq 0, |r_i| \triangleq r_i^+ + r_i^-,$$

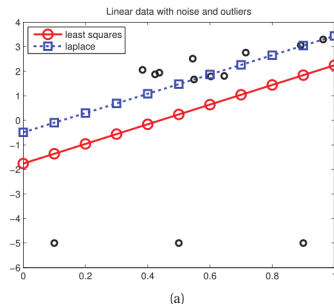
and then we impose the linear inequality constraints that  $r_i^+ \geq 0, r_i^- \geq 0$ . Now the constrained objective becomes

$$\min_{\mathbf{w}, \mathbf{r}^+, \mathbf{r}^-} \sum_i (r_i^+ + r_i^-) \text{ s.t. } r_i^+ \geq 0, r_i^- \geq 0, \boldsymbol{\beta}^T \mathbf{x}_i + r_i^+ - r_i^- = y_i.$$

This is an example of a linear program with  $p + 2N$  unknowns and  $3N$  constraints. It is a **convex optimization problem**, it has a unique solution solved by any Linear Programming solver.

# Outlying Y Observations

If we have outliers in our data, this can result in a poor fit.

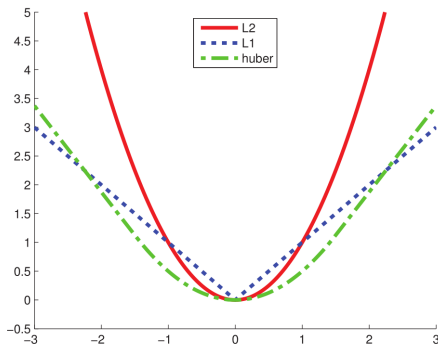


This is because squared error penalizes deviations quadratically, so points far from the line have more effect on the fit than points near the line.

# Robust linear regression

An alternative to using NLL under a Laplace likelihood is to minimize the Huber loss function (Huber 1964), defined as follows:

$$L_H(r, \delta) = \begin{cases} r^2/2 & \text{if } |r| \leq \delta \\ \delta|r| - \delta^2/2 & \text{if } |r| > \delta \end{cases}$$



# R example

```
> library(MASS)
> names(Boston)
[1] "crim"      "zn"        "indus"     "chas"      "nox"       "rm"        "age"       "dis"       "rad"
[10] "tax"       "ptratio"   "black"     "lstat"     "medv"
> lm.fit = lm(medv ~ ., data = Boston)
> summary(lm.fit)
```

Call:  
lm(formula = medv ~ ., data = Boston)

Residuals:

Min	1Q	Median	3Q	Max
-15.595	-2.730	-0.518	1.777	26.199

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	3.646e+01	5.103e+00	7.144	3.28e-12	***
crim	-1.080e-01	3.286e-02	-3.287	0.001087	**
zn	4.642e-02	1.373e-02	3.382	0.000778	***
indus	2.056e-02	6.150e-02	0.334	0.738288	
chas	2.687e+00	8.616e-01	3.118	0.001925	**
nox	-1.777e+01	3.820e+00	-4.651	4.25e-06	***
rm	3.810e+00	4.179e-01	9.116	< 2e-16	***
age	6.922e-04	1.321e-02	0.052	0.958229	
dis	-1.476e+00	1.995e-01	-7.398	6.01e-13	***
rad	3.060e-01	6.635e-02	4.613	5.07e-06	***
tax	-1.233e-02	3.760e-03	-3.280	0.001112	**
ptratio	-9.527e-01	1.308e-01	-7.283	1.31e-12	***
black	9.312e-03	2.686e-03	3.467	0.000573	***
lstat	-5.248e-01	5.072e-02	-10.347	< 2e-16	***

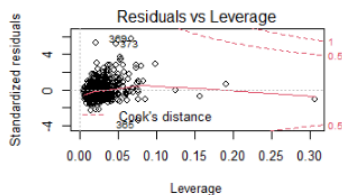
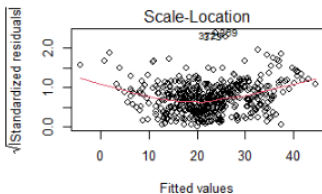
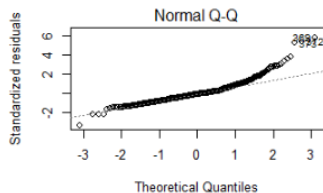
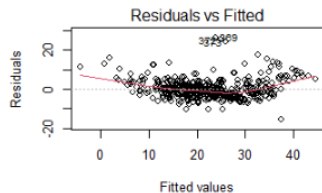
---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 4.745 on 492 degrees of freedom  
Multiple R-squared: 0.7406, Adjusted R-squared: 0.7338  
F-statistic: 108.1 on 13 and 492 DF, p-value: < 2.2e-16

# R example

```
{r}  
par(mfrow=c(2,2))  
plot(lm.fit)
```



# Table of Contents

- 1 Introduction
- 2 Model specification
- 3 Some properties
- 4 Robust linear regression
- 5 Subset Selection Methods**
- 6 Ridge Regression and PCA
- 7 Appendix



# Too Much Variables

Some variables does not contain enough useful information and can be redundant. **An example:** A hospital is trying to study the survival time after a liver resection. 10 variables of 54 patients are collected, but we only take 6 of them into consideration.

## Variables

V1	Blood clotting (凝血) score
V2	Prognostic index (预后指数)
V3	Enzyme function test score (酶功能测试分数)
V4	Liver function test score (肝功能测试分数)
V5	Age (年龄)
V10 (considered as Y)	Survival time (存活时间) (log-form)

## Models

Empty Model	$V_{10} \sim null$
Model 1	$V_{10} \sim V_1 + V_2 + V_3$
Model 2	$V_{10} \sim V_1 + V_2 + V_3 + V_4 + V_5$

## Error

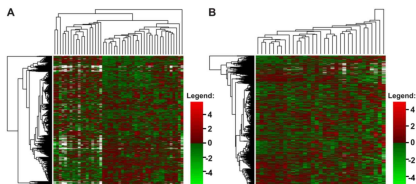
NO. of Variables	0	3	5
Predicted Error	0.147569	0.075036	0.071955

# High Dimensional Problem

Sometimes extreme cases occur and the amount of variables  $p$  can be so large and surpass sample size  $n$ . And a problem like this is unsolvable (basically because  $X^T X$  is not invertible).

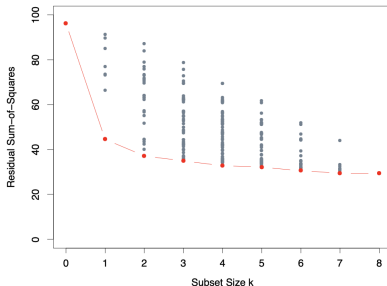
**Example:** Gene expression pattern of a certain kind of disease

- The amount of genes are simply too large but the NO. of patients included in research is limited.
- For practical use, researchers want to find a few genes that are highly correlated with the disease.



# Intuitive Thought

Intuitively, to do subset selection, assuming there are  $p$  variables in total, we can do regression with  $k$  variables for each  $k \in \{0, 1, 2, \dots, p\}$  (when  $n > p$ ),  $(k \in \{0, 1, 2, \dots, p\})$  (when  $n < p$ ) find the best one for each  $k$ . However, the problem arises when we try to define the best 'k'.



Note that we cannot use the RSS to determine the right  $k$  because it is always decreasing when  $k$  increase.

# Model Selection Methods

We change our choice of  $\lambda$  to create different criteria: Some criteria that are commonly used

$$MSE_p + \lambda p \frac{\hat{\sigma}^2}{n},$$

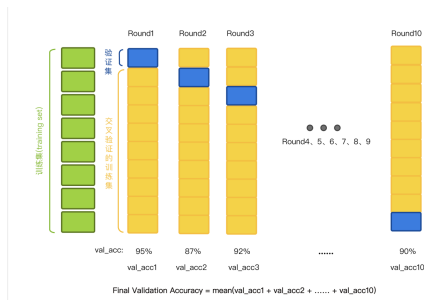
- Akaike Information Criterion(AIC), given as  $AIC = 2p - 2\ln(\hat{L})$ , under the assumption that the noise satisfies Gaussian distribution, is a coefficient 'away' of taking  $\lambda = 2$  in the above equation;
- Bayesian Information Criterion(BIC),also given as  $BIC = p \ln(n) - 2\ln(\hat{L})$ , under the assumption that the noise satisfies Gaussian distribution, is a coefficient 'away' of taking  $\lambda = \ln(n)$  in the above equation;

# Model Selection Methods

## Cross Validation

Cross-validation is a model validation techniques for assessing how the results of a statistical analysis will generalize to an independent data set.

The most commonly used CV is called k-fold cross-validation,



---

**Algorithm 6.1** *Best subset selection*

---

1. Let  $\mathcal{M}_0$  denote the *null model*, which contains no predictors. This model simply predicts the sample mean for each observation.
  2. For  $k = 1, 2, \dots, p$ :
    - (a) Fit all  $\binom{p}{k}$  models that contain exactly  $k$  predictors.
    - (b) Pick the best among these  $\binom{p}{k}$  models, and call it  $\mathcal{M}_k$ . Here *best* is defined as having the smallest RSS, or equivalently largest  $R^2$ .
  3. Select a single best model from among  $\mathcal{M}_0, \dots, \mathcal{M}_p$  using cross-validated prediction error,  $C_p$  (AIC), BIC, or adjusted  $R^2$ .
-

# Subset Selection Methods

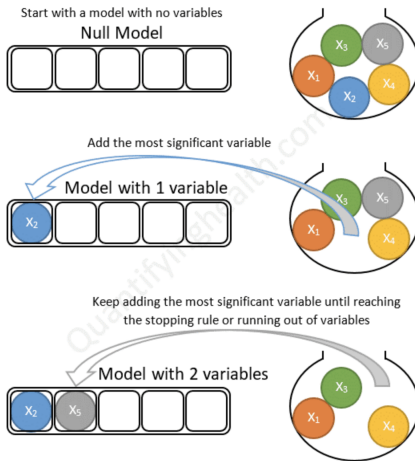
Best subset regression finds for each  $k \in \{0, 1, 2, \dots, p\}$  the subset of size  $k$  that gives smallest residual sum of squares. And then use methods such as AIC or BIC to determine the right  $k$ .

## Remark

Best subset selection suffers from computational limitations. In general, there are  $2^p$  models that involve subsets of  $p$  predictors. If  $p = 20$ , there are approximately million possibilities. As for  $p \geq 40$ , best subset selection becomes literally infeasible.

# Forward Stepwise Selection

Forward stepwise selection example with 5 variables:





---

**Algorithm 6.2** *Forward stepwise selection*

---

1. Let  $\mathcal{M}_0$  denote the *null* model, which contains no predictors.
  2. For  $k = 0, \dots, p - 1$ :
    - (a) Consider all  $p - k$  models that augment the predictors in  $\mathcal{M}_k$  with one additional predictor.
    - (b) Choose the *best* among these  $p - k$  models, and call it  $\mathcal{M}_{k+1}$ . Here *best* is defined as having smallest RSS or highest  $R^2$ .
  3. Select a single best model from among  $\mathcal{M}_0, \dots, \mathcal{M}_p$  using cross-validated prediction error,  $C_p$  (AIC), BIC, or adjusted  $R^2$ .
-

# Forward Stepwise Selection

- To choose a stopping rule, we can stop when all remaining variables to consider have a p-value larger than some threshold if added to the model.

The threshold can be:

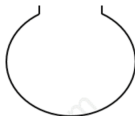
1. A fixed value (for instance: 0.05 or 0.2 or 0.5)
2. Determined by AIC: AIC chooses the threshold according to how many degrees of freedom the variable under consideration has, generally speaking, the more degrees of freedom a variable has, the lower the threshold will be.
3. Determined by BIC: BIC chooses the threshold according to the effective sample size  $n$ .

# Backward Stepwise Selection

## Backward stepwise selection example with 5 variables:

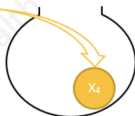
Start with a model that contains all the variables

Full Model



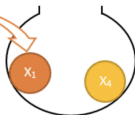
Remove the least significant variable

Model with 4 variables



Keep removing the least significant variable until reaching the stopping rule or running out of variables

Model with 3 variables



# Backward Stepwise Selection

---

## Algorithm 6.3 *Backward stepwise selection*

---

1. Let  $\mathcal{M}_p$  denote the *full* model, which contains all  $p$  predictors.
  2. For  $k = p, p - 1, \dots, 1$ :
    - (a) Consider all  $k$  models that contain all but one of the predictors in  $\mathcal{M}_k$ , for a total of  $k - 1$  predictors.
    - (b) Choose the *best* among these  $k$  models, and call it  $\mathcal{M}_{k-1}$ . Here *best* is defined as having smallest RSS or highest  $R^2$ .
  3. Select a single best model from among  $\mathcal{M}_0, \dots, \mathcal{M}_p$  using cross-validated prediction error,  $C_p$  (AIC), BIC, or adjusted  $R^2$ .
-

# Backward Stepwise Selection

- To choose a stopping rule, we can stop when all remaining variables to consider have a p-value smaller than some threshold if added to the model.

The threshold can be:

1. A fixed value (for instance: 0.05 or 0.2 or 0.5)
2. Determined by AIC
3. Determined by BIC

# Comparing Backward with Forward Stepwise Selection

- Where forward stepwise is better: It can be applied in settings where the number of variables under consideration is larger than the sample size.
- Stepwise selection is more effective , but it is not guaranteed to select the best possible combination of variables. In fact the selection of variables using a stepwise regression will be highly unstable, especially when we have a small sample size.

---

**Algorithm 3.4** *Incremental Forward Stagewise Regression— $FS_\epsilon$ .*

---

1. Start with the residual  $\mathbf{r}$  equal to  $\mathbf{y}$  and  $\beta_1, \beta_2, \dots, \beta_p = 0$ . All the predictors are standardized to have mean zero and unit norm.
  2. Find the predictor  $\mathbf{x}_j$  most correlated with  $\mathbf{r}$
  3. Update  $\beta_j \leftarrow \beta_j + \delta_j$ , where  $\delta_j = \epsilon \cdot \text{sign}[\langle \mathbf{x}_j, \mathbf{r} \rangle]$  and  $\epsilon > 0$  is a small step size, and set  $\mathbf{r} \leftarrow \mathbf{r} - \delta_j \mathbf{x}_j$ .
  4. Repeat steps 2 and 3 many times, until the residuals are uncorrelated with all the predictors.
-

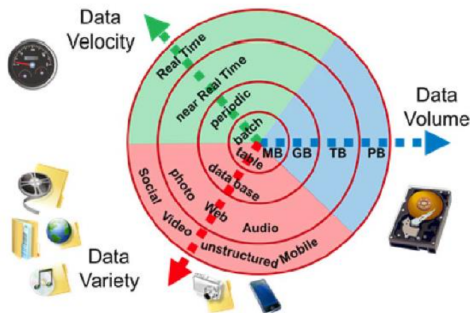
# Comparing Stagewise with Stepwise

- Greediness is counterbalanced by the small step size  $\epsilon > 0$ ;
- The learning process in stagewise is slower: In stepwise we increase the coefficient of  $X_i$  by large amount and we might change several coefficients each time; in stagewise a coefficient is increased by  $\epsilon$  each time.
- The stagewise learning process is computationally cheap.



# Introduction to Sure Independent Screening

Background: **Ultrahigh dimensional data** in which the number of features  $p$  can be much larger than the number of observations  $n$  is common but hard to deal with. Challenges includes scalability, high collinearity, spurious correlation and noise accumulation.



# Introduction to Sure Independent Screening

The idea of feature screening became important for

- Computational efficiency
- Statistical efficiency: by alleviated noise accumulation through dimension reduction.

# Introduction to Sure Independent Screening

Let us consider the linear regression model:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \epsilon$$

where  $\mathbf{y}$  is an  $n$ -dimensional vector,  $\mathbf{X}$  is an  $n \times p$  matrix,  $\boldsymbol{\beta}$  is a  $p$ -dimensional vector and  $\epsilon$  is a  $p$ -dimensional error vector. Firstly, the SIS ranks all the  $p$  features using the marginal utilities based on the marginal correlations  $\hat{corr}(\mathbf{x}_j, \mathbf{y})$  of  $\mathbf{x}_j$ s with the response of  $\mathbf{y}$  and retains the top  $d$  covariates with the largest absolute correlations collected in the set  $\hat{\mathcal{M}}$ , that is

$$\hat{\mathcal{M}} = \{1 \leq j \leq p : |\hat{corr}(\mathbf{x}_j, \mathbf{y})| \text{ is among the top } d \text{ largest ones}\}$$

# Introduction to Sure Independent Screening

The next step of SIS: Allow user's to use any favourite regularization method such as Lasso, SCAD AND Dantzig selector on the features in  $\hat{\mathcal{M}}$ .

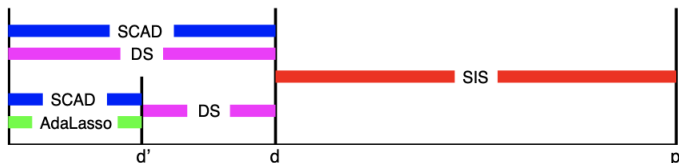


Figure 2: Methods of model selection with ultra high dimensionality.

The property this equation enjoys:  $\mathbb{P}\{\mathcal{M} \in \hat{\mathcal{M}}\} \rightarrow 1$ . It is called sure screening property and it is crucial to the  $2^{nd}$  step of refined variable selection.

# Table of Contents

- 1 Introduction
- 2 Model specification
- 3 Some properties
- 4 Robust linear regression
- 5 Subset Selection Methods
- 6 Ridge Regression and PCA**
- 7 Appendix

## Multicollinearity

Collinearity (or multicollinearity or ill-conditioning) occurs when independent variables in a regression are so highly correlated that it becomes difficult or impossible to distinguish their individual effects on the dependent variable.

The classical linear problem is unsolvable(i.e not having a unique solution) if there is an exact linear relationship between variables. Even if variables are highly correlated, the following problems arises:

- Confounding variables
- Computational problem
- Large variance

# Confounding Variables

For simplicity, let's assume complete linear relationship between variables. In the following example, we assume  $X_1 = X_2$ , and try to fit the following model (Note that when we do regression with  $X_1$  only we get  $Y = 20 - 3X_1$ ) :

$$Y = \alpha + \beta_1 X_1 + \beta_2 X_2,$$

all the above solutions are equivalent

$$Y = 20 - 5X_1 + 2X_2$$

$$Y = 20 + 6X_1 - 3X_2$$

$$Y = 20 + 3X_2$$

But the meaning of each equation can be significantly different since the coefficients of each one implies a different relation

# Computational Problem

Definition of condition number: Let  $e$  be the error in  $b$ . Assuming that  $A$  is a nonsingular matrix, the error in the solution of  $AX = b$  is  $A^{-1}e$ . The maximum value of the ratio of the relative error

$$\begin{aligned} \max_{e, b \neq 0} \left\{ \frac{\|A^{-1}e\|}{\|e\|} \frac{\|b\|}{\|A^{-1}b\|} \right\} &= \max_{e \neq 0} \left\{ \frac{\|A^{-1}e\|}{\|e\|} \right\} \max_{b \neq 0} \left\{ \frac{\|b\|}{\|A^{-1}b\|} \right\} \\ &= \|A^{-1}\| \|A\| \end{aligned}$$

Take  $\|\cdot\|$  to be  $\|\cdot\|_2$ , then:

$$\kappa(A) = \frac{\sigma_{\max}(A)}{\sigma_{\min}(A)}$$

## Remark

The matrix  $X^T X$  is 'close' to singular. This means its smallest eigenvalue is close to zero, thus its condition number is large and the solution is sensitive to disturbance.



# Large Variance

We have learned previously that  $\text{Var}(\hat{\beta}) = (X^T X)^{-1} \sigma^2$ , with

$$\text{Var}(\hat{\beta}_p) = \frac{\sigma^2}{\langle z_p, z_p \rangle} = \frac{\sigma^2}{\|z_p\|^2}$$

being the variance of a individual coefficient( $z_p$  being the variable we get from Orthogonalization).

For collinearity problem, if two columns in  $\mathbf{X}$ , say  $X_i, X_j$  is highly correlated, then the length of  $z_j$  must be rather small as most of the information in  $z_j$  has already been explained by  $z_i$ , making the variance very large.

# Examples of Colinearity

In reality, we give an example of generating highly correlated variables:

- Suppose  $(x, y)$  is generated as follows:

$$z_1, z_2 \sim \mathcal{N}(0, 1) \text{ (independent)}$$

$$\varepsilon_0, \varepsilon_1, \dots, \varepsilon_6 \sim \mathcal{N}(0, 1) \text{ (independent)}$$

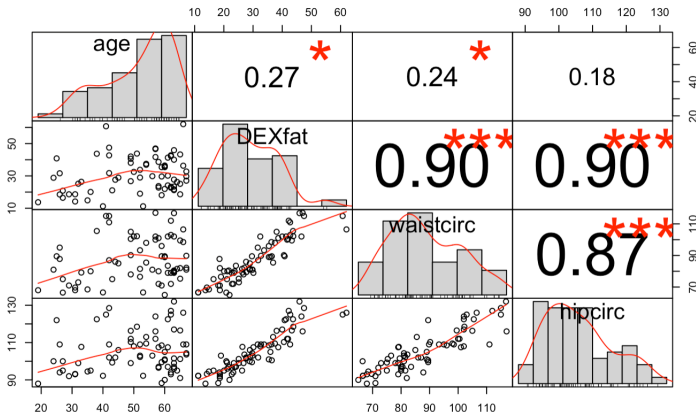
$$y = 3z_1 - 1.5z_2 + 2\varepsilon_0$$

$$x_j = \begin{cases} z_1 + \varepsilon_j/5 & \text{for } j = 1, 2, 3 \\ z_2 + \varepsilon_j/5 & \text{for } j = 4, 5, 6 \end{cases}$$

- Generated a sample of  $(x, y)$  pairs of size 100.
- Correlations within the groups of  $x$ 's were around 0.97.

# Examples of Colinearity

The following is a dataset called 'bodyfat'.



DEXfat:body fat measured by DXA, response variable;  
waistcirc:waist circumference; hipcirc: hip circumference.

# Shrinkage Method

The idea here is to perform a linear regression, while shrinking the regression coefficients  $\hat{\beta}$  toward 0.

- Wildly large positive/negative coefficients on one variables became less likely to occur.
- Discrete process(i.e variables are either retained or discarded) often exhibits high variance.
- Shrinkage methods are more continuous and don't suffer from high variability.
- It can reduce high variance of regression coefficient caused by collinearity between explanatory variables.

# Ridge Regression Model

Ridge regression shrinks the regression coefficients by imposing a penalty on their size. The ridge coefficients minimize a penalized residual sum of squares,

$$\hat{\beta}^{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}$$

An equivalent way to write the ridge problem is

$$\hat{\beta}^{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2, \text{ subject to } \sum_{j=1}^p \beta_j^2 \leq t$$

# Ridge Regression Solution

Writing in matrix form,

$$\text{RSS}(\lambda) = (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta$$

The regression solution is

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

Proof: let

$$X_* = \begin{pmatrix} X \\ \sqrt{\lambda} I \end{pmatrix}$$

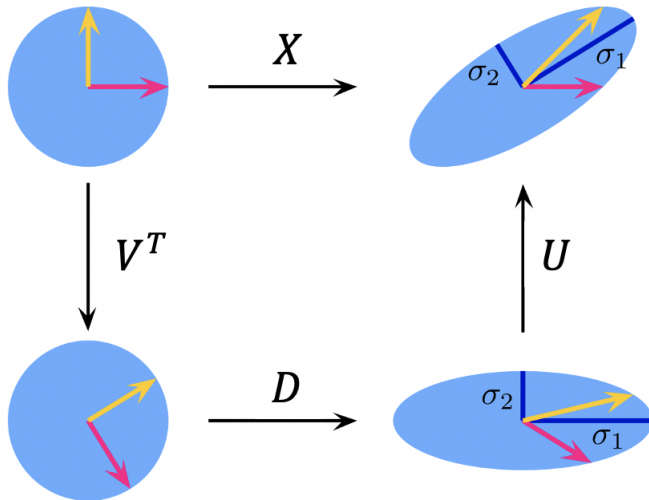
Solving ridge estimator for  $X$  is equivalent to solving OLS for  $X_*$

$$\begin{aligned} (y_* - X_* \beta)' (y_* - X_* \beta) &= (y - X\beta)' (y - X\beta) + \lambda \beta' \beta \\ (X_*^T X_*) \beta &= X_*' y_* \end{aligned}$$

Plug in  $X_*$ , and we get

$$(X^T X + \lambda I) \beta = X^T y$$

# SVD Decmoposition



$$X = UDV^T$$

# Bias of Ridge Estimator

Under the commonly assumption about the noise:

$$\mathbb{E}[\varepsilon \mid \mathbf{X}] = 0$$

$$\text{Var}[\varepsilon \mid \mathbf{X}] = \sigma^2 I$$

The bias of  $\mathbf{X}\hat{\beta}_\lambda$  satisfies:

$$\text{Bias}^2(\mathbf{X}\hat{\beta}_\lambda) = \|\mathbf{X}(\mathbb{E}_\varepsilon \beta_\lambda - \beta)\|^2 \leq \lambda \|\beta\|_2^2$$

First consider  $\mathbb{E}_\varepsilon \beta$  is the minimizer of  $\|\mathbf{X}\beta - \mathbf{X}\gamma\|^2 + \lambda \|\gamma\|_2^2$

$$\begin{aligned} \|\mathbf{X}(\mathbb{E}_\varepsilon \beta_\lambda - \beta)\|^2 &\leq \|\mathbf{X}(\mathbb{E}_\varepsilon \beta_\lambda - \beta)\|^2 + \lambda \|\hat{\beta}_\lambda\|_2^2 \\ &\leq \|\mathbf{X}\beta - \mathbf{X}\gamma\|^2 + \lambda \|\gamma\|_2^2. \end{aligned}$$

This suggests that the solution of ridge regression is actually a biased estimator of  $\beta$  and bias increase with  $\lambda$ .



# Variance of Ridge Estimator

$$\text{Var} \left[ X \hat{\beta}_\lambda \right] = \sigma^2 \sum_{i=1}^n \mu_i^2, \text{ where } \mu_i = \frac{\sigma_i^2}{\sigma_i^2 + \lambda}$$

Proof:

$$\begin{aligned} & \frac{1}{n} \mathbb{E}_\varepsilon \left\| X \left( \mathbb{E} \hat{\beta}_\lambda - \hat{\beta}_\lambda \right) \right\|_2^2 \\ &= \frac{1}{n} \mathbb{E}_\varepsilon \left\| X \left( X^T X + \lambda I \right)^{-1} X^T X \beta - X \left( X^T X + \lambda I \right)^{-1} X^T (X \beta + \varepsilon) \right\|^2 \\ &= \frac{1}{n} \mathbb{E}_\varepsilon \| W \varepsilon \|^2 = \frac{\sigma^2}{n} \sum_{i=1}^n \mu_i^2, \end{aligned}$$

where  $\mu_i$  is the eigenvalues of  $W$ , and  $W = X(X^T X + \lambda I)^{-1} X^T$ . Let  $X = UDV^T$ , and

$$W = UDV^T \left( VD^2 V^T + \lambda I \right)^{-1} = UD(D^2 + \lambda I)^{-1} DU^T,$$

thus  $\mu_i = \frac{\sigma_i^2}{\sigma_i^2 + \lambda}$ , which suggests  $\text{Var} \left[ X \hat{\beta}_\lambda \right] < \text{Var} \left[ X \hat{\beta} \right]$ .

# SVD Decomposition

The singular value decomposition (SVD) of the centered input matrix  $X$  is

$$X = U_{n \times p} D_{p \times p} V_{p \times p}^T.$$

Here  $U$  and  $V$  are orthogonal matrices, and  $D$  is a diagonal matrix, with diagonal entries  $d_1 \geq d_2 \geq \cdots d_p \geq 0$ . Using the SVD of  $X$  we have:

$$\begin{aligned} X\hat{\beta}^{1s} &= X \left( X^T X \right)^{-1} X^T y \\ &= U U^T y = \sum_{j=1}^n u u_j^T y \end{aligned}$$

$$\begin{aligned} X\hat{\beta}^{\text{ridge}} &= X \left( X^T X + \lambda I \right)^{-1} X^T y \\ &= U D (D^2 + \lambda I)^{-1} D U^T y \\ &= \sum_{j=1}^p u_j \frac{d_j^2}{d_j^2 + \lambda} u_j^T y \end{aligned}$$

# SVD Decomposition

Note that for  $\lambda \geq 0$ , we have  $\frac{d_j^2}{d_j^2 + \lambda} \leq 1$

- Like linear regression, ridge regression computes the coordinates of  $y$  w.r.t the orthonormal basis  $U$ .
- It shrinks these coordinates by the factor  $\frac{d_j^2}{d_j^2 + \lambda}$ .
- A greater amount of shrinkage is applied to the coordinates of basis vectors with smaller  $d_j^2$ .
- We define the **effective degrees of freedom** as follows:

$$\text{df}(\lambda) = \sum_{j=1}^D \frac{d_j^2}{d_j^2 + \lambda}$$

Note that  $\text{df}(\lambda) = p$  when  $\lambda = 0$  (no regularization);  
and  $\text{df}(\lambda) \rightarrow 0$  as  $\lambda \rightarrow \infty$ .

# PCA as a Variance-Maximization Technique

**Principal component analysis (PCA)** was introduced as a technique for deriving a reduced set of orthogonal linear projections of a single collection of correlated variables,  $X = (X_1, \dots, X_p)$ , where the projections are ordered by decreasing variances.

**Variance** is an important measurement of the amount of information in the variable.

PCA can be viewed as an optimization problem as following:

$$\begin{aligned} & \max \text{Var}(Xb_j) \\ \text{s.t. } & b_j^T b_j = 1, \text{Var}(Xb_i, Xb_j) = b_i^T \Sigma_{XX} b_j = 0, i < j, \\ & \text{notice that } (Xb_i, Xb_j) = b_i^T \Sigma_{XX} b_j \end{aligned}$$

# PCA as a Variance-Maximization Technique

The objective function is:

$$f(b_1) = b_1^T \Sigma_{XX} b_1 - \lambda_1 (1 - b_1^T b_1),$$

$\lambda_1$  is a Lagrange multiplier. Differentiating  $f(b_1)$  with respect to  $b_1$  and setting the result equal to zero for a maximum to yield.

$$\frac{\partial f(b_1)}{\partial b_1} = 2(\Sigma_{XX} - \lambda_1 I_r) b_1 = 0$$

If  $b_1 \neq 0$ , then  $\lambda_1$  must be chosen to satisfy the determinantal equation

$$|\Sigma_{XX} - \lambda_1 I| = 0$$

Thus,  $\lambda_1$  has to be the largest eigenvalue of  $\Sigma_{XX}$ , and  $b_1$  be the eigenvector associated with  $\lambda_1$ .

# PCA as a Variance-Maximization Technique

The second principal component,  $\xi_2$ , is then obtained by choosing a second set of coefficients,  $b_2$ , for the next linear projection,  $\xi_2$ , so that the variance of  $\xi_2$  is largest among all linear projections of  $X$  that are also uncorrelated with  $\xi_1$  above.

The variance of  $\xi_2$  is  $\text{Var}(\xi_2) = b_2^T \Sigma_{XX} b_2$ , and this has to be maximized subject to the normalization constraint  $b_2^T b_2 = 1$  and orthogonality constraint  $b_1^T b_2 = 0$ .

$$f(b_2) = b_2^T \Sigma_{XX} b_2 - \lambda_2 (1 - b_2^T b_2) + \mu b_1^T b_2$$

Differentiating  $f(b_2)$  with respect to  $b_2$  and setting the result equal to zero for a maximum yields

$$\frac{\partial f(b_2)}{\partial b_2} = 2(\Sigma_{XX} - \lambda_2 I_r) b_2 + \mu b_1 = 0$$

# PCA as a Variance-Maximization Technique

Premultiplying this derivative by  $b_1^T$  and using the orthogonality and normalization constraints, we have that

$$2b_1^T \Sigma_{XX} b_2 + \mu = 0.$$

Premultiplying the equation  $(\Sigma_{XX} - \lambda_1 I_r) b_1 = 0$  by  $b_2^T$  yields

$$b_2^T \Sigma_{XX} b_1 = 0$$

whence  $\mu = 0$ .

This means that  $\lambda_2$  is the second largest eigenvalue of  $\Sigma_{XX}$ , and the coefficient vector  $b_2$  for the second principal component is the eigenvector associated with  $\lambda_2$ .

# PCA as a Variance-Maximization Technique

In this sequential manner, we obtain the remaining sets of coefficients for the principal components  $\xi_3, \xi_4, \dots, \xi_r$ , where the  $i$ th principal component  $\xi_i$  is obtained by choosing the set of coefficients,  $b_i$ , for the linear projection  $\xi_i$  so that  $\xi_i$  has the largest variance among all linear projections of  $X$  that are also uncorrelated with  $\xi_1, \xi_2, \dots, \xi_{i-1}$ .

The coefficients of these linear projections are given by the ordered sequence of eigenvectors  $\{b_j\}$  associated with the  $j$ th largest eigenvalue  $\lambda_j$ , of  $\Sigma_{XX}$ .



# Principle Components of $X$

We use the SVD of the centered matrix  $X$  to express the principal components of the variables in  $X$ . The sample covariance matrix is given by

$$S = X^T X / (N - 1).$$

And the SVD of  $S$  is

$$S = X^T X / (N - 1) = V D^2 V^T / (N - 1).$$

The first eigenvectors  $v_j$  has the property that:  $z_1 = X v_1$  has the largest variance among all  $z_j = X v_j$ :

$$\text{Var}(z_1) = \text{Var}(X v_1) = \frac{d_1^2}{N - 1},$$

In fact,  $z_1 = X v_1 = u_1 d_1$ , hence  $u_1$  is the normalized first PC.

# PCA with Least-Square Optimality

Let  $B$  be a  $(t \times p)$  matrix  $B = (b_1, \dots, b_t)^T$ . Let  $\xi = BX$ , we want to find a  $p$ -vector  $\mu$  ( $p \times 1$ ) matrix  $A$ .

$$E \left\{ (X - \mu - A\xi)^T (X - \mu - A\xi) \right\},$$

we use least-squares error criterion as our measure of how well we can reconstruct  $X$  by linear projection on  $\xi$ .

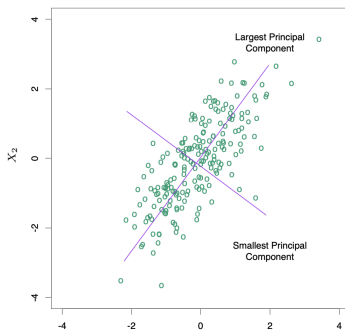
More specifically, our goal is to choose  $A$ ,  $B$  and  $\mu$  to minimize

$$E \left\{ (X - \mu - ABX)^T (X - \mu - ABX) \right\}.$$

# PCR

The PCR method may be broadly divided into three major steps:

1. Obtain the principal components following the above step;
2. Regress the observed vector of outcomes on the selected first  $k$  principal components as covariates, get OLS. (notice we can take different value of  $k$ )
3. Decide for the best  $k$  (Here we present one rule: We find the least  $k$  that satisfies:  $\sum d_k \geq 0.8$ )



# Compare PCR with Ridge Regression

Let  $W_k = XV_k = [Xv_1, \dots, Xv_k]$ ,  $\hat{\gamma}_k = (W_k^T W_k)^{-1} W_k^T Y \in \mathbb{R}^k$  we here give the solution of PCR with  $k$  principle components:

$$X\hat{\beta}_k = V_k \hat{\gamma}_k = \sum_{j=1}^k u_j u_j^T y$$

It is similar to ridge regression in the way that it also **reduced variance** of the coefficients:

$$\begin{aligned}\text{Var}(\hat{\beta}_k) &= \sigma^2 V_k (W_k^T W_k)^{-1} V_k^T \\ &= \sigma^2 V_k \text{diag}(\lambda_1^{-1}, \dots, \lambda_k^{-1}) V_k^T \\ &= \sigma^2 \sum_{j=1}^k \frac{v_j v_j^T}{\lambda_j}\end{aligned}$$

$$\text{Var}(\hat{\beta}_{\text{ols}}) - \text{Var}(\hat{\beta}_k) = \sigma^2 \sum_{j=k+1}^p \frac{v_j v_j^T}{\lambda_j}$$

# Compare PCR with Ridge Regression

Also, **multicollinearity** can be addressed by both PCR and ridge regression. We have known from the above analysis of the variance of small estimator that it is the small eigenvalue that has the maximum inflation effect on the variance of the least squares estimator.

In fact, **PCR** eliminates the effect on those eigenvectors by simply exclude them from the regression model  $\sum_{j=1}^k u_j u_j^T y$ .

On the other hand, **ridge regression** eliminates the effect by multiplying a coefficient  $\frac{d_j^2}{d_j^2 + \lambda}$  for each term  $\sum_{j=1}^p u_j \frac{d_j^2}{d_j^2 + \lambda} u_j^T y$ . Ridge regression projects  $y$  onto every principle components, but shrinks the coefficients of the low variance components more than high ones since lower ones contains less information.

# Table of Contents

- 1 Introduction
- 2 Model specification
- 3 Some properties
- 4 Robust linear regression
- 5 Subset Selection Methods
- 6 Ridge Regression and PCA
- 7 Appendix**

# Appendix

(Cochran's Theorem) Let  $X_1, X_2, \dots, X_n$  be independent  $N(\mu_i, \sigma^2)$ , i.e.  $X \sim N(\mu, \sigma^2 I)$ , and

$$\sum_{i=1}^n X_i^2 = Q_1 + Q_2 + \dots + Q_k,$$

where  $Q_i = X' A_i X$ ,  $A_1, A_2, \dots, A_k$  are symmetric and idempotent  $n \times n$  matrices with  $\text{rank}(A_i) = r_i$ ,  $i = 1, 2, \dots, k$ .

Then:

- (1)  $Q_1, Q_2, \dots, Q_k$  are independent.
- (2)  $\frac{Q_i}{\sigma^2} \sim \chi^2(r_i, \delta_i)$  with  $\delta_i = \mu' A_i \mu / \sigma^2$ .

Corollary: Let  $X \sim N(\mu, \sigma^2 I)$ ,  $A$  is symmetric with  $\text{rank}(A) = r$ ,  $\delta = \mu' A \mu / \sigma^2$ . Then

$$X' A X / \sigma^2 \sim \chi^2(r, \delta) \Leftrightarrow A \text{ is idempotent}.$$

## (Cook's Distance)

$$\hat{\beta}_{(i)} = \left( \mathbf{X}_{(i)}^T \mathbf{X}_{(i)} \right)^{-1} \mathbf{X}_{(i)}^T \mathbf{y}_{(i)},$$

$$\mathbf{X}_{(i)}^T \mathbf{X}_{(i)} = \mathbf{X}^T \mathbf{X} - \mathbf{x}_i \mathbf{x}_i^T,$$

$$\left( \mathbf{X}_{(i)}^T \mathbf{X}_{(i)} \right)^{-1} = \left( \mathbf{X}^T \mathbf{X} \right)^{-1} + \frac{\left( \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{x}_i \mathbf{x}_i^T \left( \mathbf{X}^T \mathbf{X} \right)^{-1}}{1 - \mathbf{x}_i^T \left( \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{x}_i}.$$

$$\hat{\mathbf{y}}_{(i)} = \mathbf{X} \hat{\beta}_{(i)},$$

$$\hat{\beta}_{(i)} = \left( \mathbf{X}_{(i)}^T \mathbf{X}_{(i)} \right)^{-1} \mathbf{X}_{(i)}^T \mathbf{y}_{(i)},$$

$$h_i = H_{ii} = \mathbf{x}_i^T \left( \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{x}_i,$$

$$\hat{\beta}_{(i)} = \left( \left( \mathbf{X}^T \mathbf{X} \right)^{-1} + \frac{\left( \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{x}_i \mathbf{x}_i^T \left( \mathbf{X}^T \mathbf{X} \right)^{-1}}{1 - h_i} \right) \left( \mathbf{X}^T \mathbf{y} - y_i \mathbf{x}_i \right).$$



## (Cook's Distance)

$$\hat{\beta} = (X^T X)^{-1} X^T y,$$

$$\hat{\beta}_{(i)} = \hat{\beta} + \frac{(X^T X)^{-1} x_i x_i^T \hat{\beta}}{1 - h_i} - y_i (X^T X)^{-1} x_i - \frac{h_i y_i}{1 - h_i} (X^T X)^{-1} x_i.$$

$$y_i = x_i^T \hat{\beta}, \hat{\beta}_{(i)} = \hat{\beta} - \frac{y_i - \hat{y}_i}{1 - h_i} (X^T X)^{-1} x_i = \hat{\beta} - \frac{\hat{\epsilon}_i}{1 - h_i} (X^T X)^{-1} x_i,$$

$$y_i - \hat{y}_{(i)} = y_i - x_i^T \hat{\beta}_{(i)} \hat{\beta}_{(i)},$$

$$y_i - \hat{y}_{(i)} = y_i - x_i^T \left( \hat{\beta} - \frac{\hat{\epsilon}_i}{1 - h_i} (X^T X)^{-1} x_i \right),$$

$$y_i - \hat{y}_{(i)} = \hat{\epsilon}_i + \frac{\hat{\epsilon}_i h_i}{1 - h_i} = \frac{\hat{\epsilon}_i}{1 - h_i}.$$