

STAT 540: Data Analysis and Regression

Wen Zhou

<http://www.stat.colostate.edu/~riczw/>

Email: riczw@stat.colostate.edu

Department of Statistics
Colorado State University

Fall 2015

Contents

- 1 Multiple Linear Regression Model
- 2 Inference on Multiple Regression
- 3 Inference about Regression Parameters
- 4 Estimation and Prediction
- 5 Geometric View of Regression and Linear Models
- 6 Estimating estimable function of coefficient

Multiple Linear Regression I

- Multiple linear regression model
 - ① Multiple linear regression model in matrix terms
 - ② Estimation of regression coefficients
- Inference
 - ① ANOVA results
 - ② Inference about regression parameters
 - ③ Estimation of mean response and prediction of new observation
- Inference about regression parameters
- Estimation and prediction
- Geometric interpretation of linear model and regression
- Estimating estimable function of regression or linear coefficient β

- 1 Multiple Linear Regression Model
- 2 Inference on Multiple Regression
- 3 Inference about Regression Parameters
- 4 Estimation and Prediction
- 5 Geometric View of Regression and Linear Models
- 6 Estimating estimable function of coefficient

Multiple Linear Regression

- Example: # of predictor variables = 2.

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \epsilon_i, \quad \epsilon_i \sim \text{iid } N(0, \sigma^2),$$

for $i = 1, \dots, n$.

- Response surface:

$$\mathbb{E}(Y_i) =$$

- Example:
 - ▶ Y = Pine bark beetle density
 - ▶ X_1 = Temperature
 - ▶ X_2 = Tree species

Interpretation of Coefficients

- β_0 : Intercept. When the model scope includes $X_1 = X_2 = 0$.
 - ▶ β_0 is interpreted as the mean response $E(Y)$ at $X_1 = X_2 = 0$.
- β_j : Slope in the direction of X_j (effect).
 - ▶ $\partial \mathbb{E}(Y) / \partial X_j =$
 - ▶ $\mathbb{E}_{Y|X=(X_1, X_2)}(Y) - \mathbb{E}_{Y|X=(X'_1, X_2)}(Y) =$

- Interpreted as the change in the mean response $\mathbb{E}(Y)$ per unit increase in X_j , **when X_{-j} are held constant.**
- What if X_j is qualitative?

Multiple Linear Regression

- A “general” linear regression model is, for $i = 1, \dots, n$

$$Y_i = \beta_0 + \sum_{j=1}^p X_{ij}\beta_j + \epsilon_i, \quad \epsilon_i \sim \text{iid } N(0, \sigma^2).$$

- Response surface:

$$\mathbb{E}(Y_i) = \beta_0 + \sum_{j=1}^p X_{ij}\beta_j$$

- Regression coefficients: $\beta_0, \beta_1, \dots, \beta_{p-1}, \beta_p$.
- Predictor variables: X_1, \dots, X_p are known constants/values.
- The model is linear in the parameters, **not necessarily in the shape of the response surface.**

Response Surface Examples

- Polynomial regression

$$\mathbb{E}(Y) = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3.$$

- Transformed variables

$$\mathbb{E}(\log(Y)) = \beta_0 + \beta_1 X_1 + \beta_2 \sqrt{X_2}.$$

- Interaction effects

$$\mathbb{E}(Y) = \beta_0 + \beta_1 X_1 + \beta_2 \sqrt{X_2} + \beta_3 X_1 X_2.$$

- ▶ The change in the mean response corresponding to a unit change in X_1 depends on X_2 and vice versa.
- ▶ Testing whether $\beta_3 = 0$ or not is very challenging in high-dimensional ($n = o(p)$).

Qualitative Predictor Variables

- Example: Let Y = length of hospital stay, X_1 = age, and X_2 = gender: 0 for male and 1 for female.

- ▶ An additive model is
- ▶ Thus the response surface for males is

and for females is

- β_2 is _____
- This kind of model sometimes is called ANVOCA model.

Qualitative Predictor Variables

- Interaction: the relationship between X_1 and Y for a fixed value of $X_2 = x_2$ depends on x_2 .
- An interaction model is

- Thus the response surface for males is

and for females is

Notation

- n observations, 1 response variable, $p - 1$ β 's with predictors (i.e. β_0 is the p th).
- Response variable: $\mathbf{Y}_{n \times 1} = (Y_1, Y_2, \dots, Y_n)^T$.
- The predictors are arranged in the **design matrix**

$$\mathbf{X}_{n \times p} = \begin{bmatrix} 1 & X_{11} & X_{12} & \cdots & X_{1,p-1} \\ 1 & X_{21} & X_{22} & \cdots & X_{2,p-1} \\ \vdots & \vdots & \vdots & & \\ 1 & X_{n1} & X_{n2} & \cdots & X_{n,p-1} \end{bmatrix}$$

- Random error: $\boldsymbol{\epsilon}_{n \times 1} = (\epsilon_1, \epsilon_2, \dots, \epsilon_n)^T$.
- Regression coefficients: $\boldsymbol{\beta}_{p \times 1} = (\beta_0, \beta_1, \dots, \beta_{p-1})^T$.

Multiple Linear Regression Model in Matrix Terms

- The multiple linear regression model can be written as

where as we have seen before

$$\mathbb{E}(\boldsymbol{\epsilon}) = \mathbf{0}_{n \times 1}, \quad \text{Var}\{\boldsymbol{\epsilon}\} = \sigma^2 \mathbf{I}_{n \times n}.$$

- Thus,

and

$$\mathbf{Y} \sim$$

Least Squares Estimation

- Consider the criterion:

$$Q = \sum_{i=1}^n (Y_i - \beta_0 - \sum_{j=1}^{p-1} \beta_j X_{ij})^2 =$$

- The least squares estimate of β is

assuming that $\mathbf{X}^T \mathbf{X}$ is invertible.

- ▶ This is also the MLE.
- ▶ What condition on \mathbf{X} do we need to have $\mathbf{X}^T \mathbf{X}$ invertible?

- ▶ What if $\mathbf{X}^T \mathbf{X}$ is not invertible?

Fitted Values and Residuals

- Fitted values: $\hat{Y} =$

where the hat matrix is

- Residuals: $e =$

- 1 Multiple Linear Regression Model
- 2 Inference on Multiple Regression**
- 3 Inference about Regression Parameters
- 4 Estimation and Prediction
- 5 Geometric View of Regression and Linear Models
- 6 Estimating estimable function of coefficient

Sums of Squares

- We have sums of squares in matrix forms that

$$SSR = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 =$$

$$SSE = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 =$$

$$SSTO = \sum_{i=1}^n (Y_i - \bar{Y})^2 =$$

- Partitioning of total sum of squares and particularly the **df** are

$$\underbrace{SSTO}_{df=n-1} = \underbrace{SSR}_{df=p-1} + \underbrace{SSE}_{df=n-p} .$$

Mean Squares

- Define mean squares

$$MSR = \frac{SSR}{p-1}, \quad MSE = \frac{SSE}{n-p}.$$

- It can be shown that $\mathbb{E}(MSE) =$
- Also can be shown that

$$\mathbb{E}(MSR) \begin{cases} = \sigma^2 & \text{if } \beta_j = 0 \text{ for } \forall j \\ > \sigma^2 & \text{otherwise} \end{cases}.$$

ANOVA Table

The ANOVA table is

Source	SS	df	MS	F
Regression	SSR		MSR	$F = MSR/MSE$
Error	SSE		MSE	
Total	$SSTO$			

- If _____ then

$$\mathbb{E}(MSE) = \mathbb{E}(MSR) = \sigma^2$$

in which case $MSR/MSE \approx 1$.

Overall F Test for Regression Relation

- Test

H_0 : _____ v.s. H_a : _____.

- ▶ It can be shown that under H_0 ,

$$F^* = \frac{MSR}{MSE} \sim$$

- ▶ Thus we can perform an F -test at level α by the decision rule:

- Conditional on H_0 being rejected, we may want to find

$$S = \{j \mid \beta_j \neq 0\}$$

(or a.s.)– Identification/Selection.

Coefficient of Multiple Determination, R^2

- The coefficient of multiple determination is denoted by R^2 and is defined as

$$R^2 = \frac{SSR}{SSTO} = 1 - \frac{SSE}{SSTO}$$

- Interpretation: The proportion of variation in the Y_i 's explained by the regression relation.

More on R^2

- As more predictors are added to the model ($p \uparrow$), R^2 **must increase**. Why?

- ▶ Recall

$$SSTO = SSR + SSE$$

$SSTO$ is fixed for \mathbf{Y} while SSE is a minimum of the unconstrained convex optimization problem $\beta = \arg \min SSE(\beta_0, \dots, \beta_{p-1})$.

- ▶ Suppose we consider an extra predictor and thus consider $SSE(\beta_0, \dots, \beta_p)$. The β that minimizes this SSE cannot be inferior to the previous minimizer because $\beta_p = 0$ is a special case within the new minimization problem that incorporates the previous one.

Adjusted R^2

- R^2 depends on p (even for $p \ll n$), how to remove that dependence?
- The **adjusted coefficient of multiple determination** is denoted by R_a^2 and is defined as

$$R_a^2 = 1 - \frac{SSE/n - p}{SSTO/n - 1} = 1 - \left(\frac{n - 1}{n - p} \right) \frac{SSE}{SSTO}.$$

- The adjusted coefficient of multiple determination R_a^2 may decrease when more predictors are in the model.
- Many other statistics such as AIC, BIC, Mallows's C_p , etc. will be discussed and they are superior over R_a^2 .

- 1 Multiple Linear Regression Model
- 2 Inference on Multiple Regression
- 3 Inference about Regression Parameters**
- 4 Estimation and Prediction
- 5 Geometric View of Regression and Linear Models
- 6 Estimating estimable function of coefficient

Estimation of Regression Coefficients

- Mean satisfies

$$\mathbb{E}(\hat{\beta}) = \beta.$$

That is, the LS estimate $\hat{\beta}$ is an unbiased estimate of β .

- Variance-covariance matrix:

$$\Sigma_{\beta} := \text{Var}\{\hat{\beta}\} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}.$$

▶ $(\Sigma_{\beta})_{kk} =$

▶ $(\Sigma_{\beta})_{kl} =$

Inference about Regression Coefficients

- The estimated variance-covariance matrix.

$$\hat{\Sigma}_{\beta} := s^2\{\hat{\beta}\} = MSE \cdot (\mathbf{X}^T \mathbf{X})^{-1}$$
$$:= \begin{bmatrix} s^2\{\hat{\beta}_0\} & s\{\hat{\beta}_0, \hat{\beta}_1\} & \cdots & s\{\hat{\beta}_0, \hat{\beta}_{p-1}\} \\ s\{\hat{\beta}_1, \hat{\beta}_0\} & s^2\{\hat{\beta}_1\} & \cdots & s\{\hat{\beta}_1, \hat{\beta}_{p-1}\} \\ \vdots & \vdots & \vdots & \vdots \\ s\{\hat{\beta}_{p-1}, \hat{\beta}_0\} & s\{\hat{\beta}_{p-1}, \hat{\beta}_1\} & \cdots & s^2\{\hat{\beta}_{p-1}\} \end{bmatrix}$$

- Under the multiple linear regression model, we have

$$\frac{\hat{\beta}_k - \beta_k}{s\{\hat{\beta}_k\}} \sim$$

for $k = 0, 1, \dots, p-1$.

Inference about Regression Coefficients

- Thus the $(1 - \alpha)$ confidence interval for β_k is

$$\hat{\beta}_k \pm t_{1-\alpha/2; n-p} s\{\hat{\beta}_k\}.$$

- Test $H_0 : \beta_k = \beta_{k0}$ versus $H_a : \beta_k \neq \beta_{k0}$.
- Under H_0 , we have

$$t^* = \frac{\hat{\beta}_k - \beta_{k0}}{s\{\hat{\beta}_k\}} \sim t_{n-p}$$

- Thus we can perform a t -test at level α by the decision rule:

- 1 Multiple Linear Regression Model
- 2 Inference on Multiple Regression
- 3 Inference about Regression Parameters
- 4 Estimation and Prediction**
- 5 Geometric View of Regression and Linear Models
- 6 Estimating estimable function of coefficient

Estimation of Mean Response–Hidden Extrapolation

- Define $\mathbf{X}_h = (1, X_{h1}, \dots, X_{h,p-1})^T$.

- **Caution about hidden extrapolations.**

- ▶ The region (with respect to \mathbf{X}_0) defined by

$$d(\mathbf{X}_0) = \mathbf{X}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_0 \leq h_{max}$$

where $h_{max} = \max_i h_{ii}$, is an ellipsoid enclosing all data points inside the “regressor variable hull” (RVH).

- ▶ Predictions for any \mathbf{X}_0 outside the RVH (i.e., $d(\mathbf{X}_0) > h_{max}$) is hidden extrapolation, at least to some degree.

Estimation of Mean Response

- The estimated mean response corresponding to $\mathbf{X}_h =$
- Mean $\mathbb{E}(\hat{Y}_h) =$
- Variance $\text{Var}\{\hat{Y}_h\} =$
- Estimated variance is $s^2\{\hat{Y}_h\} =$

Confidence Intervals for Mean Response

- The $(1 - \alpha)$ confidence interval for $\mathbb{E}(Y_h)$ is

$$\hat{Y}_h \pm t_{1-\alpha/2; n-p} s\{\hat{Y}_h\}$$

- The Working-Hotelling $(1 - \alpha)$ confidence band for the regression surface is

$$\hat{Y}_h \pm W s\{\hat{Y}_h\}$$

where $W^2 = pF(1 - \alpha; p, n - p)$.

- The Bonferroni $(1 - \alpha)$ joint confidence intervals for g mean responses are

$$\hat{Y}_h \pm B s\{\hat{Y}_h\}$$

where $B = t_{1-\alpha/(2g); n-p}$.

Prediction of New Observation

- The predicted new observation corresponding to \mathbf{X}_h is $\hat{Y}_h = \mathbf{X}_h^T \hat{\boldsymbol{\beta}}$, and

- ▶ Mean $\mathbb{E}(\hat{Y}_h) = \mathbf{X}_h^T \boldsymbol{\beta} = \mathbb{E}(Y_{h(\text{new})})$.

- ▶ Prediction error variance

$$\sigma_{\text{pred}}^2 = \text{Var}(\hat{Y}_h - Y_{h(\text{new})}) =$$

- ▶ Estimated prediction error variance is

$$s^2\{\text{pred}\} =$$

Prediction Intervals for New Observation

- The $(1 - \alpha)$ prediction interval for $Y_{h(\text{new})}$ is

$$\hat{Y}_h \pm t_{1-\alpha/2; n-p} s\{\text{pred}\}$$

- The Scheffé $(1 - \alpha)$ joint confidence intervals for g new observations are

$$\hat{Y}_h \pm S s\{\text{pred}\}$$

where $S^2 = gF(1 - \alpha; g, n - p)$.

- The Bonferroni $(1 - \alpha)$ joint confidence intervals for g new observations are

$$\hat{Y}_h \pm B s\{\text{pred}\}$$

where $B = t_{1-\alpha/(2g); n-p}$.

- 1 Multiple Linear Regression Model
- 2 Inference on Multiple Regression
- 3 Inference about Regression Parameters
- 4 Estimation and Prediction
- 5 Geometric View of Regression and Linear Models**
- 6 Estimating estimable function of coefficient

Geometric Viewpoint: The Column Space of the Design Matrix

- $\mathbf{X}\boldsymbol{\beta}$ is a linear combination of the columns of \mathbf{X}

$$\mathbf{X}\boldsymbol{\beta} = [\mathbf{x}_1, \dots, \mathbf{x}_p] \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} = \beta_1 \mathbf{x}_1 + \dots + \beta_p \mathbf{x}_p$$

- The set of all possible linear combinations of the columns of \mathbf{X} is called the column space of \mathbf{X} and is denoted by

$$\mathcal{C}(X) = \{\mathbf{X}\mathbf{a} : \mathbf{a} \in \mathbb{R}^p\}$$

- The Gauss-Markov linear model says \mathbf{y} is a random vector whose mean is in the column space of \mathbf{X} and whose variance is $\sigma^2 \mathbf{I}$ for some positive real number σ^2 , i.e.

$$\mathbb{E}(\mathbf{y}) \in \mathcal{C}(\mathbf{X}) \quad \text{and} \quad \text{Var}(\mathbf{y}) = \sigma^2 \mathbf{I}, \quad \sigma^2 \in \mathbb{R}^+$$

An Example Column Space

$$\begin{aligned}\mathbf{X} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} &\Rightarrow \mathcal{C}(\mathbf{X}) = \{\mathbf{X}\mathbf{a} : \mathbf{a} \in \mathbb{R}^p\} \\ &= \left\{ \begin{bmatrix} 1 \\ 1 \end{bmatrix} [a_1] : a_1 \in \mathbb{R} \right\} \\ &= \left\{ a_1 \begin{bmatrix} 1 \\ 1 \end{bmatrix} : a_1 \in \mathbb{R} \right\} \\ &= \left\{ \begin{bmatrix} a_1 \\ a_1 \end{bmatrix} : a_1 \in \mathbb{R} \right\}\end{aligned}$$

Another Example Column Space

$$\begin{aligned}\mathbf{X} = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{bmatrix} &\Rightarrow \mathcal{C}(\mathbf{X}) = \left\{ \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} : \mathbf{a} \in \mathbb{R}^2 \right\} \\ &= \left\{ a_1 \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \end{bmatrix} + a_2 \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \end{bmatrix} : a_1, a_2 \in \mathbb{R} \right\} \\ &= \left\{ \begin{bmatrix} a_1 \\ a_1 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ a_2 \\ a_2 \end{bmatrix} : a_1, a_2 \in \mathbb{R} \right\} \\ &= \left\{ \begin{bmatrix} a_1 \\ a_1 \\ a_2 \\ a_2 \end{bmatrix} : a_1, a_2 \in \mathbb{R} \right\}\end{aligned}$$

Another Example Column Space

$$\mathbf{X}_1 = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{X}_2 = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix}$$

$$\mathbf{x} \in \mathcal{C}(\mathbf{X}_1) \Rightarrow \mathbf{x} = \mathbf{X}_1 \mathbf{a} \text{ for some } \mathbf{a} \in \mathbb{R}^2$$

$$\Rightarrow \mathbf{x} = \mathbf{X}_2 \begin{bmatrix} 0 \\ \mathbf{a} \end{bmatrix} \text{ for some } \mathbf{a} \in \mathbb{R}^2$$

$$\Rightarrow \mathbf{x} = \mathbf{X}_2 \mathbf{b} \text{ for some } \mathbf{b} \in \mathbb{R}^3$$

$$\Rightarrow \mathbf{x} \in \mathcal{C}(\mathbf{X}_2)$$

Thus

$$\mathcal{C}(\mathbf{X}_1) \subset \mathcal{C}(\mathbf{X}_2)$$

Another Example Column Space (continued)

$$\mathbf{x} \in \mathcal{C}(\mathbf{X}_2) \Rightarrow \mathbf{x} = \mathbf{X}_2 \mathbf{a} \text{ for some } \mathbf{a} \in \mathbb{R}^3$$

$$\Rightarrow \mathbf{x} = a_1 \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} + a_2 \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \end{bmatrix} + a_3 \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \end{bmatrix} \text{ for some } \mathbf{a} \in \mathbb{R}^3$$

$$\Rightarrow \mathbf{x} = \begin{bmatrix} a_1 + a_2 \\ a_1 + a_2 \\ a_1 + a_3 \\ a_1 + a_3 \end{bmatrix} \text{ for some } a_1, a_2, a_3 \in \mathbb{R}$$

$$\Rightarrow \mathbf{x} = \mathbf{X}_1 \begin{bmatrix} a_1 + a_2 \\ a_1 + a_3 \end{bmatrix} \text{ for some } a_1, a_2, a_3 \in \mathbb{R}$$

Another Example Column Space (continued)

$$\Rightarrow \mathbf{x} = \mathbf{X}_1 \begin{bmatrix} a_1 + a_2 \\ a_1 + a_3 \end{bmatrix} \text{ for some } a_1, a_2, a_3 \in \mathbb{R}$$

$$\Rightarrow \mathbf{x} = \mathbf{X}_1 \mathbf{b} \text{ for some } \mathbf{b} \in \mathbb{R}^2$$

$$\Rightarrow \mathbf{x} \in \mathcal{C}(\mathbf{X}_1)$$

Thus, $\mathcal{C}(\mathbf{X}_2) \subset \mathcal{C}(\mathbf{X}_1)$, as we have shown $\mathcal{C}(\mathbf{X}_1) \subset \mathcal{C}(\mathbf{X}_2)$. It follows that $\mathcal{C}(\mathbf{X}_1) = \mathcal{C}(\mathbf{X}_2)$.

Estimation of $\mathbb{E}(\mathbf{y})$

- A fundamental goal of linear model analysis is to estimate $\mathbb{E}(\mathbf{y})$
- We could, of course, use \mathbf{y} to estimate $\mathbb{E}(\mathbf{y})$
- \mathbf{y} is obviously an unbiased estimator of $\mathbb{E}(\mathbf{y})$, but it is often not a very sensible estimator.
- For example, suppose

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \mu + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \end{bmatrix}, \quad \text{and we observe } \mathbf{y} = [6.1, 2.3]'$$

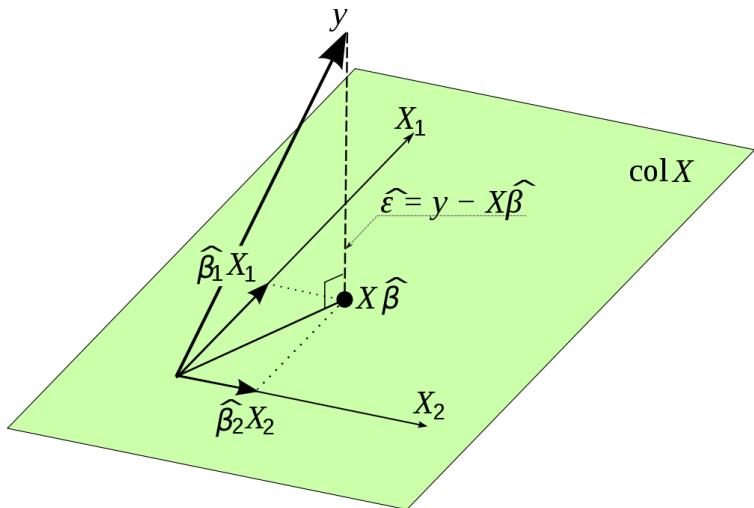
Should we estimate $\mathbb{E}(\mathbf{y}) = [\mu, \mu]'$ by $\mathbf{y} = [6.1, 2.3]'$?

Estimation of $\mathbb{E}(\mathbf{y})$

- The Gauss-Markov linear models says that $\mathbb{E}(\mathbf{y}) \in \mathcal{C}(\mathbf{X})$, so we should use that information when estimating $\mathbb{E}(\mathbf{y})$
- Consider estimating $\mathbb{E}(\mathbf{y})$ by the point in $\mathcal{C}(\mathbf{X})$ that is closest to \mathbf{y} (as measured by the usual Euclidean distance).
- This unique point is called the orthogonal projection of \mathbf{y} onto $\mathcal{C}(\mathbf{X})$ and denoted by $\hat{\mathbf{y}}$ (although it could be argued that $\widehat{\mathbb{E}(\mathbf{y})}$ might be better notation).
- By definition, $\|\mathbf{y} - \hat{\mathbf{y}}\| = \min_{\mathbf{z} \in \mathcal{C}(\mathbf{X})} \|\mathbf{y} - \mathbf{z}\|$ where $\|\mathbf{a}\| = \sqrt{\sum_{i=1}^n a_i^2}$

Geometric Viewpoint on Multiple Regression (and LM)

- Geometrically, how to minimize the distance between \mathbf{Y} and $\mathcal{C}(\mathbf{X})$?
 - ▶ That point is _____
 - ▶ The vector between \mathbf{Y} and $\mathbf{X}\hat{\boldsymbol{\beta}}$ is _____, and the distance is _____
- For R^2 : if we add another predictor, $\mathcal{C}(\mathbf{X})$ gains 1 more dimension, so $\|e\|$ can only decrease. $\mathcal{C}(\mathbf{X})$
 - ▶ Note: if $\dim(\mathcal{S}) = n$ then _____



Orthogonal Projection Matrices

It can be shown that, as we did for least square estimators

- $\forall \mathbf{y} \in \mathbb{R}^n$, $\hat{\mathbf{y}} = \mathbb{P}_{\mathbf{X}}\mathbf{y}$ is the optimal one, i.e.

$\hat{\mathbf{y}} = \mathbb{P}_{\mathbf{X}}\mathbf{y}$ is the best estimator of $\mathbb{E}(\mathbf{y})$ in the class of linear unbiased estimators

for the unique matrix $\mathbb{P}_{\mathbf{X}} = \mathbf{H}$, the hat matrix, and is called **orthogonal projection matrix**

- $\mathbf{H}\mathbf{H} = \mathbf{H}$, idempotent
- $\mathbf{H} = \mathbf{H}'$, symmetric
- $\mathbf{H}\mathbf{X} = \mathbf{X}$ and $\mathbf{X}'\mathbf{H} = \mathbf{X}'$ (Why? Intuitively...)
- If $(\mathbf{X}'\mathbf{X})$ is not invertible, we use its generalized inverse $(\mathbf{X}'\mathbf{X})^-$ where $\mathbf{A}\mathbf{A}^- \mathbf{A} = \mathbf{A}$.
- The \mathbf{H} is invariant to the choice of $(\mathbf{X}'\mathbf{X})^-$, which is itself not unique
- $\hat{\mathbf{y}}$ and $\mathbf{y} - \hat{\mathbf{y}}$ are orthogonal (Why?)

An Example Orthogonal Projection

Suppose $\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \mu + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \end{bmatrix}$, and we observe $y = [6.1, 2.3]$. Then

$$\begin{aligned} \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' &= \begin{bmatrix} 1 \\ 1 \end{bmatrix} \left(\begin{bmatrix} 1 \\ 1 \end{bmatrix}' \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right)^{-1} \begin{bmatrix} 1 \\ 1 \end{bmatrix}' \\ &= \begin{bmatrix} 1 \\ 1 \end{bmatrix} [2]^{-1} [1 \ 1] \\ &= \begin{bmatrix} 1 \\ 1 \end{bmatrix} \left[\frac{1}{2} \right] [1 \ 1] \\ &= \frac{1}{2} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \\ &= \begin{bmatrix} 1/2 & 1/2 \\ 1/2 & 1/2 \end{bmatrix} \end{aligned}$$

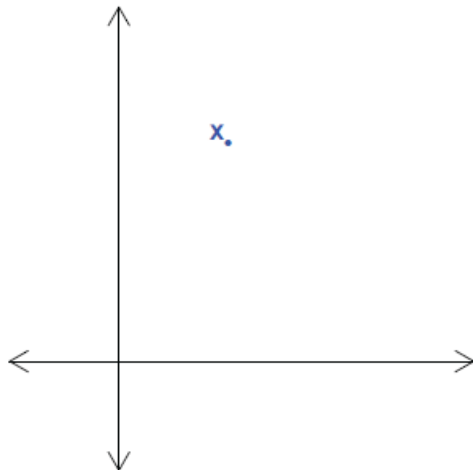
An Example Orthogonal Projection

Thus, the orthogonal projection of $y = [6.1, 2.3]$ onto the column space of $\mathbf{X} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ is

$$\mathbb{P}_{\mathbf{X}}\mathbf{y} = \mathbf{H}\mathbf{y} = \begin{bmatrix} 1/2 & 1/2 \\ 1/2 & 1/2 \end{bmatrix} \begin{bmatrix} 6.1 \\ 2.3 \end{bmatrix} = \begin{bmatrix} 4.2 \\ 4.2 \end{bmatrix}$$

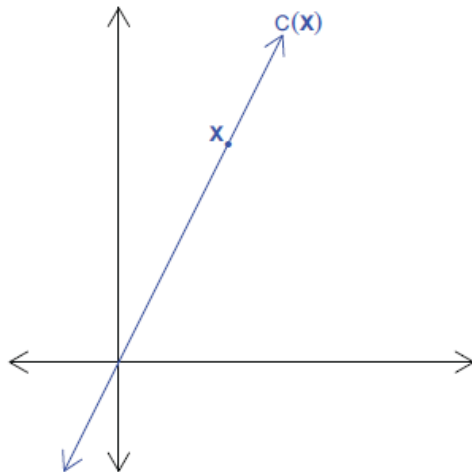
Geometric illustration

Suppose $\mathbf{x} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$ and $\mathbf{y} = \begin{bmatrix} 2 \\ 3/4 \end{bmatrix}$



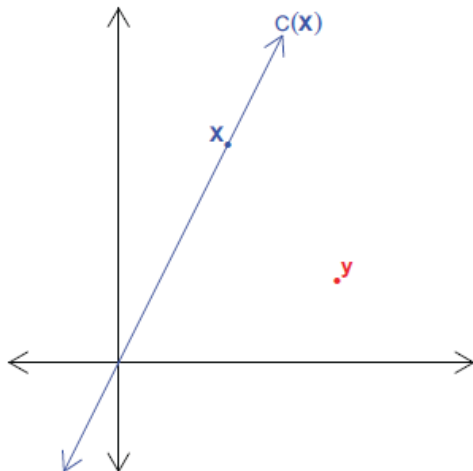
Geometric illustration

Suppose $\mathbf{x} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$ and $\mathbf{y} = \begin{bmatrix} 2 \\ 3/4 \end{bmatrix}$



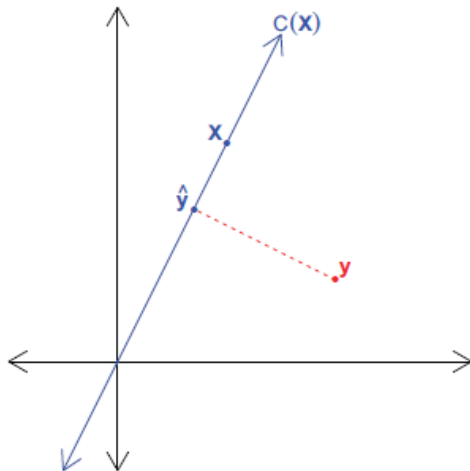
Geometric illustration

Suppose $\mathbf{x} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$ and $\mathbf{y} = \begin{bmatrix} 2 \\ 3/4 \end{bmatrix}$



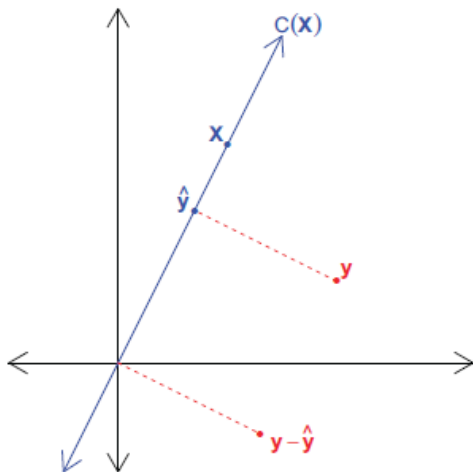
Geometric illustration

Suppose $\mathbf{X} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$ and $\mathbf{y} = \begin{bmatrix} 2 \\ 3/4 \end{bmatrix}$



Geometric illustration

The angle between \hat{y} and residual $y - \hat{y}$ is 90. So, “orthogonal projection”.



- 1 Multiple Linear Regression Model
- 2 Inference on Multiple Regression
- 3 Inference about Regression Parameters
- 4 Estimation and Prediction
- 5 Geometric View of Regression and Linear Models
- 6 Estimating estimable function of coefficient**

What if \mathbf{X} is not full column rank?

- $\mathbf{X}^T \mathbf{X}$ is not invertible, then $(\mathbf{X}^T \mathbf{X})^{-1}$ has to be defined based on the **generalized inverse matrix**.
- If \mathbf{X} is not of full column rank, then there are infinitely many vectors in the set $\{\mathbf{b} : \mathbf{X}\mathbf{b} = \mathbf{X}\boldsymbol{\beta}\}$ for any fixed value of $\boldsymbol{\beta}$.
- Thus, no matter what the value of $\mathbb{E}(\mathbf{y})$, there will be infinitely many vectors \mathbf{b} such that $\mathbf{X}\mathbf{b} = \mathbb{E}(\mathbf{y})$ when \mathbf{X} is not of full column rank.
- Our response vector \mathbf{y} can help us learn about $\mathbb{E}(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$, but when \mathbf{X} is NOT of full column rank, there is NO hope of learning about $\boldsymbol{\beta}$ alone unless additional information about $\boldsymbol{\beta}$ is available.
- How, **we could estimate estimable function of $\boldsymbol{\beta}$**

Treatment Effects Model

Researchers randomly assigned a total of six experimental units to two treatments and measured a response of interest.

$$y_{ij} = \mu + \tau_i + \epsilon_{ij}, \quad i = 1, 2; \quad j = 1, 2, 3$$

$$\begin{bmatrix} y_{11} \\ y_{12} \\ y_{13} \\ y_{21} \\ y_{22} \\ y_{23} \end{bmatrix} = \begin{bmatrix} \mu + \tau_1 \\ \mu + \tau_1 \\ \mu + \tau_1 \\ \mu + \tau_2 \\ \mu + \tau_2 \\ \mu + \tau_2 \end{bmatrix} + \begin{bmatrix} \epsilon_{11} \\ \epsilon_{12} \\ \epsilon_{13} \\ \epsilon_{21} \\ \epsilon_{22} \\ \epsilon_{23} \end{bmatrix}$$

Question: what is X , β ?

Treatment Effects Model (continued)

- In this case, it makes no sense to estimate $\beta = [\mu, \tau_1, \tau_2]'$ because there are multiple (infinitely many, in fact) choices of β that define the same mean for \mathbf{y} .
- For example

$$\begin{bmatrix} \mu \\ \tau_1 \\ \tau_2 \end{bmatrix} = \begin{bmatrix} 5 \\ -1 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 4 \\ 6 \end{bmatrix}, \begin{bmatrix} 999 \\ -995 \\ -993 \end{bmatrix}$$

all yield same $\mathbf{X}\beta = \mathbb{E}(\mathbf{y})$.

- When multiple values for β define the same $\mathbb{E}(\mathbf{y})$, we say that β is *non-estimable*.

Estimable Functions of β

- A linear function of β , $C\beta$, is said to be **estimable** if there is a linear function of \mathbf{y} , say $A\mathbf{y}$, that is an **unbiased** estimator for $C\beta$. Otherwise, nonexistence of such linear function implies that $C\beta$ is *non-estimable*.
- Note that $A\mathbf{y}$ is an unbiased estimator of $C\beta$ if and only if

$$\mathbb{E}(A\mathbf{y}) = C\beta, \text{ for } \forall \beta \in \mathbb{R}^p \Leftrightarrow \mathbf{A}\mathbf{X}\beta = C\beta \Leftrightarrow \mathbf{A}\mathbf{X} = C$$

- This says that we can estimate $C\beta$ as long as $C\beta = \mathbf{A}\mathbf{X}\beta = \mathbf{A}\mathbb{E}(\mathbf{y})$ for some \mathbf{A} , i.e. as long as $C\beta$ is a linear function of $\mathbb{E}(\mathbf{y})$
- The bottom line is **that we can always estimate $\mathbb{E}(\mathbf{y})$ and all linear functions of $\mathbb{E}(\mathbf{y})$; all other linear functions of β are non-estimable**

Treatment Effects Model (continued)

$$\mathbf{X}\boldsymbol{\beta} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} \mu \\ \tau_1 \\ \tau_2 \end{bmatrix} = \begin{bmatrix} \mu + \tau_1 \\ \mu + \tau_1 \\ \mu + \tau_1 \\ \mu + \tau_1 \\ \mu + \tau_2 \\ \mu + \tau_2 \\ \mu + \tau_2 \\ \mu + \tau_2 \end{bmatrix}$$

so that

$$[1, 0, 0, 0, 0, 0] \mathbf{X}\boldsymbol{\beta} = [1, 1, 0] \boldsymbol{\beta} = \mu + \tau_1$$

$$[0, 0, 0, 1, 0, 0] \mathbf{X}\boldsymbol{\beta} = [1, 0, 1] \boldsymbol{\beta} = \mu + \tau_2$$

$$[1, 0, 0, -1, 0, 0] \mathbf{X}\boldsymbol{\beta} = [0, 1, 1] \boldsymbol{\beta} = \tau_1 - \tau_2$$

are estimable functions of $\boldsymbol{\beta}$

Estimating Estimable Functions of β

- If $C\beta$ is estimable, then there exists a matrix A such that $C = AX$ and $C\beta = AX\beta = A\mathbb{E}(\mathbf{y})$ for any $\beta \in \mathbb{R}^p$
- It makes sense to estimate $C\beta$ by

$$\begin{aligned}A\widehat{\mathbb{E}(\mathbf{y})} &= A\hat{\mathbf{y}} = AP_X\mathbf{y} = AX(X'X)^{-1}X'\mathbf{y} \\ &= AX(X'X)^{-1}X'X\hat{\beta} = AP_XX\hat{\beta} = AX\hat{\beta} = C\hat{\beta}\end{aligned}$$

- $C\hat{\beta}$ is called an Ordinary Least Squares (OLS) estimator of $C\beta$
- Note that although the “hat” is on β , it is $C\beta$ that we are estimating
- **Invariance of $C\hat{\beta}$ to the choice of $\hat{\beta}$:** Although there are infinitely many solutions to the normal equations when X is not of full column rank, $C\hat{\beta}$ is the same for all normal equation solutions $\hat{\beta}$ whenever $C\beta$ is estimable (STAT 640)

Treatment Effects Model (continued)

- Suppose our aim is to estimate $\tau_1 - \tau_2$
- As noted before

$$\mathbf{X}\boldsymbol{\beta} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} \mu \\ \tau_1 \\ \tau_2 \end{bmatrix} = \begin{bmatrix} \mu + \tau_1 \\ \mu + \tau_1 \\ \mu + \tau_1 \\ \mu + \tau_1 \\ \mu + \tau_2 \\ \mu + \tau_2 \\ \mu + \tau_2 \\ \mu + \tau_2 \end{bmatrix}, \text{ so that}$$

$$[1, 0, 0, -1, 0, 0] \mathbf{X}\boldsymbol{\beta} = [0, 1, 1] \boldsymbol{\beta} = \tau_1 - \tau_2$$

- Thus, we can compute the OLS estimator of $\tau_1 - \tau_2$ as

$$[1, 0, 0, -1, 0, 0] \hat{\mathbf{y}} = [0, 1, 1] \hat{\boldsymbol{\beta}}$$

where $\hat{\boldsymbol{\beta}}$ is **any** solution to the normal equations.

Treatment Effects Model (continued)

The normal equation in this case is

$$\begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix}' \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix}' \begin{bmatrix} y_{11} \\ y_{12} \\ y_{13} \\ y_{21} \\ y_{22} \\ y_{23} \end{bmatrix}$$

so that

$$\begin{bmatrix} 6 & 3 & 3 \\ 3 & 3 & 0 \\ 3 & 0 & 3 \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} = \begin{bmatrix} y_{..} \\ y_{1.} \\ y_{2.} \end{bmatrix}$$

Treatment Effects Model (continued)

$\hat{\beta}_1 = \begin{bmatrix} \bar{y}_{..} \\ \bar{y}_{1.} - \bar{y}_{..} \\ \bar{y}_{2.} - \bar{y}_{..} \end{bmatrix}$ and $\hat{\beta}_2 = \begin{bmatrix} 0 \\ \bar{y}_{1.} \\ \bar{y}_{2.} \end{bmatrix}$ are both solutions to the normal equation
(Check this).

Thus, the OLS estimator of $C\beta = [0, 1, -1]\beta = \tau_1 - \tau_2$ is

$$C\hat{\beta}_1 = [0, 1, -1] \begin{bmatrix} \bar{y}_{..} \\ \bar{y}_{1.} - \bar{y}_{..} \\ \bar{y}_{2.} - \bar{y}_{..} \end{bmatrix} = \bar{y}_{1.} - \bar{y}_{2.} = [0, 1, -1] \begin{bmatrix} 0 \\ \bar{y}_{1.} \\ \bar{y}_{2.} \end{bmatrix} = C\hat{\beta}_2$$

HW: Can you find two different generalized inverse of $(X'X)$, A_1 and A_2 that $(X'X)A_i(X'X) = (X'X)$ so that $A_i = (X'X)^-$ for each i , and they will give you $\hat{\beta}_1$ and $\hat{\beta}_2$, respectively?

The Gauss-Markov Theorem

Under the Gauss-Markov Linear Model, the OLS estimator $\mathbf{c}'\hat{\boldsymbol{\beta}}$ of an estimable linear function $\mathbf{c}'\boldsymbol{\beta}$ is the unique Best Linear Unbiased Estimator (BLUE) in the sense that $\text{Var}(\mathbf{c}'\hat{\boldsymbol{\beta}})$ is strictly less than the variance of any other linear unbiased estimator of $\mathbf{c}'\boldsymbol{\beta}$ for all $\boldsymbol{\beta} \in \mathbb{R}^p$ and all $\sigma^2 \in \mathbb{R}^+$.

- The Gauss-Markov Theorem says that if we want to estimate an estimable linear function $\mathbf{c}'\boldsymbol{\beta}$ using a linear estimator that is unbiased, we should always use the OLS estimator.
- In our simple example of the treatment effects model, we could have used $y_{11}y_{21}$ to estimate $\tau_1\tau_2$. It is easy to see that $y_{11}y_{21}$ is a linear estimator that is unbiased for $\tau_1\tau_2$, but its variance is clearly larger than the variance of the OLS estimator $\bar{y}_1\bar{y}_2$. (as guaranteed by the Gauss-Markov Theorem).