

Estimation of the Response Mean

The Gauss-Markov Linear Model

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

- \mathbf{y} is an $n \times 1$ random vector of responses.
- \mathbf{X} is an $n \times p$ matrix of constants with columns corresponding to explanatory variables. \mathbf{X} is sometimes referred to as the *design matrix*.
- $\boldsymbol{\beta}$ is an unknown parameter vector in \mathbb{R}^p .
- $\boldsymbol{\epsilon}$ is an $n \times 1$ random vector of errors.
- $E(\boldsymbol{\epsilon}) = \mathbf{0}$ and $\text{Var}(\boldsymbol{\epsilon}) = \sigma^2\mathbf{I}$, where σ^2 is an unknown parameter in \mathbb{R}^+ .

The Column Space of the Design Matrix

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

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

- The set of all possible linear combinations of the columns of X is called the *column space* of X and is denoted by

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

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

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

An Example Column Space

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

Another Example Column Space

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

Different Matrices with the Same Column Space

$$W = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{bmatrix} \quad X = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix}$$

$$\begin{aligned} \mathbf{x} \in \mathcal{C}(W) &\implies \mathbf{x} = W\mathbf{a} \text{ for some } \mathbf{a} \in \mathbb{R}^2 \\ &\implies \mathbf{x} = X \begin{bmatrix} 0 \\ \mathbf{a} \end{bmatrix} \text{ for some } \mathbf{a} \in \mathbb{R}^2 \\ &\implies \mathbf{x} = X\mathbf{b} \text{ for some } \mathbf{b} \in \mathbb{R}^3 \\ &\implies \mathbf{x} \in \mathcal{C}(X) \end{aligned}$$

Thus, $\mathcal{C}(W) \subseteq \mathcal{C}(X)$.

$$\mathbf{W} = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix}$$

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

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

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

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

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

$$\implies \mathbf{x} = \mathbf{W}\mathbf{b} \text{ for some } \mathbf{b} \in \mathbb{R}^2$$

$$\implies \mathbf{x} \in \mathcal{C}(\mathbf{W})$$

Thus, $\mathcal{C}(\mathbf{X}) \subseteq \mathcal{C}(\mathbf{W})$.

We previously showed that $\mathcal{C}(\mathbf{W}) \subseteq \mathcal{C}(\mathbf{X})$.

Thus, it follows that $\mathcal{C}(\mathbf{W}) = \mathcal{C}(\mathbf{X})$.

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

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

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

Should we estimate $E(\mathbf{y}) = \begin{bmatrix} \mu \\ \mu \end{bmatrix}$ by $\mathbf{y} = \begin{bmatrix} 6.1 \\ 2.3 \end{bmatrix}$?

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

- The Gauss-Markov linear models says that $E(\mathbf{y}) \in \mathcal{C}(\mathbf{X})$, so we should use that information when estimating $E(\mathbf{y})$.
- Consider estimating $E(\mathbf{y})$ 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 $\widehat{E(\mathbf{y})}$ might be better notation).
- By definition, $\|\mathbf{y} - \hat{\mathbf{y}}\| = \min_{\mathbf{z} \in \mathcal{C}(\mathbf{X})} \|\mathbf{y} - \mathbf{z}\|$, where $\|\mathbf{a}\| \equiv \sqrt{\sum_{i=1}^n a_i^2}$.

Orthogonal Projection Matrices

In Homework Assignment 2, we will formally prove the following:

- 1 $\forall \mathbf{y} \in \mathbb{R}^n$, $\hat{\mathbf{y}} = \mathbf{P}_X \mathbf{y}$, where \mathbf{P}_X is a unique $n \times n$ matrix known as an *orthogonal projection matrix*.
- 2 \mathbf{P}_X is idempotent: $\mathbf{P}_X \mathbf{P}_X = \mathbf{P}_X$.
- 3 \mathbf{P}_X is symmetric: $\mathbf{P}_X = \mathbf{P}_X'$.
- 4 $\mathbf{P}_X \mathbf{X} = \mathbf{X}$ and $\mathbf{X}' \mathbf{P}_X = \mathbf{X}'$.
- 5 $\mathbf{P}_X = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'$, where $(\mathbf{X}'\mathbf{X})^{-}$ is any generalized inverse of $\mathbf{X}'\mathbf{X}$.

Why Does $P_X X = X$?

$$\begin{aligned} P_X X &= P_X [\mathbf{x}_1, \dots, \mathbf{x}_p] \\ &= [P_X \mathbf{x}_1, \dots, P_X \mathbf{x}_p] \\ &= [\mathbf{x}_1, \dots, \mathbf{x}_p] \\ &= X. \end{aligned}$$

Generalized Inverses

- G is a *generalized inverse* of a matrix A if $AGA = A$.
- We usually denote a generalized inverse of A by A^- .
- If A is nonsingular, i.e., if A^{-1} exists, then A^{-1} is the one and only generalized inverse of A .

$$AA^{-1}A = AI = IA = A$$

- If A is singular, i.e., if A^{-1} does not exist, then there are infinitely many generalized inverses of A .

An Algorithm for Finding a Generalized Inverse of a Matrix A

- 1 Find any $r \times r$ nonsingular submatrix of A where $r = \text{rank}(A)$. Call this matrix W .
- 2 Invert and transpose W , ie., compute $(W^{-1})'$.
- 3 Replace each element of W in A with the corresponding element of $(W^{-1})'$.
- 4 Replace all other elements in A with zeros.
- 5 Transpose the resulting matrix to obtain G , a generalized inverse for A .

Invariance of $\mathbf{P}_X = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'$ to Choice of $(\mathbf{X}'\mathbf{X})^{-}$

- If $\mathbf{X}'\mathbf{X}$ is nonsingular, then $\mathbf{P}_X = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ because the only generalized inverse of $\mathbf{X}'\mathbf{X}$ is $(\mathbf{X}'\mathbf{X})^{-1}$.
- If $\mathbf{X}'\mathbf{X}$ is singular, then $\mathbf{P}_X = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'$ and the choice of the generalized inverse $(\mathbf{X}'\mathbf{X})^{-}$ does not matter because $\mathbf{P}_X = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'$ will turn out to be the same matrix no matter which generalized inverse of $\mathbf{X}'\mathbf{X}$ is used.
- To see this, suppose $(\mathbf{X}'\mathbf{X})_1^{-}$ and $(\mathbf{X}'\mathbf{X})_2^{-}$ are any two generalized inverses of $\mathbf{X}'\mathbf{X}$. Then

$$\mathbf{X}(\mathbf{X}'\mathbf{X})_1^{-}\mathbf{X}' = \mathbf{X}(\mathbf{X}'\mathbf{X})_2^{-}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})_1^{-}\mathbf{X}' = \mathbf{X}(\mathbf{X}'\mathbf{X})_2^{-}\mathbf{X}'.$$

An Example Orthogonal Projection Matrix

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 $\mathbf{y} = \begin{bmatrix} 6.1 \\ 2.3 \end{bmatrix}$.

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

An Example Orthogonal Projection

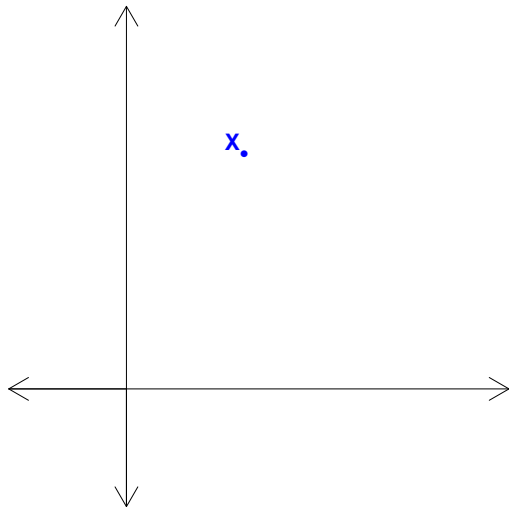
Thus, the orthogonal projection of $\mathbf{y} = \begin{bmatrix} 6.1 \\ 2.3 \end{bmatrix}$

onto the column space of $\mathbf{X} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

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

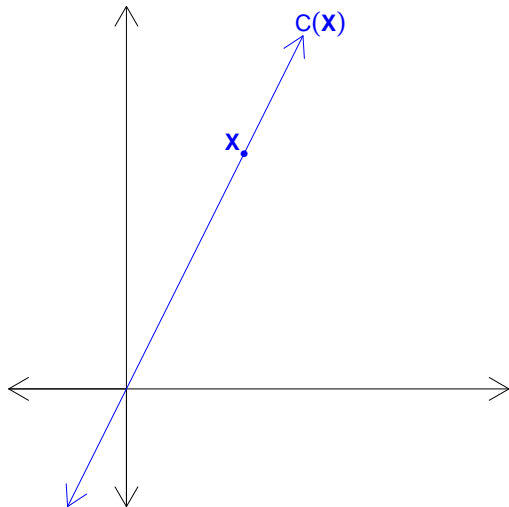
Why is P_X called an *orthogonal* projection matrix?

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



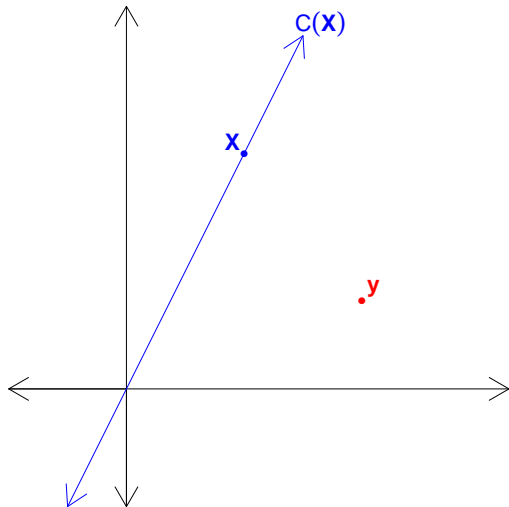
Why is P_X called an *orthogonal* projection matrix?

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



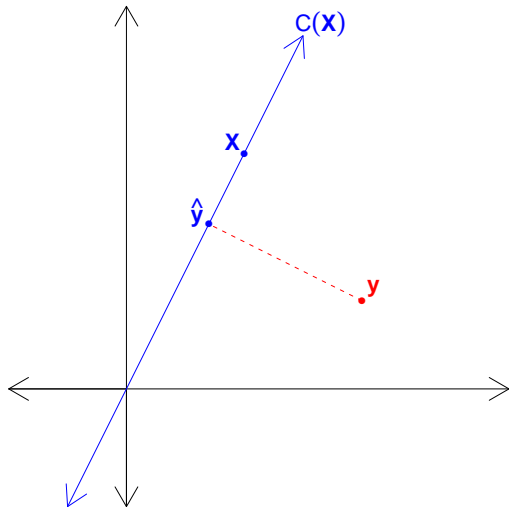
Why is P_X called an *orthogonal* projection matrix?

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



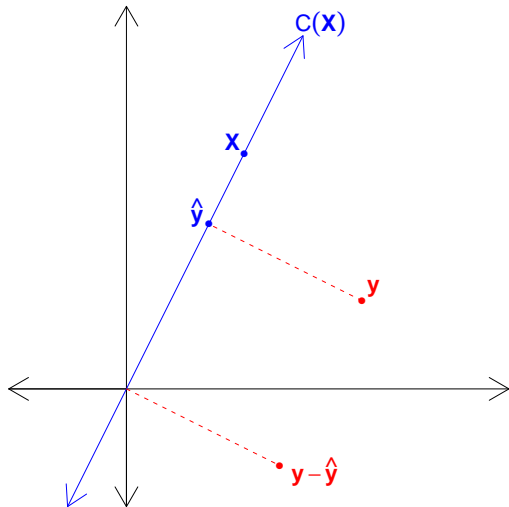
Why is P_X called an *orthogonal* projection matrix?

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



Why is P_X called an *orthogonal* projection matrix?

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



Why is P_X called an *orthogonal* projection matrix?

- The angle between \hat{y} and $y - \hat{y}$ is 90° .
- The vectors \hat{y} and $y - \hat{y}$ are orthogonal.

$$\begin{aligned}\hat{y}'(y - \hat{y}) &= \hat{y}'(y - P_X y) = \hat{y}'(I - P_X)y \\ &= (P_X y)'(I - P_X)y = y'P_X'(I - P_X)y \\ &= y'P_X(I - P_X)y = y'(P_X - P_X P_X)y \\ &= y'(P_X - P_X)y = 0.\end{aligned}$$

Optimality of $\hat{\mathbf{y}}$ as an Estimator of $E(\mathbf{y})$

- $\hat{\mathbf{y}}$ is an unbiased estimator of $E(\mathbf{y})$:

$$E(\hat{\mathbf{y}}) = E(\mathbf{P}_X \mathbf{y}) = \mathbf{P}_X E(\mathbf{y}) = \mathbf{P}_X \mathbf{X} \boldsymbol{\beta} = \mathbf{X} \boldsymbol{\beta} = E(\mathbf{y}).$$

- It can be shown that $\hat{\mathbf{y}} = \mathbf{P}_X \mathbf{y}$ is the best estimator of $E(\mathbf{y})$ in the class of linear unbiased estimators, i.e., estimators of the form $\mathbf{M} \mathbf{y}$ for \mathbf{M} satisfying

$$E(\mathbf{M} \mathbf{y}) = E(\mathbf{y}) \quad \forall \boldsymbol{\beta} \in \mathbb{R}^p \iff \mathbf{M} \mathbf{X} \boldsymbol{\beta} = \mathbf{X} \boldsymbol{\beta} \quad \forall \boldsymbol{\beta} \in \mathbb{R}^p \iff \mathbf{M} \mathbf{X} = \mathbf{X}.$$

- Under the Gauss-Markov Linear Model, $\hat{\mathbf{y}} = \mathbf{P}_X \mathbf{y}$ is best among all unbiased estimators of $E(\mathbf{y})$.

Ordinary Least Squares (OLS) Estimation of $E(\mathbf{y})$

- OLS: Find a vector $\mathbf{b}^* \in \mathbb{R}^p$ such that

$$Q(\mathbf{b}^*) \leq Q(\mathbf{b}) \quad \forall \mathbf{b} \in \mathbb{R}^p, \text{ where } Q(\mathbf{b}) \equiv \sum_{i=1}^n (y_i - \mathbf{x}'_{(i)}\mathbf{b})^2.$$

- Note that

$$Q(\mathbf{b}) = \sum_{i=1}^n (y_i - \mathbf{x}'_{(i)}\mathbf{b})^2 = (\mathbf{y} - \mathbf{X}\mathbf{b})'(\mathbf{y} - \mathbf{X}\mathbf{b}) = \|\mathbf{y} - \mathbf{X}\mathbf{b}\|^2.$$

- To minimize this sum of squares, we need to choose $\mathbf{b}^* \in \mathbb{R}^p$ such $\mathbf{X}\mathbf{b}^*$ will be the point in $\mathcal{C}(\mathbf{X})$ that is closest to \mathbf{y} .
- In other words, we need to choose \mathbf{b}^* such that $\mathbf{X}\mathbf{b}^* = \mathbf{P}_X\mathbf{y} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$.
- Clearly, choosing $\mathbf{b}^* = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ will work.

Ordinary Least Squares and the Normal Equations

- It can be shown that $Q(\mathbf{b}^*) \leq Q(\mathbf{b}) \forall \mathbf{b} \in \mathbb{R}^p$ if and only if \mathbf{b}^* is a solution to the *normal equations*:

$$\mathbf{X}'\mathbf{X}\mathbf{b} = \mathbf{X}'\mathbf{y}.$$

- If $\mathbf{X}'\mathbf{X}$ is nonsingular, multiplying both sides of the normal equations by $(\mathbf{X}'\mathbf{X})^{-1}$ shows that the only solution to the normal equations is $\mathbf{b}^* = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$.
- If $\mathbf{X}'\mathbf{X}$ is singular, there are infinitely many solutions that include $(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{y}$ for all choices of generalized inverse of $\mathbf{X}'\mathbf{X}$.

$$\mathbf{X}'\mathbf{X}[(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{y}] = \mathbf{X}'[\mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}']\mathbf{y} = \mathbf{X}'\mathbf{P}_X\mathbf{y} = \mathbf{X}'\mathbf{y}$$

- Henceforth, we will use $\hat{\beta}$ to denote any solution to the normal equations.

Ordinary Least Squares Estimator of $E(\mathbf{y}) = \mathbf{X}\beta$

- We call $\mathbf{X}\hat{\beta} = \mathbf{P}_X\mathbf{X}\hat{\beta} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\hat{\beta} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \mathbf{P}_{X\mathbf{y}} = \hat{\mathbf{y}}$ the OLS estimator of $E(\mathbf{y}) = \mathbf{X}\beta$.
- It might be more appropriate to use $\widehat{\mathbf{X}}\hat{\beta}$ rather than $\mathbf{X}\hat{\beta}$ to denote our estimator because we are estimating $\mathbf{X}\beta$ rather than pre-multiplying an estimator of β by \mathbf{X} .
- As we shall soon see, it does not make sense to estimate β when $\mathbf{X}'\mathbf{X}$ is singular.
- However, it does make sense to estimate $E(\mathbf{y}) = \mathbf{X}\beta$ whether $\mathbf{X}'\mathbf{X}$ is singular or nonsingular.