

# More Properties of Least Squares Estimation<sup>1</sup>

STA 302 Fall 2020

---

<sup>1</sup>See last slide for copyright information.

# Overview

- 1 Unbiased Estimation
- 2 Gauss-Markov Theorem
- 3 Projections

## Reading in In Rencher and Schaalje's *Linear Models In Statistics*

Much of this material is in Section 7.3.2 (pp. 145-149), except

- The Gauss-Markov Theorem is done better here.
- They discuss projections *briefly* in Chapter 9.

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

where

$\mathbf{X}$  is an  $n \times (k + 1)$  matrix of observed constants with linearly independent columns.

$\boldsymbol{\beta}$  is a  $(k + 1) \times 1$  matrix of unknown constants (parameters).

$\boldsymbol{\epsilon}$  is an  $n \times 1$  random vector with  $E(\boldsymbol{\epsilon}) = \mathbf{0}$  and  $\text{cov}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}_n$ .

$\sigma^2$  is an unknown constant.

Least squares estimator of  $\boldsymbol{\beta}$  is

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

# Unbiased Estimation

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

$$\begin{aligned} E\{\widehat{\boldsymbol{\beta}}\} &= E\{(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}\} \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'E\{\mathbf{y}\} \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\boldsymbol{\beta} \\ &= \boldsymbol{\beta} \end{aligned}$$

for any  $\boldsymbol{\beta} \in \mathbb{R}^{k+1}$ , so  $\widehat{\boldsymbol{\beta}}$  is an unbiased estimator of  $\boldsymbol{\beta}$ .

## Covariance matrix

Using  $\text{cov}(\mathbf{A}\mathbf{w}) = \mathbf{A}\text{cov}(\mathbf{w})\mathbf{A}'$

$$\begin{aligned}\text{cov}(\widehat{\boldsymbol{\beta}}) &= \text{cov}((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}) \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\text{cov}(\mathbf{y})((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')' \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\sigma^2\mathbf{I}_n\mathbf{X}''(\mathbf{X}'\mathbf{X})^{-1'} \\ &= \sigma^2(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \\ &= \sigma^2(\mathbf{X}'\mathbf{X})^{-1}\end{aligned}$$

# What are we estimating when we estimate $\beta$ ?

Human resources example:  $y = \beta_0 + \beta_1x_1 + \beta_2x_2 + \beta_3x_3 + \epsilon$

- $x_1 =$  University GPA.
- $x_2 =$  Job interview score.
- $x_3 =$  Test score.
- $y =$  Percent salary increase after one year.
  
- $E(y) = \beta_0 + \beta_1x_1 + \beta_2x_2 + \beta_3x_3$ .
- $\beta_1, \beta_2$  and  $\beta_3$  are *links* between predictor variables and (expected) response variable value.
- $\beta_0$  is for curve fitting – no interpretation in this example.
- Question: Holding interview and test scores constant, how much does GPA matter?

$$E(y) = \beta_0 + \beta_2x_2 + \beta_3x_3 + \beta_1x_1.$$

Estimating linear combinations of  $\beta$  values

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \epsilon$$

$$\ell_0 \beta_0 + \ell_1 \beta_1 + \cdots + \ell_k \beta_k$$

$x_1$  = University GPA,  $x_2$  = Interview score,  $x_3$  = Test score.  
For fixed job interview score and test score, what's the connection between GPA and salary increase?

$$\ell' \beta = (0 \quad 1 \quad 0 \quad 0) \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix} = \beta_1$$



## Another linear combination

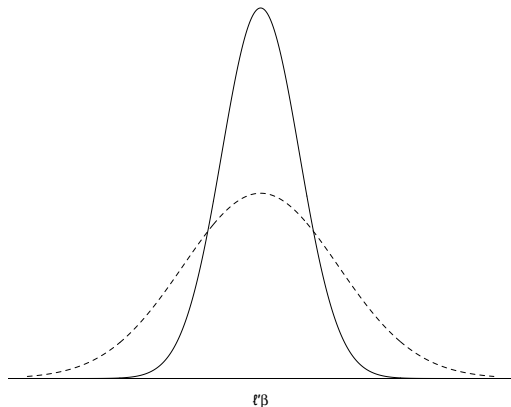
What's the expected salary increase for a job candidate with a university GPA of 2.5, an interview score of 80% and a test score of 70%?

$$\ell' \boldsymbol{\beta} = (1 \quad 2.5 \quad 80 \quad 70) \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix}$$

Estimated expected value is often used for prediction.

# Natural Estimator

- Natural Estimator of  $\ell' \beta$  is  $\ell' \hat{\beta}$ .
- It's unbiased:  $E\{\ell' \hat{\beta}\} = \ell' E\{\hat{\beta}\} = \ell' \beta$
- Small variance in an unbiased estimator is good. It's the variance of the *sampling distribution*.



# Linear Combination

- The natural estimator of  $\ell'\beta$  is a linear combination of the  $y_i$  values.

$$\ell'\hat{\beta} = \ell'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \mathbf{a}'_0\mathbf{y}$$

- Let  $L = a_1y_1 + a_2y_2 + \cdots + a_ny_n$  be another linear combination of  $y_i$  with  $E(L) = \ell'\beta$  for every  $\beta \in \mathbb{R}^{k+1}$ .
- If we can find  $L$ , unbiased, with  $Var(L) < Var(\ell'\hat{\beta})$ , we should use that  $L$  to estimate  $\ell'\beta$  instead of  $\ell'\hat{\beta}$ .

# A Serious $L = \mathbf{a}'\mathbf{y}$

$$\widehat{\boldsymbol{\beta}}_w = (\mathbf{X}'\mathbf{W}\mathbf{X})^{-1}\mathbf{X}'\mathbf{W}\mathbf{y}$$

where  $\mathbf{W}$  is an  $n \times n$  matrix of rank at least  $k + 1$ .

$$\begin{aligned} E\{\widehat{\boldsymbol{\beta}}_w\} &= E\{(\mathbf{X}'\mathbf{W}\mathbf{X})^{-1}\mathbf{X}'\mathbf{W}\mathbf{y}\} \\ &= (\mathbf{X}'\mathbf{W}\mathbf{X})^{-1}\mathbf{X}'\mathbf{W}E\{\mathbf{y}\} \\ &= (\mathbf{X}'\mathbf{W}\mathbf{X})^{-1}\mathbf{X}'\mathbf{W}\mathbf{X}\boldsymbol{\beta} \\ &= \boldsymbol{\beta} \end{aligned}$$

- Let  $L = \boldsymbol{\ell}'\widehat{\boldsymbol{\beta}}_w$ .
- Then  $E\{L\} = \boldsymbol{\ell}'E\{\widehat{\boldsymbol{\beta}}_w\} = \boldsymbol{\ell}'\boldsymbol{\beta}$ .
- Should we seek  $\mathbf{W}$  with  $Var(\boldsymbol{\ell}'\widehat{\boldsymbol{\beta}}_w) < Var(\boldsymbol{\ell}'\widehat{\boldsymbol{\beta}})$ ?
- The Gauss-Markov Theorem says don't bother.

# The Gauss-Markov Theorem

For the general linear model  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ , etc., let  $E(\mathbf{a}'\mathbf{y}) = \boldsymbol{\ell}'\boldsymbol{\beta}$  for all  $\boldsymbol{\beta} \in \mathbb{R}^{k+1}$ .

Then  $Var(\boldsymbol{\ell}'\hat{\boldsymbol{\beta}}) \leq Var(\mathbf{a}'\mathbf{y})$ , with equality only when  $\mathbf{a} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\boldsymbol{\ell}$  (in which case  $\mathbf{a}'\mathbf{y} = \boldsymbol{\ell}'\hat{\boldsymbol{\beta}}$ ).

# Proof of the Gauss-Markov-Theorem

- The impressive part.
- The rest of the proof (just a calculation).

## The impressive part

$$\begin{aligned} E(\mathbf{a}'\mathbf{y}) &= \mathbf{a}'E(\mathbf{y}) \\ &= \mathbf{a}'\mathbf{X}\boldsymbol{\beta} \\ &= \boldsymbol{\ell}'\boldsymbol{\beta} \end{aligned}$$

For all  $\boldsymbol{\beta} \in \mathbb{R}^{k+1}$ .

- This implies  $\mathbf{a}'\mathbf{X} = \boldsymbol{\ell}'$ .
- But *not* by cancelling  $\boldsymbol{\beta}$ !

$$\mathbf{a}'\mathbf{X}\boldsymbol{\beta} = \boldsymbol{\ell}'\boldsymbol{\beta} \text{ for all } \boldsymbol{\beta} \in \mathbb{R}^{k+1}$$

- $\mathbf{a}'\mathbf{X} = \mathbf{v}'$  is  $1 \times (k+1)$ .
- $\mathbf{v}' = (v_0, v_1, \dots, v_k)$ .
- $\mathbf{v}'\boldsymbol{\beta} = \boldsymbol{\ell}'\boldsymbol{\beta}$ .
- For all  $\boldsymbol{\beta} \in \mathbb{R}^{k+1}$ , meaning even for very funny  $\boldsymbol{\beta}$  vectors.

$$\mathbf{v}'\boldsymbol{\beta} = (v_0 \quad v_1 \quad \cdots \quad v_k) \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = (\ell_0 \quad \ell_1 \quad \cdots \quad \ell_k) \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = \boldsymbol{\ell}'\boldsymbol{\beta}$$

So  $v_0 = \ell_0$ .



$$\mathbf{v}'\boldsymbol{\beta} = \boldsymbol{\ell}'\boldsymbol{\beta}$$

For all  $\boldsymbol{\beta} \in \mathbb{R}^{k+1}$

$$\begin{aligned}\mathbf{v}'\boldsymbol{\beta} &= (v_0 \quad v_1 \quad v_2 \quad \cdots \quad v_k) \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \\ &= \boldsymbol{\ell}'\boldsymbol{\beta} \\ &= (\ell_0 \quad \ell_1 \quad \ell_2 \quad \cdots \quad \ell_k) \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}\end{aligned}$$

So  $v_1 = \ell_1$ .

$$\mathbf{v}'\boldsymbol{\beta} = \boldsymbol{\ell}'\boldsymbol{\beta}$$

For all  $\boldsymbol{\beta} \in \mathbb{R}^{k+1}$

$$\begin{aligned}\mathbf{v}'\boldsymbol{\beta} &= \begin{pmatrix} v_0 & v_1 & v_2 & \cdots & v_k \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 1 \\ \vdots \\ 0 \end{pmatrix} \\ &= \boldsymbol{\ell}'\boldsymbol{\beta} \\ &= \begin{pmatrix} \ell_0 & \ell_1 & \ell_2 & \cdots & \ell_k \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 1 \\ \vdots \\ 0 \end{pmatrix}\end{aligned}$$

So  $v_2 = \ell_2$ .

# Continuing ...

• • •

$$\mathbf{v}'\boldsymbol{\beta} = \boldsymbol{\ell}'\boldsymbol{\beta}$$

For all  $\boldsymbol{\beta} \in \mathbb{R}^{k+1}$

$$\begin{aligned}\mathbf{v}'\boldsymbol{\beta} &= (v_0 \quad v_1 \quad v_2 \quad \cdots \quad v_k) \begin{pmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix} \\ &= \boldsymbol{\ell}'\boldsymbol{\beta} \\ &= (\ell_0 \quad \ell_1 \quad \ell_2 \quad \cdots \quad \ell_k) \begin{pmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix}\end{aligned}$$

So  $v_k = \ell_k$ .

# Conclusion

$$\mathbf{a}'\mathbf{X} = \boldsymbol{\ell}' \Leftrightarrow \boldsymbol{\ell} = \mathbf{X}'\mathbf{a}$$

- This condition is both necessary and sufficient for  $\mathbf{a}'\mathbf{y}$  to be an unbiased estimator of  $\boldsymbol{\ell}'\boldsymbol{\beta}$ .
- We have proved necessary.

# Calculation part of the Proof

Using  $\boldsymbol{\ell}' = \mathbf{a}'\mathbf{X} \Leftrightarrow \boldsymbol{\ell} = \mathbf{X}'\mathbf{a}$

$$\begin{aligned} \text{Var}(\mathbf{a}'\mathbf{y}) - \text{Var}(\boldsymbol{\ell}'\widehat{\boldsymbol{\beta}}) &= \text{cov}(\mathbf{a}'\mathbf{y}) - \text{cov}(\boldsymbol{\ell}'\widehat{\boldsymbol{\beta}}) \\ &= \mathbf{a}'\text{cov}(\mathbf{y})\mathbf{a} - \boldsymbol{\ell}'\text{cov}(\widehat{\boldsymbol{\beta}})\boldsymbol{\ell} \\ &= \mathbf{a}'\sigma^2\mathbf{I}_n\mathbf{a} - \boldsymbol{\ell}'\sigma^2(\mathbf{X}'\mathbf{X})^{-1}\boldsymbol{\ell} \\ &= \sigma^2(\mathbf{a}'\mathbf{I}_n\mathbf{a} - \boldsymbol{\ell}'(\mathbf{X}'\mathbf{X})^{-1}\boldsymbol{\ell}) \\ &= \sigma^2(\mathbf{a}'\mathbf{I}_n\mathbf{a} - \mathbf{a}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{a}) \\ &= \sigma^2\mathbf{a}'(\mathbf{I}_n - \mathbf{H})\mathbf{a} \\ &= \sigma^2\mathbf{a}'(\mathbf{I}_n - \mathbf{H})'(\mathbf{I}_n - \mathbf{H})\mathbf{a} \\ &= \sigma^2((\mathbf{I}_n - \mathbf{H})\mathbf{a})'(\mathbf{I}_n - \mathbf{H})\mathbf{a} \\ &= \sigma^2\mathbf{z}'\mathbf{z} = \sigma^2\sum_{i=1}^n z_i^2 \\ &\geq 0. \end{aligned}$$

## Continuing

And using  $\ell' = \mathbf{a}'\mathbf{X} \Leftrightarrow \ell = \mathbf{X}'\mathbf{a}$  again

- Have  $Var(\mathbf{a}'\mathbf{y}) - Var(\ell'\hat{\boldsymbol{\beta}}) = \sigma^2 \mathbf{z}'\mathbf{z} \geq 0$ ,
- Where  $\mathbf{z} = (\mathbf{I}_n - \mathbf{H})\mathbf{a}$ .
- Variances are the same if and only if  $\mathbf{z} = \mathbf{0}$ .

$$\Rightarrow (\mathbf{I}_n - \mathbf{H})\mathbf{a} = \mathbf{0}$$

$$\Rightarrow \mathbf{a} = \mathbf{H}\mathbf{a}$$

$$\Rightarrow \mathbf{a} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{a}$$

$$\Rightarrow \mathbf{a} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\ell$$

$$\Rightarrow \mathbf{a}'\mathbf{y} = \ell'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \ell'\hat{\boldsymbol{\beta}}$$

So  $\ell'\hat{\boldsymbol{\beta}}$  is the *unique* minimum variance linear unbiased estimator of  $\ell'\boldsymbol{\beta}$ . ■

## BLUE

Sometimes we say that  $\hat{\beta}$  is the

Best

Linear

Unbiased

Estimator.

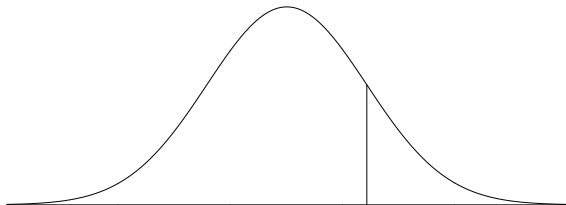


# Projections

- Let  $\mathcal{V} = \{\mathbf{v} \in \mathbb{R}^n : \mathbf{v} = \mathbf{X}\mathbf{b}, \mathbf{b} \in \mathbb{R}^{k+1}\}$
- The space *spanned* by the columns of  $\mathbf{X}$ .
- All linear combinations of the columns of  $\mathbf{X}$ . The elements of  $\mathbf{b}$  are the coefficients of the linear combination.
- Some important vectors are in  $\mathcal{V}$ .
  - $E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$ :  $\boldsymbol{\beta}$  is a vector  $\mathbf{b}$ .
  - $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$ :  $\hat{\boldsymbol{\beta}}$  is a vector  $\mathbf{b}$ .
  - Every column of  $\mathbf{X}$  is in  $\mathcal{V}$ .
  - Is  $\mathbf{y} \in \mathcal{V}$ ?

Is  $\mathbf{y} \in \mathcal{V} = \{\mathbf{v} \in \mathbb{R}^n : \mathbf{v} = \mathbf{X}\mathbf{b}, \mathbf{b} \in \mathbb{R}^{k+1}\}$ ?

- The  $k + 1$  linearly independent columns of  $\mathbf{X}$  span  $\mathcal{V}$ .
- So  $\mathcal{V}$  is of dimension  $k + 1 < n$ .
- And  $\mathcal{V}$  is a set of volume zero in  $\mathbb{R}^n$ .
- If  $\epsilon_i$  have a continuous distribution (with a density), then the distribution of the random vector  $\mathbf{y}$  is also continuous.
- And the probability that  $\mathbf{y}$  will fall into a set of volume zero is equal to zero:  $P\{\mathbf{y} \in \mathcal{V}\} = 0$ .



# What point $\mathbf{p} \in \mathcal{V}$ is closest to $\mathbf{y}$ ?

Euclidean distance is

$$\sqrt{(y_1 - p_1)^2 + (y_2 - p_2)^2 + \cdots + (y_n - p_n)^2}$$

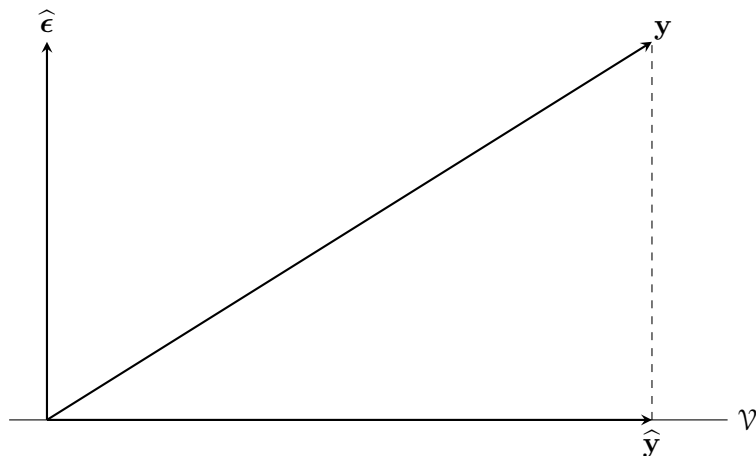
where  $\mathbf{p} = \mathbf{X}\mathbf{b}$ , some  $\mathbf{b} \in \mathbb{R}^{k+1}$ . To find it, minimize

$$(\mathbf{y} - \mathbf{p})'(\mathbf{y} - \mathbf{p}) = (\mathbf{y} - \mathbf{X}\mathbf{b})'(\mathbf{y} - \mathbf{X}\mathbf{b})$$

over all  $\mathbf{b} \in \mathbb{R}^{k+1}$ .

- We've already done this!
- The answer is  $\mathbf{b} = \hat{\boldsymbol{\beta}}$ .
- $\mathbf{p} = \mathbf{X}\hat{\boldsymbol{\beta}} = \hat{\mathbf{y}}$ .
- The closest point in  $\mathcal{V}$  to  $\mathbf{y}$  is  $\hat{\mathbf{y}}$ .

Projection:  $\hat{\mathbf{y}}$  is the shadow of  $\mathbf{y}$  on  $\mathcal{V}$



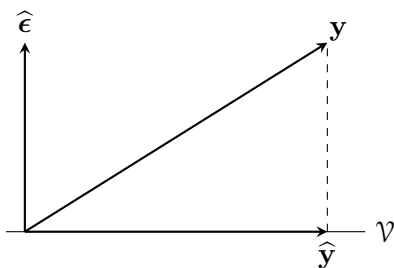
$$\hat{\mathbf{y}} + \hat{\boldsymbol{\epsilon}} = \hat{\mathbf{y}} + (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y}$$

# Projection Operator

$$\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$$

- $\hat{\mathbf{y}}$  is the projection of  $\mathbf{y}$  onto  $\mathcal{V}$ .
- $\mathbf{H}$  is the projection operator:  $\mathbf{H}\mathbf{y} = \hat{\mathbf{y}}$ .
- $\mathbf{H}$  sends any point in  $\mathbb{R}^n$  to  $\mathcal{V}$ .  
$$\mathbf{H}\mathbf{p} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{p} = \mathbf{X}\mathbf{b}.$$
- The projection is the closest point.
- If  $\mathbf{p} \in \mathcal{V}$  already,  $\mathbf{H}\mathbf{p} = \mathbf{p}$ .  
$$\mathbf{H}\mathbf{p} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\mathbf{b} = \mathbf{X}\mathbf{b} = \mathbf{p}.$$

Picture suggests  $\hat{\boldsymbol{\epsilon}} \perp \hat{\mathbf{y}}$

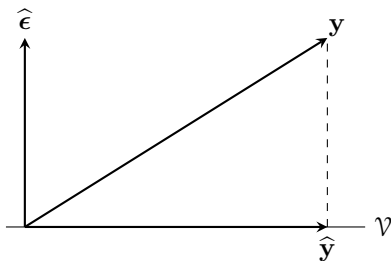


- In fact,  $\hat{\boldsymbol{\epsilon}} \perp \mathbf{v}$  for all  $\mathbf{v} \in \mathcal{V}$ .

$$\begin{aligned} \mathbf{v}'\hat{\boldsymbol{\epsilon}} &= (\mathbf{X}\mathbf{b})'\hat{\boldsymbol{\epsilon}} \\ &= \mathbf{b}'\mathbf{X}'\hat{\boldsymbol{\epsilon}} \\ &= \mathbf{b}'\mathbf{0} = 0 \end{aligned}$$

- $\mathbf{v} \in \mathcal{V}$  includes
  - $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$ .
  - $E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$ .
  - Every column of  $\mathbf{X}$ .

## Another way to arrive at the normal equations



- Least squares task is to minimize  $Q = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$ .
- Find the  $\mathbf{X}\boldsymbol{\beta}$  point in  $\mathcal{V}$  that is closest to  $\mathbf{y}$ . Call it  $\mathbf{X}\hat{\boldsymbol{\beta}}$ .
- Drop a perpendicular (normal) from  $\mathbf{y}$  to  $\mathcal{V}$ .

- This perpendicular is parallel to  $\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\epsilon}}$ .
- So  $\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}$  is at right angles to all basis vectors of  $\mathcal{V}$ . Inner products are all zero.
- That is,  $\mathbf{X}'(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \mathbf{0}$ .  
 $\Rightarrow \mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}'\mathbf{y}$ .
- These are the “normal equations.”
- Wikipedia says “In geometry, a normal is an object such as a line, ray, or vector that is perpendicular to a given object.”

## Copyright Information

This slide show was prepared by **Jerry Brunner**, Department of Statistical Sciences, University of Toronto. It is licensed under a **Creative Commons Attribution - ShareAlike 3.0 Unported License**. Use any part of it as you like and share the result freely. The L<sup>A</sup>T<sub>E</sub>X source code is available from the course website:

<http://www.utstat.toronto.edu/~brunner/oldclass/302f20>