



# Modern Machine Learning

## Gauss-Markov Theorem

Kenneth E. Barner

Department of Electrical and Computer Engineering  
University of Delaware



**Objective:** Prove that the least-squares determine linear model parameters,  $\hat{\beta}_{LS}$ , have the **smallest variance among all linear unbiased estimates**

## Gauss-Markov Theorem

Theorem Assumptions:

- $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ ,
  - where  $\mathbf{X}$  is fixed (non-random) and  $\boldsymbol{\epsilon}$  components are iid  $N(0, \sigma^2)$
- $\hat{\boldsymbol{\beta}} = \mathbf{C}\mathbf{y}$  be a linear unbiased estimator of  $\boldsymbol{\beta}$

Then if for all  $\mathbf{a} \in \mathbb{R}^{(p+1)}$

$$MSE(\mathbf{a}^T \hat{\boldsymbol{\beta}}_{LS}) \leq MSE(\mathbf{a}^T \hat{\boldsymbol{\beta}}) \quad [\text{MSE} = \text{Mean Squared Error}]$$

We say that  $\hat{\boldsymbol{\beta}}_{LS}$  is the **Best Linear Unbiased Estimator** (BLUE) of  $\boldsymbol{\beta}$



To prove the result consider first the: **bias-variance decomposition of the MSE**

Let  $\mathbf{Z} = \mathbf{a}^T \boldsymbol{\beta}$  and  $\hat{\mathbf{Z}} = \mathbf{a}^T \hat{\boldsymbol{\beta}}$  then

$$\begin{aligned} \text{MSE}(\mathbf{a}^T \hat{\boldsymbol{\beta}}) &= E \left[ \left( \mathbf{a}^T (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) \right)^2 \right] = E \left[ (\mathbf{Z} - \hat{\mathbf{Z}})^2 \right] \\ &= E \left[ \mathbf{Z}^2 - 2\mathbf{Z}\hat{\mathbf{Z}} + \hat{\mathbf{Z}}^2 \right] \\ &= E[\mathbf{Z}^2] - 2E[\mathbf{Z}\hat{\mathbf{Z}}] + E[\hat{\mathbf{Z}}^2] \\ &= \mathbf{Z}^2 - 2\mathbf{Z}E[\hat{\mathbf{Z}}] + \text{Var}[\hat{\mathbf{Z}}] + E[\hat{\mathbf{Z}}]^2 \quad (*) \\ &= (\mathbf{Z} - E[\hat{\mathbf{Z}}])^2 + \text{Var}[\hat{\mathbf{Z}}] \end{aligned}$$

$\text{bias}^2 := (\mathbf{Z} - E[\hat{\mathbf{Z}}])^2$ , variance:  $= \text{Var}[\hat{\mathbf{Z}}]$

Therefore, if  $\hat{\boldsymbol{\beta}}$  is unbiased then  $\text{MSE}(\mathbf{a}^T \hat{\boldsymbol{\beta}}) = \text{Var}[\hat{\mathbf{Z}}]$

**NOTE:** derivation (\*) above utilizes known result  $\text{Var}[\hat{\mathbf{Z}}] = E[\hat{\mathbf{Z}}^2] - E[\hat{\mathbf{Z}}]^2$



To prove the Gauss-Markov theorem it suffices to show that

$$\text{Var}(\mathbf{a}^T \hat{\boldsymbol{\beta}}_{LS}) \leq \text{Var}(\mathbf{a}^T \hat{\boldsymbol{\beta}}), \forall \mathbf{a} \in \mathbb{R}^{(p+1)}$$

**Proof:** Let  $\hat{\boldsymbol{\beta}} = \mathbf{C}\mathbf{y}$ , where  $\mathbf{C} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T + \mathbf{D}$ , for some  $\mathbf{D} \in \mathbb{R}^{(p+1) \times n}$  then

$$\begin{aligned} E[\hat{\boldsymbol{\beta}}] &= E\left[\left((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T + \mathbf{D}\right) \mathbf{y}\right] \\ &= E\left[\left((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T + \mathbf{D}\right) (\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon})\right] \\ &= (\mathbf{I} + \mathbf{D}\mathbf{X})\boldsymbol{\beta} \end{aligned}$$

**Observation:** It *must* hold that  $\mathbf{D}\mathbf{X} = \mathbf{0}$  for  $\hat{\boldsymbol{\beta}}$  to be unbiased

Now calculate  $\text{Var}(\mathbf{a}^T \hat{\boldsymbol{\beta}})$



Then we have

$$\begin{aligned}\text{Var}(\hat{\boldsymbol{\beta}}) &= \text{Var}(\mathbf{C}\mathbf{y}) \\ &= \mathbf{C} \text{Var}(\mathbf{y})\mathbf{C}^T = \sigma^2 \mathbf{C}\mathbf{C}^T \quad [\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \Rightarrow \text{Var}(\mathbf{y}) = \sigma^2 \mathbf{I}] \\ &= \sigma^2 \left( (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T + \mathbf{D} \right) \left( (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T + \mathbf{D} \right)^T \\ &= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} + \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{D}^T \\ &\quad + \sigma^2 \mathbf{D} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} + \sigma^2 \mathbf{D} \mathbf{D}^T \\ &= \sigma^2 \left[ (\mathbf{X}^T \mathbf{X})^{-1} + \mathbf{D} \mathbf{D}^T \right] \quad [\text{since } \mathbf{D} \mathbf{X} = \mathbf{0}]\end{aligned}$$

**Note:**  $\mathbf{X}^T \mathbf{X}$  and  $\mathbf{D} \mathbf{D}^T$  are positive semidefinite

$$\Rightarrow \mathbf{a}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{a} \geq 0, \quad \mathbf{a}^T (\mathbf{D}^T \mathbf{D}) \mathbf{a} \geq 0$$



Therefore

$$\begin{aligned}\text{Var}(\mathbf{a}^T \hat{\beta}) &= \mathbf{a}^T \sigma^2 [(\mathbf{X}^T \mathbf{X})^{-1} + \mathbf{D} \mathbf{D}^T] \mathbf{a} \geq \mathbf{a}^T \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{a} \\ &\geq \text{Var}(\mathbf{a}^T \hat{\beta}_{LS})\end{aligned}$$

$$\text{Since } \text{Var}(\mathbf{a}^T \hat{\beta}_{LS}) = \text{Var}(\mathbf{a}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{y}) = \mathbf{a}^T \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{a}$$

This concludes the proof

**Result:** The Least Squares estimate  $\hat{\beta}_{LS}$  is unbiased and has the smallest weight variance  $\Rightarrow$  it is the **Best Linear Unbiased Estimate (BLUE)**