# Some Key Linear Models Results

Ulrike Genschel

January 23, 2025

# A General Linear Model (GLM)

Suppose

$$\boldsymbol{y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \qquad \text{where} \tag{1}$$

- $\boldsymbol{y} \in \mathbb{R}^n$ is the response vector,

- $\boldsymbol{X}$ is an $n \times p$ matrix of known/fixed constants,

  } known

- $\boldsymbol{\beta} \in \mathbb{R}^p$ is an unknown parameter vector, and

- $\boldsymbol{\epsilon}$ is a vector of unobserved random "errors" satisfying $\mathrm{E}(\boldsymbol{\epsilon}) = \boldsymbol{0}$ and $\mathrm{Cov}(\boldsymbol{\epsilon}) = \boldsymbol{\Sigma}$.

( unknown

The model is called a linear model because the mean of the response vector $\boldsymbol{y}$ is linear in the unknown parameter vector $\boldsymbol{\beta}$. ($\mathrm{E}(\boldsymbol{y}) = \boldsymbol{X}\boldsymbol{\beta}$)

# A General Linear Model

- This GLM says simply that $y$ is a random vector with expectation $\mathrm{E}(y) = X\beta$ for some $\beta \in \mathbb{R}^p$.

- The distribution of $y$ is left unspecified but generally depends on the distribution of $\epsilon$.

- Goal: estimate $\mathrm{E}(y)$

- Available: <u>observed</u> values of $y$ and $X$,

- Estimate $X\beta$, which by definition corresponds to the mean of $y$, i.e., $\mathrm{E}(y)$.

# Examples

There are many special cases of (1) depending on the distribution of $\epsilon$, the structure of the $\mathbf{\Sigma}$, and the rank and the structure of $\mathbf{X}$.

We will start out by considering the following two cases generally known as the Gauss-Markov Model:

*G hh* (handwritten)

*GMNE* (handwritten)

1. the distribution of $\epsilon$ is Normal with $\mathrm{E}(\epsilon) = \mathbf{0}$ and $\mathrm{Cov}(\epsilon) = \mathbf{\Sigma}_\epsilon = \sigma^2 \mathbf{I}$, where $\sigma^2 > 0$ is unknown; $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$

2. the distribution of $\epsilon$ is unknown with $\mathrm{E}(\epsilon) = \mathbf{0}$ and $\mathrm{Cov}(\epsilon) = \mathbf{\Sigma}_\epsilon = \sigma^2 \mathbf{I}$, where $\sigma^2 > 0$ is unknown

We will later relax the form of $\mathrm{Cov}(\epsilon) = \mathbf{\Sigma}_\epsilon$ to allow for more flexibility, e.g., $\mathrm{Cov}(\epsilon) = \mathbf{\Sigma}_\epsilon = \sigma^2 \mathbf{V}$, where $\mathbf{V}$ is known and $\sigma^2 > 0$ is unknown. This model is known as the Aitken model.

# Ordinary Least Squares (OLS) Estimation

Suppose $\quad \boldsymbol{y} = \boldsymbol{X\beta} + \boldsymbol{\epsilon}, \qquad \mathrm{E}(\boldsymbol{\epsilon}) = 0, \qquad \mathrm{Cov}(\boldsymbol{\epsilon}) = \sigma^2 \boldsymbol{I}$

- $\mathrm{E}(\boldsymbol{y}) = \boldsymbol{X\beta} \in \mathcal{C}(\boldsymbol{X})$ with $\boldsymbol{\beta}$ unknown, $\boldsymbol{X}$ is full-rank

  *Column space of $X$*

- To estimate $\mathrm{E}(\boldsymbol{y})$, consider $\boldsymbol{X}\widehat{\boldsymbol{\beta}}$.

  *estimate: $\widehat{y}$*

- To estimate $\mathrm{E}(\boldsymbol{y})$, find the vector in $\mathcal{C}(\boldsymbol{X})$ that is closest to $\boldsymbol{y}$.

- Let $\mathcal{N}(\boldsymbol{X}^\top)$ denote the null space of $\boldsymbol{X}^\top$ and note that $\mathcal{N}(\boldsymbol{X}^\top)$ and $\mathcal{C}(\boldsymbol{X})$ are orthogonal to each other, i.e., $\mathcal{N}(\boldsymbol{X}^\top) \perp \mathcal{C}(\boldsymbol{X})$

*residuals $\widehat{e}$ are in the null space*

The null space of a matrix $\boldsymbol{A}$, denoted by $\mathcal{N}(\boldsymbol{A})$, is given as   *of $X^\top$*

$\mathcal{N}(\boldsymbol{A}) = \{\boldsymbol{x} : \boldsymbol{xA} = \boldsymbol{0}.\}$

*allowing for orthogonal decomposition of $y$ into $\widehat{y} + \widehat{e}$*

# Ordinary Least Squares (OLS) Estimation

An estimate $\widehat{\boldsymbol{\beta}}$ is a **least squares estimate** (LSE) of $\boldsymbol{\beta}$ if $\boldsymbol{X}\widehat{\boldsymbol{\beta}}$ is the

*Ordinary*

vector in $\mathcal{C}(\boldsymbol{X})$ that is closes to $\boldsymbol{y}$

$$\widehat{\boldsymbol{\beta}} = \min_{\boldsymbol{\beta} \in \mathbb{R}^p} (\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta})^\top (\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}).$$

Method of least squares identifies the value of $\boldsymbol{\beta}$ for which the squared Euclidean norm of the residual vector, i.e., **error sum of squares**

$$\mathcal{Q}(\boldsymbol{\beta}) = \|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}\|_2^2 = (\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta})^\top (\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta})$$

is minimized.

# Ordinary Least Squares (OLS) Estimation

There exist two distinct ways to identify the LSE:

- algebraically: normal equations

- geometrically: orthogonal projection of $y$ onto $\mathcal{C}(X)$

# OLS Estimation: Normal Equations

Recall that the method of least squares seeks the $\boldsymbol{\beta}$ that minimizes the Euclidean norm of the residual vector

$$
\begin{aligned}
\mathcal{Q}(\boldsymbol{\beta}) &= \|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}\|_2^2 = (\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta})^\top(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}) \\
&= \boldsymbol{y}^\top\boldsymbol{y} - 2\boldsymbol{\beta}^\top\boldsymbol{X}^\top\boldsymbol{y} + \boldsymbol{\beta}^\top\boldsymbol{X}^\top\boldsymbol{X}\boldsymbol{\beta}.
\end{aligned}
$$

To find the minimum, we take the derivative and set the gradient equal to the null vector

$$
\nabla\mathcal{Q}(\boldsymbol{\beta}) = -2\boldsymbol{X}^\top\boldsymbol{y} + 2\boldsymbol{X}^\top\boldsymbol{X}\boldsymbol{\beta} = \boldsymbol{0}
$$

leading to the **normal equations**

$$
\boldsymbol{X}^\top\boldsymbol{X}\boldsymbol{\beta} = \boldsymbol{X}^\top\boldsymbol{y}. \tag{2}
$$

The normal equations

$$\boldsymbol{X}^\top \boldsymbol{X} \boldsymbol{\beta} = \boldsymbol{X}^\top \boldsymbol{y}$$

have $(\boldsymbol{X}^\top \boldsymbol{X})^{-1} \boldsymbol{X}^\top \boldsymbol{y}$ as the **unique** solution for $\boldsymbol{\beta}$ if $rank(\boldsymbol{X}) = p$.

$(X^\top X)^{-1}$ unique inverse when $X$ is full rank

The normal equations have infinitely many solutions for $\boldsymbol{\beta}$ if $rank(\boldsymbol{X}) < p$.

generalized inverse

While $\widehat{\boldsymbol{\beta}} = (\boldsymbol{X}^\top \boldsymbol{X})^- \boldsymbol{X}^\top \boldsymbol{y}$ may not always be a unique solution, $\boldsymbol{X}\widehat{\boldsymbol{\beta}} = \widehat{\boldsymbol{y}}$ will be unique.

# OLS Estimation: Geometric Approach

Let $P_X$ denote the underlined orthogonal projection matrix onto $\mathcal{C}(X)$

$$P_X = X(X^\top X)^- X^\top.$$

Properties:

- $P_X$ is idempotent, (i.e., $P_X P_X = P_X$)
- $P_X$ projects onto $\mathcal{C}(X)$
- $P_X$ is invariant to the choice of $(X^\top X)^-$, i.e., it is the same matrix for all generalized inverses $(X^\top X)^-$ of $X^\top X$
- $P_X$ is symmetric (i.e., $P_X = P_X^\top$) and unique
- $P_X X = X$ and $X^\top P_X = X^\top$.      *trace*
- $rank(X) = rank(P_X) = tr(P_X)$.

# OLS Estimation: Geometric Approach

An estimate $\widehat{\boldsymbol{\beta}}$ is a least squares estimate if and only if

$$X\widehat{\boldsymbol{\beta}} = \boxed{P_X y}.$$

*projecting $y$ onto column space of $X$*

The OLS Estimator of $\mathrm{E}(\boldsymbol{y})$ is thus given by

$$P_X y = X\widehat{\boldsymbol{\beta}} \equiv \widehat{y} \; = \; \widehat{\mathcal{E}(y)} \tag{3}$$

because $P_X y \in \mathcal{C}(X)$ and

$$||\boldsymbol{y} - P_X y||^2 < ||\boldsymbol{y} - \boldsymbol{z}||^2 \; \forall \; \boldsymbol{z} \in \mathcal{C}(X) \setminus \{P_X y\}.$$

Even when $\widehat{\boldsymbol{\beta}}$ is not unique, $P_X y = X\widehat{\boldsymbol{\beta}} \equiv \widehat{y}$ always will.

# OLS Estimation: Fitted Values

$\widehat{y} = P_X y$ is the vector of fitted values. Recall that geometrically, $\widehat{y}$ is the point in $\mathcal{C}(X)$ that is closest to $y$. Now, note that $I - P_X$ is the perpendicular projection matrix onto $\mathcal{N}(X^\top)$ and

$$(I - P_X)y = y - P_X y = y - \widehat{y} \equiv \widehat{e}.$$

$\widehat{e}$ is the vector of **residuals** and $\widehat{e} \in \mathcal{N}(X^\top)$. Because $\mathcal{C}(X)$ and $\mathcal{N}(X^\top)$ are orthogonal complements, we can uniquely decompose $y$ as

$$y = \widehat{y} + \widehat{e}.$$

We know that $\widehat{y}$ and $\widehat{e}$ are orthogonal vectors. Thus,

*Projects onto $C(x)$*

*Projects onto $N(x^\top)$*

*adding zero*

$$
\begin{aligned}
y^\top y = y^\top I y &= y^\top (P_X + (I - P_X)) y \\
&= y^\top P_X y + y^\top (I - P_X) y \\
&= y^\top P_X P_X y + y^\top (I - P_X)(I - P_X) y \\
&= \widehat{y}^\top \widehat{y} + \widehat{e}^\top \widehat{e},
\end{aligned}
$$

$P_X$

since $P_X$ and $(I - P_X)$ are both symmetric and idempotent.

# Orthogonal Decomposition of $\boldsymbol{y}^\top \boldsymbol{y}$ & ANOVA Table

This orthogonal decomposition of $\boldsymbol{y}^\top \boldsymbol{y}$ is often given in a tabular display called an analysis of variance (ANOVA) table.

Suppose $\boldsymbol{y}$ is $n \times 1$, $\boldsymbol{X}$ is $n \times p$ with rank $r \leq p$, $\boldsymbol{\beta}$ is $p \times 1$, and $\boldsymbol{\epsilon}$ is $n \times 1$. We assume the the model given in (1): $\boldsymbol{y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$. Then, the ANOVA table looks as follows

| Source | df | Sum of Squares |
|:---:|:---:|:---:|
| Model | $r$ | $\widehat{\boldsymbol{y}}^\top \widehat{\boldsymbol{y}} = \boldsymbol{y}^\top \boldsymbol{P_X} \boldsymbol{y}$ |
| Residual | $n - r$ | $\widehat{\mathbf{e}}^\top \widehat{\mathbf{e}} = \boldsymbol{y}^\top (\boldsymbol{I} - \boldsymbol{P_X})\boldsymbol{y}$ |
| Total | $n - 1$ | $\boldsymbol{y}^\top \boldsymbol{y} = \boldsymbol{y}^\top \boldsymbol{I} \boldsymbol{y}$ |

Table: ANOVA Table

# The OLS Estimator of a Linear Function of $\mathrm{E}(\boldsymbol{y})$

For any $q \times n$ matrix $\boldsymbol{A}$, $\boldsymbol{A}\mathrm{E}(\boldsymbol{y})$ is a linear function of $\mathrm{E}(\boldsymbol{y})$.

For any $q \times n$ matrix $\boldsymbol{A}$, the OLS Estimator of $\boldsymbol{A}\mathrm{E}(\boldsymbol{y}) = \boldsymbol{A}\boldsymbol{X}\boldsymbol{\beta}$ is

$$
\begin{aligned}
\boldsymbol{A}\,[\text{OLS Estimator of } \mathrm{E}(\boldsymbol{y})] &= \boldsymbol{A}\widehat{\boldsymbol{y}} = \boldsymbol{A}\boldsymbol{P}_{\boldsymbol{X}}\boldsymbol{y} \\
&= \boldsymbol{A}\boldsymbol{X}(\boldsymbol{X}^{\top}\boldsymbol{X})^{-}\boldsymbol{X}^{\top}\boldsymbol{y}.
\end{aligned}
$$

- $\boldsymbol{A}\mathrm{E}(\boldsymbol{y}) = \boldsymbol{A}\boldsymbol{X}\boldsymbol{\beta}$ is automatically a linear function of $\boldsymbol{\beta}$ of the form $\boldsymbol{C}\boldsymbol{\beta}$, where $\boldsymbol{C} = \boldsymbol{A}\boldsymbol{X}$.

- If $\boldsymbol{C}$ is any $q \times p$ matrix, we say that the linear function of $\boldsymbol{\beta}$ given by $\boldsymbol{C}\boldsymbol{\beta}$ is estimable if and only if $\boldsymbol{C} = \boldsymbol{A}\boldsymbol{X}$ for some matrix $q \times n$ matrix $\boldsymbol{A}$.

- The OLS Estimator of an estimable linear function $\boldsymbol{C}\boldsymbol{\beta}$ is $\boldsymbol{C}(\boldsymbol{X}^{\top}\boldsymbol{X})^{-}\boldsymbol{X}^{\top}\boldsymbol{y}$.

# Uniqueness of the OLS Estimator of an Estimable $C\beta$

If $C\beta$ is estimable, then $C\widehat{\beta}$ is the same for all solutions $\widehat{\beta}$ to the Normal Equations.

In particular, the unique OLS Estimator of $C\beta$ is

$$C\widehat{\beta} = C(X^\top X)^- X^\top y = AX(X^\top X)^- X^\top y = AP_X y,$$

where $C = AX$.

# The OLS Estimator is a Linear Unbiased Estimator

If $C\beta$ is estimable, then $C\widehat{\beta}$ is a linear unbiased estimator of $C\beta$.

The OLS Estimator is a linear estimator because it is a linear function of $y$:

$$C\widehat{\beta} = C(X^\top X)^- X^\top y = My, \text{ where } M = C(X^\top X)^- X^\top.$$

The OLS Estimator is unbiased because, for all $\beta \in \mathbb{R}^p$,

$$
\begin{aligned}
\mathrm{E}(C\widehat{\beta}) &= \mathrm{E}(C(X^\top X)^- X^\top y) = C(X^\top X)^- X^\top \mathrm{E}(y) \\
&= AX(X^\top X)^- X^\top \mathrm{E}(y) = AP_X \mathrm{E}(y) \\
&= AP_X X\beta = AX\beta = C\beta.
\end{aligned}
$$

$My$

Constant

$C = AX$

$= X\beta$

Slide 10

# The Gauss-Markov Model (GMM)

Suppose $y = X\beta + \epsilon$, where

- $y \in \mathbb{R}^n$ is the response vector,

- $X$ is an $n \times p$ matrix of known constants,

- $\beta \in \mathbb{R}^p$ is an unknown parameter vector, and

- $\epsilon$ is a vector of random "errors" satisfying $\mathrm{E}(\epsilon) = \mathbf{0}$ and $\mathrm{Var}(\epsilon) = \sigma^2 I$ for some unknown variance parameter $\sigma^2 \in \mathbb{R}^+$.

# The GMM is a Special Case of the GLM

The GMM is a special case of the GLM presented previously.

We have added the assumption $\mathrm{Var}(\boldsymbol{\epsilon}) = \sigma^2 \boldsymbol{I}$; i.e., we assume the errors are uncorrelated and have constant variance.

All the results presented for the GLM hold for the GMM.

# The Gauss-Markov Theorem

For the GMM, we have an additional result provided by the
*Gauss-Markov Theorem*:

## The Gauss-Markov Theorem

The OLS Estimator of an estimable function $C\beta$ is the

$$\textit{Best Linear Unbiased Estimator (BLUE)} \text{ of } C\beta$$

in the sense that the OLS Estimator $C\hat{\beta}$ has the smallest variance
among all linear unbiased estimators of $C\beta$.

*End lecture 2*