

Linear combinations. If $\mathbf{a} = (a_0, \dots, a_p)^\top \in \mathbb{R}^{p+1}$ is a fixed vector, then

$$\theta := \mathbf{a}^\top \boldsymbol{\beta} = \sum_{j=0}^p a_j \beta_j \in \mathbb{R} \quad (19)$$

is called a *linear combination* (of $\boldsymbol{\beta}$). Consider estimating a linear combination (19). A natural estimator is

$$\hat{\theta} = \mathbf{a}^\top \hat{\boldsymbol{\beta}} = \mathbf{a}^\top \left(\mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{X}^\top \mathbf{Y} = \mathbf{c}^\top \mathbf{Y}, \quad (20)$$

where

$$\mathbf{c} := \mathbf{X} \left(\mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{a}, \quad (21)$$

and note that $\mathbf{c} \in \mathbb{R}^n$ whereas $\mathbf{a} \in \mathbb{R}^{p+1}$. This estimator is a linear function of \mathbf{Y} , and we can calculate its mean and variance under the linear model,

$$\mathbb{E}[\hat{\theta}] = \mathbb{E}[\mathbf{a}^\top \hat{\boldsymbol{\beta}}] = \mathbf{a}^\top \mathbb{E}[\hat{\boldsymbol{\beta}}] = \mathbf{a}^\top \boldsymbol{\beta} = \theta$$

and

$$V(\hat{\theta}) = V(\mathbf{a}^\top \hat{\boldsymbol{\beta}}) = V(\mathbf{c}^\top \mathbf{Y}) = \mathbf{c}^\top \text{cov}(\mathbf{Y}) \mathbf{c} = \mathbf{c}^\top [\sigma^2 \mathbf{I}_n] \mathbf{c} = \sigma^2 \mathbf{c}^\top \mathbf{c}.$$

Thus $\hat{\theta}$ is a *linear, unbiased* estimator of θ with variance $\sigma^2 \mathbf{c}^\top \mathbf{c}$. Is there any *better* linear unbiased estimator of θ ? First we need to define “better”. The mean squared error (MSE) of an estimator $\hat{\theta}$ of θ is

$$\text{MSE}(\hat{\theta}) := \mathbb{E}_\theta[(\hat{\theta} - \theta)^2],$$

and notice that this generally depends on the true value θ . We will say that an estimator $\hat{\theta}$ of θ is *better* than another estimator $\tilde{\theta}$ if

$$\text{MSE}(\hat{\theta}) \leq \text{MSE}(\tilde{\theta}) \quad \text{for all } \theta.$$

For any estimator $\hat{\theta}$, we have

$$\begin{aligned} \text{MSE}(\hat{\theta}) &= \mathbb{E}[(\hat{\theta} - \theta)^2] = \mathbb{E}[(\hat{\theta} - \mathbb{E}\hat{\theta}) + (\mathbb{E}\hat{\theta} - \theta)^2] = \\ &= \mathbb{E}[(\hat{\theta} - \mathbb{E}\hat{\theta})^2] + \mathbb{E}[(\mathbb{E}\hat{\theta} - \theta)^2] + 2 \underbrace{\mathbb{E}[(\hat{\theta} - \mathbb{E}\hat{\theta})(\mathbb{E}\hat{\theta} - \theta)]}_{=0} = \\ &= \underbrace{\mathbb{E}[(\hat{\theta} - \mathbb{E}\hat{\theta})^2]}_{V(\hat{\theta})} + \underbrace{\mathbb{E}[(\mathbb{E}\hat{\theta} - \theta)^2]}_{(\text{bias}(\hat{\theta}))^2} \end{aligned}$$

where we used that fact that $\mathbb{E}[(\hat{\theta} - \mathbb{E}\hat{\theta})] = \mathbb{E}\hat{\theta} - \mathbb{E}\hat{\theta} = 0$.

We conclude from the general decomposition (5) that an unbiased estimator has

$$\text{MSE}(\hat{\theta}) = V(\hat{\theta}).$$

Hence, restricting attention to unbiased estimators, the estimator $\hat{\theta}$ is better than another estimator $\tilde{\theta}$ if

$$V(\hat{\theta}) \leq V(\tilde{\theta}) \quad \forall \theta.$$

The following theorem, maybe the most famous result in all of linear regression, says that, under the linear model (15), the LS estimator $\hat{\boldsymbol{\beta}}$ is the *best linear unbiased estimator* (BLUE) of $\boldsymbol{\beta}$.

Theorem 1 (Gauss-Markov). . Let $\theta := \mathbf{a}^\top \boldsymbol{\beta}$ be a linear combination, and assume the linear model (15). Denote by $\hat{\theta}$ the LS estimator of θ in (20), and consider another linear unbiased estimator $\tilde{\theta}$ of θ

$$\tilde{\theta} = \mathbf{d}^\top \mathbf{Y}, \quad \mathbb{E}[\tilde{\theta}] = \theta \quad \forall \theta.$$

Then

$$V(\hat{\theta}) \leq V(\tilde{\theta}) \quad \forall \theta$$

Proof. For \mathbf{c} defined in (21), write

$$\mathbf{d} = \mathbf{c} + \boldsymbol{\Delta}, \quad \boldsymbol{\Delta} := \mathbf{d} - \mathbf{c} \in \mathbb{R}^n.$$

$\tilde{\theta}$ is unbiased, hence for all $\boldsymbol{\beta}$ we have

$$\begin{aligned} \theta = \mathbb{E}\tilde{\theta} &= \mathbb{E}[\mathbf{d}^\top \mathbf{Y}] = \mathbb{E}[(\mathbf{c} + \boldsymbol{\Delta})^\top \mathbf{Y}] = \mathbb{E}[(\mathbf{c} + \boldsymbol{\Delta})^\top \mathbf{Y}] = \mathbb{E}[\mathbf{c}^\top \mathbf{Y}] + \mathbb{E}[(\boldsymbol{\Delta}^\top \mathbf{Y})] \\ &= \theta + \boldsymbol{\Delta}^\top \mathbb{E}[\mathbf{Y}] = \theta + \boldsymbol{\Delta}^\top \mathbf{X}\boldsymbol{\beta}, \end{aligned}$$

where the second-to-last equality is due to unbiasedness of $\hat{\theta}$. Comparing the two extreme sides of the sequence of the equality, we get

$$\boldsymbol{\Delta}^\top \mathbf{X}\boldsymbol{\beta} = 0 \quad \forall \boldsymbol{\beta} \quad \Rightarrow \quad \boldsymbol{\Delta}^\top \mathbf{X} = \mathbf{0},$$

so

$$\boldsymbol{\Delta}^\top \mathbf{c} = \underbrace{\boldsymbol{\Delta}^\top \mathbf{X}}_{=0} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a} = \mathbf{0}.$$

We then calculate

$$\begin{aligned} V(\tilde{\theta}) &= V(\mathbf{d}^\top \mathbf{Y}) = V[(\mathbf{c} + \boldsymbol{\Delta})^\top \mathbf{Y}] \\ &= \text{cov}[(\mathbf{c} + \boldsymbol{\Delta})^\top \mathbf{Y}] \\ &= (\mathbf{c} + \boldsymbol{\Delta})^\top \text{cov}[\mathbf{Y}] (\mathbf{c} + \boldsymbol{\Delta}) \\ &= (\mathbf{c} + \boldsymbol{\Delta})^\top \sigma^2 [\mathbf{I}_n] (\mathbf{c} + \boldsymbol{\Delta}) \\ &= \sigma^2 (\mathbf{c} + \boldsymbol{\Delta})^\top (\mathbf{c} + \boldsymbol{\Delta}) \\ &= \sigma^2 (\mathbf{c}^\top \mathbf{c} + \boldsymbol{\Delta}^\top \boldsymbol{\Delta}) \\ &\geq \sigma^2 \mathbf{c}^\top \mathbf{c} \\ &= V(\hat{\theta}). \end{aligned}$$

□

We have considered point estimation of a scalar $\theta = \theta(\beta)$, more specifically unbiased estimation of a linear function of β . We now want to move on to other inferential tasks, for example we'll want to use the LS estimator $\hat{\beta}$ to construct a confidence interval for β , or to test whether a particular coordinate β_j is equal to zero. For this we will need some further assumptions on the linear model.

Review of multivariate distributions. All the concepts presented here generalize naturally beyond the two dimensional case. If Z_1, Z_2 are two random variables, then $Z = (Z_1, Z_2)^\top$ is a random vector of dimension 2. The joint cumulative distribution function (CDF) of Z is

$$F_Z(z_1, z_2) := P(Z_1 \leq z_1, Z_2 \leq z_2),$$

which is always defined and determines the distribution of Z . The variables Z_1 and Z_2 are (statistically) independent if

$$F_Z(z_1, z_2) = P(Z_1 \leq z_1) P(Z_2 \leq z_2) \quad \text{for all } z_1, z_2 \in \mathbb{R}.$$

If the derivative

$$f_Z(z_1, z_2) = \frac{\partial^2}{\partial z_1 \partial z_2} F_Z(z_1, z_2)$$

exists (for all except maybe a subset of \mathbb{R}^2 of probability zero), we call f_Z the *joint density* of Z , and we have the relation

$$F_Z(z_1, z_2) = \int_{-\infty}^{z_1} \int_{-\infty}^{z_2} f_Z(u_1, u_2) du_1 du_2.$$

Of course, the derivative is in that case an equivalent characterization of the distribution of Z .

The multivariate Normal distribution.

Definition 4. We say that a random vector $W = (W_1, \dots, W_k)^\top$ has a multivariate normal distribution if there exists a representation

$$W \stackrel{d}{=} \mu + AZ \tag{22}$$

where $\mu \in \mathbb{R}^k$ and $A \in \mathbb{R}^{k \times l}$ are constant (nonrandom) and where $Z = (Z_1, \dots, Z_l)^\top$ is a random vector whose components Z_i are i.i.d. $\mathcal{N}(0, 1)$ random variables (“ $\stackrel{d}{=}$ ” means “equal in distribution”).

Properties of the multivariate Normal distribution.

1. If W has a multivariate normal distribution, then

$$\begin{aligned} \mathbb{E}[W] &= \mathbb{E}[\mu + AZ] = \mathbb{E}[\mu] + \mathbb{E}[AZ] = \mu + A\mathbb{E}[Z] = \mu \\ \text{cov}(W) &= \text{cov}(\mu + AZ) = \text{cov}(AZ) = A \text{cov}(Z) A^\top = AA^\top \end{aligned}$$

Therefore, if there is another representation $W \stackrel{d}{=} \mu' + A'Z$, then necessarily $\mu' = \mu$ and $A'A'^\top = AA^\top$ (this, in turn, can be shown to hold if and only if $A' = AU^\top$ for an orthogonal matrix U – try to prove this).

2. In (22) suppose that $l = k$, and if $A_{k \times k}$ has linearly independent columns, and denote $V := AA^\top$. Then

$$f_W(w) = (2\pi)^{-m/2} |V|^{-1/2} \exp \left[- (w - \mu)^\top V^{-1} (w - \mu) / 2 \right], \quad w \in \mathbb{R}^m$$

Hence, the distribution of W in (22) is completely determined by μ and V .

We write

$$W \sim \mathcal{N}_k(\mu, V)$$

for the multivariate distribution with mean μ and covariance matrix V (this notation applies whether or not $A_{k \times k}$ has linearly independent columns).

3. It is a consequence of 2 that if $\mathbf{W}^{(1)} = \boldsymbol{\mu} + \mathbf{A}^{(1)}\mathbf{Z}^{(1)}$ and $\mathbf{W}^{(2)} = \boldsymbol{\mu} + \mathbf{A}^{(2)}\mathbf{Z}^{(2)}$, and if $\mathbf{A}^{(1)}\mathbf{A}^{(1)\top} = \mathbf{A}^{(2)}\mathbf{A}^{(2)\top}$, then $\mathbf{W}^{(1)} \stackrel{d}{=} \mathbf{W}^{(2)} \sim \mathcal{N}_k(\boldsymbol{\mu}, \mathbf{V})$.

4. From the previous properties, if $\mathbf{c} \in \mathbb{R}^k$ is a constant vector, then

$$\mathbf{c}^\top \mathbf{W} \sim \mathcal{N}(\mathbf{c}^\top \boldsymbol{\mu}, \mathbf{c}^\top \mathbf{V} \mathbf{c})$$

In words, a linear combination of a multivariate normal vector has a univariate normal distribution. In particular, if we take $\mathbf{c} = (\underbrace{0, \dots, 0}_{j-1}, 1, \underbrace{0, \dots, 0}_{k-j})^\top$, then

$$W_j = \mathbf{c}^\top \mathbf{W} \sim \mathcal{N}(\mu_j, \mathbf{V}_{jj})$$

5. If for a random vector \mathbf{W} it holds that $\mathbf{c}^\top \mathbf{W} \sim \mathcal{N}(\mathbf{c}^\top \boldsymbol{\mu}, \mathbf{c}^\top \mathbf{V} \mathbf{c}) \forall \mathbf{c} \in \mathbb{R}^m$, where $\boldsymbol{\mu}$ and \mathbf{V} denote the mean and covariance of \mathbf{W} , then \mathbf{W} has a multivariate normal distribution. Combined with property 4, this says

$$\mathbf{c}^\top \mathbf{W} \sim \mathcal{N}(\mathbf{c}^\top \boldsymbol{\mu}, \mathbf{c}^\top \mathbf{V} \mathbf{c}) \quad \forall \mathbf{c} \in \mathbb{R}^m \quad \Longleftrightarrow \quad \mathbf{W} \sim \mathcal{N}_m(\boldsymbol{\mu}, \mathbf{V}).$$

Thus, Property 4 is in fact a *defining property* of the multivariate normal distribution.

5. If $\mathbf{C} \in \mathbb{R}^{m \times k}$ constant matrix then $\mathbf{C}\mathbf{W} \sim \mathcal{N}_m(\mathbf{C}\boldsymbol{\mu}, \mathbf{C}\mathbf{V}\mathbf{C}^\top)$.

6. If $\mathbf{W}^{(j)} \sim \mathcal{N}_k(\boldsymbol{\mu}^{(j)}, \mathbf{V}^{(j)})$, $j = 1, \dots, p$, independent, and if d_j are scalar constants, then

$$\sum_{j=1}^p d_j \mathbf{W}^{(j)} \sim \mathcal{N}_k\left(\sum_{j=1}^p d_j \boldsymbol{\mu}^{(j)}, \sum_{j=1}^p d_j^2 \mathbf{V}^{(j)}\right)$$

7. Let $\mathbf{W} \sim \mathcal{N}_k(\boldsymbol{\mu}, \mathbf{V})$ and $\mathcal{J}_1, \mathcal{J}_2 \subseteq \{1, \dots, k\}$ disjoint subsets of indices. If $\text{Cov}(W_i, W_j) = 0 \quad \forall i \in \mathcal{J}_1, j \in \mathcal{J}_2$, then the vectors

$$\mathbf{W}^{(1)} = (W_l : l \in \mathcal{J}_2) \in \mathbb{R}^{|\mathcal{J}_2|}, \quad \mathbf{W}^{(2)} = (W_k : k \in \mathcal{J}_1) \in \mathbb{R}^{|\mathcal{J}_1|}$$

are statistically independent.

Distributions related to the normal.

Definition 5 (Chi-square distribution). If $Z_1, Z_2, \dots, Z_k \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$, then the distribution of

$$Q = \sum_{j=1}^k Z_j^2$$

is called the Chi-square distribution with k degrees of freedom, and we denote $Q \sim \chi_k^2$ (in R: `pchisq()`, `qchisq()`, `rchisq()`).

Definition 6 (t -distribution). If $Z \sim \mathcal{N}(0, 1)$, $V \sim \chi_k^2$, are independent random variables, then the distribution of

$$T = \frac{Z}{\sqrt{V/k}}$$

is called the t -distribution with k degrees of freedom, and we denote $T \sim t_k$ (in R: `pt()`, `qt()`, `rt()`).

Definition 7 (*F distribution*). If $V_1 \sim \chi_{k_1}^2$, $V_2 \sim \chi_{k_2}^2$, are independent random variables, the distribution of

$$F = \frac{V_1/k_1}{V_2/k_2}$$

is called the *F-distribution with k_1 and k_2 (numerator and denominator, respectively) degrees of freedom*, and we denote $F \sim F_{k_1, k_2}$.

Proposition. If $Q \sim \chi_k^2$, then $\mathbb{E}Q = k$.

Proof. For $Z_i \sim \mathcal{N}(0, 1)$, iid for $i = 1, \dots, k$, we can write $Q \stackrel{d}{=} \sum_{i=1}^k Z_i^2$, where “ $\stackrel{d}{=}$ ” means “equal in distribution”. Then $\mathbb{E}Q \stackrel{d}{=} \mathbb{E} \sum_{i=1}^k Z_i^2 = \sum_{i=1}^k \mathbb{E}Z_i^2 = \sum_{i=1}^k \mathbb{E}(Z_i) = k$. \square

Proposition. Let $\mathbf{Z} \sim \mathcal{N}_n(\mathbf{0}, \mathbf{I})$, and \mathbf{P} be a square symmetric ($\mathbf{P}^\top = \mathbf{P}$) and idempotent ($\mathbf{P}^2 = \mathbf{P}$) matrix with $\text{rank}(\mathbf{P}) = r$. Then $\|\mathbf{P}\mathbf{Z}\|^2 \sim \chi_r^2$.

Proof. From a previous lemma, since $\mathbb{E}[\mathbf{P}\mathbf{Z}] = \mathbf{P}\mathbb{E}[\mathbf{Z}] = \mathbf{0}$, we have $\mathbb{E}\|\mathbf{P}\mathbf{Z}\|^2 = \text{tr}(\text{cov}[\mathbf{P}\mathbf{Z}]) = \text{tr}(\mathbf{P}\mathbf{I}\mathbf{P}^\top) = \text{tr}(\mathbf{P}) = r$, where the last equality is because \mathbf{P} is similar to a diagonal matrix with r nonzero elements on its diagonal. \square

Inference under the normal linear model.

Recall:

$$\text{The linear model:} \quad \mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \mathbb{E}\boldsymbol{\epsilon} = \mathbf{0}, \quad \text{cov}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}_n$$

We will now make the additional assumption that the error term $\boldsymbol{\epsilon}$ has a *multivariate normal* distribution. In other words, we will assume The normal linear model:

$$\text{The normal linear model:} \quad \mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

The additional normality assumption will enable us to address inferential tasks beyond point estimation, e.g., to construct a confidence interval for a linear combination of $\hat{\boldsymbol{\beta}}$. Indeed, if we assume $\boldsymbol{\epsilon}$ has a multivariate normal distribution, then we can derive exact distributions of $\hat{\boldsymbol{\beta}}$, $\hat{\sigma}^2$, and their joint.

Distribution of $\hat{\boldsymbol{\beta}}$. Recall that, for $\mathbf{A} := (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \in \mathbb{R}^{(p+1) \times n}$, we have

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= \mathbf{A}\mathbf{Y} = \mathbf{A}(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}) \\ &= \mathbf{A}\mathbf{X}\boldsymbol{\beta} + \mathbf{A}\boldsymbol{\epsilon} \\ &= \boldsymbol{\beta} + \mathbf{A}\boldsymbol{\epsilon} \\ &\stackrel{d}{=} \boldsymbol{\beta} + (\sigma \mathbf{A})\mathbf{Z} \end{aligned}$$

where $\mathbf{Z} \sim \mathcal{N}_n(\mathbf{0}, \mathbf{I})$. Hence, by definition, $\hat{\boldsymbol{\beta}}$ has a multivariate normal distribution. We have already calculated the moments of $\hat{\boldsymbol{\beta}}$,

$$\mathbb{E}\hat{\boldsymbol{\beta}} = \boldsymbol{\beta}, \quad \text{cov}(\hat{\boldsymbol{\beta}}) = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1},$$

so in conclusion we have

$$\hat{\boldsymbol{\beta}} \sim \mathcal{N}_{p+1}\left(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}\right)$$

Distribution of $\hat{\sigma}^2$. Recall that $\mathbf{e} = \mathbf{Q}\boldsymbol{\epsilon}$, where \mathbf{Q} is the $n \times n$ projection matrix onto the orthogonal complement of $\text{Im}(\mathbf{X})$. By a previous result, $\|\mathbf{e}\|^2 \sim \sigma^2 \chi_{n-p-1}^2$. This gives

$$\frac{n-p-1}{\sigma^2} \hat{\sigma}^2 \sim \chi_{n-p-1}^2 \iff \frac{\hat{\sigma}^2}{\sigma^2} \sim \frac{\chi_{n-p-1}^2}{n-p-1} \sim t_{n-p-1}.$$

Joint distribution of $\hat{\boldsymbol{\beta}}$ and $\hat{\sigma}^2$. For $\mathbf{A} := (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \in \mathbb{R}^{(p+1) \times n}$, first note that

$$\hat{\boldsymbol{\beta}} = \mathbf{A}\mathbf{Y} = \mathbf{A}(\mathbf{P}_M \mathbf{Y} + \mathbf{P}_{M^\perp} \mathbf{Y}) = \mathbf{P} \mathbf{P} \mathbf{P}_M \mathbf{Y} + \mathbf{A} \mathbf{P}_{M^\perp} \mathbf{Y} = \mathbf{A} \mathbf{P}_M \mathbf{Y}$$

Then,

$$\begin{aligned} \text{cov}(\hat{\boldsymbol{\beta}}, \mathbf{e}) &= \text{cov}(\mathbf{A} \mathbf{P}_M \mathbf{Y}, (\mathbf{I}_n - \mathbf{P}_M) \mathbf{Y}) = \mathbf{A} \mathbf{P}_M \text{cov}(\mathbf{Y}) (\mathbf{I}_n - \mathbf{P}_M)^\top \\ &= \sigma^2 \mathbf{A} \mathbf{P}_M (\mathbf{I}_n - \mathbf{P}_M) = \mathbf{0} \end{aligned} \quad (23)$$

Moreover,

$$\begin{bmatrix} \hat{\boldsymbol{\beta}} \\ \mathbf{e} \end{bmatrix} = \begin{bmatrix} \mathbf{A} \\ \mathbf{I}_n - \mathbf{P}_M \end{bmatrix} \mathbf{Y} \stackrel{d}{=} \begin{bmatrix} \mathbf{A} \\ \mathbf{I}_n - \mathbf{P}_M \end{bmatrix} (\mathbf{X}\boldsymbol{\beta} + \sigma \mathbf{Z}) \quad (24)$$

i.e., $\begin{bmatrix} \hat{\boldsymbol{\beta}} \\ \mathbf{e} \end{bmatrix}$ has a multivariate normal distribution. Together, (23) and (24) imply that $\hat{\boldsymbol{\beta}}$ and $\hat{\sigma}^2$ are statistically independent (because uncorrelated=independent under joint normality).