

Statistics for Data Science: Week 9

Myrto Limnios and Rajita Chandak

Institute of Mathematics – EPFL

`rajita.chandak@epfl.ch`, `myrto.limnios@epfl.ch`



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Linear Algebra Intermezzo- *continued*

Linear Subspaces, Orthogonal Projections, Gaussian Vectors

Definition (Multivariate Gaussian Distribution)

A random vector \mathbf{Y} in \mathbb{R}^d has the multivariate normal distribution if and only if $\beta^\top \mathbf{Y}$ has the univariate normal distribution, $\forall \beta \in \mathbb{R}^d$.

How can we use this definition to determine basic properties?

Recall that the *moment generating function* (MGF) of a random vector \mathbf{W} in \mathbb{R}^d is defined as

$$M_{\mathbf{W}}(\boldsymbol{\theta}) = \mathbb{E}[e^{\boldsymbol{\theta}^\top \mathbf{W}}], \quad \boldsymbol{\theta} \in \mathbb{R}^d,$$

provided the expectation exists. When the MGF exists *it characterises the distribution of the random vector*. Furthermore, two random vectors are independent if and only if their joint MGF is the product of their marginal MGF's.

Most important facts about Gaussian vectors:

- ❶ Moment generating function of $\mathbf{Y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Omega})$:

$$M_{\mathbf{Y}}(\mathbf{u}) = \exp \left(\mathbf{u}^\top \boldsymbol{\mu} + \frac{1}{2} \mathbf{u}^\top \boldsymbol{\Omega} \mathbf{u} \right).$$

- ❷ $\mathbf{Y} \sim \mathcal{N}(\boldsymbol{\mu}_{p \times 1}, \boldsymbol{\Omega}_{p \times p})$ and given $\mathbf{B}_{n \times p}$ and $\boldsymbol{\theta}_{n \times 1}$, then
 $\boldsymbol{\theta} + \mathbf{B}\mathbf{Y} \sim \mathcal{N}(\boldsymbol{\theta} + \mathbf{B}\boldsymbol{\mu}, \mathbf{B}\boldsymbol{\Omega}\mathbf{B}^\top).$

- ❸ $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Omega})$ density, assuming $\boldsymbol{\Omega}$ nonsingular:

$$f_{\mathbf{Y}}(\mathbf{y}) = \frac{1}{(2\pi)^{p/2} |\boldsymbol{\Omega}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \boldsymbol{\mu})^\top \boldsymbol{\Omega}^{-1} (\mathbf{y} - \boldsymbol{\mu}) \right\}.$$

- ❹ Constant density isosurfaces are ellipsoidal
❺ Marginals of Gaussian are Gaussian (converse NOT true).
❻ $\boldsymbol{\Omega}$ diagonal \Leftrightarrow independent coordinates Y_j .
❼ If $\mathbf{Y} \sim \mathcal{N}(\boldsymbol{\mu}_{p \times 1}, \boldsymbol{\Omega}_{p \times p})$,

$$\mathbf{A}\mathbf{Y} \text{ independent of } \mathbf{B}\mathbf{Y} \iff \mathbf{A}\boldsymbol{\Omega}\mathbf{B}^\top = \mathbf{0}.$$

Proposition (Property 1: Moment Generating Function)

The moment generating function of $\mathbf{Y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Omega})$ is

$$M_{\mathbf{Y}}(\mathbf{u}) = \exp\left(\mathbf{u}^\top \boldsymbol{\mu} + \frac{1}{2} \mathbf{u}^\top \boldsymbol{\Omega} \mathbf{u}\right)$$

Proof (*).

Let $\mathbf{u} \in \mathbb{R}^p$ be arbitrary. Then $\mathbf{u}^\top \mathbf{Y}$ is Gaussian with mean $\mathbf{u}^\top \boldsymbol{\mu}$ and variance $\mathbf{u}^\top \boldsymbol{\Omega} \mathbf{u}$. Hence it has moment generating function:

$$M_{\mathbf{u}^\top \mathbf{Y}}(t) = \mathbb{E}\left(e^{t \mathbf{u}^\top \mathbf{Y}}\right) = \exp\left\{t(\mathbf{u}^\top \boldsymbol{\mu}) + \frac{t^2}{2}(\mathbf{u}^\top \boldsymbol{\Omega} \mathbf{u})\right\}.$$

Now take $t = 1$ and observe that

$$M_{\mathbf{u}^\top \mathbf{Y}}(1) = \mathbb{E}\left(e^{\mathbf{u}^\top \mathbf{Y}}\right) = M_{\mathbf{Y}}(\mathbf{u}).$$

Combining the two, we conclude that

$$M_{\mathbf{Y}}(\mathbf{u}) = \exp\left(\mathbf{u}^\top \boldsymbol{\mu} + \frac{1}{2} \mathbf{u}^\top \boldsymbol{\Omega} \mathbf{u}\right), \quad \mathbf{u} \in \mathbb{R}^p.$$

□

Proposition (Property 2: Affine Transformation)

For $\mathbf{Y} \sim \mathcal{N}(\boldsymbol{\mu}_{p \times 1}, \boldsymbol{\Omega}_{p \times p})$ and given $\mathbf{B}_{n \times p}$ and $\boldsymbol{\theta}_{n \times 1}$, we have

$$\boldsymbol{\theta} + \mathbf{B}\mathbf{Y} \sim \mathcal{N}(\boldsymbol{\theta} + \mathbf{B}\boldsymbol{\mu}, \mathbf{B}\boldsymbol{\Omega}\mathbf{B}^\top)$$

Proof (*).

$$\begin{aligned} M_{\boldsymbol{\theta} + \mathbf{B}\mathbf{Y}}(\mathbf{u}) &= \mathbb{E} \left[\exp \{ \mathbf{u}^\top (\boldsymbol{\theta} + \mathbf{B}\mathbf{Y}) \} \right] = \exp \{ \mathbf{u}^\top \boldsymbol{\theta} \} \mathbb{E} \left[\exp \{ (\mathbf{B}^\top \mathbf{u})^\top \mathbf{Y} \} \right] \\ &= \exp \{ \mathbf{u}^\top \boldsymbol{\theta} \} M_{\mathbf{Y}}(\mathbf{B}^\top \mathbf{u}) \\ &= \exp \{ \mathbf{u}^\top \boldsymbol{\theta} \} \exp \left\{ (\mathbf{B}^\top \mathbf{u})^\top \boldsymbol{\mu} + \frac{1}{2} \mathbf{u}^\top \mathbf{B} \boldsymbol{\Omega} \mathbf{B}^\top \mathbf{u} \right\} \\ &= \exp \left\{ \mathbf{u}^\top \boldsymbol{\theta} + \mathbf{u}^\top (\mathbf{B}\boldsymbol{\mu}) + \frac{1}{2} \mathbf{u}^\top \mathbf{B} \boldsymbol{\Omega} \mathbf{B}^\top \mathbf{u} \right\} \\ &= \exp \left\{ \mathbf{u}^\top (\boldsymbol{\theta} + \mathbf{B}\boldsymbol{\mu}) + \frac{1}{2} \mathbf{u}^\top \mathbf{B} \boldsymbol{\Omega} \mathbf{B}^\top \mathbf{u} \right\} \end{aligned}$$

And this last expression is the MGF of a $\mathcal{N}(\boldsymbol{\theta} + \mathbf{B}\boldsymbol{\mu}, \mathbf{B}\boldsymbol{\Omega}\mathbf{B}^\top)$ distribution. □

Proposition (Property 3: Density Function)

Let $\Omega_{p \times p}$ be nonsingular. The density of $\mathcal{N}(\mu_{p \times 1}, \Omega_{p \times p})$ is

$$f_{\mathbf{Y}}(\mathbf{y}) = \frac{1}{(2\pi)^{p/2} |\Omega|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \boldsymbol{\mu})^\top \Omega^{-1} (\mathbf{y} - \boldsymbol{\mu}) \right\}$$

Proof (*).

Let $\mathbf{Z} = (Z_1, \dots, Z_p)^\top$ be a vector of iid $\mathcal{N}(0, 1)$ random variables. Then, because of independence,

(a) the density of \mathbf{Z} is

$$f_{\mathbf{Z}}(\mathbf{z}) = \prod_{i=1}^p f_{Z_i}(z_i) = \prod_{i=1}^p \frac{1}{\sqrt{2\pi}} \exp \left(-\frac{1}{2} z_i^2 \right) = \frac{1}{(2\pi)^{p/2}} \exp \left(-\frac{1}{2} \mathbf{z}^\top \mathbf{z} \right).$$

(b) The MGF of \mathbf{Z} is

$$M_{\mathbf{Z}}(\mathbf{u}) = \mathbb{E} \left\{ \exp \left(\sum_{i=1}^p u_i Z_i \right) \right\} = \prod_{i=1}^p \mathbb{E} \{ \exp(u_i Z_i) \} = \exp(\mathbf{u}^\top \mathbf{u} / 2),$$

which is the MGF of a p -variate $\mathcal{N}(0, \mathbf{I})$ distribution.

$\stackrel{(a)+(b)}{\implies}$ the $\mathcal{N}(0, \mathbf{I})$ density is $f_{\mathbf{Z}}(\mathbf{z}) = \frac{1}{(2\pi)^{p/2}} \exp\left(-\frac{1}{2}\mathbf{z}^\top \mathbf{z}\right)$.

By the spectral theorem, $\mathbf{\Omega}$ admits a square root, $\mathbf{\Omega}^{1/2}$. Furthermore, since $\mathbf{\Omega}$ is non-singular, so is $\mathbf{\Omega}^{1/2}$.

Now observe that from our Property 2, we have $\mathbf{Y} \stackrel{d}{=} \mathbf{\Omega}^{1/2}\mathbf{Z} + \boldsymbol{\mu} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{\Omega})$.

By the change of variables formula,

$$\begin{aligned} f_{\mathbf{Y}}(\mathbf{y}) &= f_{\mathbf{\Omega}^{1/2}\mathbf{Z} + \boldsymbol{\mu}}(\mathbf{y}) \\ &= |\mathbf{\Omega}^{-1/2}| f_{\mathbf{Z}}\{\mathbf{\Omega}^{-1/2}(\mathbf{y} - \boldsymbol{\mu})\} \\ &= \frac{1}{(2\pi)^{p/2} |\mathbf{\Omega}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^\top \mathbf{\Omega}^{-1}(\mathbf{y} - \boldsymbol{\mu})\right\}. \end{aligned}$$

[Recall that to obtain the density of $\mathbf{W} = g(\mathbf{X})$ at \mathbf{w} , we need to evaluate $f_{\mathbf{X}}$ at $g^{-1}(\mathbf{w})$ but also multiply by the Jacobian determinant of g^{-1} at \mathbf{w} .]



Proposition (Property 4: Isosurfaces)

The isosurfaces of a $\mathcal{N}(\boldsymbol{\mu}_{p \times 1}, \boldsymbol{\Omega}_{p \times p})$ are $(p - 1)$ -dimensional ellipsoids centred at $\boldsymbol{\mu}$, with principal axes given by the eigenvectors of $\boldsymbol{\Omega}$ and with anisotropies given by the ratios of the square roots of the corresponding eigenvalues of $\boldsymbol{\Omega}$.

Proof (*).

Exercise: Use Property 3, and the spectral theorem. □

Proposition (Property 5: Coordinate Distributions)

Let $\mathbf{Y} = (Y_1, \dots, Y_p)^\top \sim \mathcal{N}(\boldsymbol{\mu}_{p \times 1}, \boldsymbol{\Omega}_{p \times p})$. Then $Y_j \sim \mathcal{N}(\mu_j, \Omega_{jj})$.

Proof (*).

Observe that $Y_j = (0, 0, \dots, \underbrace{1}_{j\text{th position}}, \dots, 0, 0) \mathbf{Y}$ and use Property 2. □

Proposition (Property 6: Diagonal $\Omega \iff$ Independence)

Let $\mathbf{Y} = (Y_1, \dots, Y_p)^\top \sim \mathcal{N}(\boldsymbol{\mu}_{p \times 1}, \boldsymbol{\Omega}_{p \times p})$. Then the Y_i are mutually independent if and only if $\boldsymbol{\Omega}$ is diagonal.

Proof (*).

Suppose that the Y_j are independent. Property 5 yields $Y_j \sim \mathcal{N}(\mu_j, \sigma_j^2)$ for some $\sigma_j > 0$. Thus the density of \mathbf{Y} is

$$\begin{aligned} f_{\mathbf{Y}}(\mathbf{y}) &= \prod_{j=1}^p f_{Y_j}(y_j) = \prod_{j=1}^p \frac{1}{\sigma_j \sqrt{2\pi}} \exp \left\{ -\frac{1}{2} \frac{(y_j - \mu_j)^2}{\sigma_j^2} \right\} \\ &= \frac{1}{(2\pi)^{p/2} |\text{diag}(\sigma_1^2, \dots, \sigma_p^2)|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \boldsymbol{\mu})^\top \text{diag}(\sigma_1^{-2}, \dots, \sigma_p^{-2}) (\mathbf{y} - \boldsymbol{\mu}) \right\}. \end{aligned}$$

Hence $\mathbf{Y} \sim \mathcal{N}\{\boldsymbol{\mu}, \text{diag}(\sigma_1^2, \dots, \sigma_p^2)\}$, i.e. the covariance $\boldsymbol{\Omega}$ is diagonal.

Conversely, assume $\boldsymbol{\Omega}$ is diagonal, say $\boldsymbol{\Omega} = \text{diag}(\sigma_1^2, \dots, \sigma_p^2)$. Then we can reverse the steps of the first part to see that the joint density $f_{\mathbf{Y}}(\mathbf{y})$ can be written as a product of the marginal densities $f_{Y_j}(y_j)$, thus proving independence.



Proposition (Property 7: \mathbf{AY}, \mathbf{BY} indep $\iff \mathbf{A}\Omega\mathbf{B}^\top = 0$)

If $\mathbf{Y} \sim \mathcal{N}(\boldsymbol{\mu}_{p \times 1}, \Omega_{p \times p})$, and $\mathbf{A}_{m \times p}$, $\mathbf{B}_{d \times p}$ be real matrices. Then,

$$\mathbf{AY} \text{ independent of } \mathbf{BY} \iff \mathbf{A}\Omega\mathbf{B}^\top = 0.$$

Proof (*). [wlog assuming $\boldsymbol{\mu} = 0$ (simplifies the algebra)]

First assume $\mathbf{A}\Omega\mathbf{B}^\top = 0$. Let $\mathbf{W}_{(m+d) \times 1} = \begin{pmatrix} \mathbf{AY} \\ \mathbf{BY} \end{pmatrix}$ and $\boldsymbol{\theta}_{(m+d) \times 1} = \begin{pmatrix} \mathbf{u}_{m \times 1} \\ \mathbf{v}_{d \times 1} \end{pmatrix}$.

$$\begin{aligned} M_{\mathbf{W}}(\boldsymbol{\theta}) &= \mathbb{E}[\exp\{\mathbf{W}^\top \boldsymbol{\theta}\}] = \mathbb{E}[\exp\{\mathbf{Y}^\top \mathbf{A}^\top \mathbf{u} + \mathbf{Y}^\top \mathbf{B}^\top \mathbf{v}\}] \\ &= \mathbb{E}[\exp\{\mathbf{Y}^\top (\mathbf{A}^\top \mathbf{u} + \mathbf{B}^\top \mathbf{v})\}] = M_{\mathbf{Y}}(\mathbf{A}^\top \mathbf{u} + \mathbf{B}^\top \mathbf{v}) \\ &= \exp\left\{\frac{1}{2}(\mathbf{A}^\top \mathbf{u} + \mathbf{B}^\top \mathbf{v})^\top \Omega (\mathbf{A}^\top \mathbf{u} + \mathbf{B}^\top \mathbf{v})\right\} \end{aligned}$$

$$= \exp\left\{\frac{1}{2}\left(\mathbf{u}^\top \mathbf{A}\Omega\mathbf{A}^\top \mathbf{u} + \mathbf{v}^\top \mathbf{B}\Omega\mathbf{B}^\top \mathbf{v} + \underbrace{\mathbf{u}^\top \mathbf{A}\Omega\mathbf{B}^\top \mathbf{v}}_{=0} + \underbrace{\mathbf{v}^\top \mathbf{B}\Omega\mathbf{A}^\top \mathbf{u}}_{=0}\right)\right\}$$

$$= M_{\mathbf{AY}}(\mathbf{u})M_{\mathbf{BY}}(\mathbf{v}) \quad (\text{joint MGF} = \text{product of marginal MGFs, thus independence})$$

For the converse, assume that $\mathbf{A}\mathbf{Y}$ and $\mathbf{B}\mathbf{Y}$ are independent. Then, $\forall \mathbf{u}, \mathbf{v}$,

$$M_{\mathbf{W}}(\boldsymbol{\theta}) = M_{\mathbf{A}\mathbf{Y}}(\mathbf{u})M_{\mathbf{B}\mathbf{Y}}(\mathbf{v}), \quad \forall \mathbf{u}, \mathbf{v},$$

$$\implies \exp \left\{ \frac{1}{2} (\mathbf{u}^\top \mathbf{A}\boldsymbol{\Omega}\mathbf{A}^\top \mathbf{u} + \mathbf{v}^\top \mathbf{B}\boldsymbol{\Omega}\mathbf{B}^\top \mathbf{v} + \mathbf{u}^\top \mathbf{A}\boldsymbol{\Omega}\mathbf{B}^\top \mathbf{v} + \mathbf{v}^\top \mathbf{B}\boldsymbol{\Omega}\mathbf{A}^\top \mathbf{u}) \right\}$$

$$= \exp \left\{ \frac{1}{2} \mathbf{u}^\top \mathbf{A}\boldsymbol{\Omega}\mathbf{A}^\top \mathbf{u} \right\} \exp \left\{ \frac{1}{2} \mathbf{v}^\top \mathbf{B}\boldsymbol{\Omega}\mathbf{B}^\top \mathbf{v} \right\}$$

$$\implies \exp \left\{ \frac{1}{2} \times 2 \mathbf{v}^\top \mathbf{A}\boldsymbol{\Omega}\mathbf{B}^\top \mathbf{u} \right\} = 1$$

$$\implies \mathbf{v}^\top \mathbf{A}\boldsymbol{\Omega}\mathbf{B}^\top \mathbf{u} = 0, \quad \forall \mathbf{u}, \mathbf{v},$$

$$\implies \mathbf{A}\boldsymbol{\Omega}\mathbf{B}^\top = \mathbf{0}.$$



Reminder:

Definition (χ^2 distribution)

Let $\mathbf{Z} \sim \mathcal{N}(0, \mathbf{I}_{p \times p})$. Then $\|\mathbf{Z}\|^2 = \sum_{j=1}^p Z_j^2$ is said to have the chi-square (χ^2) distribution with p degrees of freedom; we write $\|\mathbf{Z}\|^2 \sim \chi_p^2$.

[Thus, χ_p^2 is the distribution of the sum of squares of p real independent standard Gaussian random variates.]

Definition (F distribution)

Let $V \sim \chi_p^2$ and $W \sim \chi_q^2$ be independent random variables. Then $(V/p)/(W/q)$ is said to have the F distribution with p and q degrees of freedom; we write $(V/p)/(W/q) \sim F_{p,q}$.

Proposition (Gaussian Quadratic Forms)

- ❶ If $\mathbf{Z} \sim \mathcal{N}(0_{p \times 1}, \mathbf{I}_{p \times p})$ and \mathbf{H} is a projection of rank $r \leq p$,

$$\mathbf{Z}^\top \mathbf{H} \mathbf{Z} \sim \chi_r^2.$$

- ❷ $\mathbf{Y} \sim \mathcal{N}(\boldsymbol{\mu}_{p \times 1}, \boldsymbol{\Omega}_{p \times p})$ with $\boldsymbol{\Omega}$ nonsingular \implies

$$(\mathbf{Y} - \boldsymbol{\mu})^\top \boldsymbol{\Omega}^{-1} (\mathbf{Y} - \boldsymbol{\mu}) \sim \chi_p^2.$$

Gaussian Linear Regression: Likelihood and Geometry

General formulation:

$$Y_i | x_i \stackrel{\text{ind}}{\sim} \text{Distribution}\{g(x_i)\}, \quad i = 1, \dots, n.$$

Simple Normal Linear Regression:

$$\begin{cases} \text{Distribution} = \mathcal{N}\{g(x), \sigma^2\} \\ g(x) = \beta_0 + \beta_1 x \end{cases}$$

Resulting Model:

$$Y_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\beta_0 + \beta_1 x_i, \sigma^2)$$

$$\Updownarrow$$

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad \varepsilon_i \stackrel{\text{ind}}{\sim} \mathcal{N}(0, \sigma^2)$$

Jargon: Y is *response variable* and x is *explanatory variable* (or *covariate*)

Linearity: Linearity is in the *parameters*, not the *explanatory variable*.

Example: Flexibility in what we define as explanatory:

$$Y_j = \beta_0 + \beta_1 \underbrace{\sin(x_j)}_{x_j^*} + \varepsilon_j, \quad \varepsilon_j \stackrel{iid}{\sim} \text{Normal}(0, \sigma^2).$$

Example: Sometimes a transformation may be required:

$$Y_j = \beta_0 e^{\beta_1 x_j} \eta_j, \quad \eta_j \stackrel{iid}{\sim} \text{Lognormal}$$

$$\log(\cdot) \downarrow \quad \uparrow \exp(\cdot)$$

$$\log Y_j = \log \beta_0 + \beta_1 x_j + \log \eta_j, \quad \log \eta_j \stackrel{iid}{\sim} \text{Normal}$$

Data Structure:

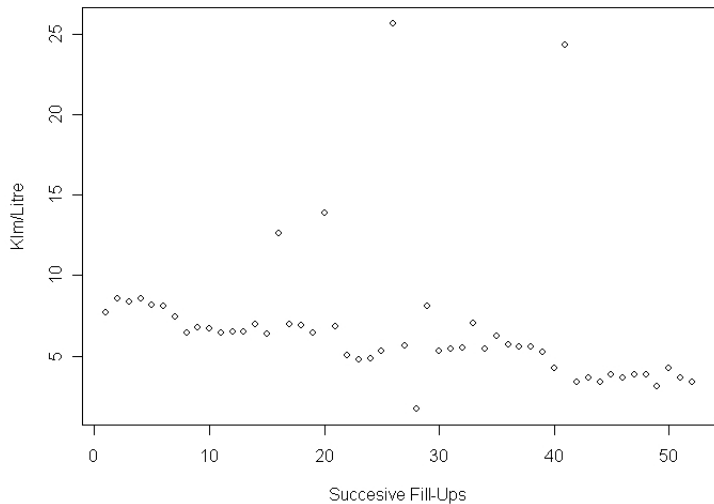
For $i = 1, \dots, n$, pairs

$$(x_i, y_i) \xrightarrow{\text{red arrow}} \begin{cases} x_i \text{ fixed values of } x \\ Y_i \text{ random output } Y_i \text{ when input is } x_i \end{cases}$$

Example: Professor's Van

Fillup	Km/L
1	7.72
2	8.54
3	8.35
4	8.55
5	8.16
6	8.12
7	7.46
8	6.43
9	6.74
10	6.72

Example: Professor's Van



Instead of $x_i \in \mathbb{R}$ could have $\mathbf{x}_i^\top \in \mathbb{R}^q$:

$$Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_q x_{iq} + \varepsilon_i, \quad \varepsilon_i \stackrel{\text{ind}}{\sim} \mathcal{N}(0, \sigma^2).$$

Letting $p = q + 1$, this can be summarised via matrix notation:

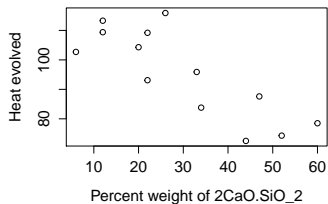
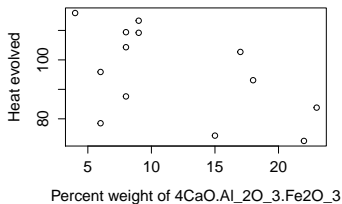
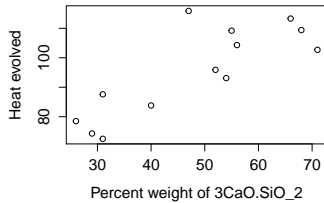
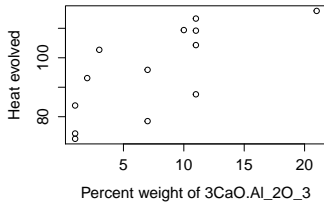
$$\underbrace{\begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix}}_{\mathbf{Y}} = \underbrace{\begin{pmatrix} 1 & x_{11} & \dots & x_{1q} \\ 1 & x_{21} & & x_{2q} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \dots & x_{nq} \end{pmatrix}}_{\mathbf{X}} \underbrace{\begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_q \end{pmatrix}}_{\boldsymbol{\beta}} + \underbrace{\begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}}_{\boldsymbol{\varepsilon}}$$

$$\implies \underset{n \times 1}{\mathbf{Y}} = \underset{n \times p}{\mathbf{X}} \underset{p \times 1}{\boldsymbol{\beta}} + \underset{n \times 1}{\boldsymbol{\varepsilon}}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}_n(0, \sigma^2 \mathbf{I}_{n \times n})$$

\mathbf{X} is called the **design matrix**.

Example: Cement Heat Evolution

Case	$3\text{CaO} \cdot \text{Al}_2\text{O}_3$	$3\text{CaO} \cdot \text{SiO}_2$	$4\text{CaO} \cdot \text{Al}_2\text{O}_3 \cdot \text{Fe}_2\text{O}_3$	$2\text{CaO} \cdot \text{SiO}_2$	Heat
1	7.00	26.00	6.00	60.00	78.50
2	1.00	29.00	15.00	52.00	74.30
3	11.00	56.00	8.00	20.00	104.30
4	11.00	31.00	8.00	47.00	87.60
5	7.00	52.00	6.00	33.00	95.90
6	11.00	55.00	9.00	22.00	109.20
7	3.00	71.00	17.00	6.00	102.70
8	1.00	31.00	22.00	44.00	72.50
9	2.00	54.00	18.00	22.00	93.10
10	21.00	47.00	4.00	26.00	115.90
11	1.00	40.00	23.00	34.00	83.80
12	11.00	66.00	9.00	12.00	113.30
13	10.00	68.00	8.00	12.00	109.40



Model is:

$$Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_q x_{iq} + \varepsilon_i, \quad \varepsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$$

$$\Updownarrow$$

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}_n(0, \sigma^2 \mathbf{I}_{n \times n})$$

Observe: $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$ for given fixed design matrix \mathbf{X} , i.e.:

$$(Y_1, x_{11}, \dots, x_{1q}), \dots, (Y_i, x_{i1}, \dots, x_{iq}), \dots, (Y_n, x_{n1}, \dots, x_{nq})$$

Likelihood and Loglikelihood

$$L(\boldsymbol{\beta}, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \right\}$$

$$\ell(\boldsymbol{\beta}, \sigma^2) = -\frac{1}{2} \left\{ n \log 2\pi + n \log \sigma^2 + \frac{1}{\sigma^2} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \right\}$$

Whatever the value of σ , the log-likelihood is maximised when $(\mathbf{Y} - \mathbf{X}\beta)^\top (\mathbf{Y} - \mathbf{X}\beta)$ is minimised. Hence, the MLE of β is:

$$\hat{\beta} = \arg \max_{\beta} \{ -(\mathbf{Y} - \mathbf{X}\beta)^\top (\mathbf{Y} - \mathbf{X}\beta) \} = \arg \min_{\beta} (\mathbf{Y} - \mathbf{X}\beta)^\top (-\mathbf{X}\beta)$$

Obtain minimum by solving:

$$0 = \frac{\partial}{\partial \beta} (\mathbf{Y} - \mathbf{X}\beta)^\top (\mathbf{Y} - \mathbf{X}\beta)$$

$$0 = \frac{\partial (\mathbf{Y} - \mathbf{X}\beta)}{\partial \beta} \frac{\partial (\mathbf{Y} - \mathbf{X}\beta)^\top (\mathbf{Y} - \mathbf{X}\beta)}{\partial (\mathbf{Y} - \mathbf{X}\beta)} \quad (\text{chain rule})$$

$$0 = \mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\beta) \quad (\text{normal equations})$$

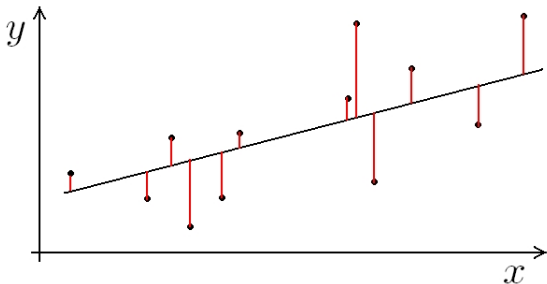
$$\mathbf{X}^\top \mathbf{X} \beta = \mathbf{X}^\top \mathbf{Y}$$

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} \quad (\text{if } \mathbf{X} \text{ has rank } p)$$

The MLE $\hat{\beta}$ is called the **least squares estimator** because it is a result of minimising

$$(\mathbf{Y} - \mathbf{X}\beta)^\top (\mathbf{Y} - \mathbf{X}\beta) = \underbrace{\sum_{i=1}^n (Y_i - \beta_0 - \beta_1 x_{i1} - \beta_2 x_{i2} - \cdots - \beta_q x_{iq})^2}_{\text{sum of squares}}.$$

Thus we are trying to find the β that gives the hyperplane with minimum sum of squared vertical distances from our observations.



Residuals: $\mathbf{e} = \mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}$, so that $\mathbf{e} = (\mathbf{e}_1, \dots, \mathbf{e}_n)^\top$, with

$$e_i = Y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \hat{\beta}_2 x_{i2} - \dots - \hat{\beta}_q x_{iq}$$

“Regression Line” is such that $\sum e_i^2$ is minimised over all $\boldsymbol{\beta}$.

Fitted Values: $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$, so that $\hat{\mathbf{Y}} = (\hat{Y}_1, \dots, \hat{Y}_n)^\top$, with

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_q x_{iq}$$

Since the MLE of β is $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$ for all values of σ^2 , we have

$$\begin{aligned}\hat{\sigma}^2 &= \arg \max_{\sigma^2} \left\{ \max_{\beta} \ell(\beta, \sigma^2) \right\} \\ &= \arg \max_{\sigma^2} \ell(\hat{\beta}, \sigma^2) \\ &= \arg \max_{\sigma^2} -\frac{1}{2} \left\{ n \log \sigma^2 + \frac{1}{\sigma^2} (\mathbf{Y} - \mathbf{X}\hat{\beta})^\top (\mathbf{Y} - \mathbf{X}\hat{\beta}) \right\}.\end{aligned}$$

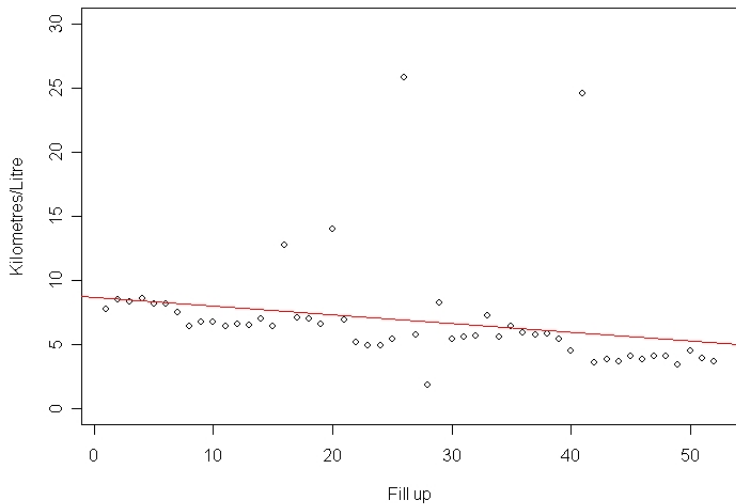
Differentiating and setting equal to zero yields

$$\hat{\sigma}^2 = \frac{1}{n} (\mathbf{Y} - \mathbf{X}\hat{\beta})^\top (\mathbf{Y} - \mathbf{X}\hat{\beta}).$$

We will soon see that a better (unbiased) estimator is

$$S^2 = \frac{1}{n-p} (\mathbf{Y} - \mathbf{X}\hat{\beta})^\top (\mathbf{Y} - \mathbf{X}\hat{\beta}).$$

Example: Professor's Van



$$\hat{\beta}_0 = 8.6 \quad \hat{\beta}_1 = -0.068 \quad S^2 = 17.4$$

There are two dual geometrical viewpoints that one may adopt:

$$\begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1q} \\ 1 & x_{21} & x_{22} & & x_{2q} \\ \vdots & \vdots & & \vdots & \\ 1 & x_{(n-1)1} & x_{(n-1)2} & \cdots & x_{(n-1)q} \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nq} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_q \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

- Row geometry: focus on the n OBSERVATIONS
- Column geometry: focus on the p covariates

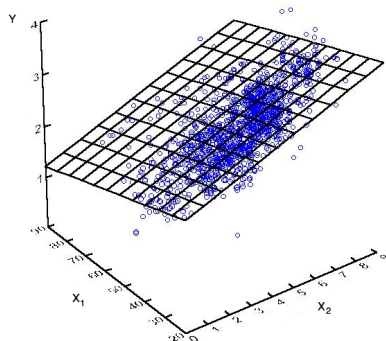
Both are useful, usually for different things:

- Row geometry useful for exploratory analysis.
- Column geometry useful for theoretical analysis.

Both geometries give useful, but different, intuitive interpretations of the least squares estimators.

Corresponds to the “scatterplot geometry” – (data space)

- n points in \mathbb{R}^p
- each corresponds to an observation
- least squares parameters give parametric equation for a hyperplane
- hyperplane has property that it minimizes the sum of squared vertical distances of observations from the plane itself over all possible hyperplanes



- Fitted values are vertical projections (NOT orthogonal projections!) of observations onto plane, residuals are signed vertical distances of observations from plane.

Adopt the dual perspective:

- Consider the entire vector \mathbf{Y} as a **single** point living in \mathbb{R}^n
- Then consider each variable (column of \mathbf{X}) as a point also in \mathbb{R}^n

What is the interpretation of the p -dimensional vector $\hat{\beta}$, and the n -dimensional vectors $\hat{\mathbf{Y}}$ and \mathbf{e} in this dual space?

Turns out there is another important plane here: the plane spanned by the variable vectors (the column vectors of \mathbf{X}).

Recall that this is the *column space* of \mathbf{X} , denoted by $\mathcal{M}(\mathbf{X})$.

Recall: $\underbrace{\mathcal{M}(\mathbf{X})}_{\text{Column Space}} := \{\mathbf{X}\boldsymbol{\gamma} : \boldsymbol{\gamma} \in \mathbb{R}^p\}$

Q: What does $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \varepsilon$ mean?

A: \mathbf{Y} is [some element of $\mathcal{M}(\mathbf{X})$] + [Gaussian disturbance].

Any realisation of \mathbf{Y} will lie outside $\mathcal{M}(\mathbf{X})$ (almost surely). MLE estimates $\boldsymbol{\beta}$ by minimising

$$(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) = \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2$$

Thus we search for a $\boldsymbol{\beta}$ giving the element of $\mathcal{M}(\mathbf{X})$ with the minimum distance from \mathbf{Y} .

Hence $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$ is the projection of \mathbf{Y} onto $\mathcal{M}(\mathbf{X})$:

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}} := \underbrace{\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top}_{\mathbf{H}} \mathbf{Y} = \mathbf{H}\mathbf{Y}.$$

\mathbf{H} is the **hat matrix** (because it puts a hat on \mathbf{Y} !)

Leads to geometric derivation of the MLE of β :

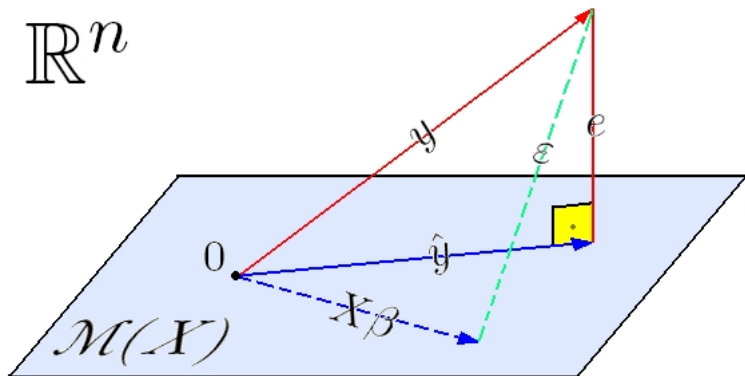
- Choose $\hat{\beta}$ to minimise $(\mathbf{Y} - \mathbf{X}\beta)^\top (\mathbf{Y} - \mathbf{X}\beta) = \|\mathbf{Y} - \mathbf{X}\beta\|^2$, so

$$\hat{\beta} = \arg \min \|\mathbf{Y} - \mathbf{X}\beta\|^2.$$

- $\min_{\beta \in \mathbb{R}^p} \|\mathbf{Y} - \mathbf{X}\beta\|^2 = \min_{\gamma \in \mathcal{M}(\mathbf{X})} \|\mathbf{Y} - \gamma\|^2$
- But the unique γ that yields $\min_{\gamma \in \mathcal{M}(\mathbf{X})} \|\mathbf{Y} - \gamma\|^2$ is $\gamma = \mathbf{P}\mathbf{Y}$.
- Here \mathbf{P} is the projection onto the column space of \mathbf{X} , $\mathcal{M}(\mathbf{X})$.
- Since \mathbf{X} is of full rank, $\mathbf{P} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$.
- So $\gamma = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$
- $\hat{\beta}$ will now be the unique (since \mathbf{X} non-singular) vector of coordinates of γ with respect to the basis of columns of \mathbf{X} .
- So

$$\mathbf{X}\hat{\beta} = \gamma = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y},$$

which implies that $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$



Important facts that will repeatedly be made use of:

- 1 $\mathbf{e} = (\mathbf{I} - \mathbf{H})\mathbf{Y} = (\mathbf{I} - \mathbf{H})\boldsymbol{\varepsilon}$.
- 2 $\hat{\mathbf{Y}}$ and \mathbf{e} are orthogonal, i.e. $\hat{\mathbf{Y}}^\top \mathbf{e} = 0$
- 3 Pythagoras: $\mathbf{Y}^\top \mathbf{Y} = \hat{\mathbf{Y}}^\top \hat{\mathbf{Y}} + \mathbf{e}^\top \mathbf{e} = \mathbf{Y}^\top \mathbf{H} \mathbf{Y} + \boldsymbol{\varepsilon}^\top (\mathbf{I} - \mathbf{H}) \boldsymbol{\varepsilon}$

Derivation:

- 1 $\mathbf{e} = \mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{Y} - \mathbf{H}\mathbf{Y} = (\mathbf{I} - \mathbf{H})\mathbf{Y} = (\mathbf{I} - \mathbf{H})(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) = (\mathbf{I} - \mathbf{H})\mathbf{X}\boldsymbol{\beta} + (\mathbf{I} - \mathbf{H})\boldsymbol{\varepsilon} = (\mathbf{I} - \mathbf{H})\boldsymbol{\varepsilon}$
- 2 $\mathbf{e} = \mathbf{Y} - \hat{\mathbf{Y}} = (\mathbf{I} - \mathbf{H})\mathbf{Y} \implies \hat{\mathbf{Y}}^\top \mathbf{e} = \mathbf{Y}^\top \mathbf{H}^\top (\mathbf{I} - \mathbf{H})\mathbf{Y} = 0$
- 3 $\mathbf{Y}^\top \mathbf{Y} = (\mathbf{H}\mathbf{Y} + (\mathbf{I} - \mathbf{H})\mathbf{Y})^\top (\mathbf{H}\mathbf{Y} + (\mathbf{I} - \mathbf{H})\mathbf{Y}) = \hat{\mathbf{Y}}^\top \hat{\mathbf{Y}} + \mathbf{e}^\top \mathbf{e} + \underbrace{2\mathbf{Y}^\top \mathbf{H} (\mathbf{I} - \mathbf{H}) \mathbf{Y}}_{=0}$.

Could also assume slightly different model:

$$Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_q x_{iq} + \frac{\varepsilon_i}{\sqrt{w_i}}, \quad \varepsilon_i \stackrel{\text{ind}}{\sim} \mathcal{N}(0, \sigma^2), \quad w_i > 0$$

\Updownarrow

$$Y_i \stackrel{\text{ind}}{\sim} \mathcal{N} \left(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_q x_{iq}, \frac{\sigma^2}{w_i} \right).$$

With the w_j *known* weights (example: each Y_j is an average of w_j measurements).

Arises often in practice (e.g., in sample surveys), but also arises in theory (will see in GLM).

Transformation:

$$\mathbf{Y}^* = \mathbf{W}^{1/2} \mathbf{Y}, \quad \mathbf{X}^* = \mathbf{W}^{1/2} \mathbf{X}$$

with

$$\mathbf{W}_{n \times n} = \text{diag}(w_1, \dots, w_n)$$

Leads to usual scenario. In this notation we obtain:

$$\begin{aligned} \hat{\beta} &= [(\mathbf{X}^*)^\top \mathbf{X}^*]^{-1} (\mathbf{X}^*)^\top \mathbf{Y}^* \\ &= (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W} \mathbf{Y} \end{aligned}$$

Similarly:

$$S^2 = \frac{1}{n-p} \mathbf{Y}^\top [\mathbf{W} - \mathbf{W} \mathbf{X} (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W}] \mathbf{Y}$$