

# Statistics for Data Science: Week 8

Myrto Limnios and Rajita Chandak

Institute of Mathematics – EPFL

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



ÉCOLE POLYTECHNIQUE  
FÉDÉRALE DE LAUSANNE

# Regression

In the beginning we distinguished between:

- ① **Marginal Inference.** Here  $(Y_1, \dots, Y_n)^\top$  has i.i.d. entries each from the same distribution  $F(y; \theta)$  with the same parameter  $\theta$ .
  - In other words, all observations were obtained under identical experimental conditions, and thus depend in the same way on the same unknown  $\theta$ .
- ② **Regression.** Here  $(Y_1, \dots, Y_n)^\top$  has independent entries, each with distribution  $F(y; \theta_i)$  of the same family but with different parameters.
  - Each observation was generated under slightly different experimental conditions. They depend in a similar way on different  $\theta_i$ .
  - These  $\theta_i$  correspond to different experimental conditions, say  $x_i$ .
  - Each  $x_i$  is called a covariate/feature, and is an input that the experimenter can vary. They are known. The index  $i$  reminds us that it corresponds to the  $i$ th observation  $Y_i$ .
  - Usually  $\theta_i$  is postulated to have a special relationship to  $x_i$ , for example  $\theta_i = \exp\{\alpha + \beta x_i\}$ , for  $(\alpha, \beta)$  unknown parameters.
  - The point here is to understand the effect of varying the covariate/feature on the distribution of the observable. (✓)

# What is a Regression Model?

Statistical model for:

$Y$  (random output) <sup>whose law is influenced by</sup>  $x$  (non-random input)

Aim: understand the effect of  $x$  on the distribution of random variable  $Y$

General formulation<sup>1</sup>:

$$\underline{Y_i} \overset{\text{independent}}{\sim} \text{Distribution} \underbrace{\{g(x_i)\}}_{=\theta_i}, \quad \underline{i = 1, \dots, n.}$$

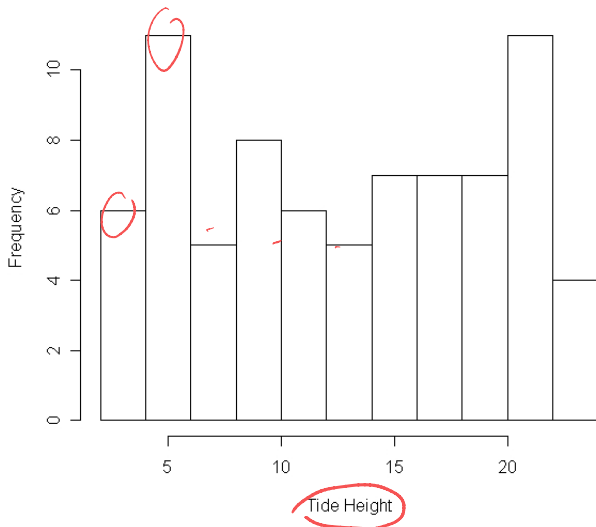
*Handwritten notes: A red question mark is above the tilde symbol. A red arrow points from the tilde to the distribution. A red arrow points from the distribution to the index i. A red arrow points from the index i to the expression i = 1, ..., n.*

**Statistical Problem:** Estimate (learn)  $g(\cdot)$  from data  $\{(x_i, Y_i)\}_{i=1}^n$ . Use for:

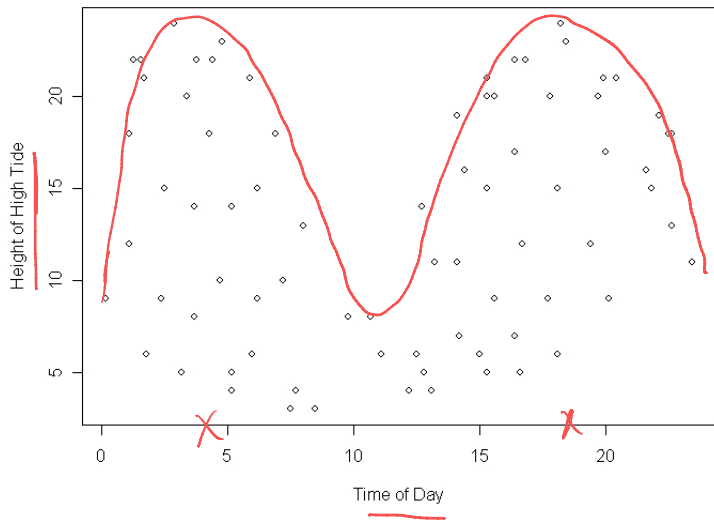
- Inference
- Prediction
- Data compression (parsimonious representations)

<sup>1</sup>Sometimes we write  $\underline{Y_i | x_i} \overset{\text{independent}}{\sim} \text{Distribution} \{g(x_i) = \theta_i\}$  to highlight that the distribution of  $Y$  depends on  $x$ , but without meaning that  $(X, Y)$  are jointly random; such an assumption is unnecessary (e.g., in a designed experiment we choose values for  $x$ ).

## Example: How to model the height of Honolulu tides throughout the day - Histogram



## Example: Height of Honolulu tides as function of the time of day



A **bewildering variety** of models can be captured by the general specification

$$Y_i \overset{\text{independent}}{\sim} \text{Distribution} \{ \underbrace{g(x_i)}_{=\theta_i} \}, \quad i = 1, \dots, n.$$

$x_i$  can be:

(Gaussian)

students

grades

$Y = \text{salary}$

$X = \{ \text{age, sex, } \text{cat.} \}$

geographical region,

cat level of studies, personal wealth, continuous

deterministic.

- continuous, discrete, categorical, vector ...
- arrive randomly, or be chosen by experimenter, or both
- however  $x$  arises, we treat it as constant in the analysis

Distribution can be:

- Gaussian, Laplace, Bernoulli, Poisson, gamma, general exponential family, ...

Function  $g(\cdot)$  can be:

- $g(x) = \beta_0 + \beta_1 x$ ,  $g(x) = \sum_{k=-K}^K \beta_k e^{-ikx}$ , cubic spline, neural net...

Table: A coarse classification of regression models we will consider

Distribution / Function $g$	$g(x_i^\top) = x_i^\top \beta$	$g$ nonparametric
Gaussian	Linear Regression	Smoothing
Exponential Family	GLM	GAM

GLM: Generalized Linear Model and GAM: Generalized Additive Model

We start with a very standard model: Linear Regression with  $Y|x$  being Gaussian.



- $Y, x \in \mathbb{R}, g(x) = \beta_0 + \beta_1 x$

$$Y | x \sim \mathcal{N}(\beta_0 + \beta_1 x, \sigma^2)$$

$\Updownarrow$

$$Y = \beta_0 + \beta_1 x + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

The second version is useful for mathematical work, but is puzzling statistically, since we don't observe  $\epsilon$ .

- Also, covariate could be vector ( $Y, \beta_0 \in \mathbb{R}, \mathbf{x} \in \mathbb{R}^D, \beta \in \mathbb{R}^D$ ):

$$Y | \mathbf{x} \sim \mathcal{N}(\beta_0 + \beta^\top \mathbf{x}, \sigma^2)$$

$\Updownarrow$

$$Y = \beta_0 + \beta^\top \mathbf{x} + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$\tilde{\beta} = (\beta_0, \beta)$$

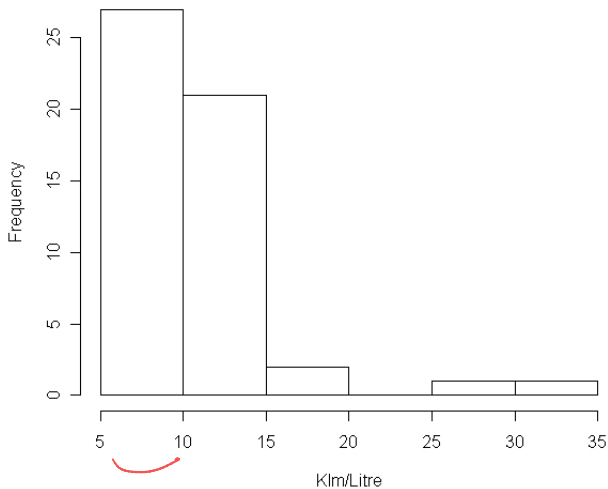
$$\mathbf{x} = (1, \mathbf{x})$$

$$Y | \mathbf{x} \sim \mathcal{N}(\tilde{\beta}^\top \mathbf{x}, \sigma^2)$$

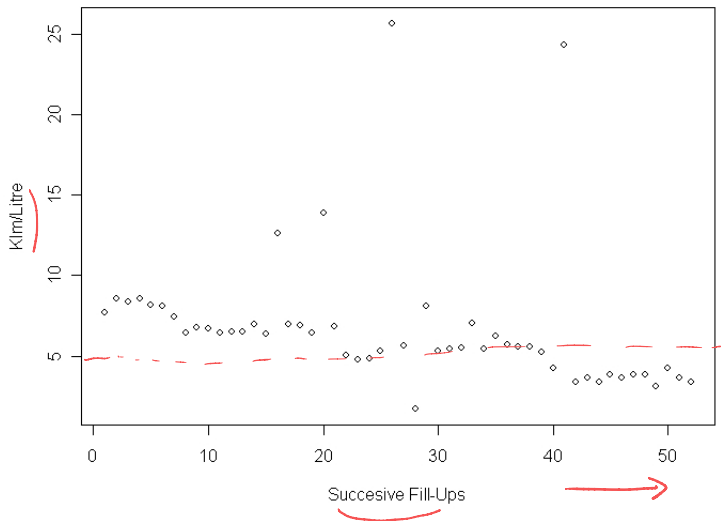
## Example: How to model my van's consumption of gas



## Example: Histogram of consumption of gas (km/L)



## Example: Gas consumption as function of successive fill-ups



Start from **Gaussian linear regression** then gradually generalise ...

Obviously: important features of Gaussian linear model are

- Gaussian distribution
- Linearity

These two **combine well** and give **geometric insights** to solve the estimation problem. Thus we need to revise some probabilistic linear algebra...

- Subspaces and projection matrices .
- Multivariate Gaussian Distribution .
- Optimal dimension reduction
- Random quadratic forms

# Linear Algebra Intermezzo

Linear Subspaces, Orthogonal Projections, Gaussian Vectors

$$Q = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 2 & 0 \end{pmatrix} \quad n=3, p=2 \quad \mathcal{M}(Q) = \left\{ \begin{pmatrix} 1 \\ 0 \\ 2 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \right\}$$

$$\rightarrow a_1 \begin{pmatrix} 1 \\ 0 \\ 2 \end{pmatrix} + a_2 \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} a_1 \\ a_2 \\ 2a_1 \end{pmatrix} \quad \mathcal{M}(Q) = \{ y = (y_1, y_2, y_3), y_3 = 2y_1 \}$$

If  $Q$  is an  $n \times p$  real matrix, we define the column space (or range) of  $Q$  to be the set spanned by its columns:

$$\mathcal{M}(Q) = \{ \underline{y} \in \mathbb{R}^n : \exists \underline{\beta} \in \mathbb{R}^p, \underline{y} = Q\underline{\beta} \}.$$

- Recall that  $\mathcal{M}(Q)$  is a subspace of  $\mathbb{R}^n$ .
- The columns of  $Q$  provide a coordinate system for the subspace  $\mathcal{M}(Q)$
- If  $Q$  is of full column rank ( $p$ ), then the coordinates  $\beta$  corresponding to a  $y \in \mathcal{M}(Q)$  are unique.
- Allows interpretation of system of linear equations

$$Q\underline{\beta} = \underline{y}.$$

[existence of solution  $\leftrightarrow$  is  $\underline{y}$  an element of  $\mathcal{M}(Q)$ ?

[uniqueness of solution  $\leftrightarrow$  is there a unique coordinate vector  $\beta$ ?

Two further important subspaces associated with a real  $n \times p$  matrix  $\mathbf{Q}$ :

- the null space (or kernel),  $\ker(\mathbf{Q})$ , of  $\mathbf{Q}$  is the subspace defined as

$$\ker(\mathbf{Q}) = \{\mathbf{x} \in \mathbb{R}^p : \mathbf{Q}\mathbf{x} = \mathbf{0}\};$$

- the orthogonal complement of  $\mathcal{M}(\mathbf{Q})$ ,  $\mathcal{M}^\perp(\mathbf{Q})$ , is the subspace defined as

$$\begin{aligned}\mathcal{M}^\perp(\mathbf{Q}) &= \{\mathbf{y} \in \mathbb{R}^n : \mathbf{y}^\top \mathbf{Q}\mathbf{x} = 0, \forall \mathbf{x} \in \mathbb{R}^p\} \\ &= \{\mathbf{y} \in \mathbb{R}^n : \mathbf{y}^\top \mathbf{v} = 0, \forall \mathbf{v} \in \mathcal{M}(\mathbf{Q})\}.\end{aligned}$$

$\langle \mathbf{y}, \mathbf{v} \rangle = 0$

The orthogonal complement may be defined for arbitrary subspaces by using the second equality.



## Theorem (Spectral Theorem)

A  $p \times p$  matrix  $\mathbf{Q}$  is symmetric if and only if there exists a  $p \times p$  orthogonal matrix<sup>a</sup>  $\mathbf{U}$  and a diagonal matrix  $\mathbf{\Lambda}$  such that

$$\mathbf{Q} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$$

$$\mathbf{Q} = \begin{pmatrix} * & & \\ & \ddots & \\ & & * \end{pmatrix} \quad \mathbf{\Lambda} = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_p \end{pmatrix}$$

In particular:

- 1 the columns of  $\mathbf{U} = (\mathbf{u}_1 \cdots \mathbf{u}_p)$  are eigenvectors of  $\mathbf{Q}$ , i.e. there exist  $\lambda_j$  such that

$$\mathbf{Q}\mathbf{u}_j = \lambda_j \mathbf{u}_j, \quad j = 1, \dots, p;$$

- 2 the entries of  $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_p)$  are the corresponding eigenvalues of  $\mathbf{Q}$ , which are real; and
- 3 the rank of  $\mathbf{Q}$  is the number of non-zero eigenvalues.

---

<sup>a</sup>A matrix is orthogonal if  $\mathbf{U}\mathbf{U}^\top = \mathbf{U}^\top \mathbf{U} = \mathbf{I}_p$

Note: if the eigenvalues are distinct, the eigenvectors are unique (up to changes in signs).

## Theorem (Singular Value Decomposition)

Any  $n \times p$  real matrix can be factorised as

$$\underset{n \times p}{Q} = \underset{n \times n}{U} \underset{n \times p}{\Sigma} \underset{p \times p}{V}^T,$$

where  $U$  and  $V^T$  are orthogonal with columns called left singular vectors and right singular vectors, respectively, and  $\Sigma$  is diagonal with real entries called singular values.

- 1 The left singular vectors are eigenvectors of  $QQ^T$ .<sup>2</sup>
- 2 The right singular vectors are eigenvectors of  $Q^T Q$ .
- 3 The squares of the singular values are eigenvalues of both  $QQ^T$  and  $Q^T Q$ .
- 4 The left singular vectors corresponding to non-zero singular values form an orthonormal basis for  $\mathcal{M}(Q)$ .
- 5 The left singular vectors corresponding to zero singular values form an orthonormal basis for  $\mathcal{M}^\perp(Q)$ .

<sup>2</sup>hint: compute  $QQ^T U_i = \lambda_i^2 U_i$  for all  $i \leq p$ . And similarly with  $Q^T Q V_i = \lambda_i^2 V_i$

A matrix  $\mathbf{Q}$  is called **idempotent** if  $\mathbf{Q}^2 = \mathbf{Q}$ .

An **orthogonal projection** (henceforth **projection**) onto a subspace  $\mathcal{V}$  is a symmetric idempotent matrix  $\mathbf{H}$  such that  $\mathcal{M}(\mathbf{H}) = \mathcal{V}$ , i.e. the column space of  $\mathbf{H}$  coincides with the subspace  $\mathcal{V}$ .

### Proposition

*The only possible eigenvalues of a projection matrix are 0 and 1.*

## Proposition

Let  $\mathcal{V}$  be a subspace and  $\mathbf{H}$  be a projection onto  $\mathcal{V}$ . Then  $\mathbf{I} - \mathbf{H}$  is the projection matrix onto  $\mathcal{V}^\perp$ .

## Proof (\*).

We first prove that  $\mathbf{I} - \mathbf{H}$  is a projection matrix (idempotent and symmetric).

$$(\mathbf{I} - \mathbf{H})^\top = \mathbf{I} - \mathbf{H}^\top = \mathbf{I} - \mathbf{H} \text{ since } \mathbf{H} \text{ is symmetric and,}$$

$$(\mathbf{I} - \mathbf{H})^2 = \mathbf{I}^2 - 2\mathbf{H} + \mathbf{H}^2 = \mathbf{I} - \mathbf{H}.$$

It remains to identify the column space of  $\mathbf{I} - \mathbf{H}$ . Let  $\mathbf{H} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$  be the spectral decomposition of  $\mathbf{H}$ .

$$\text{Then } \mathbf{I} - \mathbf{H} = \mathbf{U}\mathbf{U}^\top - \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top = \mathbf{U}(\mathbf{I} - \mathbf{\Lambda})\mathbf{U}^\top.$$

Hence the column space of  $\mathbf{I} - \mathbf{H}$  is spanned by the eigenvectors of  $\mathbf{H}$  corresponding to zero eigenvalues of  $\mathbf{H}$ , which coincides with  $\mathcal{M}^\perp(\mathbf{H}) = \mathcal{V}^\perp$ .  $\square$

$$\begin{aligned} & \text{diag}(1, 0, 0, 0, 1, \dots) \\ \mathbf{I} &= \text{diag}(1, 1, \dots, 1) \end{aligned}$$

## Proposition

Let  $\mathcal{V}$  be a subspace and  $H$  be a projection onto  $\mathcal{V}$ . Then  $Hy = y$  for all  $y \in \mathcal{V}$ .

## Proposition

If  $P$  and  $Q$  are projection matrices onto a subspace  $\mathcal{V}$ , then  $P = Q$ .

## Proposition

If  $x_1, \dots, x_p$  are linearly independent<sup>a</sup> and are such that  $\text{span}(x_1, \dots, x_p) = \mathcal{V}$ , then the projection onto  $\mathcal{V}$  can be represented as

$$H = X(X^T X)^{-1} X^T$$

where  $X$  is a matrix with columns  $x_1, \dots, x_p$ .

---

<sup>a</sup> $\sum_{i=1}^p a_i x_i = 0$  iff.  $a_i = 0$ , for all  $i \leq p$

## Proposition

Let  $\mathcal{V}$  be a subspace of  $\mathbb{R}^n$  and  $H$  be a projection onto  $\mathcal{V}$ . Then

$$\overset{\text{distance}}{\|x - \underline{Hx}\|} \leq \|x - v\|, \quad \forall v \in \mathcal{V}.$$

## Proof (\*).

Let  $H = U\Lambda U^\top$  be the spectral decomposition of  $H$ ,  $U = (u_1 \cdots u_n)$  and  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ . Letting  $p = \dim(\mathcal{V})$ , then

by assumption of  $H$

- ①  $\lambda_1 = \cdots = \lambda_p = 1$  and  $\lambda_{p+1} = \cdots = \lambda_n = 0$ , (by definition of a projection matrix s19)
- ②  $u_1, \dots, u_n$  is an orthonormal basis of  $\mathbb{R}^n$ ,
- ③  $u_1, \dots, u_p$  is an an orthonormal basis of  $\mathcal{V}$ .

Let's us it in the following computations

$$\begin{aligned}
\| \underline{x - Hx} \|^2 &= \sum_{i=1}^n (\underline{x^\top u_i} - \underline{(Hx)^\top u_i})^2 \quad [\text{orthonormal basis}] \\
&= \sum_{i=1}^n (\underline{x^\top u_i} - \underline{x^\top H u_i})^2 \quad [H \text{ is symmetric}] \\
&= \sum_{i=1}^n (\underline{x^\top u_i} - \underline{\lambda_i x^\top u_i})^2 \quad [u\text{'s are eigenvectors of } H] \\
&= 0 + \sum_{i=p+1}^n (x^\top u_i)^2 \quad [\text{eigenvalues 0 or 1}] \\
&\leq \sum_{i=1}^p (x^\top u_i - v^\top u_i)^2 + \sum_{i=p+1}^n (x^\top u_i)^2 \quad \forall v \in \mathcal{V} \\
&= \underline{\|x - v\|^2}.
\end{aligned}$$



## Proposition

Let  $\mathcal{V}_1 \subseteq \mathcal{V} \subseteq \mathbb{R}^n$  be two nested linear subspaces. If  $\mathbf{H}_1$  is the projection onto  $\mathcal{V}_1$  and  $\mathbf{H}$  is the projection onto  $\mathcal{V}$ , then

$$\mathbf{H}\mathbf{H}_1 = \mathbf{H}_1 = \mathbf{H}_1\mathbf{H}.$$

## Proof (\*).

First we show that  $\mathbf{H}\mathbf{H}_1 = \mathbf{H}_1$ , and then that  $\mathbf{H}_1\mathbf{H} = \mathbf{H}\mathbf{H}_1$ . For all  $\mathbf{y} \in \mathbb{R}^n$  we have  $\mathbf{H}_1\mathbf{y} \in \mathcal{V}_1$ . But then  $\mathbf{H}_1\mathbf{y} \in \mathcal{V}$ , since  $\mathcal{V}_1 \subseteq \mathcal{V}$ .

Therefore  $\mathbf{H}\mathbf{H}_1\mathbf{y} = \mathbf{H}_1\mathbf{y}$ . We have shown that  $(\mathbf{H}\mathbf{H}_1 - \mathbf{H}_1)\mathbf{y} = \mathbf{0}$  for all  $\mathbf{y} \in \mathbb{R}^n$ , so that  $\mathbf{H}\mathbf{H}_1 - \mathbf{H}_1 = \mathbf{0}$ , as its kernel is all  $\mathbb{R}^n$ . Hence  $\mathbf{H}\mathbf{H}_1 = \mathbf{H}_1$ .

To prove that  $\mathbf{H}_1\mathbf{H} = \mathbf{H}\mathbf{H}_1$ , note that symmetry of projection matrices and the first part of the proof give

$$\mathbf{H}_1\mathbf{H} = \mathbf{H}_1^\top \mathbf{H}^\top = (\mathbf{H}\mathbf{H}_1)^\top = (\mathbf{H}_1)^\top = \mathbf{H}_1 = \mathbf{H}\mathbf{H}_1.$$





## Definition (Quadratic Form Definition)

A  $p \times p$  real symmetric matrix  $\Omega$  is called **non-negative definite** (written  $\Omega \succeq 0$ ) if and only if  $\mathbf{x}^\top \Omega \mathbf{x} \geq 0$  for all  $\mathbf{x} \in \mathbb{R}^p$ . If  $\mathbf{x}^\top \Omega \mathbf{x} > 0$  for all  $\mathbf{x} \in \mathbb{R}^p \setminus \{0\}$ , then we call  $\Omega$  **positive definite** (written  $\Omega \succ 0$ ).

## Definition (Spectral Definition)

A  $p \times p$  real symmetric matrix  $\Omega$  is called **non-negative definite** (written  $\Omega \succeq 0$ ) if and only if the eigenvalues of  $\Omega$  are non-negative. If the eigenvalues of  $\Omega$  are strictly positive, then  $\Omega$  is called **positive definite** (written  $\Omega \succ 0$ ).

$$\Omega = U \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_p \end{pmatrix} U^\top$$

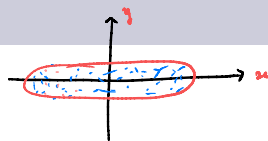
### Lemma (Little exercise)

*The two definitions are equivalent.*

### Proposition (Non-Negative and Covariance Matrices)

*Let  $\Omega$  be a real symmetric matrix.*

*Then  $\Omega$  is non-negative definite iff.  $\Omega$  is the covariance matrix of some random vector  $\mathbf{Y}$ .*



We want to find the subspace that explains the most a random vector  $\mathbf{Y}$  in  $\mathbb{R}^d$  with covariance matrix  $\mathbf{\Omega}$ .

- Step  $j = 1$ : Find direction  $\mathbf{v}_1 \in \mathbb{S}^{d-1}$  <sup>*unit sphere*</sup> such that the projection of  $\mathbf{Y}$  onto  $\mathbf{v}_1$  has maximal variance.
- Steps  $j = 2, 3, \dots, d$ : Find direction  $\mathbf{v}_j \perp \{\mathbf{v}_1, \dots, \mathbf{v}_{j-1}\}$  such that projection of  $\mathbf{Y}$  onto  $\mathbf{v}_j$  has maximal variance.

- First, by Proposition s26,  $\Omega$  is symmetric, non-negative definite of size  $d \times d$ .
- Step  $j = 1$ :** Maximise  $\text{var}(\mathbf{v}_1^\top \mathbf{Y}) = \mathbf{v}_1^\top \Omega \mathbf{v}_1$  over  $\|\mathbf{v}_1\| = 1$
- $\mathbf{v}_1^\top \Omega \mathbf{v}_1 = \mathbf{v}_1^\top \underbrace{\mathbf{U} \Lambda \mathbf{U}^\top}_{\text{diag}(\lambda_i)_{i=1}^d} \mathbf{v}_1 = \|\underbrace{\Lambda^{1/2} \mathbf{U}^\top \mathbf{v}_1}_{\substack{\uparrow \\ p_i}}\|^2 = \sum_{i=1}^d \underbrace{\lambda_i (\underbrace{\mathbf{u}_i^\top \mathbf{v}_1}_{p_i})^2}_{p_i}$  [change of basis]

Now  $\sum_{i=1}^d (\mathbf{u}_i^\top \mathbf{v}_1)^2 = \|\mathbf{v}_1\|^2 = 1$  so we have a convex combination of  $\{\lambda_j\}_{j=1}^d$ ,

$$\sum_{i=1}^d p_i \lambda_i, \quad \sum_i p_i = 1, \quad p_i \geq 0, \quad i = 1, \dots, d.$$

But  $\lambda_1 \geq \lambda_i \geq 0$  so clearly this sum is maximised when  $p_1 = 1$  and  $p_j = 0$   $\forall j \neq 1$ , i.e.  $\mathbf{v}_1 = \pm \mathbf{u}_1$ .

- Steps  $j = 2, 3, \dots, d$ :** Iteratively,  $\mathbf{v}_j = \pm \mathbf{u}_j$ , i.e. principal components are eigenvectors of  $\Omega$ .

$$\Omega = U \Lambda U^T \quad \Lambda = \text{diag}(\lambda_1, \dots, \lambda_k, \lambda_{k+1}, \dots)$$

## Theorem (Optimal (Linear) Dimension Reduction Theorem)

Let  $\mathbf{Y}$  be a mean-zero random variable in  $\mathbb{R}^d$  with  $d \times d$  covariance  $\Omega$ . Let  $\mathbf{H}$  be the projection matrix onto the span of the first  $k$  eigenvectors of  $\Omega$ . Then

$$\mathbb{E} \|\mathbf{Y} - \mathbf{H}\mathbf{Y}\|^2 \leq \mathbb{E} \|\mathbf{Y} - \mathbf{Q}\mathbf{Y}\|^2$$

for any  $d \times d$  projection matrix  $\mathbf{Q}$  or rank at most  $k$ .

**Intuitively:** if you want to approximate a mean-zero random variable taking values  $\mathbb{R}^d$  by a random variable that ranges over a subspace of dimension at most  $k \leq d$ , the optimal choice is the projection of the random variable onto the space spanned by its first  $k$  principal components (eigenvectors of the covariance).

“Optimal” is with respect to the mean squared error.

For the proof, use lemma below (follows immediately from spectral decomposition)

### Lemma

$Q$  is a rank  $k$  projection matrix iff. there exist orthonormal vectors  $\{\mathbf{v}_j\}_{j=1}^k$  such that  $Q = \sum_{j=1}^k \mathbf{v}_j \mathbf{v}_j^\top$

## Proof of Optimal Linear Dimension Reduction (\*).

Write  $\mathbf{Q} = \sum_{j=1}^k \mathbf{v}_j \mathbf{v}_j^\top$  for some orthonormal  $\{\mathbf{v}_j\}_{j=1}^k$ . Then

$$\mathbb{E} \|\mathbf{Y} - \mathbf{QY}\|^2 =$$

MSE( $\mathbf{Y}, \mathbf{QY}$ )

$$= \mathbb{E} [\mathbf{Y}^\top (\mathbf{I} - \mathbf{Q})^\top (\mathbf{I} - \mathbf{Q}) \mathbf{Y}] = \mathbb{E} [\text{tr}\{(\mathbf{I} - \mathbf{Q}) \mathbf{Y} \mathbf{Y}^\top (\mathbf{I} - \mathbf{Q})^\top\}]$$

$$= \text{tr}\{(\mathbf{I} - \mathbf{Q}) \mathbb{E} [\mathbf{Y} \mathbf{Y}^\top] (\mathbf{I} - \mathbf{Q})^\top\} = \text{tr}\{(\mathbf{I} - \mathbf{Q})^\top (\mathbf{I} - \mathbf{Q}) \Omega\}$$

$$= \text{tr}\{(\mathbf{I} - \mathbf{Q}) \Omega\} = \text{tr}\{\Omega\} - \text{tr}\{\mathbf{Q} \Omega\} = \sum_{i=1}^d \lambda_i - \text{tr}\left\{\sum_{j=1}^k \mathbf{v}_j \mathbf{v}_j^\top \Omega\right\}$$

$$= \sum_{i=1}^d \lambda_i - \sum_{j=1}^k \text{tr}\{\mathbf{v}_j \mathbf{v}_j^\top \Omega\} = \sum_{i=1}^d \lambda_i - \sum_{j=1}^k \mathbf{v}_j^\top \Omega \mathbf{v}_j$$

$$= \sum_{i=1}^d \lambda_i - \sum_{j=1}^k \text{var}[\mathbf{v}_j^\top \mathbf{Y}]$$

If we can minimise this expression over all  $\{\mathbf{v}_j\}_{j=1}^k$  with  $\mathbf{v}_i^\top \mathbf{v}_j = \mathbf{1}\{i=j\}$ , then we're done. By PCA, this is done by choosing the top  $k$  eigenvectors of  $\Omega$ .  $\square$

Recall that for any matrices  $\mathbf{A}, \mathbf{B}, \mathbf{C}$ , we have  $\text{tr}(\mathbf{ABC}) = \text{tr}(\mathbf{BCA}) = \text{tr}(\mathbf{CAB})$  under conditions (cf A3W8).

## Corollary (Deterministic Version)

Let  $\{\mathbf{x}_1, \dots, \mathbf{x}_p\} \subset \mathbb{R}^d$  be such that  $\mathbf{x}_1 + \dots + \mathbf{x}_p = 0$ , and let  $\mathbf{X}$  be the matrix with columns  $\{\mathbf{x}_i\}_{i=1}^p$ . The best approximating  $k$ -hyperplane to the points  $\{\mathbf{x}_1, \dots, \mathbf{x}_p\}$  is given by the span of the first  $k$  eigenvectors of the matrix  $\mathbf{X}\mathbf{X}^\top$ , i.e. if  $\mathbf{H}$  is the projection onto this span, it holds that

$$\sum_{i=1}^p \|\mathbf{x}_i - \mathbf{H}\mathbf{x}_i\|^2 \leq \sum_{i=1}^p \|\mathbf{x}_i - \mathbf{Q}\mathbf{x}_i\|^2$$

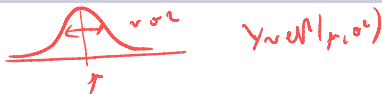
$(\mathbf{Q} = \mathbb{E}[\mathbf{Y}\mathbf{Y}^\top])$

for any  $d \times d$  projection operator  $\mathbf{Q}$  or rank at most  $k$ .

## Proof.

Define the discrete random vector  $\mathbf{Y}$  by  $\mathbb{P}[\mathbf{Y} = \mathbf{x}_i] = 1/p$ , and use optimal linear dimension reduction as stated earlier.  $\square$





## 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.

$$M_{(\mathbf{W}, \mathbf{W}')}(\boldsymbol{\theta}) \stackrel{!}{=} M_{\mathbf{W}}(\boldsymbol{\theta}) M_{\mathbf{W}'}(\boldsymbol{\theta})$$

*joint density*

## 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 if  $p=1$ ,  $1-\Omega=0^2$   $\frac{(\mathbf{y}-\boldsymbol{\mu})^2}{2\sigma^2}$
- ❺ Marginals of Gaussian are Gaussian (converse NOT true).
- ❻  $\boldsymbol{\Omega}$  diagonal  $\Leftrightarrow$  independent coordinates  $Y_i$ .
- 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  $\mathbf{u} \in \mathbb{R}^d$  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

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

□

## 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  $\mathcal{N}(\boldsymbol{\mu}_{p \times 1}, \Omega_{p \times p})$  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,  $\forall \mathbf{u}, \mathbf{v}$ ,

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

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

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

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

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

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



## Reminder:

Definition ( $\chi^2$  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 ( $\chi^2$ ) distribution with  $p$  degrees of freedom; we write  $\|\mathbf{Z}\|^2 \sim \chi_p^2$ .

[Thus,  $\chi_p^2$  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  $\implies$

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