# Random vectors and the multivariate normal distribution,

## conditional expectation

# 1.1. Random vectors

For the purposes of statistics we will understand random vectors as columns, or $n \times 1$ matrices. We will write

$$\underline{X} = \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{pmatrix}$$

Definition: The expected value of a random vector is the vector of expected values of components. In symbols

$$E(\underline{X}) = \begin{pmatrix} E(X_1) \\ E(X_2) \\ \vdots \\ E(X_n) \end{pmatrix}$$

Definition: An $m \times n$ away of random variables $X_{i,j}$ is called a random matrix.

Similarly as before we define the expectation of a random matrix as the matrix of expected values. In symbols

$$E(\underline{x}) = \begin{pmatrix} E(X_{11}) & \cdots & E(X_{1n}) \\ & \vdots & \\ E(X_{m1}) & \cdots & E(X_{mn}) \end{pmatrix}$$

Theorem 1.1: Let $\underline{A}, \underline{B}$ be fixed matrices and $\underline{x}$ a random matrix. We have

$$E(\underline{A}\,\underline{x}\,\underline{B}) = \underline{A}\, E(\underline{x}) \cdot \underline{B}$$

Comment: we assume $\underline{A}, \underline{B}$ are of right dimension.

**Proof:** We have

$$(\underline{A} \times \underline{B})_{i,j} = \sum_{k,\ell} a_{ik} \cdot X_{k,\ell} \, b_{\ell j}.$$

The proof follows from linearity of expected value.

We need to extend the notion of variance and covariance to random vectors. In one dimension we have

$$\text{cov}(x,y) = E(x \cdot y) - E(x) \cdot E(y).$$

The analogy for vectors would be

$$\text{cov}(\underline{x},\underline{y}) = E(\underline{x} \cdot \underline{y}^T) - E(\underline{x}) \cdot E(\underline{y})^T$$

**Comment:** Obviously

$$E(\underline{y}^T) = E(\underline{y})^T.$$

If we write it componentwise and say $\underline{c} = \text{cov}(\underline{x}, \underline{y})$ we get

$$\underline{c}_{ij} = E(x_i \cdot y_j) - E(x_i)E(y_j)$$

$$= \text{cov}(x_i, y_j)$$

The covariance of $\underline{x}, \underline{y}$ is an array or matrix of covariances between components of $\underline{x}$ and $\underline{y}$.

$$\text{cov}(\underline{x}, \underline{y}) = \begin{pmatrix} \text{cov}(x_1, y_1), \ldots, \text{cov}(x_1, y_n) \\ \vdots \\ \text{cov}(x_m, y_1) \ldots \text{cov}(x_m, y_n) \end{pmatrix}$$

why is this definition good? It makes many computations elegant.

Theorem 1.2 :  Let  $\underline{A}, \underline{B}$  be

fixed matrices . We have

$$\text{cov} \left( \underline{A}\underline{x}, \underline{B}\underline{y} \right) = \underline{A} \, \text{cov} \left( \underline{x}, \underline{y} \right) \underline{B}^T$$

Proof :  We compute

$$E \left( (\underline{A}\underline{x})(\underline{B}\underline{y})^T \right)$$

$$= E \left( \underline{A} \, \underline{x} \cdot \underline{y}^T \underline{B}^T \right)$$

$$= \underline{A} \, E \left( \underline{x} \, \underline{y}^T \right) \underline{B}^T$$

and  we  know

$$E(\underline{A}\underline{x}) = \underline{A} \cdot E(\underline{x})$$
$$E(\underline{y}^T \cdot \underline{B}^T) = E(\underline{y}^T) \underline{B}^T$$

We  need  to  define  var$(\underline{x})$.

By analogy we have

$$\text{var}(\underline{x}) = \text{cov}(\underline{x}, \underline{x}).$$

For variances Theorem 1.2 becomes

$$\text{var}(\underline{A}\,\underline{x}) = \underline{A}\,\text{var}(\underline{x})\underline{B}^T$$

Example: Let $\underline{x}^T = (x_1, x_2, \ldots, x_n)$ where all $x_k$ are independent and $x_k \sim N(\mu, \sigma^2)$ for $k = 1, 2, \ldots, n$. What is the variance of

$$\underline{y} = \begin{pmatrix} x_1 - \overline{x} \\ x_2 - \overline{x} \\ \vdots \\ x_n - \overline{x} \end{pmatrix}.$$

Let $\underline{H} = \underline{I} - \frac{1}{n}\underline{1}\cdot\underline{1}^T$ where

$$\underline{1} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}.$$ We have

$$\underline{y} = \underline{H}\cdot\underline{x} \quad \text{and}$$

$$\underline{H}^T = \underline{H}, \quad \text{and.}$$

$$\underline{H}^2 = \left(\underline{I} - \frac{1}{n} \cdot \underline{1} \cdot \underline{1}^T\right)\left(\underline{I} - \frac{1}{n} \underline{1} \cdot \underline{1}^T\right)$$

$$= \underline{I} - \frac{2}{n} \underline{1} \cdot \underline{1}^T + \frac{1}{n^2} \underline{1} \underbrace{\underline{1}^T \underline{1}}_{= n} \cdot \underline{1}^T$$

$$= \underline{I} - \frac{1}{n} \underline{1} \cdot \underline{1}^T$$

$$= \underline{H} .$$

It follows that

$$\text{var}(\underline{Y}) = \text{var}(\underline{H} \cdot \underline{x})$$

$$= \underline{H} \cdot \text{var}(\underline{x}) \cdot \underline{H}^T$$

$$= \underline{H} \cdot \sigma^2 \underline{I} \cdot \underline{H}^T$$

$$= \sigma^2 \underline{H} \cdot \underline{H}^T$$

$$= \sigma^2 \cdot \underline{H}^2$$

$$= \sigma^2 \left(\underline{I} - \frac{1}{n} \underline{1} \cdot \underline{1}^T\right) .$$

Theorem 1.3 : If $\underline{x}$ and $\underline{y}$
are random vectors of the same
dimension then

$$\text{var}(\underline{x} + \underline{y}) = \text{var}(\underline{x}) + \text{var}(\underline{y})$$
$$+ \text{cov}(\underline{x}, \underline{y})$$
$$+ \text{cov}(\underline{y}, \underline{x})$$

Proof : This follows from
definitions :

$$(\underline{x} + \underline{y})(\underline{x} + \underline{y})^T$$

$$= \underline{x} \cdot \underline{x}^T + \underline{x} \cdot \underline{y}^T$$
$$+ \underline{y} \cdot \underline{x}^T + \underline{y} \cdot \underline{y}^T$$
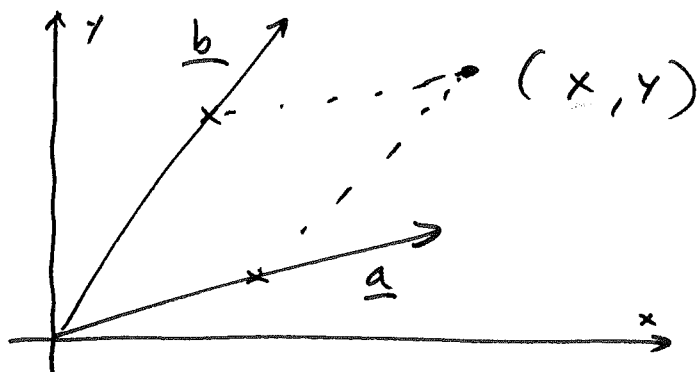
Take expectations and collect
terms.

# 1.2. Multivariate normal distribution

Idea: We want to create a random point in the plane or a random point in space.

Take two vectors $\underline{a}$ and $\underline{b}$ and create a linear combination with random components

Figure:



We can take only one vector as well but then all the points will be on a line.

One possibility is to take coefficients as independent random variables. We can specialize and say the coefficients $z_1, z_2, \ldots$ are independent normal random variables. Denote

$$\underline{z}^T = (z_1, \ldots, z_p) .$$

We can create a point in $\mathbb{R}^n$ by taking $p$ vectors $\underline{a}_1, \underline{a}_2, \ldots, \underline{a}_p$ and say

$$\underline{x} = z_1 \cdot \underline{a}_1 + \cdots + z_p \underline{a}_p + \underline{\mu} .$$

If we collect the vectors in a matrix

$$\underline{A} = \begin{pmatrix} \underline{a}_1, \underline{a}_2, \ldots, \underline{a}_p \end{pmatrix}$$

in matrix notation we have

$$\underline{x} = \underline{A} \cdot \underline{z} + \underline{\mu}$$

Comment: If $p < n$ the random points created will be in a hyperplane.

From 1.1 we have that

$$E(\underline{x}) = \underline{A} \cdot E(\underline{z}) + \underline{\mu} = \underline{\mu}$$

and

$$var(\underline{x}) = \underline{A} \cdot var(\underline{z}) \cdot \underline{A}^T$$

$$= \underline{A} \cdot \underline{I} \cdot \underline{A}^T$$

$$= \underline{A} \cdot \underline{A}^T.$$

We will say later that $\underline{X}$ has multivariate normal distribution with parameters $\underline{\mu}$ and $\underline{\Sigma} = \underline{A}\underline{A}^T$. But we must answer another question first. Suppose

$$\underline{Y} = \underline{B} \cdot \underline{Z} + \underline{\mu}$$

and $\underline{A}\underline{A}^T = \underline{B}\underline{B}^T$. Do $\underline{X}$ and $\underline{Y}$ have the same distribution? This means $P(\underline{X} \in u) = P(\underline{Y} \in \underline{u})$ for all reasonable sets $u$. The idea of the proof is the following: we will show that

$$\underline{X} = \underline{M}\underline{Z} + \underline{\mu} \quad \text{and} \quad \underline{Y} = \underline{M}\underline{\tilde{Z}} + \underline{\mu}$$

where $\underline{z}$ and $\bar{\underline{z}}$ will have the same distribution. Once we have that we will know that $\underline{A}\underline{A}^T$ and $\hat{\underline{\mu}}$ uniquely determine the distribution of $\underline{x}$.

We need the singular value decomposition for matrices: if $\underline{A}$ is a $p \times q$ matrix there are orthogonal matrices $\underline{S}$ ($p \times p$) and $T$ ($q \times q$) such that

$$\underline{A} = \underline{S} \cdot \underline{D} \cdot T^T$$

where

(i) $\underline{D}$ is of the form

$$\underline{D} = \begin{pmatrix} \sqrt{\lambda_1} & & & 0 \\ 0 & \ddots & \sqrt{\lambda_r} & \\ & & & 0 \end{pmatrix}$$

where $\lambda_1, \lambda_2, \ldots, \lambda_r$ are non-zero eigenvalues of matrices $\underline{A}\underline{A}^T$ or $\underline{A}^T\underline{A}$ (they are the same).

(ii) The columns of $\underline{S}$ are orthogonal eigenvectors of $\underline{A}\underline{A}^T$, the first $r$ belonging to eigenvalues $\lambda_1, \ldots, \lambda_r$

(iii) The columns of $\underline{T}$ are eigenvectors of $\underline{A}^T\underline{A}$; the first $r$ belong to eigenvalues $\lambda_1, \lambda_2, \ldots, \lambda_r$.

If $A A^T = B B^T$ this means that we can write

$$A = S \cdot D_1 \cdot T_1^T$$

$$B = S \cdot D_2 \cdot T_2^T$$

Fact from probability: if

$X$ has density $f_X(x)$ and

$A$ is an invertible matrix

the vector

$$Y = A \cdot X + v$$ has

density

$$f_Y(y) = \frac{1}{|\det(A)|} \times$$

$$\times f_X\left(A^{-1}(y-v)\right).$$

If $\underline{z}^T = (z_1, z_2, \ldots, z_n)$ has independent normal components the density is the product i.e.

$$f_{\underline{z}}(\underline{z}) = \prod_{k=1}^{n} f_{z_k}(z_k)$$

$$= \left(\frac{1}{\sqrt{2\pi}}\right)^n \cdot e^{-\frac{z_1^2 + \cdots + z_n^2}{2}}$$

$$= \frac{1}{(2\pi)^{n/2}} \, e^{-\frac{\underline{z}^T \cdot \underline{z}}{2}}$$

If $\underline{A}$ is invertible the density of $\underline{w} = \underline{A}\,\underline{z} + \underline{v}$ is

$$f_{\underline{w}}(\underline{w}) = \frac{1}{(2\pi)^{n/2} |\det A|} \times$$

$$\times \; e^{-\frac{1}{2}[\underline{A}^{-1}(\underline{w}-\underline{v})]^T [\underline{A}^{-1}(\underline{w}-\underline{v})]}$$

$$= \qquad (*)$$

$$(*) = \frac{1}{(2\pi)^{n/2} |\det(\underline{A})|}$$

$$\times \ e^{-\frac{1}{2}(w-\underline{\nu})^T (\underline{A}^{-1})^T A^{-1}(w-\underline{\nu})}$$

$$= \frac{1}{(2\pi)^{n/2} |\det(\underline{A})|}$$

$$\times \ e^{-\frac{1}{2}(w-\underline{\nu})^T (\underline{A}\underline{A}^T)^{-1}(w-\underline{\nu})}$$

This means that

$$T_1^T \underline{z}_1 = \underline{w}_1 \quad \text{has}$$

independent normal $N(0,1)$

components because $\underline{\nu} = 0$

and $T_1^T T_1 = \underline{I}$. The same

is true for $\underline{w}_2 = T_2^T \cdot \underline{z}_2$.

Putting these facts together gives

$$\underline{X} = \underline{S} \cdot \underline{D_1} \underline{W_1} + \underline{\mu}$$

$$\underline{Y} = \underline{S} \cdot \underline{D_2} \underline{W_2} + \underline{\mu}$$

But $\underline{W_1}$ and $\underline{W_2}$ have the $N(0, \underline{I})$ density. This implies $\underline{X}$ and $\underline{Y}$ have the same distribution. because the components are the same linear combinations of indep. normals.

<u>Theorem 1.4</u> : Let $\underline{Z}^T = (z_1, \ldots, z_n)$ have independent standard normal distribution and let $\underline{A}$ be a matrix and $\underline{\mu} \in \mathbb{R}^m$ a vector. The distribution of

$$\underline{X} = \underline{A} \underline{Z} + \underline{\mu}$$

is uniquely determined by $\underline{\mu}$ and $\underline{\Sigma} = \underline{A} \underline{A}^T$.

Proof: Done already.

This justifies the following definition:

Definition: Let $\underline{z}^T = (z_1, \ldots, z_n)$ be a vector of independent standard normal variables. The distribution of the vector

$$\underline{X} = \underline{A}\,\underline{z} + \underline{\mu}$$

is called multivariate normal with parameters $\underline{\mu}$ and $\underline{\Sigma} = \underline{A}\underline{A}^T$. Symbol: $\underline{X} \sim N(\underline{\mu}, \underline{\Sigma})$

Comments:

(i) just as in the one-dimensional case we have

$$\underline{\mu} = E(\underline{X}) \quad \text{and} \quad \underline{\Sigma} = \text{var}(\underline{X}).$$

(ii) If $\underline{A}$ is invertible $\underline{x}$ has a density

Theorem A.5 : If $\underline{x} \sim N(\underline{\mu}, \underline{\Sigma})$ then :

(i) all marginal distributions are normal.

(ii) If $\underline{Y} = \underline{B}\underline{x} + \underline{\nu}$ then $\underline{Y}$ is multivariate normal with parameter $\underline{\Sigma}' = \underline{B} \underline{\Sigma} \underline{B}^T$ and $\underline{\mu}' = \underline{B}\underline{\mu} + \underline{\nu}$.

Proof: Follows directly from definitions.

Suppose $A$ is invertible and

$$A A^T = \begin{pmatrix} \underline{\Sigma}_{11} & 0 \\ 0 & \underline{\Sigma}_{22} \end{pmatrix} \begin{array}{l} \} \, p \\ \} \, q \end{array}$$

We assume $\underline{\Sigma}_{11}$ and $\underline{\Sigma}_{22}$ are invertible and square. If we write $\underline{X} = \begin{pmatrix} \underline{X}_1 \\ \underline{X}_2 \end{pmatrix} \begin{array}{l} \} \, p \\ \} \, q \end{array}$ then the density has the form

$$f_{\underline{X}}(\underline{x}) = f_{\underline{X}_1}(\underline{x}_1) \cdot f_{\underline{X}_2}(\underline{x}_2).$$

This means that $\underline{X}_1$ and $\underline{X}_2$ are independent. But we have

$$\text{cov}(\underline{X}_1, \underline{X}_2) = 0.$$

In this particular case

we have that uncorrelated vectors are independent. This is not generally true.

Theorem 1.6 : Let $\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \sim N(\underline{\mu}, \underline{\Sigma})$. If $\text{cov}(\underline{x}_1, \underline{x}_2) = 0$ then $\underline{x}_1$ and $\underline{x}_2$ are independent.

Proof : By Cholesky we can find a matrix $\underline{A}$ such that $\underline{\Sigma} = \underline{A} \cdot \underline{A}^T$. Write

$$\underline{x} = \underline{A} \cdot \underline{z} + \underline{\mu}$$

$$= \underline{S}\,\underline{D} \cdot \underline{T}^T \cdot \underline{z} + \underline{\mu} \quad \text{by SVD.}$$

$$= \underline{S} \cdot \underline{D} \cdot \hat{\underline{z}} + \underline{\mu}$$

But linear algebra gives us that $\underline{S}$ must be of the form

$$\underline{S} = \begin{pmatrix} \underline{S}_{11} & 0 \\ 0 & \underline{S}_{22} \end{pmatrix} \begin{matrix} \} p \\ \} 2 \end{matrix}$$

This means that $X_1$ depends on the first $p$ components of $\tilde{\underline{z}}$ and $X_2$ on the last $p$ components of $\hat{\underline{z}}$. This implies independence.

Example: Let $\underline{X} = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \sim N(\underline{\mu}, \underline{\Sigma})$

with $\underline{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$ and $\underline{\Sigma} = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$.

Assume that $\Sigma_{11}$ is invertible.

Is there a matrix $\underline{C}$ such that

$\underline{X}_1$ and $\underline{X}_2 - \underline{C}\underline{X}_1$ are independent?

We just need to compute

$\text{cov}(\underline{X}_1, \underline{X}_2 - \underline{C}\underline{X}_1)$ because

we know that $\begin{pmatrix} \underline{x}_1 \\ \underline{x}_2 - C\underline{x}_1 \end{pmatrix}$

is multivariate normal.

We have

$$\text{cov}(\underline{x}_1, \underline{x}_2 - C\underline{x}_1)$$

$$= \text{cov}(\underline{x}_1, \underline{x}_2) - \text{cov}(\underline{x}_1, C\underline{x}_1)$$

$$= \Sigma_{12} - \Sigma_{11} \cdot C^T$$

If we want this to be 0

we need

$$C^T = \Sigma_{11}^{-1} \Sigma_{12} \quad \text{or}$$

$$C = \Sigma_{21} \Sigma_{11}^{-1}.$$

## 1.3. Quadratic forms

Let us recall a few classical definitions.

(i) Let $z_1, \ldots, z_r$ be independent standard normal. Let

$$u = z_1^2 + z_2^2 + \cdots + z_r^2.$$

We say that $u$ has the $\chi^2$-distribution with $v$ degrees of freedom. Symbol: $u \sim \chi^2(r)$.

Comment: We know that

$$\chi^2(r) = \Gamma\left(\frac{r}{2}, \frac{1}{2}\right).$$

(ii) Let $z \sim N(0,1)$ independent of $u \sim \chi^2(v)$. The random variable

$$T = \frac{z}{\sqrt{u/n}}$$

has the $t$-distribution with

$r$ degrees of freedom.

Symbol:   $T \sim t_r$

(iii)   Let   $U \sim X^2(m)$ independent of

$V \sim X^2(n)$.   Then

$$F = \frac{U/m}{V/n}$$

has the F - distribution with

$m, n$ degrees of freedom.

Symbol:   $F \sim F_{m,n}$.

In linear algebra a quadratic

form is the expression

$$Q = \sum_{k,\ell = 1}^{n} a_{k\ell} x_k x_\ell .$$

If we collect the constants

$a_{k,\ell}$ into a matrix $\underline{A}$ and define

$\underline{x}^T = (x_1, x_2, .., x_n)$   we can write

$$Q = \underline{a}^T \underline{A} \, \underline{a}$$

Without loss of generality $A$ can be assumed to be symmetric.

In statistics $\underline{x}$ will be replaced by a random vector $\underline{X}$ so that $Q$ will be a random variable. Such expressions arise frequently. We will be interested in such random variables and their distributions.

We will need some facts from linear algebra.

<u>Definition</u> A symmetric matrix $\underline{H}$ is called idempotent if $\underline{H}^2 = \underline{H}$.

If $\underline{x}$ is an eigenvector of $\underline{H}$ we have

$$\underline{H}^2 \underline{x} = \underline{H} \cdot (\underline{H} \underline{x})$$

$$= \underline{H} \cdot \lambda \underline{x}$$

$$= \lambda^2 \cdot \underline{x}$$

$$= \lambda \underline{x}$$

If $\underline{x} \neq o$ this means $\lambda^2 = \lambda$. The eigenvalues of $\underline{H}$ can only be in $\{0, 1\}$. Because every symmetric matrix can be diagonalized we have

$$\underline{H} = \underline{Q} \, \text{diag} \{1, \ldots, 1, 0, \ldots, 0\} \underline{Q}^T$$

for an orthogonal matrix. It is clear that the rank of $\underline{H}$ in equal to the number of non-zero eigenvalues.

**Definition:** If $\underline{A}$ is a square ($n \times n$) matrix its trace is

$$\text{Tr}(\underline{A}) = \sum_{k=1}^{n} a_{kk}.$$

From linear algebra we will borrow the fact that for matrices $A\,(p \times q)$ and $B\,(q \times p)$ we have

$$\text{Tr}(\underline{A} \cdot \underline{B}) = \text{Tr}(\underline{B} \cdot \underline{A}).$$

From this we have for idempotent $\underline{H}$:

$$\text{Tr}(\underline{H}) = \text{Tr}\left(\underline{Q}\,\text{diag}(1, .., 1, 0, 0)\,\underline{Q}^T\right)$$

$$= \text{Tr}\left(\text{diag}(1, .., 1, 0..0)\,\underbrace{\underline{Q}^T\underline{Q}}_{=I}\right)$$

$$= \text{rank}(\underline{H}).$$

Let $\underline{z}$ be a vector with
independent standard normal
components. In the notation
of the multivariate normal
random vectors we have

$\underline{z} \sim N(\underline{0}, \underline{I})$. Let $\underline{H}$ be idempotent.

Define

$$U = \underline{z}^T \underline{H} \underline{z}$$

Rewrite

$$\underline{U} = \underline{z}^T \underline{Q}^T \, \text{diag} \langle 1, 1, \ldots, 1, 0, 0 \rangle \, \underline{Q} \underline{z} .$$

We know that $\hat{\underline{z}} = \underline{Q} \underline{z} \sim N(\underline{0}, \underline{Q} \underline{I} \underline{Q}^T)$.

But $\underline{Q} \underline{I} \underline{Q}^T = \underline{I}$ so $\hat{\underline{z}} \sim N(\underline{0}, \underline{I})$.

It follows

$$U = \hat{z}_1^2 + \cdots + \hat{z}_r^2$$

where $r = \text{rank}(\underline{H})$.

The result is still true if

$\underline{z} \sim N(\underline{\mu}, \underline{\underline{I}})$ and $\underline{\underline{H}}\underline{\mu} = 0$.

We simply rewrite

$$u = \underline{z}^T \underline{\underline{H}} \underline{z}$$

$$= (\underline{z} - \underline{\mu})^T \cdot \underline{\underline{H}} \underbrace{(\underline{z} - \underline{\mu})}_{N(\underline{0}, \underline{\underline{I}})}.$$

We can take this further.

Theorem 1.7 (Cochran) Let

$\underline{\underline{H}}_1, \ldots, \underline{\underline{H}}_s$ be idempotent

matrices such that $\underline{\underline{H}}_k \cdot \underline{\underline{H}}_l = 0$

for $k \neq l$. Then the

random variables

$$u_k = \underline{z}^T \underline{\underline{H}}_k \underline{z} \quad \text{are independent}$$

and $u_k \sim \chi^2(\text{rank}(\underline{\underline{H}}_k))$.

**Proof:**   The   vector

$$\begin{pmatrix} \underline{H_1 z_k} \\ \underline{H_2 z} \\ \\ \underline{H_3 z} \end{pmatrix}$$

is   multivariate

normal   because   its   components
are   linear   combinations   of
a   multivariate   normal   vector.
But

$$\text{cov}\left( \underline{H_k z}, \underline{H_e z} \right) = \underline{H_k} \cdot H_e^T$$

$$= \underline{H_k} \cdot \underline{H_e}$$

$$= 0$$

by   assumption.   So   all

$$\underline{H_1 z}, \dots, \underline{H_s z}$$   are   independent.

But

$$U_k = \underline{z}^T \underline{H_k} \, \underline{z}$$

$$= \underline{z}^T \underline{H_k}^T \cdot \underline{H_k} \, \underline{z}$$

$$= \left( \underline{H_e z} \right)^T \left( \underline{H_k} \cdot \underline{z} \right)$$

This implies that all $U_k$ are independent. The fact that $U_k \sim \chi^2(\text{rank}(\underline{H}_k))$ we have already proved.

Example: If $\underline{H}$ is idempotent so is $\underline{I} - \underline{H}$. But

$$\underline{H}(\underline{I} - \underline{H}) = \underline{H} - \underline{H}^2$$
$$= \underline{H} - \underline{H}$$
$$= 0.$$

So for $\underline{z} \sim N(0, \underline{I})$ we have that $U = \underline{z}^T \underline{H} \underline{z}$ and $V = \underline{z}^T(\underline{I} - \underline{H})\underline{z}$ are independent and $\chi^2$-distributed.

Comment: The independence assertion of Theorem 1.7 is true even if $\underline{z} \sim N(\mu, \underline{I})$.

**Example:** In regression we assume

$$\underline{Y} = \underline{X} \cdot \beta + \underline{\varepsilon}$$

with $E(\underline{\varepsilon}) = 0$ and $\text{Var}(\underline{\varepsilon}) = \delta^2 \cdot \underline{I}$.

Assume further that $\underline{\varepsilon} \sim N(0, \delta^2 \cdot \underline{I})$.

The best unbiased linear estimator of $\beta$ is

$$\hat{\beta} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \cdot \underline{Y}$$

We define residuals as

$$\hat{\underline{\varepsilon}} = (\underline{I} - \underline{X}(\underline{X}^T \underline{X})^{-1} \underline{X}^T) \underline{Y}$$

Let $\underline{H} = \underline{X}(\underline{X}^T \underline{X})^{-1} \underline{X}^T$. We have that $\underline{H}$ is symmetric and

$$\underline{H}^2 = \underline{X}(\underline{X}^T\underline{X})^{-1} \underline{X}^T \underline{X} (\underline{X}^T\underline{X})^{-1} \underline{X}^T = \underline{H}.$$

The matrix $\underline{H}$ is idempotent. In fact it is the projection onto the subspace spanned by the columns of $X$ and hence

$$\text{rank}(\underline{H}) = \text{rank}(\underline{x}) = m.$$

Since $\begin{pmatrix} \hat{\underline{\beta}} \\ \hat{\underline{\varepsilon}} \end{pmatrix}$ is multivariate normal we have that

$$\text{cov}(\hat{\underline{\beta}}, \hat{\underline{\varepsilon}}) = \text{cov}((\underline{x}^T\underline{x})^{-1}\underline{x}^T Y, (\underline{I} - \underline{H})\underline{Y})$$

$$= (\underline{x}^T\underline{x})^{-1}\underline{x}^T \cdot \sigma^2 \cdot \underline{I} \left(\underline{I} - x(\underline{x}^T\underline{x})^{-1}\underline{x}^T\right)$$

$$= \sigma^2 \left((\underline{x}^T x)^{-1}\underline{x}^T - (\underline{x}^T\underline{x})(\underline{x}^T\underline{x})^{-1}(\underline{x}^T\underline{x})^{-1}\underline{x}^T\right)$$

$$= 0.$$

This means that $\hat{\underline{\beta}}$ and $\hat{\underline{\varepsilon}}$ are independent.

Further

$$\hat{\varepsilon}^T \hat{\varepsilon} = [(\underline{I} - \underline{H})\underline{Y}]^T [(\underline{I} - \underline{H})\underline{Y}]$$

$$= \underline{Y}^T (\underline{I} - \underline{H}) \underline{Y}$$

$$= (\underline{Y} - \underline{X}\beta)^T (\underline{I} - \underline{H})(\underline{Y} - \underline{X}\beta)$$

because $(\underline{I} - \underline{H})\underline{X} = (I - \underline{X}(\underline{x}^T\underline{x})^{-1}\underline{x}^T\underline{x})$

$$= I - I$$

$$= 0.$$

But $\underline{Y} - \underline{X}\beta \sim N(0, \delta^2 I)$. This

implies that $\quad \underline{\varepsilon}^T \sim \delta$

$$\hat{\varepsilon}^T \underline{\varepsilon} \sim X^2(\ n-m)$$

where $\quad m = \text{rank}(\underline{x}).$

If $\hat{\beta}_k$ is the $k$-th component of $\hat{\beta}$ it is independent of $\hat{\underline{\varepsilon}}$.

We know that $\hat{\beta}_k \sim N(\beta_k, \sigma^2 \cdot c_{kk})$ where $\underline{C} = (\underline{X}^T \underline{X})^{-1}$ and $c_{kk}$ is the diagonal element. The expression

$$T = \frac{\dfrac{\hat{\beta}_k - \beta_k}{\sqrt{c_{kk}}}}{\sqrt{\hat{\underline{\varepsilon}}^T \underline{\varepsilon} / (n-m)}}$$

is a quotient of a $N(0, \sigma^2)$ random variable and the square root of a $\sigma^2 \chi^2 (n-m)/n-m$ random variable hence

$$T \sim t_m.$$

This is what you see on regression printouts.

## 2.  Abstract expected values.

**Example :**   Players A and B each get 5 cards from a well shuffled deck of cards. Let X be the number of aces of player A and Y be the number of aces of player B. From elementary probability we know that

$$E(Y \mid X = k) = 5 \cdot \frac{4-k}{47}$$

for $k = 0, 1, 2, 3, 4$. The right side is a function of $k$. Call it $\Psi(k)$.

Suppose we ask about the

conditional expectation before
the cards are dealt. At that
time $t$ is unknown and is
a random variable. But then
the conditional expectation is
also a random variable! Which
one? Obviously $\Psi(X)$. This is
Kolmogorov's idea of a random
variable. that plays the role of
abstract conditional expectation
$E(Y|X) = \Psi(X)$.

Let $X, Y$ be a pair of discrete
random variables and denote
$E(Y|X=x) = \Psi(x)$.

Assume $E|Y| < \infty$ and let $g$ be a bounded function. We compute

$$E[\,\psi(x)\,g(x)\,]$$

$$= \sum_x \psi(x)\,g(x)\,P(X = x)$$

$$= \sum_x \underbrace{\sum_y y \cdot P(Y = y \mid X = x)}_{\psi(x)}\, g(x)\, P(X = x)$$

$$= \sum_{x,y} y\, g(x)\, P(Y = y \mid X = x)\, P(X = x)$$

$$= \sum_{x,y} y\, g(x)\, P(X = x, Y = y)$$

$$= E[\,Y \cdot g(x)\,].$$

This expression uniquely determines the function $\psi$ because we can take $g(x) = I_{A \times y}$.

This is the idea for the general mathematical definition of $E(Y|X)$.

## Definition :

(i) Let $Y$ be a random variable with $E|Y| < \infty$. The conditional expectation of $Y$ with respect of $X$ is the function $\psi(x)$ such that for any bounded $g$ we have

$$E(Y g(x)) = E(\psi(x) g(x)).$$

(ii) Let $Y$ be a random variable with $E(|Y|) < \infty$. The conditional expectation of $Y$ given $X_1, X_2, \ldots, X_n$ is a function $\psi(X_1, X_2, \ldots, X_n)$

such that for any bounded $g: \mathbb{R}^n \to \mathbb{R}$ we have

$$E\left(Y\, g(x_1, \dots, x_n)\right) = E\left[Y(x_1 \dots x_n)\, g(x_1, \dots, x_n)\right]$$

### Theorem 2.1 (Radon - Kolmogorov)

The conditional expectation of $Y$ given $X_1, X_2, \dots, X_n$ for $Y$ with $E|Y| < \infty$ exists and is uniquely determined.

Proof: R. Durrett, Probability: Theory and Examples, 2nd Ed., Duxbury 1995

Let us look at elementary properties of conditional expectation. We will write

$$E(Y \mid X_1, \dots, X_n) = E(Y \mid \underline{x}).$$

**Theorem 2.2** ( Linearity ) If $Y_1, Y_2$ are random variables with $E|Y_1| < \infty$ and $E|Y_2| < \infty$. Then

$$E(\alpha Y_1 + \beta Y_2 \mid \underline{x})$$

$$= \alpha E(Y_1 \mid \underline{x}) + \beta E(Y_2 \mid \underline{x})$$

**Proof:** The right side should satisfy the definition. We compute

$$E[(\alpha E(Y_1 \mid \underline{x}) + \beta E(Y_2 \mid \underline{x})) \cdot g(\underline{x})]$$

$$= \alpha E[E(Y_1 \mid \underline{x}) g(\underline{x})]$$

$$+ \beta E[E(Y_2 \mid \underline{x}) g(\underline{x})]$$

$$= \alpha E(Y_1 g(\underline{x})) + \beta E(Y_2 g(x))$$

$$= E[(\alpha Y_1 + \beta Y_2) g(\underline{x})].$$

**Theorem 2.3** : (Tower property).

Let $E|Y| < \infty$ and $m < n$.

Then

$$E[Y | X_1, \dots, X_m]$$

$$= E[E(Y | X_1, \dots, X_n) | X_1, \dots, X_m].$$

**Proof** : We compute

$$E\left[E[E(Y | X_1, \dots, X_n) | X_1, \dots, X_m] \, g(X_1, \dots, X_m)\right]$$

$$\overset{def.}{=} E\left[E(Y | X_1, \dots X_n) \, g(X_1, \dots, X_m)\right]$$

$$\overset{def}{=} E(Y \, g(X_1, \dots, X_m))$$

This concludes the proof.

**Theorem 2.4**: Let $f : \mathbb{R}^n \to \mathbb{R}$ be a function such that $E|f(\underline{x})Y| < \infty$. Then

$$E(f(\underline{x})Y \mid \underline{x}) = f(\underline{x}) E(Y \mid \underline{x}).$$

**Proof**: Left to the reader.

In a similar way we can think about variances. They are functions of $\underline{x}$. To get a formal definition we just replace expectations by conditional expectation.

**Definitions** : The (abstract) conditional variance of $Y$ given $\underline{x}$ is defined by

$$\text{var}(Y \mid \underline{x}) = E(Y^2 \mid \underline{x}) - E(Y \mid \underline{x})^2.$$

**Examples** :

(i) Let $X_1, X_2, \ldots, X_u$ be independent and equally distributed. Let $S_u = X_1 + X_2 + \cdots + X_u$. What is $E(X_1 \mid S_u)$ ?

By symmetry $(X_k, S_u)$ have the same distribution and so

$$E(X_k \, g(S_u)) = E[Y(S_u) \, g(S_u)]$$

for the same $Y$.

By linearity

$$E(X_1 | S_n) + \cdots + E(X_n | S_n)$$

$$= E(S_n | S_n)$$

$$= S_n$$

But all the terms on the left are equal to $\psi(S_n)$. So

$$E(X_1 | S_n) = \frac{S_n}{n}$$

(ii) Let $X, Z$ be discrete and independent. Let $f$ be a function such that $E|f(X, Z)| < \infty$. Define $\psi(x) = E[f(x, Z)]$. We claim that

$$E[f(X, Z) | X] = \psi(X).$$

We compute

$$E[f(x) g(x)]$$

$$= \sum_x f(x) g(x) P(X = x)$$

$$= \sum_x \sum_z f(x, z) P(Z = z) P(X = x) g(x)$$

(indep)
$$= \sum_{x, z} f(x, z) P(X = x, Z = z) g(x)$$

$$= E[f(x, z) g(x)].$$

In general this is also true but slightly more difficult to prove. A more general version for vectors is also true. If $\underline{x}, \underline{z}$ are independent we have

$$E[f(\underline{x}, \underline{z}) | \underline{z}] = f(\underline{x}) \quad \text{where}$$

$$f(\underline{x}) = E[f(\underline{x}, \underline{z})].$$

(iii)    Compute

$$E[E(Y|\underline{x})]$$

- $E[E(Y|\underline{x}) \cdot g(\underline{x})]$    $g \equiv 1$

$$= E(Y)$$

So

$$\boxed{E[E(Y|\underline{x})] = E(Y)}$$

For variances we get

$$E[var(Y|\underline{x})]$$

$$= E[E(Y^2|\underline{x}) - E(Y|\underline{x})^2]$$

$$= E(Y^2) - E[E(Y|\underline{x})^2]$$

$$= E(Y^2) - E(Y)^2$$

$$- \left( E[E(Y|\underline{x})^2] \right.$$

$$\left. - E(Y)^2 \right)$$

$$= \text{var}(y) - \text{var}(E(y|\underline{x}))$$

Rearrange to get

$$\boxed{\begin{aligned} \text{var}(y) &= E[\text{var}(y|\underline{x})] \\ &\quad + \text{var}(E(y|x)) \end{aligned}}$$

This is a well known variance decomposition formula.

(iv) Let $\underline{x}$ be multivariate normal. Write

$$\underline{x} \sim N\left( \begin{pmatrix} \alpha_1 \\ \underline{\mu}^2 \end{pmatrix}, \begin{pmatrix} \delta_{11} & \Sigma_{12} \\ \Sigma_{22} & \Sigma_{22} \end{pmatrix} \right).$$

What is $E(x_1 | x_2, \ldots, x_n)$ ?

$\text{var}(x_1 | x_2, \ldots, x_n)$.

Preliminary calculation: Let

$\underline{X}, Y$ be independent. In

this case $E(Y|\underline{X}) = E(Y)$ and

$\text{var}(Y|\underline{X}) = \text{var}(Y)$. We know

that $X_1 - \underline{\Sigma}_{12} \underline{\Sigma}_{22}^{-1} \underline{X}^2$ is

independent of $\underline{X}^2$ so

$$E\left( X_1 - \underline{\Sigma}_{12} \underline{\Sigma}_{22}^{-1} \underline{X}^2 \mid \underline{X}^2 \right)$$

$$= E\left( X_1 - \underline{\Sigma}_{12} \underline{\Sigma}_{22}^{-1} \underline{X}^2 \right)$$

$$= \mu_1 - \underline{\Sigma}_{12} \underline{\Sigma}_{22} \underline{\mu}^2$$

By linearity, however,

$$E\left( \underline{\Sigma}_{12} \underline{\Sigma}_{22}^{-1} \underline{X}^2 \mid \underline{X}^2 \right) = \underline{\Sigma}_{12} \underline{\Sigma}_{22}^{-1} \underline{X}^2.$$

Putting all the pieces together

$$E(X_1 | \underline{X}^2) = \mu_1 + \underline{\Sigma}_{12} \underline{\Sigma}_{11}^{-1} \left( \underline{X}^2 - \underline{\mu}^2 \right)$$

For any function $h$ we have

$$\text{var}(Y + h(\underline{x}) \mid \underline{x}) = \text{var}(Y \mid \underline{x}).$$

The reader can check that. So

$$\text{var}(X_1 \mid \underline{X}^2)$$

$$= \text{var}(X_1 - \Sigma_{12} \Sigma_{22} \underline{X}^2 \mid \underline{X}^2)$$

(indep)
$$= \text{var}(X_1 - \Sigma_{12} \Sigma_{22}^{-1} \underline{X}^2)$$

$$= \text{var}(X_1) + \text{var}(\Sigma_{12} \Sigma_{22}^{-1} \underline{X}^2)$$

$$+ 2 \text{cov}(X_1, \Sigma_{12} \Sigma_{22}^{-1} \underline{X}^2)$$

$$= \Sigma_{11} + \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{22} \Sigma_{22}^{-1} \Sigma_{21}$$

$$+ 2 \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

$$= \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

(iv)   What is the best predictor of $Y$ given $X_1, X_2, \ldots, X_n$?

A predictor is a function $f(x_1, x_2, \ldots, x_n)$. The best predictor is the conditional expectation. We prove this by direct calculation:

$$E\left[\left(f(x_1, \ldots, x_n) - Y\right)^2\right]$$

$$= E\left[\left(f(\underline{x}) - \Psi(\underline{x}) + \Psi(\underline{x}) - Y\right)^2\right]$$

$$= E\underbrace{\left[\left(f(\underline{x}) - \Psi(\underline{x})\right)^2\right]}_{> 0.}$$

$$+ E\left[\left(\Psi(\underline{x}) - Y\right)^2\right]$$

$$+ 2\, E\left[\left(f(\underline{x}) - \Psi(\underline{x})\right)\left(\Psi(\underline{x}) - \underline{Y}\right)\right]$$

Side calculation:

$$E\left[ \left( f(\underline{x}) - \gamma(\underline{x}) \right) \left( \gamma(\underline{x}) - y \right) \right]$$

$$= E\left[ E\left[ " - | \underline{x} \right] \right]$$

$$\overset{Thm\ 2.3}{=} E\left[ \left( f(\underline{x}) - \gamma(\underline{x}) \right) \underbrace{E\left[ \gamma(\underline{x}) - y | \underline{x} \right]}_{= 0} \right]$$

$$= 0$$
by def.

$$= 0.$$

So $\gamma(\underline{x})$ is the best predictor!