

We turn now to a technical result, which is important in reducing n -dimensional problems to one-dimensional ones.

Theorem (Cramér-Wold device). The distribution of a random n -vector \mathbf{X} is completely determined by the set of all one-dimensional distributions of linear combinations $\mathbf{t}^T \mathbf{X} = \sum_i t_i X_i$, where \mathbf{t} ranges over all fixed n -vectors.

Proof. $Y := \mathbf{t}^T \mathbf{X}$ has CF

$$\phi_Y(t) := E \exp\{itY\} = E \exp\{it\mathbf{t}^T \mathbf{X}\}.$$

If we know the distribution of each Y , we know its CF $\phi_Y(t)$. In particular, taking $t = 1$, we know $E \exp\{i\mathbf{t}^T \mathbf{X}\}$. But this is the CF of $\mathbf{X} = (X_1, \dots, X_n)^T$ evaluated at $\mathbf{t} = (t_1, \dots, t_n)^T$. But this determines the distribution of \mathbf{X} . //

Thus by the Cramér-Wold device, to define an n -dimensional distribution it suffices to define the distributions of *all linear combinations*.

The Cramér-Wold device suggests a way to *define* the multivariate normal distribution. The definition below seems indirect, but it has the advantage of handling the full-rank and singular cases together ($\rho = \pm 1$ as well as $-1 < \rho < 1$ for the bivariate case).

Definition. An n -vector \mathbf{X} has an n -variate normal distribution iff $\mathbf{a}^T \mathbf{X}$ has a univariate normal distribution for all constant n -vectors \mathbf{a} .

Proposition. (i) Any linear transformation of a multinormal n -vector is multinormal,

(ii) Any vector of elements from a multinormal n -vector is multinormal. In particular, the components are univariate normal.

Proof. (i) If $\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{c}$ (\mathbf{A} an $m \times n$ matrix, \mathbf{c} an m -vector) is an m -vector, and \mathbf{b} is any m -vector,

$$\mathbf{b}^T \mathbf{Y} = \mathbf{b}^T (\mathbf{A}\mathbf{X} + \mathbf{c}) = (\mathbf{b}^T \mathbf{A})\mathbf{X} + \mathbf{b}^T \mathbf{c}.$$

If $\mathbf{a} = \mathbf{A}^T \mathbf{b}$ (an n -vector), $\mathbf{a}^T \mathbf{X} = \mathbf{b}^T \mathbf{A}\mathbf{X}$ is univariate normal as \mathbf{X} is multinormal. Adding the constant $\mathbf{b}^T \mathbf{c}$, $\mathbf{b}^T \mathbf{Y}$ is univariate normal. This holds for

all \mathbf{b} , so \mathbf{Y} is m -variate normal.

(ii) Take a suitable matrix \mathbf{A} of 1s and 0s to pick out the required sub-vector.

Theorem 1. If \mathbf{X} is n -variate normal with mean μ and covariance matrix Σ , its CF is

$$\phi(\mathbf{t}) := E \exp\{i\mathbf{t}^T \mathbf{X}\} = \exp\{i\mathbf{t}^T \mu - \frac{1}{2}\mathbf{t}^T \Sigma \mathbf{t}\}.$$

Proof. By Proposition 1, $Y := \mathbf{t}^T \mathbf{X}$ has mean $\mathbf{t}^T \mu$ and variance $\mathbf{t}^T \Sigma \mathbf{t}$. By definition of multinormality, $Y = \mathbf{t}^T \mathbf{X}$ is univariate normal. So Y is $N(\mathbf{t}^T \mu, \mathbf{t}^T \Sigma \mathbf{t})$, so Y has CF

$$\phi_Y(t) := E \exp\{itY\} = E \exp\{itt^T \mathbf{X}\} = \exp\{itt^T \mu - \frac{1}{2}t^2 \mathbf{t}^T \Sigma \mathbf{t}\}.$$

Taking $t = 1$ (as in the proof of the Cramér-Wold device),

$$E \exp\{i\mathbf{t}^T \mathbf{X}\} = \exp\{i\mathbf{t}^T \mu - \frac{1}{2}\mathbf{t}^T \Sigma \mathbf{t}\}. \quad //$$

Corollary. The components of \mathbf{X} are independent iff Σ is diagonal.

Proof. The components are independent iff the joint CF factors into the product of the marginal CFs. This factorization takes place, into $\prod_j \exp\{i\mu_j t_j - \frac{1}{2}\sigma_{jj}t_j^2\}$, in the diagonal case only. //

Recall that a covariance matrix Σ is always

- (a) symmetric ($\sigma_{ij} = \sigma_{ji}$, as $\sigma_{ij} = \text{cov}(X_i, X_j)$),
- (b) non-negative definite, written $\Sigma \geq 0$: $\mathbf{a}^T \Sigma \mathbf{a} \geq 0$ for all n -vectors \mathbf{a} .

Suppose that Σ is, further, *positive definite*, written $\Sigma > 0$:

$$\mathbf{a}^T \Sigma \mathbf{a} > 0 \quad \text{unless} \quad \mathbf{a} = \mathbf{0}.$$

The Multinormal Density.

If \mathbf{X} is n -variate normal, $N(\mu, \Sigma)$, its density (in n dimensions) need not exist (e.g. the singular case $\rho = \pm 1$ with $n = 2$). But if $\Sigma > \mathbf{0}$ (so Σ^{-1} exists), \mathbf{X} has a density. The link between the multinormal density below and the multinormal MGF above is due to the English statistician F. Y. Edgeworth (1845-1926) in 1893.

Theorem (Edgeworth). If μ is an n -vector, $\Sigma > \mathbf{0}$ a symmetric positive definite $n \times n$ matrix, then

(i)

$$f(\mathbf{x}) := \frac{1}{(2\pi)^{\frac{1}{2}n} |\Sigma|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)\right\}$$

is an n -dimensional probability density function (of a random n -vector \mathbf{X} , say),

(ii) \mathbf{X} has CF $\phi(\mathbf{t}) = \exp\{i\mathbf{t}^T \mu - \frac{1}{2}\mathbf{t}^T \Sigma \mathbf{t}\}$,

(iii) \mathbf{X} is multinormal $N(\mu, \Sigma)$.

Proof. Write $\mathbf{Y} := \Sigma^{-\frac{1}{2}} \mathbf{X}$ ($\Sigma^{-\frac{1}{2}}$ exists as $\Sigma > \mathbf{0}$, by above). Then \mathbf{Y} has covariance matrix $\Sigma^{-\frac{1}{2}} \Sigma (\Sigma^{-\frac{1}{2}})^T$. Since $\Sigma = \Sigma^T$ and $\Sigma = \Sigma^{\frac{1}{2}} \Sigma^{\frac{1}{2}}$, \mathbf{Y} has covariance matrix \mathbf{I} (the components Y_i of \mathbf{Y} are uncorrelated).

Change variables as above, with $\mathbf{y} = \Sigma^{-\frac{1}{2}} \mathbf{x}$, $\mathbf{x} = \Sigma^{\frac{1}{2}} \mathbf{y}$. The Jacobian is (taking $\mathbf{A} = \Sigma^{-\frac{1}{2}}$) $J = \partial \mathbf{x} / \partial \mathbf{y} = \det(\Sigma^{\frac{1}{2}}) = (\det \Sigma)^{\frac{1}{2}}$ by the product theorem for determinants. Substituting, the integrand is

$$\exp\left\{-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)\right\} = \exp\left\{-\frac{1}{2}(\Sigma^{\frac{1}{2}} \mathbf{y} - \Sigma^{\frac{1}{2}}(\Sigma^{-\frac{1}{2}} \mu))^T \Sigma^{-1}(\Sigma^{\frac{1}{2}} \mathbf{y} - \Sigma^{\frac{1}{2}}(\Sigma^{-\frac{1}{2}} \mu))\right\}.$$

Writing $\nu := \Sigma^{-\frac{1}{2}} \mu$, this is

$$\exp\left\{-\frac{1}{2}(\mathbf{y} - \nu)^T \Sigma^{\frac{1}{2}} \Sigma^{-1} \Sigma^{\frac{1}{2}} (\mathbf{y} - \nu)\right\} = \exp\left\{-\frac{1}{2}(\mathbf{y} - \nu)^T (\mathbf{y} - \nu)\right\}.$$

So by the change-of-density formula, \mathbf{Y} has density

$$g(\mathbf{y}) = \frac{1}{(2\pi)^{\frac{1}{2}n} |\Sigma|^{\frac{1}{2}}} \cdot |\Sigma|^{\frac{1}{2}} \cdot \exp\left\{-\frac{1}{2}(\mathbf{y} - \nu)^T (\mathbf{y} - \nu)\right\}.$$

This factorises as

$$\prod_{i=1}^n \frac{1}{(2\pi)^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}(y_i - \nu_i)^2\right\}.$$

So the components Y_i of \mathbf{Y} are independent $N(\nu_i, 1)$. So \mathbf{Y} is multinormal, $N(\nu, \mathbf{I})$.

(i) Taking $A = B = \mathbf{R}^n$, $\int_{\mathbf{R}^n} f(\mathbf{x}) d\mathbf{x} = \int_{\mathbf{R}^n} g(\mathbf{y}) d\mathbf{y} = 1$ as g is a probability density, as above. So f is also a probability density (non-negative and integrates to 1).

(ii) $\mathbf{X} = \Sigma^{\frac{1}{2}} \mathbf{Y}$ is a linear transformation of \mathbf{Y} , and \mathbf{Y} is multivariate normal,

$N(\nu, I)$. So \mathbf{X} is multivariate normal.

(iii) $E\mathbf{X} = \Sigma^{\frac{1}{2}}E\mathbf{Y} = \Sigma^{\frac{1}{2}}\nu = \Sigma^{\frac{1}{2}}\Sigma^{-\frac{1}{2}}\mu = \mu$, $cov\mathbf{X} = \Sigma^{\frac{1}{2}}cov\mathbf{Y}(\Sigma^{\frac{1}{2}})^T = \Sigma^{\frac{1}{2}}\mathbf{I}\Sigma^{\frac{1}{2}} = \Sigma$. So \mathbf{X} is multinormal $N(\mu, \Sigma)$. So its CF is

$$\phi(\mathbf{t}) = \exp\{i\mathbf{t}^T\mu - \frac{1}{2}\mathbf{t}^T\Sigma\mathbf{t}\}. \quad //$$

Note. The inverse Σ^{-1} of the covariance matrix Σ is called the *concentration matrix*, K .

Conditional independence of two components X_i, X_j of a multinormal vector given the others can be identified by vanishing of the (off-diagonal) (i, j) entry k_{ij} in the concentration matrix K . The proof needs the results on conditioning and regression in IV.6 D6 below, and the formula for the inverse of a partitioned matrix; see Problems 6.

Independence of Linear Forms

Given a normally distributed random vector $\mathbf{x} \sim N(\mu, \Sigma)$ and a matrix A , one may form the *linear form* $A\mathbf{x}$. One often encounters several of these together, and needs their joint distribution – in particular, to know when these are independent.

Theorem 3. Linear forms $A\mathbf{x}$ and $B\mathbf{x}$ with $\mathbf{x} \sim N(\mu, \Sigma)$ are independent iff

$$A\Sigma B^T = 0.$$

In particular, if A, B are symmetric and $\Sigma = \sigma^2 I$, they are independent iff

$$AB = 0.$$

Proof. The joint CF is

$$\phi(\mathbf{u}, \mathbf{v}) := E \exp\{i\mathbf{u}^T A\mathbf{x} + i\mathbf{v}^T B\mathbf{x}\} = E \exp\{i(A^T\mathbf{u} + B^T\mathbf{v})^T \mathbf{x}\}.$$

This is the CF of \mathbf{x} at argument $\mathbf{t} = A^T\mathbf{u} + B^T\mathbf{v}$, so

$$\begin{aligned} \phi(\mathbf{u}, \mathbf{v}) &= \exp\{i(\mathbf{u}^T A + \mathbf{v}^T B)\mu - \frac{1}{2}(A^T\mathbf{u} + B^T\mathbf{v})^T \Sigma (A^T\mathbf{u} + B^T\mathbf{v})\} \\ &= \exp\{i(\mathbf{u}^T A + \mathbf{v}^T B)\mu - \frac{1}{2}[\mathbf{u}^T A \Sigma A^T \mathbf{u} + \mathbf{u}^T A \Sigma B^T \mathbf{v} + \mathbf{v}^T B \Sigma A^T \mathbf{u} + \mathbf{v}^T B \Sigma B^T \mathbf{v}]\}. \end{aligned}$$

This factorises into a product of a function of \mathbf{u} and a function of \mathbf{v} iff the two cross-terms in \mathbf{u} and \mathbf{v} vanish, that is, iff $A\Sigma B^T = 0$ and $B\Sigma A^T = 0$; by symmetry of Σ , the two are equivalent.

4. Quadratic forms in normal variates

We give a brief treatment of this important material; for full detail see e.g. [BF], 3.4 – 3.6. Recall (IV.3, D5)

- (i) with $x \sim N(\mu, \Sigma)$, linear forms Ax , BX are independent iff $A\Sigma B^T = 0$;
- (ii) for a projection, $P^2 = P$ (P is *idempotent*); for a symmetric projection, $P^T P = P$.

We restrict attention, for simplicity, to $\mu = 0$, $\Sigma = \sigma^2 I$, $x \sim N(0, \sigma^2 I)$.

It turns out that the distribution theory relevant to regression depends on *quadratic forms in normal variates*, $x^T A x$ for a normally distributed random vector x , and that we can confine attention to projection matrices. For P a symmetric projection,

$$x^T P x = x^T P^T P x = (P x)^T (P x),$$

which reduces from *quadratic forms* to *linear* forms – which are much easier! So: if $x^T P_1 x$, $x^T P_2 x$ are quadratic forms in normal vectors x , with P_1, P_2 projections, $x^T P_1 x$ and $x^T P_2 x$ are independent iff

$$P_1 P_2 = 0 :$$

P_1, P_2 are *orthogonal projections*. Recall that projections P_1, P_2 are *orthogonal* if their ranges are orthogonal subspaces, i.e.

$$(P_1 x) \cdot (P_2 x) = 0 \quad \forall x : \quad x^T P_1^T P_2 x = 0 \quad \forall x; \quad P_1^T P_2 = 0 \quad \forall x; \quad P_1 P_2 = 0$$

for P_i symmetric. Note that for P a projection, $I - P$ is a projection orthogonal to it:

$$(I - P)^2 = I - 2P + P^2 = I - 2P + P = I - P; \quad P(I - P) = P - P^2 = P - P = 0.$$

If λ is an eigenvalue of A , λ^2 is an eigenvalue of A^2 (check). So if a projection P has eigenvalue λ , $\lambda^2 = \lambda$: $\lambda = 0$ or 1 . Also, the trace is the sum of the eigenvalues; for a projection, this is the number of non-zero eigenvalues; this is the rank. So:

For a projection, the eigenvalues are 0 or 1, and the trace is the rank.

By Spectral Decomposition (III.1 D4), a symmetric projection matrix P can be diagonalised by an orthogonal transformation O to a diagonal matrix D :

$$O^T P O = D, \quad P = O D O^T;$$

as above, the diagonal entries d_{ii} are 0 or 1, and we may re-order so that the 1s come first. So with $y := O^T x$,

$$x^T P x = x^T O D O^T x = y^T D y = y_1^2 + \dots + y_r^2.$$

Normality is preserved under orthogonal transformations (check!), so also $y \sim N(0, \sigma^2 I)$. So $y_1^2 + \dots + y_r^2$ is σ^2 times the sum of r independent squares of standard normal variates, and this sum is $\chi^2(r)$ (by definition of chi-square):

$$x^T P x \sim \sigma^2 \chi^2(r).$$

If P has rank r , $I - P$ has rank $n - r$ (where n is the sample size – the dimension of the vector space we are working in):

$$x^T (I - P) x \sim \sigma^2 \chi^2(n - r),$$

and the two quadratic forms are independent.

It turns out that all this can be generalised, to the sum of several projections, not just two. This result – the key to all the distribution theory in Regression – is *Cochran's theorem* (William G. COCHRAN (1909-1980) in 1934); [BF] Th. 3.27):

Theorem (Cochran's Theorem). If

$$I = P_1 + \dots + P_k$$

with each P_i a symmetric projection with rank n_i , then

- (i) the ranks sum: $n = n_1 + \dots + n_k$;
- (ii) each quadratic form $Q_i := x^T P_i x \sim \sigma^2 \chi^2(n_i)$;
- (iii) Q_1, \dots, Q_k are mutually independent;
- (iv) P_1, \dots, P_k are mutually orthogonal: $P_i P_j = 0$ for $i \neq j$.

The quadratic forms that we encounter in Statistics are called *sums of squares* (SS) – for *regression* (SSR), for *error* (SSE), for the *hypothesis* (SSH), etc.

Recall the definition of the *Fisher F-distribution* with degrees of freedom (df) m and n (note the order): $F(m, n)$ is the distribution of the ratio

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

where U, V are independent chi-square random variables with df m, n (see e.g. [BF] 2.3 for the explicit formula for the density, but we shall not need this).

Recall also (or, if you have not met these, take a look at a textbook):

(i) *Analysis of variance (ANOVA)* (see e.g. [BF] Ch. 2). Here one tests for differences between the *means* of different (normal) populations by analysing *variances*. Specifically, one looks at *within-groups* variability and *between-groups* variability, and *rejects* the null hypothesis of no difference between the group means if the second is *too big* compared to the first. As above, one forms the relevant F -statistic, and rejects if F is too big. Here one has *qualitative* factors (which group?).

(ii) *Analysis of Covariance (ANCOVA)* (see e.g. [BF] Ch. 5). Similarly for ANCOVA, where one has both qualitative factors (as with ANOVA) and quantitative ones (covariates), as with Regression.

(iii) Tests of linear hypotheses in Regression (see e.g. [BF] Ch. 6). Here we reject if SSH is too big compared to SSE.

5. Estimation theory for the multivariate normal.

Given a sample x_1, \dots, x_n from the multivariate normal $N_p(\mu, \Sigma)$, form the *sample mean* (vector) and the *sample covariance matrix* as in the one-dimensional case:

$$\bar{x} := \frac{1}{n} \sum_{i=1}^n x_i, \quad S := \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^T (x_i - \bar{x}).$$

The likelihood for a sample of size 1 is

$$L(x|\mu, \Sigma) = (2\pi)^{-p/2} |\Sigma|^{-1/2} \exp\left\{-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right\},$$

so the likelihood for a sample of size n is

$$L = (2\pi)^{-np/2} |\Sigma|^{-n/2} \exp\left\{-\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1}(x_i - \mu)\right\}.$$

Writing

$$x_i - \mu = (x_i - \bar{x}) - (\mu - \bar{x}),$$

$$\sum_1^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) = \sum_1^n (x_i - \bar{x})^T \Sigma^{-1} (x_i - \bar{x}) + n(\bar{x} - \mu)^T \Sigma^{-1} (\bar{x} - \mu)$$

(the cross-terms cancel as $\sum_1^n (x_i - \bar{x}) = 0$). The summand in the first term on the right is a scalar, so is its own trace. Since $\text{trace}(AB) = \text{trace}(BA)$ and $\text{trace}(A + B) = \text{trace}(B + A)$,

$$\begin{aligned} \text{trace}\left(\sum_1^n (x_i - \bar{x})^T \Sigma^{-1} (x_i - \bar{x})\right) &= \text{trace}\left(\Sigma^{-1} \sum_1^n (x_i - \bar{x})^T (x_i - \bar{x})\right) \\ &= \text{trace}(\Sigma^{-1} \cdot nS) = n \text{trace}(\Sigma^{-1} S). \end{aligned}$$

Combining,

$$L = (2\pi)^{-np/2} |\Sigma|^{-n/2} \exp\left\{-\frac{1}{2}n \text{trace}(\Sigma^{-1} S) - \frac{1}{2}n(\bar{x} - \mu)^T \Sigma^{-1} (\bar{x} - \mu)\right\}.$$

Write

$$V := \Sigma^{-1}$$

(‘V for variance’); then

$$\ell = \text{const} - \frac{1}{2}n \text{trace}(VS) - (\bar{x} - \mu)^T V (\bar{x} - \mu).$$

So by the Fisher-Neyman Theorem, (\bar{x}, S) is sufficient for (μ, Σ) . It is in fact minimal sufficient (Problems 2 Q2).

These natural estimators are in fact the MLEs:

Theorem. For the multivariate normal $N_p(\mu, \Sigma)$, \bar{x} and S are the maximum likelihood estimators for μ, Σ .

Proof. Write $V = (v_{ij}) := \Sigma^{-1}$. By above, the likelihood is

$$L = \text{const} \cdot |V|^{n/2} \exp\left\{-\frac{1}{2}n \text{trace}(VS) - \frac{1}{2}n(\bar{x} - \mu)^T V (\bar{x} - \mu)\right\},$$

so the log-likelihood is

$$\ell = c + \frac{1}{2}n \log |V| - \frac{1}{2}n \text{trace}(VS) - \frac{1}{2}n(\bar{x} - \mu)^T V (\bar{x} - \mu).$$