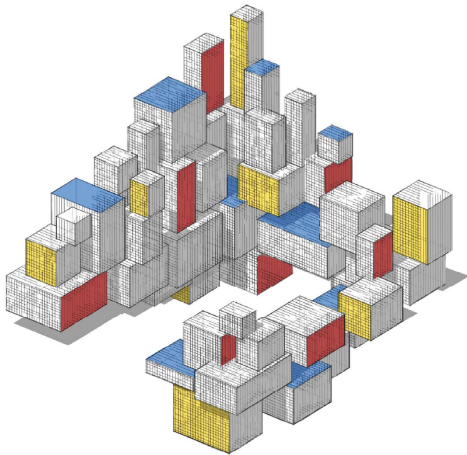


Normal Linear Models



Purpose

In this lecture we discuss:

- Properties of multivariate normal distributions
- Linear model + normal error terms = normal linear models.

Multivariate Normal Distribution

The multivariate normal (or Gaussian) distribution plays a central role in data science and machine learning.

Let Z_1, \dots, Z_n be independent and standard normal random variables. The joint pdf of $\mathbf{Z} = [Z_1, \dots, Z_n]^\top$ is given by

$$f_{\mathbf{Z}}(\mathbf{z}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z_i^2} = (2\pi)^{-\frac{n}{2}} e^{-\frac{1}{2}\mathbf{z}^\top \mathbf{z}}, \quad \mathbf{z} \in \mathbb{R}^n.$$

We write $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, where \mathbf{I} is the identity matrix. Consider the affine transformation

$$\mathbf{X} = \boldsymbol{\mu} + \mathbf{B} \mathbf{Z}$$

for some $m \times n$ matrix \mathbf{B} and m -dimensional vector $\boldsymbol{\mu}$. Then \mathbf{X} has a **multivariate normal** or *multivariate Gaussian* distribution with mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$. We write $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

Affine Transforms are Again Normal

The following theorem states that any affine combination of independent multivariate normal random variables is again multivariate normal.

Theorem: Affine Transform of Normal Random Vectors

Let X_1, X_2, \dots, X_r be independent m_i -dimensional normal random vectors, with $X_i \sim \mathcal{N}(\mu_i, \Sigma_i)$, $i = 1, \dots, r$. Then, for any $n \times 1$ vector a and $n \times m_i$ matrices B_1, \dots, B_r ,

$$a + \sum_{i=1}^r B_i X_i \sim \mathcal{N}\left(a + \sum_{i=1}^r B_i \mu_i, \sum_{i=1}^r B_i \Sigma_i B_i^\top\right).$$

Proof

Denote the n -dimensional random vector in the left-hand side of (??) by \mathbf{Y} . By definition, each \mathbf{X}_i can be written as $\boldsymbol{\mu}_i + \mathbf{A}_i \mathbf{Z}_i$, where the $\{\mathbf{Z}_i\}$ are independent (because the $\{\mathbf{X}_i\}$ are independent), so that

$$\mathbf{Y} = \mathbf{a} + \sum_{i=1}^r \mathbf{B}_i (\boldsymbol{\mu}_i + \mathbf{A}_i \mathbf{Z}_i) = \mathbf{a} + \sum_{i=1}^r \mathbf{B}_i \boldsymbol{\mu}_i + \sum_{i=1}^r \mathbf{B}_i \mathbf{A}_i \mathbf{Z}_i,$$

which is an affine combination of independent standard normal random vectors. Hence, \mathbf{Y} is multivariate normal.

By linearity of the expectation, we have $\mathbb{E}\mathbf{Y} = \mathbf{a} + \sum_{i=1}^r \mathbf{B}_i \boldsymbol{\mu}_i$. And the covariance matrix is

$$\mathbb{C}\text{ov}(\mathbf{Y}) = \sum_{i=1}^r \mathbb{C}\text{ov}(\mathbf{B}_i \mathbf{A}_i \mathbf{Z}_i) = \sum_{i=1}^r \mathbf{B}_i \mathbf{A}_i \mathbb{C}\text{ov}(\mathbf{Z}_i) \mathbf{A}_i^\top \mathbf{B}_i^\top = \sum_{i=1}^r \mathbf{B}_i \boldsymbol{\Sigma}_i \mathbf{B}_i^\top.$$

Marginals are Again Normal

The next theorem shows that the distribution of a subvector of a multivariate normal random vector is again normal.

Theorem: Marginal Distributions of Normal Vectors

Let $X \sim \mathcal{N}(\mu, \Sigma)$ be an n -dimensional normal random vector. Decompose X , μ , and Σ as

$$X = \begin{bmatrix} X_p \\ X_q \end{bmatrix}, \quad \mu = \begin{bmatrix} \mu_p \\ \mu_q \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \Sigma_p & \Sigma_r \\ \Sigma_r^\top & \Sigma_q \end{bmatrix},$$

where Σ_p is the upper left $p \times p$ corner of Σ and Σ_q is the lower right $q \times q$ corner of Σ . Then, $X_p \sim \mathcal{N}(\mu_p, \Sigma_p)$.

Proof

We give a proof assuming that Σ is positive semidefinite. Let $\mathbf{B}\mathbf{B}^\top$ be the (lower) **Cholesky decomposition** of Σ . We can write

$$\begin{bmatrix} X_p \\ X_q \end{bmatrix} = \begin{bmatrix} \mu_p \\ \mu_q \end{bmatrix} + \underbrace{\begin{bmatrix} \mathbf{B}_p & \mathbf{O} \\ \mathbf{C}_r & \mathbf{C}_q \end{bmatrix}}_{\mathbf{B}} \begin{bmatrix} Z_p \\ Z_q \end{bmatrix},$$

where Z_p and Z_q are independent p - and q -dimensional standard normal random vectors. In particular, $X_p = \mu_p + \mathbf{B}_p Z_p$, which means that $X_p \sim \mathcal{N}(\mu_p, \Sigma_p)$, since $\mathbf{B}_p \mathbf{B}_p^\top = \Sigma_p$.

By relabeling the elements of X we see that the Theorem implies that *any* subvector of X has a multivariate normal distribution. For example, $X_q \sim \mathcal{N}(\mu_q, \Sigma_q)$.

Conditionals are Again Normal

Theorem: Conditional Distributions of Normal Vectors

Let $X \sim \mathcal{N}(\mu, \Sigma)$ be an n -dimensional normal random vector with $\det(\Sigma) > 0$. Decompose X as:

$$\begin{bmatrix} X_p \\ X_q \end{bmatrix} = \begin{bmatrix} \mu_p \\ \mu_q \end{bmatrix} + \underbrace{\begin{bmatrix} \mathbf{B}_p & \mathbf{O} \\ \mathbf{C}_r & \mathbf{C}_q \end{bmatrix}}_{\mathbf{B}} \begin{bmatrix} Z_p \\ Z_q \end{bmatrix},$$

Then

$$(X_q | X_p = x_p) \sim \mathcal{N}(\mu_q + \Sigma_r^\top \Sigma_p^{-1} (x_p - \mu_p), \Sigma_q - \Sigma_r^\top \Sigma_p^{-1} \Sigma_r).$$

As a consequence, X_p and X_q are *independent* if and only if they are *uncorrelated*; that is, if $\Sigma_r = \mathbf{O}$ (zero matrix).

Proof

From the decomposition we see that $\mathbf{X}_p = \boldsymbol{\mu}_p + \mathbf{B}_p \mathbf{Z}_p$ and $\mathbf{X}_q = \boldsymbol{\mu}_q + \mathbf{C}_r \mathbf{Z}_p + \mathbf{C}_q \mathbf{Z}_q$. Consequently,

$$(\mathbf{X}_q \mid \mathbf{X}_p = \mathbf{x}_p) = \boldsymbol{\mu}_q + \mathbf{C}_r \mathbf{B}_p^{-1} (\mathbf{x}_p - \boldsymbol{\mu}_p) + \mathbf{C}_q \mathbf{Z}_q,$$

where $\mathbf{Z}_q \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_q)$. It follows that \mathbf{X}_q , conditional on $\mathbf{X}_p = \mathbf{x}_p$, has a $\mathcal{N}(\boldsymbol{\mu}_q + \mathbf{C}_r \mathbf{B}_p^{-1} (\mathbf{x}_p - \boldsymbol{\mu}_p), \mathbf{C}_q \mathbf{C}_q^\top)$ distribution.

The proof is completed by observing that

$$\boldsymbol{\Sigma}_r^\top \boldsymbol{\Sigma}_p^{-1} = \mathbf{C}_r \mathbf{B}_p^\top (\mathbf{B}_p^\top)^{-1} \mathbf{B}_p^{-1} = \mathbf{C}_r \mathbf{B}_p^{-1},$$

and

$$\boldsymbol{\Sigma}_q - \boldsymbol{\Sigma}_r^\top \boldsymbol{\Sigma}_p^{-1} \boldsymbol{\Sigma}_r = \mathbf{C}_r \mathbf{C}_r^\top + \mathbf{C}_q \mathbf{C}_q^\top - \mathbf{C}_r \mathbf{B}_p^{-1} \underbrace{\boldsymbol{\Sigma}_r}_{\mathbf{B}_p \mathbf{C}_r^\top} = \mathbf{C}_q \mathbf{C}_q^\top.$$

The next few results are about the relationships between the normal, chi-squared, Student, and F distributions, defined in Table ?? . Recall that the chi-squared family of distributions, denoted by χ_n^2 , are simply $\text{Gamma}(n/2, 1/2)$ distributions, where the parameter $n \in \{1, 2, 3, \dots\}$ is called the **degrees of freedom**.

Theorem: Normal and χ^2 Distributions

If $X \sim \mathcal{N}(\mu, \Sigma)$ is an n -dimensional normal random vector with $\det(\Sigma) > 0$, then

$$(X - \mu)^\top \Sigma^{-1} (X - \mu) \sim \chi_n^2.$$

Proof

Let $\mathbf{B}\mathbf{B}^\top$ be the Cholesky decomposition of Σ , where \mathbf{B} is invertible. Since \mathbf{X} can be written as $\boldsymbol{\mu} + \mathbf{B}\mathbf{Z}$, where $\mathbf{Z} = [Z_1, \dots, Z_n]^\top$ is a vector of independent standard normal random variables, we have

$$(\mathbf{X} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{X} - \boldsymbol{\mu}) = (\mathbf{X} - \boldsymbol{\mu})^\top (\mathbf{B}\mathbf{B}^\top)^{-1} (\mathbf{X} - \boldsymbol{\mu}) = \mathbf{Z}^\top \mathbf{Z} = \sum_{i=1}^n Z_i^2.$$

Using the independence of Z_1, \dots, Z_n , the moment generating function of $Y = \sum_{i=1}^n Z_i^2$ is given by

$$\mathbb{E} e^{sY} = \mathbb{E} e^{s(Z_1^2 + \dots + Z_n^2)} = \mathbb{E} [e^{sZ_1^2} \dots e^{sZ_n^2}] = \left(\mathbb{E} e^{sZ^2} \right)^n,$$

where $Z \sim \mathcal{N}(0, 1)$. The moment generating function of Z^2 is

$$\mathbb{E} e^{sZ^2} = \int_{-\infty}^{\infty} e^{sz^2} \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{1}{2}(1-2s)z^2} dz = \frac{1}{\sqrt{1-2s}},$$

so that $\mathbb{E} e^{sY} = \left(\frac{1}{2}/(\frac{1}{2} - s)\right)^{\frac{n}{2}}$, $s < \frac{1}{2}$, which is the moment generating function of the $\text{Gamma}(n/2, 1/2)$ distribution; that is, the χ_n^2 distribution. The proof is completed using the uniqueness of the moment generating function.

Consequently, if $\mathbf{X} = [X_1, \dots, X_n]^\top \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$, then the squared length $\|\mathbf{X}\|^2 = X_1^2 + \dots + X_n^2$ has a χ_n^2 distribution.

If instead $X_i \sim \mathcal{N}(\mu_i, 1)$, $i = 1, \dots$, independently, then $\|\mathbf{X}\|^2$ is said to have a **noncentral χ_n^2 distribution**.

This distribution depends on the $\{\mu_i\}$ only through the norm $\|\boldsymbol{\mu}\|$. We write $\|\mathbf{X}\|^2 \sim \chi_n^2(\theta)$, where $\theta = \|\boldsymbol{\mu}\|$ is the **noncentrality parameter**.

Projections of Normal Vectors

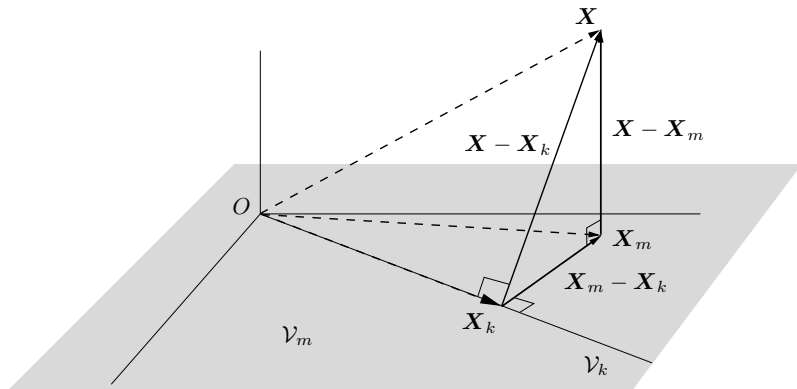
Non-central χ^2 distributions frequently occur when considering *projections* of multivariate normal random variables.

Theorem: Normal and Noncentral χ^2 Distributions

Let $X \sim \mathcal{N}(\mu, I_n)$ be an n -dimensional normal random vector and let $\mathcal{V}_k \subset \mathcal{V}_m$ be linear subspaces of dimensions k and m , respectively, with $k < m \leq n$. Let X_k and X_m be orthogonal projections of X onto \mathcal{V}_k and \mathcal{V}_m , and let μ_k and μ_m be the corresponding projections of μ . Then, the following holds.

1. The random vectors X_k , $X_m - X_k$, and $X - X_m$ are independent.
2. $\|X_k\|^2 \sim \chi_k^2(\|\mu_k\|)$, $\|X_m - X_k\|^2 \sim \chi_{m-k}^2(\|\mu_m - \mu_k\|)$, and $\|X - X_m\|^2 \sim \chi_{n-m}^2(\|\mu - \mu_m\|)$.

Projections and Pythagoras



Proof

Let $\mathbf{v}_1, \dots, \mathbf{v}_n$ be an orthonormal basis of \mathbb{R}^n such that $\mathbf{v}_1, \dots, \mathbf{v}_k$ spans \mathcal{V}_k and $\mathbf{v}_1, \dots, \mathbf{v}_m$ spans \mathcal{V}_m .

We can write the orthogonal projection matrices onto \mathcal{V}_j , as $\mathbf{P}_j = \sum_{i=1}^j \mathbf{v}_i \mathbf{v}_i^\top$, $j = k, m, n$, where \mathbf{P}_n is simply the identity matrix.

Let $\mathbf{V} := [\mathbf{v}_1, \dots, \mathbf{v}_n]$ and define $\mathbf{Z} := [Z_1, \dots, Z_n]^\top = \mathbf{V}^\top \mathbf{X}$. Recall that any orthogonal transformation such as $\mathbf{z} = \mathbf{V}^\top \mathbf{x}$ is *length preserving*; that is, $\|\mathbf{z}\| = \|\mathbf{x}\|$.

To prove the first statement of the theorem, note that $\mathbf{V}^\top \mathbf{X}_j = \mathbf{V}^\top \mathbf{P}_j \mathbf{X} = [Z_1, \dots, Z_j, 0, \dots, 0]^\top$, $j = k, m$.

It follows that $\mathbf{V}^\top (\mathbf{X}_m - \mathbf{X}_k) = [0, \dots, 0, Z_{k+1}, \dots, Z_m, 0, \dots, 0]^\top$ and $\mathbf{V}^\top (\mathbf{X} - \mathbf{X}_m) = [0, \dots, 0, Z_{m+1}, \dots, Z_n]^\top$.

Moreover, being a linear transformation of a normal random vector, \mathbf{Z} is also normal, with covariance matrix $\mathbf{V}^\top \mathbf{V} = \mathbf{I}_n$.

Proof

In particular, the $\{Z_i\}$ are *independent*. This shows that X_k , $X_m - X_k$ and $X - X_m$ are independent as well.

Next, observe that $\|X_k\| = \|V^T X_k\| = \|Z_k\|$, where $Z_k := [Z_1, \dots, Z_k]^T$. The latter vector has independent components with variances 1, and its squared norm has therefore (by definition) a $\chi_k^2(\theta)$ distribution. The noncentrality parameter is

$$\theta = \|\mathbb{E}Z_k\| = \|\mathbb{E}X_k\| = \|\mu_k\|,$$

again by the length-preserving property of orthogonal transformations. This shows that $\|X_k\|^2 \sim \chi_k^2(\|\mu_k\|)$.

The distributions of $\|X_m - X_k\|^2$ and $\|X - X_m\|^2$ follow by analogy.

Quotients of Squared Norms

The above projection theorem is frequently used in the statistical analysis of *normal linear models*. In typical situations $\boldsymbol{\mu}$ lies in the subspace \mathcal{V}_m or even \mathcal{V}_k — in which case $\|\mathbf{X}_m - \mathbf{X}_k\|^2 \sim \chi_{m-k}^2$ and $\|\mathbf{X} - \mathbf{X}_m\|^2 \sim \chi_{n-m}^2$, independently. The (scaled) quotient then turns out to have an F distribution — a consequence of the following theorem.

Theorem: Relationship Between χ^2 and F Distributions

Let $U \sim \chi_m^2$ and $V \sim \chi_n^2$ be independent. Then,

$$\frac{U/m}{V/n} \sim F(m, n).$$

Proof

For notational simplicity, let $c = m/2$ and $d = n/2$. The pdf of $W = U/V$ is given by $f_W(w) = \int_0^\infty f_U(wv) v f_V(v) dv$. Substituting the pdfs of the corresponding Gamma distributions, we have

$$\begin{aligned} f_W(w) &= \int_0^\infty \frac{(wv)^{c-1} e^{-wv/2}}{\Gamma(c) 2^c} v \frac{v^{d-1} e^{-v/2}}{\Gamma(d) 2^d} dv \\ &= \frac{\Gamma(c+d)}{\Gamma(c) \Gamma(d)} \frac{w^{c-1}}{(1+w)^{c+d}}. \end{aligned}$$

The density of $Z = \frac{n}{m} \frac{U}{V}$ is given by

$$f_Z(z) = f_W(z m/n) m/n.$$

The proof is completed by comparing the resulting expression with the pdf of the F distribution.

Normal Linear Models

Normal linear models combine the simplicity of the linear model with the tractability of the Gaussian distribution. They are the principal model for traditional statistics, and include the classic linear regression and analysis of variance models.

Definition 1: Normal Linear Model

In a **normal linear model** the response Y depends on a p -dimensional explanatory variable $\mathbf{x} = [x_1, \dots, x_p]^\top$, via the linear relationship

$$Y = \mathbf{x}^\top \boldsymbol{\beta} + \varepsilon,$$

where $\varepsilon \sim \mathcal{N}(0, \sigma^2)$.

Multivariate Normal Response Vector

Thus, a normal linear model is a linear model with normal error terms.

In our usual notation, the corresponding normal linear model for the whole training set $\{(\mathbf{x}_i, Y_i)\}$ has the form

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon},$$

where \mathbf{X} is the model matrix comprised of rows $\mathbf{x}_1^\top, \dots, \mathbf{x}_n^\top$ and $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$.

Consequently, \mathbf{Y} can be written as $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \sigma \mathbf{Z}$, where $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$, so that $\mathbf{Y} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$.

Estimating β

It follows that the joint density of \mathbf{Y} , conditional on the explanatory variables, is given by

$$g(\mathbf{y} \mid \beta, \sigma^2, \mathbf{X}) = (2\pi\sigma^2)^{-\frac{n}{2}} e^{-\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\beta\|^2}.$$

Estimation of the parameter β can be performed via

- the least-squares method, or
- the maximum likelihood method.

It is clear that for every value of σ^2 the likelihood is maximal when $\|\mathbf{y} - \mathbf{X}\beta\|^2$ is minimal.

As a consequence, the maximum likelihood estimate for β is the same as the least-squares estimate $\hat{\beta} = \mathbf{X}^+ \mathbf{y}$.

Estimating σ^2

The maximum likelihood estimate of σ^2 is equal to

$$\widehat{\sigma^2} = \frac{\|\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}\|^2}{n},$$

where $\widehat{\boldsymbol{\beta}}$ is the maximum likelihood estimate (least squares estimate in this case) of $\boldsymbol{\beta}$.