

Review: The Multivariate Normal Distribution

version: 2020-01-07 · 10:58:30

The Multivariate Normal Distribution

- ▶ **Definition:** $\mathbf{X} = (X_1, \dots, X_d)$ is multivariate normal if and only if it a linear combination of iid normals:

$$\mathbf{X} = \mathbf{C}\mathbf{Z} + \boldsymbol{\mu}, \quad \mathbf{Z} = (Z_1, \dots, Z_d), \quad Z_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1).$$

- ▶ **Mean and Variance:**

- ▶ $E[\mathbf{X}] = \mathbf{C} E[\mathbf{Z}] + \boldsymbol{\mu} = \boldsymbol{\mu}$

- ▶ $\text{var}(\mathbf{X}) = \mathbf{C} \text{var}(\mathbf{Z}) \mathbf{C}' = \mathbf{C} \mathbf{C}' := \boldsymbol{\Sigma}$ (many different \mathbf{C} give the same variance)

- ▶ **Notation:** $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

- ▶ **PDF:**

$$f(\mathbf{x}) = (2\pi)^{-d/2} \cdot \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) - \frac{1}{2} \log |\boldsymbol{\Sigma}| \right\}$$

Simulation

- ▶ To generate $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$:
 1. Find \mathbf{C} such that $\mathbf{C}\mathbf{C}' = \boldsymbol{\Sigma}$.
 2. Generate $\mathbf{Z} = (Z_1, \dots, Z_d)$ with $Z_i \sim \mathcal{N}(0, 1)$.
 3. Set $\mathbf{X} = \mathbf{C}\mathbf{Z} + \boldsymbol{\mu}$.
- ▶ To find \mathbf{C} :
 - ▶ Note that $\boldsymbol{\Sigma}$ is (i) symmetric and (ii) **positive definite**: for any vector \mathbf{a} we have $\mathbf{a}'\boldsymbol{\Sigma}\mathbf{a} \geq 0$, with equality $\iff \mathbf{a} = \mathbf{0}$.
 - ▶ Every symmetric positive-definite matrix has a **Cholesky definition**: $\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}'$, where \mathbf{L} is **lower triangular** and all diagonal elements $L_{ii} > 0$.
 - ▶ *Properties*:
 - ▶ The eigenvalues of triangular matrices are the diagonal elements
 $\implies |\boldsymbol{\Sigma}| = \prod_{i=1}^d L_{ii}^2$.
 - ▶ $\boldsymbol{\Sigma}^{-1}\mathbf{x} = (\mathbf{L}')^{-1}(\mathbf{L}^{-1}\mathbf{x})$. Solving linear systems with triangular matrices is $\mathcal{O}(d)$, as opposed to $\mathcal{O}(d^3)$ for general matrices.

Conditional Distribution

► **Block Notation:**

$$\begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{bmatrix} \sim \mathcal{N} \left\{ \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix} \right\}$$

► **Marginal Distribution:** $\mathbf{X}_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$

► **Conditional Distribution:**

$$\begin{aligned} \mathbf{X}_2 | \mathbf{X}_1 &\sim \mathcal{N}(\boldsymbol{\mu}_2^*, \boldsymbol{\Sigma}_2^*), & \boldsymbol{\mu}_2^* &= \boldsymbol{\mu}_2 + \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} (\mathbf{X}_1 - \boldsymbol{\mu}_1) \\ & & \boldsymbol{\Sigma}_2^* &= \boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12}. \end{aligned}$$

► **Verify Calculations:** $f(\mathbf{x}_1, \mathbf{x}_2) = f(\mathbf{x}_1) \times f(\mathbf{x}_2 | \mathbf{x}_1)$ for any pair $(\mathbf{x}_1, \mathbf{x}_2)$.

Can do this **analytically** (harder) or **computationally** (easier)

Parameter Inference

► **Data:** $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n)$, $\mathbf{X}_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

► **Loglikelihood function:**

$$\begin{aligned}\ell(\boldsymbol{\mu}, \boldsymbol{\Sigma} \mid \mathbf{X}) &= \log \prod_{i=1}^n f(\mathbf{X}_i \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) - \{\text{terms not involving } \boldsymbol{\mu} \text{ or } \boldsymbol{\Sigma}\} \\ &= -\frac{1}{2} \left\{ n \log |\boldsymbol{\Sigma}| + \sum_{i=1}^n (\mathbf{X}_i - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{X}_i - \boldsymbol{\mu}) \right\} \\ &= -\frac{1}{2} \left\{ n \log |\boldsymbol{\Sigma}| + \text{tr}(\boldsymbol{\Sigma}^{-1} \mathbf{S}) + n(\bar{\mathbf{X}} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{X}} - \boldsymbol{\mu}) \right\},\end{aligned}$$

where $\bar{\mathbf{X}} = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i$, $\mathbf{S}_{d \times d} = \sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{X}_i - \bar{\mathbf{X}})'$

► **MLE:** $\hat{\boldsymbol{\mu}} = \bar{\mathbf{X}}$ and $\hat{\boldsymbol{\Sigma}} = \mathbf{S}/n$.

Parameter Inference

► **Data:** $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n)$, $\mathbf{X}_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

► **Loglikelihood function:**

$$\ell(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \mathbf{X}) = -\frac{1}{2} \left\{ n \log |\boldsymbol{\Sigma}| + \text{tr}(\boldsymbol{\Sigma}^{-1} \mathbf{S}) + n(\bar{\mathbf{X}} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{X}} - \boldsymbol{\mu}) \right\},$$

where $\bar{\mathbf{X}} = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i$, $\mathbf{S}_{d \times d} = \sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{X}_i - \bar{\mathbf{X}})'$

► **Verify Calculations:** Can do this **analytically** (harder) or **computationally** (easier):

► Check that

$$\ell(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \mathbf{X}) = \log f(\mathbf{X} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \text{CONST}$$

for fixed \mathbf{X} and varying $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

Parameter Inference

► **Data:** $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n)$, $\mathbf{X}_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

► **MLE:** $\hat{\boldsymbol{\mu}} = \bar{\mathbf{X}}$ and $\hat{\boldsymbol{\Sigma}} = \mathbf{S}/n$,

where $\bar{\mathbf{X}} = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i$, $\mathbf{S}_{d \times d} = \sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{X}_i - \bar{\mathbf{X}})'$

► **Verify Calculations:** Can do this **analytically** (harder) or **computationally** (easier):

► For differentiable loglikelihood $\ell(\boldsymbol{\theta} | \mathbf{X})$ with $\boldsymbol{\theta} = (\theta_1, \dots, \theta_p)$, each component of MLE is the maximum of a 1-d function:

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta} | \mathbf{X}) \implies \hat{\theta}_i = \arg \max_{\theta_i} \ell(\theta_i, \hat{\boldsymbol{\theta}}_{[-i]} | \mathbf{X}),$$

where $\hat{\boldsymbol{\theta}}_{[-i]} = \hat{\boldsymbol{\theta}} \setminus \{\hat{\theta}_i\}$.

► Converse is **false** (as $\hat{\boldsymbol{\theta}}$ could be a saddlepoint). However, loglikelihoods are generally well-behaved, so if $\hat{\theta}_i = \arg \max_{\theta_i} \ell(\theta_i, \hat{\boldsymbol{\theta}}_{[-i]} | \mathbf{X})$ for $i = 1, \dots, p$, then very likely that $\hat{\boldsymbol{\theta}}$ is the MLE.

Applications

1. Confidence Intervals

- ▶ **Model:** $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} f(y | \boldsymbol{\theta}), \quad \boldsymbol{\theta} = (\theta_1, \dots, \theta_d).$
- ▶ **MLE:** $\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta} | \mathbf{Y}), \quad \ell(\boldsymbol{\theta} | \mathbf{Y}) = \sum_{i=1}^n \log f(Y_i | \boldsymbol{\theta}).$
- ▶ **Asymptotic Theory:** As $n \rightarrow \infty$, we have

$$\hat{\boldsymbol{\theta}} \approx \mathcal{N}(\boldsymbol{\theta}_0, \mathcal{I}_0^{-1}), \quad \text{where}$$

- ▶ $\boldsymbol{\theta}_0$ is the true parameter value
- ▶ \mathcal{I}_0 is the (expected) Fisher Information:

$$\mathcal{I}_0 = -E \left[\frac{\partial^2}{\partial \boldsymbol{\theta}^2} \ell(\boldsymbol{\theta}_0 | \mathbf{Y}) \right] = -n \int \left[\frac{\partial^2}{\partial \boldsymbol{\theta}^2} \log f(y | \boldsymbol{\theta}_0) \right] \cdot f(y | \boldsymbol{\theta}_0) dy$$

Applications

1. Confidence Intervals

- ▶ **Model:** $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} f(y | \theta)$
- ▶ **Asymptotic Theory:** As $n \rightarrow \infty$, we have $\hat{\theta} \approx \mathcal{N}(\theta_0, \mathcal{I}_0^{-1})$, where
 - ▶ θ_0 is the true parameter value
 - ▶ $\mathcal{I}_0 = -E \left[\frac{\partial^2}{\partial \theta^2} \ell(\theta_0 | \mathbf{Y}) \right]$ is the (expected) Fisher Information.

Typically \mathcal{I}_0 is impossible to calculate because (i) expectation is usually intractable and (ii) true θ_0 is unknown.

Observed Fisher Information is a consistent estimator: $\hat{\mathcal{I}} = -\frac{\partial^2}{\partial \theta^2} \ell(\hat{\theta} | \mathbf{Y}) \xrightarrow{n} \mathcal{I}_0$

- ▶ **Asymptotic Confidence Intervals:** 95% CI for each element of θ :

$$\hat{\theta}_i \pm 1.96 \cdot \text{se}(\hat{\theta}_i), \quad \text{se}(\hat{\theta}_i) = \sqrt{[\hat{\mathcal{I}}^{-1}]_{ii}}.$$

Such CI's are often valid even without iid data.

Applications

2. Linear Regression

- ▶ **Model:** $\mathbf{y} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{V})$, where
 - ▶ $\mathbf{y} = (y_1, \dots, y_n)$ is multivariate normal (random)
 - ▶ $\mathbf{X}_{n \times p}$ and $\mathbf{V}_{n \times n}$ are known (nonrandom)
 - ▶ $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)$ and σ are unknown (parameters)

- ▶ **Loglikelihood:**

$$\begin{aligned}\ell(\boldsymbol{\beta}, \sigma \mid \mathbf{y}) &= -\frac{1}{2} \left\{ (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})' [\sigma^2 \mathbf{V}]^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \log |\sigma^2 \mathbf{V}| \right\} \\ &= -\frac{1}{2} \left\{ \frac{(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})' \mathbf{X}' \mathbf{V}^{-1} \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) + n\hat{\sigma}^2}{\sigma^2} + n \log \sigma^2 \right\},\end{aligned}$$

where $\hat{\boldsymbol{\beta}} = (\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}' \mathbf{V}^{-1} \mathbf{y}$ and $\hat{\sigma}^2 = \frac{1}{n} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})' \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$.

- ▶ **Inference:** The MLE of $\boldsymbol{\theta} = (\boldsymbol{\beta}, \sigma)$ is $\hat{\boldsymbol{\theta}} = (\hat{\boldsymbol{\beta}}, \hat{\sigma})$.