

Lecture 3: Multivariate Normal Distributions

Lecturer: Sasha Rush

Scribes: Christopher Mosch, Lindsey Brown, Ryan Lapcevic

3.1 Examples

Multivariate gaussians are used for modeling in various applications, where knowing mean and variance is useful:

- radar: mean and variance of approaching objects (like invading aliens)
- weather forecasting: predicting the position of a hurricane, where the uncertainty in the storm's position increases for timepoints farther away
- tracking the likely outcome of a sports game: last year's superbowl is an example of a failure of modeling with multivariate gaussians as the Patriots still won after a large Falcons' lead

3.2 Review: Eigendecomposition

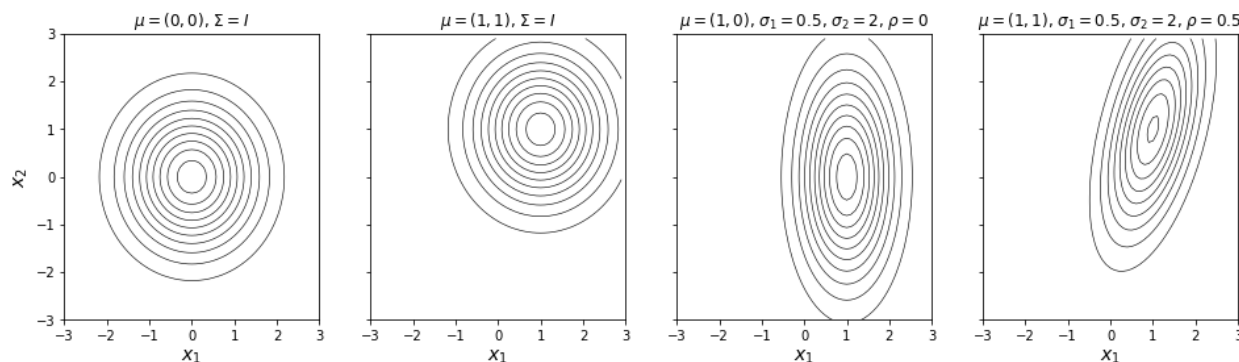
Let Σ be a square, symmetric matrix. Then its eigendecomposition is given by $\Sigma = \mathbf{U}^T \mathbf{\Lambda} \mathbf{U}$, where \mathbf{U} is an orthogonal matrix and $\mathbf{\Lambda}$ is a diagonal matrix. In the special case that Σ is positive semidefinite (as is the case for covariance matrices), denoted $\Sigma \succeq 0$, all its eigenvalues are nonnegative, $\Lambda_{ii} \geq 0$, and we can decompose its inverse as $\Sigma^{-1} = \mathbf{U}^T \mathbf{\Lambda}^{-1} \mathbf{U}$, where $\Lambda_{ii}^{-1} = 1/\Lambda_{ii}$.

3.3 Multivariate Normal Distributions (MVNs)

Let X be a D -dimensional MVN random vector with mean μ and covariance matrix Σ , denoted $X \sim \mathcal{N}(\mu, \Sigma)$. Then the pdf of X is

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

where for many problems we focus on the quadratic form $(x - \mu)^T \Sigma^{-1}(x - \mu)$ (which geometrically can be thought of as distance) and ignore the normalization factor $(2\pi)^{-D/2} |\Sigma|^{-1/2}$. The figure below plots the contours of a bivariate Normal for various μ and Σ (in the figure, ρ denotes the off-diagonal elements of Σ , given by the covariance of x_1 and x_2).



Note that we can decompose Σ as

$$\begin{aligned}
\Sigma &= (x - \mu)^T \Sigma^{-1} (x - \mu) \\
&= (x - \mu)^T \left(U^T \Lambda^{-1} U \right) (x - \mu) \\
&= (x - \mu)^T \left(\sum_d \frac{1}{\lambda_d} U_d U_d^T \right) (x - \mu) \\
&= \sum_d \frac{1}{\lambda_d} (x - \mu)^T U_d U_d^T (x - \mu),
\end{aligned}$$

where $(x - \mu)^T U_d$ can be interpreted as the projection of $(x - \mu)$ onto U_d (which can each be thought of as univariate gaussians), the eigenvector corresponding to the eigenvalue λ_d . Since Σ is the weighted sum of the dot product of such projections (with weights being given by $1/\lambda_d$, which can be thought of as the scale $1/\sigma^2$), we can describe the MVN as tiling of univariates.

3.3.1 Manipulating MVNs: Stretches, Rotations, and Shifts

Let $x \sim \mathcal{N}(0, I)$ and $y = Ax + b$. We want to consider two ways of obtaining the complete distribution of y .

- 'Overkill': We can perform a change of variables¹. Here, we have $x = A^{-1}y - A^{-1}b$ and $|dx/dy| = |A^{-1}|$, leading to

$$\begin{aligned}
p(y) &= \mathcal{N}(A^{-1}(y - b) | 0, I) |A^{-1}| \\
&= \frac{1}{z} \exp [-(A^{-1}(y - b))^T (A^{-1}(y - b))] \\
&= \frac{1}{z} \exp [-(y - b)^T (A^{-1})^T (A^{-1})(y - b)] \\
&= \mathcal{N}(y | b, AA^T),
\end{aligned}$$

where z is the normalizing constant.

- Using the properties of MVN, we know that y is also MVN, so is completely specified by its mean and covariance matrix which can easily be derived,

$$\mathbb{E}(y) = \mathbb{E}(Ax + b) = A\mathbb{E}(x) + b \quad \text{cov}(y) = AA^T.$$

Thus, we can generate MVN from $\mathcal{N}(0, I)$ via the transformation $y = Ax + b$, where we set $A = U\Lambda^{1/2}$, leading to $\Sigma_Y = U^T \Lambda U$. Then shifts are represented by b , stretches by Λ , and rotations by U .

3.3.2 Detour: MVN in High-Dimensions ($D \gg 0$)

Let x be a D -dimensional random vector, distributed as $\mathcal{N}(0, I/D)$, where I is the identity. The expected length of x is given by

$$\mathbb{E}(\|x\|^2) = \mathbb{E}\left(\sum_d x_d^2\right) = D\sigma_d^2 = 1,$$

¹A change of variables can be done in the following way: Let $y = f(x)$ and assume f is invertible so that $x = f^{-1}(y)$. Then $p(y) = p(x)|dx/dy|$. This is a technique which will be used often in this course

which means that x is expected to be on the boundary of a unit sphere centered at the origin. Moreover, the variance of the length is

$$\text{var}(\|x\|^2) = D \cdot \left(\mathbb{E}(x^4) - \mathbb{E}(x^2)^2 \right) = D \cdot (3\sigma^4 - \sigma^4) = 2D/D^2 = 2/D$$

Thus, it is not only expected that x lies on the boundary but as D increases most of its realizations will in fact fall on the boundary².

3.3.3 Key Formulas for MVN: Marginalization and Conditioning

Let $X \sim \mathcal{N}(\mu, \Sigma)$ with

$$X = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}, \quad \mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \quad \Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}.$$

Note that Σ is written in block matrix form, rather than scalar entries. It turns out that the marginals, X_1 and X_2 , are also MVN, and their mean and covariance matrices are given by μ_1 and Σ_{11} and μ_2 and Σ_{22} respectively. A sketch of the proof is provided below.

$$p(x_1) = \int_{x_2} N(x|\mu, \Sigma) dx_2,$$

which can be written as

$$0.5 \int_{x_2} \exp \left[(x_1 - \mu_1)^T \Sigma_{11}^{-1} (x_1 - \mu_1) + 2(x_1 - \mu_1)^T \Sigma_{12}^{-1} (x_2 - \mu_2) + (x_2 - \mu_2)^T \Sigma_{22}^{-1} (x_2 - \mu_2) \right] dx_2.$$

Note that this equals

$$p(x_1) \int_{x_2} p(x_2|x_1) dx_2,$$

implying that $X_1 \sim \mathcal{N}(\mu_1, \Sigma_{11})$.

While the marginals have a simple form, the conditionals are more complicated. (For a complete derivation, which requires matrix inversion lemmas, refer to Murphy.) It can be shown that $X_1|X_2 \sim \mu_{\infty|\in}, \Sigma_{\infty|\in}$ with

$$\mu_{1|2} = \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (x_2 - \mu_2), \quad \Sigma_{1|2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}.$$

3.3.4 Information Form

An alternative formulation, called information form, uses the precision matrix (inverse variance) $\Lambda = \Sigma^{-1}$. Partitioning Λ as

$$\Lambda = \begin{pmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{pmatrix},$$

the covariance matrices of the conditional distributions have a simple form. For example, the covariance matrix of X_1 given X_2 is $\Lambda_{1|2} = \Lambda_{11}$. However, the simplicity of the conditional precision comes at the cost of marginalization (which was simple when using Σ) becoming a more complicated expression (see Murphy subsection 4.3 for more details).

²It is left as an exercise to show that this formula holds. Hint: Use the fact that we assumed no covariance.