
A multivariate extension to the Exponentially-modified Gaussian distribution

Sandhya Prabhakaran*

Department of Integrated Mathematical Oncology
Moffitt Cancer Center
Florida
sandhya.prabhakaran@moffitt.org

Abstract

The exponentially-modified Gaussian (EMG) distribution is a convolution sum of a univariate Gaussian and an exponential distribution. This has been used to model univariate skewed data such as chromatographic peaks' shape, cell population dynamics from single-cell data and reaction times in neuropsychology. Currently, the EMG is only available in its univariate form. In this work, we propose a multivariate extension to the EMG, called *mvEMG*, by using an affine transformation involving rotation, translation and shearing to accommodate for the three moments (mean, variance and skew). We derive statistical properties for *mvEMG*. Although we demonstrate its performance in synthetic data compared with the multivariate skew normal distribution, we are unable to show its practical applicability, mainly due to lack of efficient sampling strategies and a viable real-world dataset.

1 Introduction and motivation

Generally, datasets are described using statistics involving mean and variance, which through the method of moments, represent the first and second order estimators. Most data-driven approaches assume the data to be generated using a Gaussian distribution where these estimators are used to address simple hypotheses about the datasets. To gain deeper insights into the data, one needs to consider higher-order estimators, such as skewness, and such estimators are rarely used.

In many real world applications, datasets do not generally follow a Gaussian distribution but are rather highly skewed with complex structures. For computational feasibility, one way to handle such data is to log normalise the data rendering it as Gaussian, but by doing so, the inherent skewness is lost. Since the tails of a probability distribution reflect the most extreme data points, it seems plausible that measures based on skewness would be useful for identifying rare events or entities such as transcriptional regulators in gene expression data or highly-variable genes that could be targeted in precision medicine studies. As a specific example, we look at multidimensional single-cell RNA-seq data where multiple genes are profiled per cell. The distribution of gene counts per cell, called library size, is highly right-skewed (Figure 1a upper panel) and log normalising renders the distribution to resemble a Gaussian distribution (Figure 1a lower panel), where the information regarding rare events triggering data skewness is lost. Skewness is important for modelling complex phenomena.

There exists the exponentially-modified Gaussian (EMG) distribution, to handle univariate skewed data. In this work, we present a multivariate extension to the EMG to cater to multidimensional skewed data whilst providing access to the first three distributional moments viz. mean, covariance *and* skew. For this, we make use of the multivariate Gaussian that provides the mean and covariance, and the multivariate exponential distribution that gives the skewness of the data. The multivariate exponential

*<http://sandhyaprabhakaran.com>

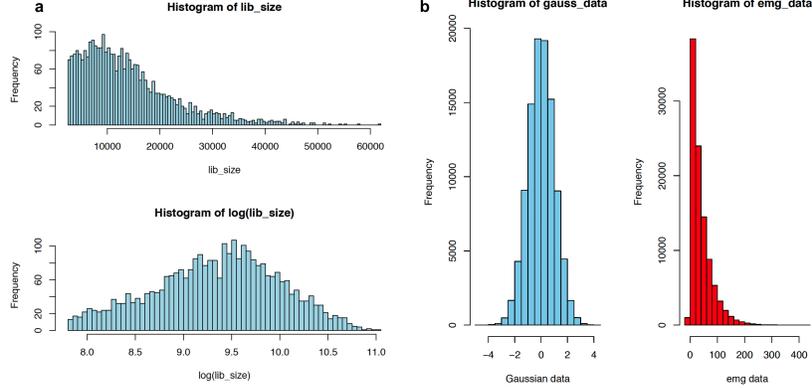


Figure 1: **a: Upper:** Histogram showing the original skewed library-size distribution i.e. the distribution of genes across cells for the *Zeisel, et al. (2015)* single-cell RNAseq dataset for 3500 cells and **a: Lower:** the same library size as above but is log-normalised depicting the shape of a Gaussian distribution. **b: Left:** Histograms showing 10K random variables drawn from a standard Gaussian distribution and **b: Right:** from a univariate EMG ($mean = 0, std = 1, skew = 0.25$).

distribution has many forms that can be derived from the univariate exponential distribution (*Esary and Marshall (1974)*). This work aims to provide one way to unify these two useful multivariate distributions into a single viable multivariate form leading to the multivariate EMG (mvEMG).

2 Univariate EMG

The exponentially-modified Gaussian (EMG) probability distribution is the convolution of a univariate Gaussian distribution and an exponential distribution which are independent of each other. Assume a r.v. $Z = Z_1 + Z_2$ with Z_1 and Z_2 as independent r.v.s where $Z_1 \sim \mathcal{N}(\mu, \sigma^2)$ and $Z_2 \sim \exp(\lambda)$ then Z is a convolution of a Gaussian and an Exponential random variable and is said to follow an exponentially-modified Gaussian (EMG) distribution with parameters (μ, σ, λ) i.e. $Z \sim \text{EMG}(\mu, \sigma, \lambda)$ and is given by: $f_Z(z) = \frac{\lambda}{2} \exp\left(\frac{\lambda(2\mu + \lambda\sigma^2 - 2z)}{2}\right) \text{erfc}\left(\frac{\mu + \lambda\sigma^2 - z}{\sqrt{(2)}\sigma}\right)$

where erfc is the complementary error function and $\text{erfc}(x) = 1 - \text{erf}(x) = 1 - \frac{2}{\sqrt{\pi}} \int_x^\infty e^{-t^2} dt$

Figure 1b depicts the histograms of r.v.s distributed according to a standard normal (left panel) and according to a univariate EMG with $mean=0, std=1$ and $skew = 0.25$ (right panel). One of the key properties of the univariate EMG distribution is its differential behavior in the right and left tails. The distribution exhibits a Gaussian-distributed left tail and an exponentially-distributed right tail. Hence, it is particularly well suited to fit empirical distributions that are right-skewed i.e. exhibit thicker right tails. This distribution has found practical applications in a variety of scientific disciplines such as chromatography (*Naish and Hartwell, 1988; Kong et al., 2005*), cellular biology (*Golubev, 2010; Tyson et al., 2012*) and microarray preprocessing (*Irizzary et al., 2003*).

The univariate skew normal (SN) distribution is also used to model right-skewed data. Since it has a normal-like right tail, it is thinner as compared to the EMG. Thus the EMG is better suited in modelling fat-tailed phenomenon consisting of outliers or data points residing in the right-tail. Additionally, the SN parameters are mean, scale and shape whereas for the EMG, the parameters are mean, variance and skew. Therefore, to access skewness in fat right-tailed distributions, we pursue building a multivariate extension to the univariate EMG.

3 Derivation of the mvEMG PDF

For each $Z_i \sim EMG(\mu_i, \sigma_i, \lambda_i)$, we have mean $\mathbb{E}(Z_i) = \mu_{Z_i} = \mu_i + \frac{1}{\lambda_i}$, $\text{var}(Z_i) = \sigma_{Z_i}^2 = \sigma_i^2 + \frac{1}{\lambda_i^2}$, and $\text{skew}(Z_i) = \lambda_{Z_i} = \frac{2}{\sigma_i^3 \lambda_i^3} \left(1 + \frac{1}{\sigma_i^2 \lambda_i^2}\right)^{-\frac{3}{2}}$ where μ_{Z_i} , $\sigma_{Z_i}^2$ and $\lambda_{Z_i} \in \mathbb{R}$. Next we assume a random vector \mathbf{Z} such that for $\mathbf{Z} \in \mathbb{R}^d$, we have mean $\mathbb{E}(\mathbf{Z}) = \mu_Z = [\mu_{Z_1}, \dots, \mu_{Z_d}]$, $\text{Cov}(\mathbf{Z}) = \Sigma_Z = \text{diag}(\sigma_{Z_i}^2) = \text{diag}(\sigma_i^2 + \frac{1}{\lambda_i^2})$, and $\text{skew}(\mathbf{Z}) = \lambda_Z = \text{diag}(\lambda_{Z_i}) = \text{diag}\left(\frac{2}{\sigma_i^3 \lambda_i^3} \left(1 + \frac{1}{\sigma_i^2 \lambda_i^2}\right)^{-\frac{3}{2}}\right)$ where $\mu_Z \in \mathbb{R}^d$, Σ_Z and $\lambda_Z \in \mathbb{R}^{d \times d}$.

Affine transformation of a d independent 0-centered univariate EMG random variable \mathbf{Z} . Given an affine transformation \mathbb{L} that comprises of (a) the Euclidean transformations of rotation via matrix A and translation via vector b followed by (b) a shear transformation via matrix S where $A, S \in \mathbb{R}^{d \times d}$, both A and S are invertible and $b \in \mathbb{R}^d$. We define the shear matrix S as $S = S^* + \mathbb{I}_{d \times d}$ where S^* is skew-symmetric i.e. $S^{*(ij)} = \lambda_{Z_i}$, $S^{*(ij)} = S^*(-ji)$, $\text{diag}(S^*) = 0$ and $\text{diag}(S) = 1$. Applying \mathbb{L} to \mathbf{Z} such that $\mathbf{Z}' = \mathbb{L}(\mathbf{Z})$ results in the image of $S(A\mathbf{Z} + b) := \mathbb{L}(\mathbf{Z}) = \mathbf{Z}'$. We note here that $\mathbb{L}(\mathbf{Z})$ still linearly transforms \mathbf{Z} .

Moments of \mathbf{Z}' The moments of this transformed random variable \mathbf{Z}' are:

1. $\mathbb{E}(\mathbf{Z}') = S(\mathbb{E}(A\mu_Z) + \mathbb{E}(b)) = S(A\mu_Z + b)$
2. $\text{Cov}(\mathbf{Z}') = \mathbb{E}[(\mathbf{Z}' - \bar{\mathbf{Z}}')(\mathbf{Z}' - \bar{\mathbf{Z}}')^T]$
 $= \mathbb{E}[(SA(\mathbf{Z} + b) - S(A\mu_Z + b))(SA(\mathbf{Z} + b) - S(A\mu_Z + b))^T] = SA\mathbb{K}_{\mathbf{Z}\mathbf{Z}^T}A^T S^T$
where $\mathbb{K}_{\mathbf{Z}\mathbf{Z}^T} = \mathbb{E}[(\mathbf{Z} - \mu_Z)(\mathbf{Z} - \mu_Z)^T] = \mathbb{E}[\mathbf{Z}\mathbf{Z}^T] - \mu_Z\mu_Z^T$
3. $\text{Skew}(\mathbf{Z}') = \mathbb{E}\left[\left(\frac{(\mathbf{Z}' - \bar{\mathbf{Z}}')}{\text{cov}(\mathbf{Z}')} \right)^3\right] = \mathbb{E}\left[\left(\frac{SA(\mathbf{Z} - \mu_Z)}{SA\mathbb{K}_{\mathbf{Z}\mathbf{Z}^T}A^T S^T} \right)^3\right]$

In our case, we will be considering \mathbb{L} given by the rotation due to Σ_Z , the translation via μ_Z and the shearing via S . Through this construction, we bring in the covariance and mean of a multivariate Gaussian distribution and the rate parameter of a multivariate exponential distribution, respectively. From the properties of mean vectors and covariance matrices for a multivariate Gaussian distribution, we have $\mathbb{E}(\mathbf{Z}) = \mu_Z$, $\text{Cov}(\mathbf{Z}) = \Sigma_Z = \text{diag}(\sigma_i^2 + \frac{1}{\lambda_i^2}) = AA^T$ where $\mu_Z \in \mathbb{R}^d$ and $\Sigma_Z \in \mathbb{R}^{d \times d}$.

Therefore, the affine transformed $\mathbf{Z}' = S(\Sigma_Z^{-\frac{1}{2}}\mathbf{Z} + \mu_Z)$ and each z'_i can be written as: $z'_i = (S(\Sigma_Z^{-\frac{1}{2}}\mathbf{Z} + \mu_Z))_i = S\sqrt{\psi_i^T}U_i \cdot \mathbf{Z} + S\mu_Z = V_i\sqrt{\psi_i^T}U_i \cdot \mathbf{Z} \cdot V_i^T + V_i\mu_Z V_i^T = \sum_j^d V_{ij}\sqrt{\psi_{ij}^T}U_{ij} \cdot Z_j \cdot V_{ij} + V_{ij}\mu_{Z_j} V_{ij}$ where $\Sigma_Z = U\Psi U^T$, ψ_i is a row of Ψ , U_i is a row of U , λ_{Z_i} is a row of λ_Z , S is Cholesky decomposed to VV^T and V_i is a row of V .

PDF of \mathbf{Z}' Suppose that \mathbf{Z} is a random vector in the subset $T \subseteq \mathbb{R}^d$ and has continuous pdf f . Suppose that $\mathbf{Z}' = r(\mathbf{Z})$ where r is a differentiable function from T to another subset $T^* \subseteq \mathbb{R}^d$, then we have the continuous pdf g for \mathbf{Z}' as $g(\mathbf{Z}') = f(\mathbf{Z})|\det(\frac{d\mathbf{Z}}{d\mathbf{Z}'})| = f(r^{-1}(\mathbf{Z}'))|\det(\frac{d\mathbf{Z}}{d\mathbf{Z}'})|$ where

$\frac{d\mathbf{Z}}{d\mathbf{Z}'}$ is the Jacobian of the inverse of r . The inverse transformation is given by $\mathbf{Z} = \Sigma_Z^{-\frac{1}{2}}(S^{-1}\mathbf{Z}' - \mu_Z)$, the Jacobian $\frac{d\mathbf{Z}}{d\mathbf{Z}'} = \Sigma_Z^{-\frac{1}{2}}S^{-1}$ and the determinant of the Jacobian of the inverse transformation

is $\frac{1}{\sqrt{\det(\Sigma_Z)} \det S}$. Therefore

$$\begin{aligned}
g(Z') &= f(r^{-1}(Z')) \left| \det \left(\frac{dZ}{dZ'} \right) \right| \\
&\propto \frac{\det(\Lambda)}{\sqrt{\det(\Sigma_Z)} \det S} \exp\left(\frac{\text{trace}(\Sigma^{*T} \Sigma^*)}{2}\right) \exp((-1)\text{trace}(\Sigma')) \frac{\exp\left(\frac{(-1)\text{trace}(Z^{*T} Z^*)}{2}\right)}{\det|Z^*|} \quad (1) \\
&:= \text{mvEMG}(Z' | \mu_Z, \Sigma_Z, \lambda_Z)
\end{aligned}$$

where $z_i = \sum_j^d V_{ij}^{-1} (\sqrt{\psi_{ij}^T} U_{ij})^{-1} \cdot Z'_j \cdot V_{ij}^{-T} - (\sqrt{\psi_{ij}^T} U_{ij})^{-1} \mu_{Zj}$. We ensure that Z^* is non-singular.

4 Experiments

We compare the mvEMG with the multivariate skew-normal distribution (mvSN). The mvSN was introduced by *Azzalini and Capitanio (1999)*; *Azzalini and DallaValle (1996)* and is a simple function of a multivariate Gaussian pdf and a univariate Gaussian cdf, given by: $f(\mathbf{x}) = 2\Phi(\mathbf{x}; \mathbf{0}, \Sigma)\Phi(\alpha^T \mathbf{x})$, $\mathbf{x} \in \mathcal{R}^d$ where Φ is the multivariate Gaussian density centered at $\mathbf{0}$, covariance matrix Σ and $\Phi(\cdot)$ is the CDF of the univariate spherical Gaussian, $\mathcal{N}(0, 1)$. α contains the shape for the Gaussian cdf. One convenient property of the mvSN is that the marginal distributions are scalar skew-normal distributions. The filled contour plots of the mvSN pdf are shown in Figure 2 (upper panel). Notice that when α are zeros (Figure 2 upper panel, first plot), the pdf reduces to that of a multivariate Gaussian. When α are positive, the pdfs attain right skewness. Figure 2 (lower panel) shows the filled contour plots for the mvEMG pdf generated using Equation 1. The scale parameters λ_Z constitute the off-diagonal elements of the shear matrix S that consists the skew-symmetric matrix as discussed in Section 3.2. We observe that the heights of the pdf are shorter than that of mvSN indicating that the mvEMG distributions, with the same shape and covariance parameters as that of the mvSN, are broader distributions with fatter tails.

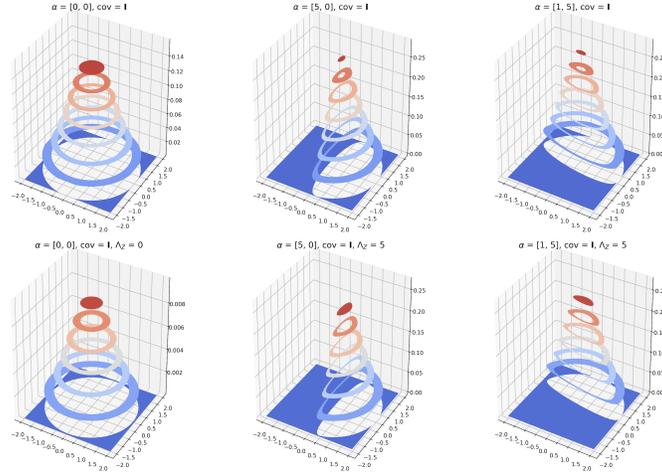


Figure 2: Filled contour plots of pdfs: **(Upper panel)**: The multivariate skew normal (mvSN) pdf for varying shape parameters α and the spherical covariance matrix Σ . **(Lower panel)**: The mvEMG pdf for the same varying shape parameters α (as above), the spherical covariance matrix Σ and varying scale parameters λ_Z that constitute the off-diagonal elements of the shear matrix S .

5 Conclusion

In this paper, we present mvEMG, the multivariate generalisation to a univariate exponentially-modified Gaussian (EMG) distribution, to handle multivariate skewed data giving access to mean,

covariance and skewness of the distribution. We make use of an affine transformation involving rotation, translation and shearing to accommodate for the three moments in mvEMG. We present the statistical properties for mvEMG. The experiments on synthetic data show that the mvEMG is potentially better suited to model fat-tailed distributions than the mvSN distribution.

Future work: We would implement an efficient sampling strategy for probabilistic inference to enable downstream analysis such as cluster identification or building graphical models. We will further validate mvEMG on a real-world dataset to depict its usefulness where skewness plays a crucial role in identifying underlying data patterns and that lead to further insights into the data.

References

- [1] Zeisel, Amit & Muñoz-Manchado, Ana B. & Codeluppi, Simone & Lönnerberg, Peter & La Manno, Gioele & Juréus, Anna & Marques, Sueli & Munguba, Hermany & He, Liquan & Betscholtz, Christer & Rolny, Charlotte & Castelo-Branco, Gonçalo & Hjerling-Lefler, Jens & Linnarsson, Sten (2015) Cell types in the mouse cortex and hippocampus revealed by single-cell RNA-seq. *Science* **347**(6226):1138-1142.
- [2] Esary, James D & Marshall, Albert (1974) Multivariate distributions with exponential minimums. *The Annals of Statistics*.
- [3] Naish, P. J. & Hartwell, S. (1998) Exponentially modified gaussian functions—a good model for chromatographic peaks in isocratic hplc? *Chromatographia*, **26**(1):285–296.
- [4] Kong, H. & Ye, F. & Lu, X. & Guo, L. & Tian, J. & Xu, G. (2005) Deconvolution of overlapped peaks based on the exponentially modified gaussian model in comprehensive two-dimensional gas chromatography. *Journal Of Chromatography A*, **1086**(1-2):160–164.
- [5] Golubev, A. (2010) Exponentially modified gaussian (emg) relevance to distributions related to cell proliferation and differentiation. *Journal of theoretical biology*, **262**(2):257–266
- [6] Tyson, D. R. & Garbett, S. P. & Frick, P. L. & Quaranta, V. Fractional proliferation: a method to deconvolve cell population dynamics from single-cell data. *Nature methods*, **9**(9):923, 2012
- [7] Irizarry, R. A. & Hobbs, B. & Collin, F. & Beazer-Barclay, Y. D. & Antonellis, K. J. & Scherf, U. & Speed, T. P. (2003) Exploration, normalization, and summaries of high density oligonucleotide array probe level data. *Biostatistics*, **4**(2):249–264.
- [8] Azzalini, A. & Capitanio, C. (1999) Statistical applications of the multivariate skew normal distribution. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **61**(3):7579–602
- [9] Azzalini, A. & DallaValle, A. (1996) The multivariate skew-normal distribution. *Biometrika*, **83**:715–726.