

# Annotated Bibliography

Cassandra Jin

## Textbooks

### Mixture Models and Applications

Bouguila, N., & Fan, W. (Eds.). (2020). Mixture models and applications. Berlin: Springer.

<https://link.springer.com/content/pdf/10.1007/978-3-030-23876-6.pdf>

- Section 5.2 of this textbook on mixture models provides some mathematical context for GMMs and a few visualizations of multivariate Gaussian distributions with different shape parameters. It also explains the component of expectation maximization and includes more specific algorithms that stem from generalized Gaussian mixture models. There are multiple applications as well, with measures of performance and comparison.

### Gaussian Mixture Models from Signals and Communication Technology

Yu, D., Deng, L. (2015). Gaussian Mixture Models. In: Automatic Speech Recognition. Signals and Communication Technology. Springer, London. [https://doi.org/10.1007/978-1-4471-5779-3\\_2](https://doi.org/10.1007/978-1-4471-5779-3_2)

[https://link.springer.com/chapter/10.1007/978-1-4471-5779-3\\_2](https://link.springer.com/chapter/10.1007/978-1-4471-5779-3_2)

- This textbook chapter lays a strong statistical groundwork for GMMs, first introducing the basic concepts of random variables and the associated distributions and then applying to Gaussian random variables and mixture-of-Gaussian random variables. It goes on to fit the model to real-world data such as speech recognition systems. This source does particularly well in discussing key advantages and weaknesses of GMMs and describing the principles of maximum likelihood and the EM algorithm in detail.

## Original Paper

### Pattern Classification and Scene Analysis

Duda, R. O., & Hart, P. E. (1973). Pattern classification and scene analysis (Vol. 3, pp. 731-739). New York: Wiley.

<https://www.svms.org/classification/DuHS95.pdf>

- This is an early paper on GMMs that begins by breaking down the restrictive assumption that the functional forms for the underlying probability densities are known. It explains that standard algorithms suffer from problems associated with parametric assumptions without providing any of the benefits of computational simplicity. So this paper provides the motivations behind learning the value of an unknown parameter vector by partitioning the data into subgroups or clusters. It then shows how this method leads to useful tools for pattern recognition problems.

## Background/Other

### What is the expectation maximization algorithm?

Do, C. B., & Batzoglou, S. (2008). What is the expectation maximization algorithm?. *Nature biotechnology*, 26(8), 897-899.

<https://www.nature.com/articles/nbt1406>

- This article explains how the expectation maximization algorithm works and shows how it arises in many computational biology applications that involve probabilistic models. This will be useful background information for GMMs, and the article also touches on parameter estimation in probabilistic models with incomplete data.

### Gaussian Mixture Models from Encyclopedia of Biometrics

Reynolds, D. A. (2009). Gaussian mixture models. *Encyclopedia of biometrics*, 741(659-663).

[https://link.springer.com/referenceworkentry/10.1007/978-1-4899-7488-4\\_196](https://link.springer.com/referenceworkentry/10.1007/978-1-4899-7488-4_196)

- This article is a digestible introduction of GMM with a few helpful, fundamental visualizations of the method.

### A Novel Gaussian Mixture Model for Classification

Wan, H., Wang, H., Scotney, B., & Liu, J. (2019, October). A novel gaussian mixture model for classification. In *2019 IEEE International Conference on Systems, Man and Cybernetics (SMC)* (pp. 3298-3303). IEEE.

<https://ieeexplore.ieee.org/document/8914215>

- This piece analyzes the performance of GMM as a classifier in comparison with other conventional classifiers such as k-nearest neighbors (KNN), support vector machine (SVM), decision tree and naive Bayes. It finds that GMMs do not do well relative to the other methods, and the authors propose a new GMM classifier, SC-GMM, which is significantly more competitive with the other classification approaches.

## Blogs

### Gaussian Mixture Models: What are they & when to use?

Kumar, A. (2022). Gaussian mixture models: What are they & when to use. Online, Apr.

[https://vitalflux.com/gaussian-mixture-models-what-are-they-when-to-use/#What\\_are\\_some\\_real-world\\_examples\\_where\\_Gaussian\\_mixture\\_models\\_can\\_be\\_used](https://vitalflux.com/gaussian-mixture-models-what-are-they-when-to-use/#What_are_some_real-world_examples_where_Gaussian_mixture_models_can_be_used)

- This blog piece explains what are GMMs are and how the foundational expectation-maximization method acts in relation to GMM. It also delves into the key steps of using Gaussian mixture models for clustering, showing the differences between GMMs and other types of clustering algorithms such as k-means (since GMMs are a generalization of k-means clustering). Lastly, it lists scenarios when Gaussian mixture models can be used and briefly explains some real-world examples of application.

### Gaussian Mixture Model Clearly Explained

<https://towardsdatascience.com/gaussian-mixture-model-clearly-explained-115010f7d4cf>

- This blog gives an elementary explanation of GMMs and does well in showing code application of the model in Python. It aligns the steps in the model process with steps in the code and includes multiple basic visualizations of how GMMs can be used to break down multi-distributional data.

## Gaussian mixture models: k-means on steroids

<https://www.r-bloggers.com/2018/12/gaussian-mixture-models-k-means-on-steroids/>

- This blog delves into the R implementation of GMMs, showing different (extreme) types of datasets containing multiple distributions and how GMM performs in capturing their nature. It also lists out the parameters up for estimation which contribute to the performance of the GMM.

## Applications

### Robust text-independent speaker identification using Gaussian mixture speaker models

Reynolds, D. A., & Rose, R. C. (1995). Robust text-independent speaker identification using Gaussian mixture speaker models. *IEEE transactions on speech and audio processing*, 3(1), 72-83.

<https://ieeexplore.ieee.org/abstract/document/365379>

- This paper introduces the use of GMMs for robust text-independent speaker identification. The application here required high identification rates using short utterance from unconstrained conversational speech and robustness to degradations produced by transmission over a telephone channel. It is a complete experimental evaluation of the Gaussian mixture speaker model, and in the process, it examines algorithmic issues (initialization, variance limiting, model order selection), spectral variability robustness techniques, large population performance, and comparisons to other speaker modeling techniques (uni-modal Gaussian, VQ codebook, tied Gaussian mixture, and radial basis functions).

### Improved performance of data-driven simulator of liquid rocket engine under varying operating conditions

Satoh, D., Omata, N., Tsutsumi, S., Hashimoto, T., Sato, M., Kimura, T., & Abe, M. (2023). Improved performance of data-driven simulator of liquid rocket engine under varying operating conditions. *Acta Astronautica*.

<https://www.sciencedirect.com/science/article/pii/S0094576523005660>

- This paper uses GMM as a data-driven approach to predict the behavior of a reusable rocket engine, which cannot be determined using a model-based approach. It proposes a clustered regression model to predict various operating conditions for the engine. Using the Gaussian mixture model, the training data were divided in advance according to the operating conditions, and multiple regression models were individually trained on each cluster. The researchers were then able to accurately predict the combustion and chill-down phases and transition between phases, compared with the prediction results of a single regression model.

### Analyzing the impact of fare-free public transport policies on crowding patterns at stations using crowdsensing data

Qing-Long, L., Mahajan, V., Cheng, L., & Antoniou, C. (2023). Analyzing the impact of fare-free public transport policies on crowding patterns at stations using crowdsensing data. Preprint.

<https://mediatum.ub.tum.de/doc/1725655/document.pdf>

- This paper explains that cities worldwide are experimenting with a variety of fare discount policies, which will lead to spatiotemporal changes in mobility patterns, including crowding in public transport (PT) stations. To better model and understand these spatiotemporal dynamics, the researchers focus on crowding patterns in PT stations using crowdsensing data and employ a GMM to cluster PT stations based on their crowding pattern deviations at different stages of policy implementation. This clustering step enables the identification of distinct station types based on their response to the policy, and the

the clustering results show three station types: unaffected, mildly stimulated, and intensely stimulated stations.

## **Deep Learning for Deep Earthquakes: Insights from OBS Observations of the Tonga Subduction Zone**

Xi, Z., Wei, S. S., Zhu, W., Beroza, G., Jie, Y., & Saloor, N. (2023). Deep Learning for Deep Earthquakes: Insights from OBS Observations of the Tonga Subduction Zone.

<https://eartharxiv.org/repository/view/6222/>

- This paper introduces an integrated deep-learning-based workflow to detect deep earthquakes recorded by ocean-bottom seismographs (OBSs) and land-based stations in the Tonga subduction zone. With the goal of detecting different wave arrivals in the time-frequency domain, the researchers use Gaussian Mixture Model Associator GaMMA-1D to find associations between predicted phases. They find that this approach suppresses false predictions and significantly improves the detection of small earthquakes.

## **R Packages**

### **Gaussian mixture models**

<https://rdr.io/cran/ClusterR/man/GMM.html>

- Contains arguments and examples of **GMM** package

### **MGMM: An R Package for fitting Gaussian Mixture Models on Incomplete Data**

<https://www.biorxiv.org/content/10.1101/2019.12.20.884551v2.full>

- This article introduces MGMM, an R package for fitting GMMs in the presence of missing data. Using three case studies on real and simulated data sets, it demonstrates that, when the underlying distribution is near to a GMM, MGMM is more effective at recovering the true cluster assignments than state of the art imputation followed by standard GMM.