# CS7641 A3: Unsupervised Learning and Dimensionality Reduction

Scott Schmidl

*sschmidl3@gatech.edu*

## I. INTRODUCTION

Clustering, Expectation Maximization (EM), Principal Component Analysis (PCA), Independent Component Analysis (ICA), Randomized Projection, and Manifold Learning represent a diverse array of techniques in the field of machine learning and data analysis. Each of these methods plays a crucial role in uncovering patterns, reducing dimensionality, and extracting meaningful information from complex datasets. As the volume and complexity of data continue to grow in various domains such as finance, healthcare, and social media, the importance of these techniques becomes increasingly evident.

In this paper, we delve into the principles, applications, and comparative analyses of these key methods. We begin by elucidating the fundamental concepts behind clustering and its variants, EM, and the role they play in grouping similar data points. Following this, we explore dimensionality reduction techniques such as PCA and ICA, which enable us to capture the essential features of high-dimensional data while preserving relevant information. Additionally, we discuss the utility of randomized projection in efficiently approximating high-dimensional data and the significance of manifold learning in uncovering the underlying geometric structure of data distributions.

Through a comprehensive examination of these techniques, we aim to provide insights into their respective strengths, weaknesses, and real-world applications. By understanding the intricacies of these methods, researchers and practitioners can make informed decisions in selecting the most suitable approach for their specific data analysis tasks, ultimately advancing the field of machine learning and data science.

## II. STEP 1

Quisque faucibus egestas fermentum.

### A. Expectation Maximization

Nulla consequat, tortor sit amet interdum tempus, ante mauris vulputate dui, et bibendum ipsum nisl vitae ante.

### B. Clustering Algorithm 2

Pellentesque efficitur magna pharetra, molestie libero vel, tempus justo.

## III. STEP 2

Quisque faucibus egestas fermentum.

### A. Principal Component Analysis

Nulla consequat, tortor sit amet interdum tempus, ante mauris vulputate dui, et bibendum ipsum nisl vitae ante.

### B. Independent Component Analysis

Pellentesque efficitur magna pharetra, molestie libero vel, tempus justo.

*1) Randomized Projections:* Pellentesque auctor eros justo, nec cursus ligula porta tincidunt.

*2) Dim. Reduction Algorithm 4:* Nulla pharetra felis ut felis auctor convallis.

## IV. STEP 3

Quisque faucibus egestas fermentum. Nulla consequat, tortor sit amet interdum tempus, ante mauris vulputate dui, et bibendum ipsum nisl vitae ante. Pellentesque efficitur magna pharetra, molestie libero vel, tempus justo.

## V. STEP 4

Quisque faucibus egestas fermentum. Nulla consequat, tortor sit amet interdum tempus, ante mauris vulputate dui, et bibendum ipsum nisl vitae ante. Pellentesque efficitur magna pharetra, molestie libero vel, tempus justo.

## VI. STEP 5

Quisque faucibus egestas fermentum. Nulla consequat, tortor sit amet interdum tempus, ante mauris vulputate dui, et bibendum ipsum nisl vitae ante. Pellentesque efficitur magna pharetra, molestie libero vel, tempus justo.

## VII. CONCLUSION

In conclusion, the methods of clustering, Expectation Maximization, PCA, ICA, randomized projection, and manifold learning stand as indispensable tools in the arsenal of data scientists and machine learning practitioners. From organizing unlabeled data points into meaningful clusters to uncovering the latent structure of high-dimensional datasets, each technique offers unique capabilities and insights into the underlying patterns within data.

As we navigate the complexities of modern datasets characterized by high dimensionality, noise, and non-linear relationships, the importance of these methods cannot be overstated. Whether it be for exploratory data analysis, feature extraction, or dimensionality reduction, understanding the principles and intricacies of these techniques empowers researchers to extract actionable insights and drive informed decision-making.

Looking ahead, further advancements in these areas hold the promise of addressing emerging challenges in fields ranging from biomedical research to natural language processing. By fostering interdisciplinary collaboration and continuing to refine these methodologies, we can unlock new avenues for knowledge discovery and innovation, propelling the frontier of machine learning and data analysis into the future.

## VIII. RESOURCES

[1] *API Reference.* Python. https://www.python.org/.
[2] *API Reference.* Pandas. https://pandas.pydata.org/.
[3] *API Reference.* Jupyter. https://jupyter.org/.
[4] *API Reference.* IPython. https://ipython.org/.
[5] *API Reference.* MatPlotLib. https://matplotlib.org/stable/.
[6] *API Reference.* NumPy. https://numpy.org/.
[7] *API Reference.* Seaborn. https://seaborn.pydata.org/index.html.
[8] *API Reference.* OpenPyXl. https://openpyxl.readthedocs.io/en/stable/tutorial.html.
[9] *API Reference.* Scikit-Learn. https://scikit-learn.org.
[10] Nakamura, K. (2023). *ML LaTeX Template.*
[11] *Data Source* Nutrition Facts Database Tools and Spreadsheet. https://tools.myfooddata.com/nutrition-facts-database-spreadsheet.php.
[12] *Data Source* Cardiovascular Disease. https://www.kaggle.com/datasets/colewelkins/cardiovascular-disease.