# UMAP(Uniform Manifold Approximation and Projection)

Uniform Manifold Approximation and Projection (UMAP) is a powerful dimension reduction technique that has gained significant traction in the fields of machine learning and data visualization. Developed by Leland McInnes, John Healy, and James Melville, UMAP is built on solid mathematical foundations, including Riemannian geometry and algebraic topology. This article delves into the technical aspects of UMAP, its underlying principles, implementation, and practical applications.

UMAP, was introduced by Leland McInnes, John Healy, and James Melville in 2018, it gained immediate popularity for its ability to address the shortcomings of traditional methods. UMAP belongs to the family of manifold learning techniques and stands out for its efficiency, scalability, and capability to capture both local and global structures in the data.

## Mathematical Foundations for UMAP

UMAP is grounded in several key mathematical concepts:

- **Riemannian Manifold**: UMAP assumes that the data is uniformly distributed on a Riemannian manifold. This means that the data points lie on a smooth, curved surface that can be locally approximated by Euclidean space.

- **Riemannian Metric**: The Riemannian metric is locally constant or can be approximated as such. This metric defines the distance between points on the manifold.

- **Topological Data Analysis**: UMAP leverages topological data analysis to capture the structure of the data. It constructs a fuzzy topological representation of the data, which is then optimized to find a low-dimensional embedding.

### 1. Graph Construction

#### a. Nearest Neighbor Search

For each data point $x_i$, UMAP identifies its k-nearest neighbors. This is typically done using a distance metric such as Euclidean distance. Let's denote the set of k-nearest neighbors of $x_i$ as $\{x_{i1}, x_{i2}, \ldots, x_{ik}\}$.

#### b. Similarity Computation

The similarity between points $x_i$ and $x_j$ is computed using a Gaussian kernel:

$$sim_{ij} = exp(-\frac{d(xi, xj)^2}{\sigma_i^2})$$

where $d(x_i, x_j)$ is the Euclidean distance between $x_i$ and $x_j$, and $\sigma_i$ is a local scaling parameter. The value of $\sigma_i$ is chosen such that the probability $p_{ij}$ of selecting $x_i$ as a neighbor of $x_j$ can be controlled, typically by ensuring that only k neighbors are considered.

The scaling parameter $\sigma_i$ is chosen so that the distance $d(x_i,x_j)$ falls within a range that makes $sim_{ij}$ meaningful for the local neighborhood structure.

## 2. Probabilistic Modeling

### a. High-Dimensional Space

For a given data point $x_i$, the probability of $x_j$ being a neighbor is given by:

$$p_{ij} = \frac{sim_{ij}}{\sum_{k \neq i} sim_{ij}}$$

This is the probability of observing $x_j$ given $x_i$ in the high-dimensional space. It ensures that the probability mass is distributed according to the similarities between $x_i$ and its neighbors.

### b. Low-Dimensional Space

In the low-dimensional space, we want to model the data points $y_i$ and $y_j$ such that the distances between them reflect the similarities from the high-dimensional space. The probability $q_{ij}$ of $y_j$ being a neighbor of $y_i$ is given by:

$$q_{ij} = \frac{\exp(-dist(y_i, y_j)2}{\sum_{k \neq i} \exp(-dist(y_i, y_k)2)}$$

where dist(yi,yj)\ is typically the Euclidean distance between yi and yj in the low-dimensional space.

## 3. Objective Function

UMAP's objective is to make the low-dimensional probability distribution qij as similar as possible to the high-dimensional probability distribution pij. To achieve this, UMAP minimizes a loss function based on the divergence between pij and qij. The typical objective function used is:

$$L = \sum_i \sum_{j \neq i} \log\left(\frac{p_{ij}}{q_{ij}}\right)p_{ij}$$

This can also be expressed in terms of Kullback-Leibler divergence:

$L=KL(P\|Q)= \sum_i \sum_{j \neq i} \log\left(\frac{p_{ij}}{q_{ij}}\right)p_{ij}$

where KL(P‖Q) measures how one probability distribution P diverges from a second, expected probability distribution Q.

## 4. Optimization

UMAP optimizes the objective function using stochastic gradient descent (SGD). The steps involved are:

### a. Initialization

- Initialize the low-dimensional coordinates {yi} randomly.

### b. Gradient Calculation

- Compute gradients of the objective function with respect to the low-dimensional coordinates. The gradient for the loss function L is:

- $\frac{\delta L}{\delta y_i} = -\sum_{j \neq i}(p_{ij} - q_{ij})\frac{\delta}{\delta y_i}dist(y_i, y_j)2$

Where $\frac{\delta}{\delta y_i}dist(y_i, y_j)2$ represents how changes in yi affect the dist(yi,yj).

## c. Update

- Update the coordinates yi using the gradients and a learning rate η:

- $y_i \leftarrow y_i - \eta\frac{\delta L}{\delta y_i}$

## 5. Embedding

Through the optimization process, the coordinates {yi} are adjusted to minimize L. The final embedding is a low-dimensional representation where the pairwise distances approximate the similarities in the high-dimensional space, effectively capturing the manifold structure.

# Applications of UMAP

UMAP has been applied in various fields, including genomics, where it has helped reveal cryptic population structures and phenotype correlations

**Classical Machine Learning:** UMAP is employed in various machine learning tasks, such as clustering, classification, and anomaly detection. Its ability to preserve the intrinsic structure of the data makes it a valuable preprocessing step for improving the performance of machine learning models.

**Image Analysis:** UMAP has proven effective in image analysis, aiding researchers and practitioners in visualizing and exploring high-dimensional image datasets. This is particularly valuable in fields such as computer vision and medical imaging.

**Social Network Analysis:** UMAP is applied to understand and visualize relationships in social networks. By reducing the dimensionality of the data while preserving relevant structures, it becomes easier to identify communities and patterns in large-scale social graphs.

In summary, UMAP stands out as a crucial tool for exploratory data analysis and preprocessing within machine learning pipelines, thanks to its efficient handling of large datasets and its adept preservation of both local and global data structures.