

고차원 데이터 임베딩은 복잡한 데이터 구조를 저차원으로 표현하는 중요한 도구다. 이런 임베딩 방식 중 주요한 두 가지 방법으로 UMAP과 PaCMAP이 있다.

UMAP, 즉 Uniform Manifold Approximation and Projection는 고차원 데이터의 기하학적 구조를 저차원으로 압축하는 알고리즘이다. 이 방법은 데이터의 로컬 및 글로벌 구조를 동시에 보존하는 데 초점을 둔다. UMAP은 고차원의 데이터 포인트 간의 topological 거리와 밀도를 저차원에서 효과적으로 반영하려고 한다. 이 때문에 UMAP은 대규모 데이터셋을 빠르게 처리할 수 있고, 다양한 데이터 유형에도 잘 맞는다.

한편, PaCMAP, Pairwise Controlled Manifold Approximation Projection는 고차원 데이터를 저차원으로 변환할 때 명시적인 거리 제약을 도입한다는 특징이 있다. 이 방법은 데이터 포인트 간의 거리를 명확하게 제어하면서, 이웃 관계를 유지하고 멀리 떨어진 포인트 간의 거리도 적절하게 반영하려고 설계되었다. 이런 접근 방식 덕분에 PaCMAP는 UMAP에 비해 더 정교한 임베딩 결과를 보여준다. 계산 효율성에서는 UMAP과 크게 다르지 않다.

결론적으로, UMAP과 PaCMAP 모두 데이터의 시각화, 차원 축소, 클러스터링 등의 다양한 응용에서 중요한 역할을 한다. 사용하려는 데이터와 목적에 따라 어떤 방법을 선택할지 결정해야 한다.