

Introduction to Machine Learning

Foundations and Applications

Paul J. Atzberger
University of California Santa
Barbara





Unsupervised Learning: Dimension Reduction

Dimension Reduction

Motivations for Unsupervised Learning:

Given data set $S = \{x_1, x_2, \dots, x_m\}$ what types of patterns or structure occur.

Insights into how features $\mathbf{x}_i = (x_i^1, x_i^2, \dots, x_i^N)$ relate to one another.

Insights into low-dimensional structure inherent in data.

Insights into groupings or clustering of data (number of classes).

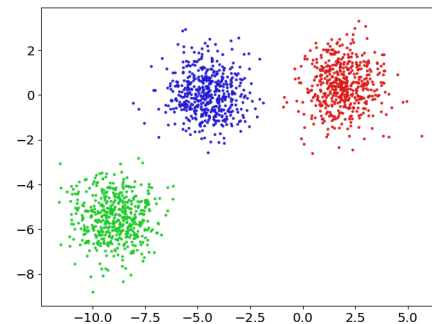
Many methods to consider depending on aims:

- **Clustering Methods (K-means)**
- **Principal Component Analysis (PCA / KPCA)**
- **Manifold Learning** + (many more methods)

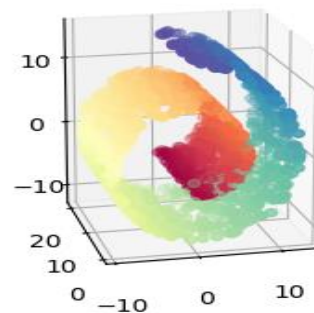
Overall methods share some common features.

Abstractly trying to learn characteristics of $\mathcal{D} \sim \mathcal{X}$.

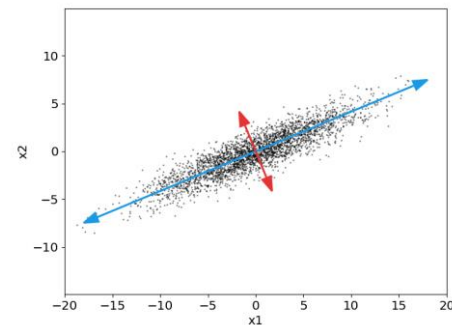
K-Means Clustering



Manifold Learning



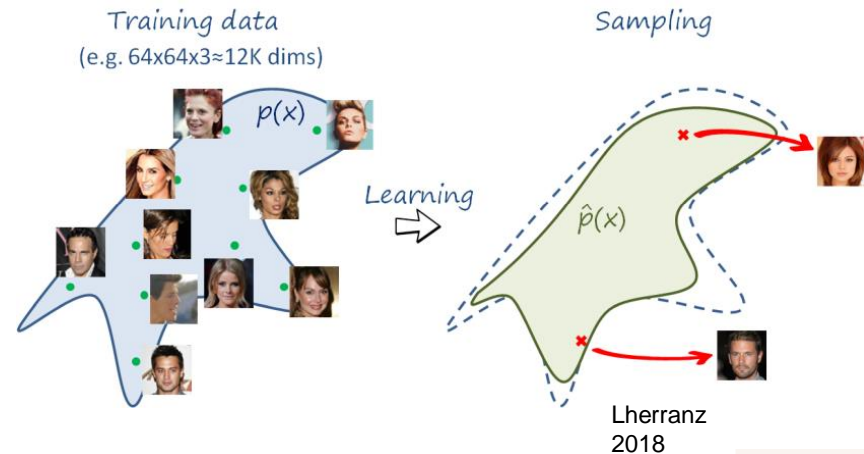
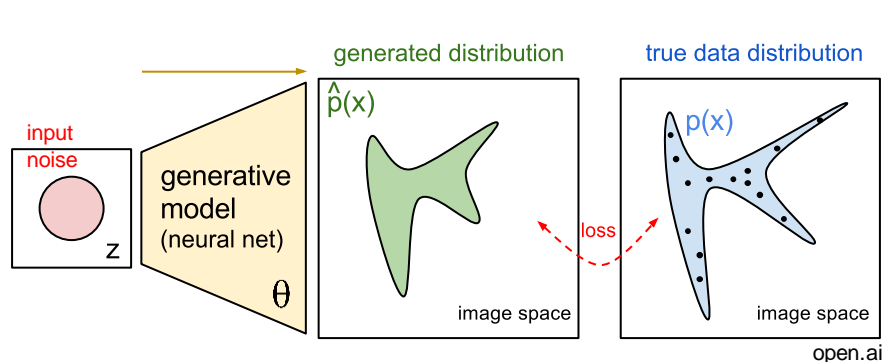
Principal Component Analysis



Motivations

Manifold-like structures in high dimensional spaces (natural images, audio, physical fields, PDE solutions).

Challenge: How to learn high dimensional probability distributions, generators $G(z)$ for sampling?

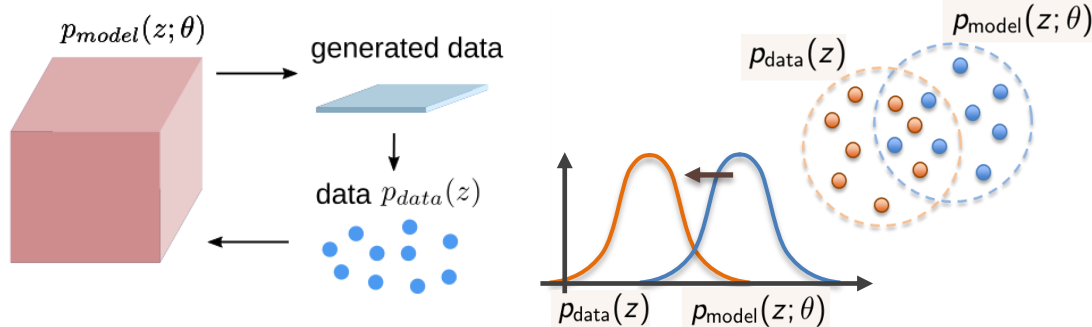



Generative Modeling

Approaches for learning models:


- Bayesian Methods
- Maximum Likelihood Estimation (MLE)
- Generative Adversarial Networks (GANs) (discussed in later lecture)
- ... and many other approaches.

Challenge: How to do this in a tractable way?





Clustering: K-Means



K-Means Clustering

Task: Given data set $S = \{x_1, x_2, \dots, x_m\}$ find partition into k sets $\Omega = \{\Omega_1, \Omega_2, \dots, \Omega_k\}$.

K-Means Clustering Optimization Problem:

$$\arg \min_{\Omega} \sum_{l=1}^k \sum_{x \in \Omega_l} \|x - \mu_l\|^2$$

Challenge: Exact solution is NP-hard (requires considering all partitions).

Need to use approximate methods.

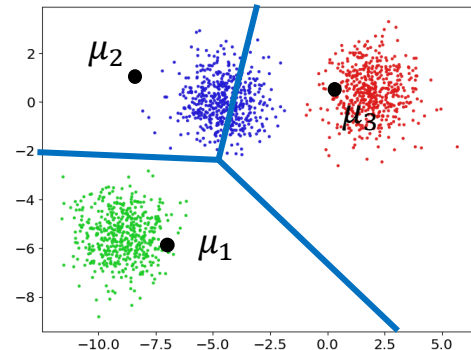
Voronoi Iteration (Lloyd's Algorithm):

- Randomly choose k initial seed points.
 - Compute Voronoi Cells of seed points and centroids.
 - Move seed points to centroid of points in each Voronoi Cell.
 - Repeat until termination criteria (seed points move less than ϵ).

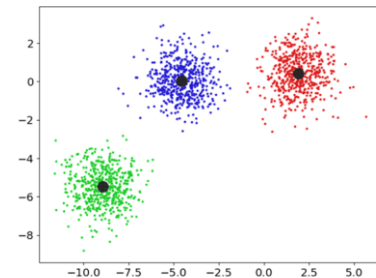
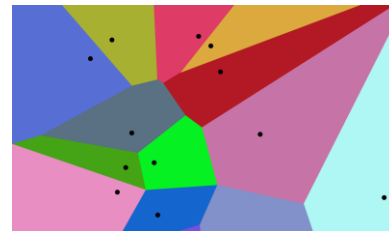
Voronoi iteration yields seeds tending toward distribution w/ more uniform cells.

Seed points tend to migrate toward cluster centers giving approximate solution.

K-Means Clustering



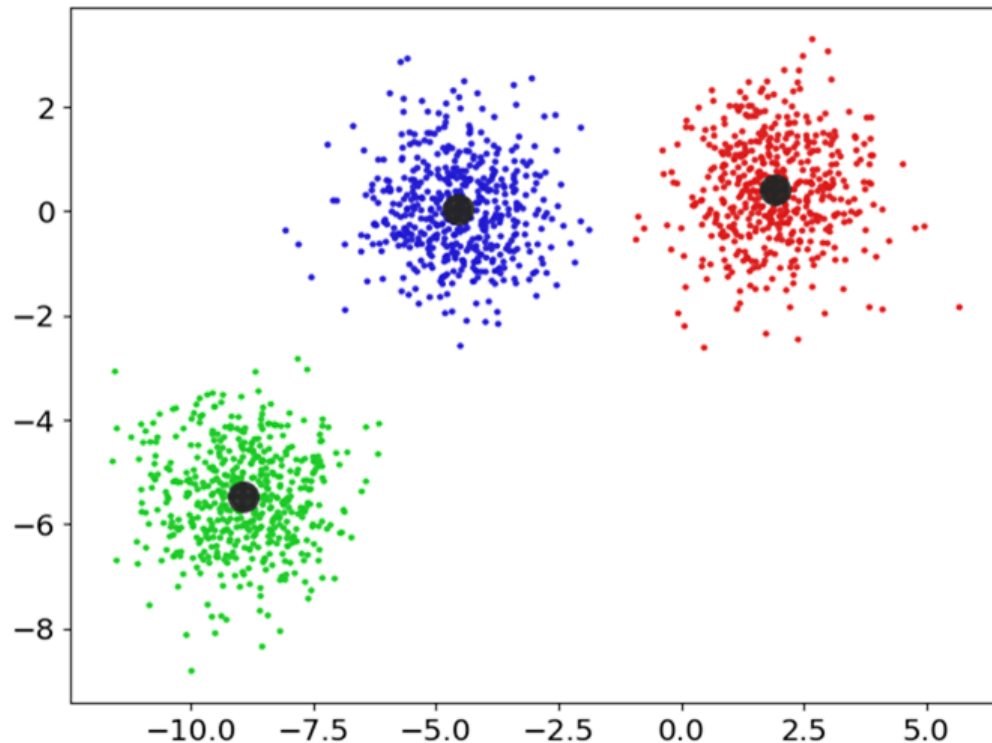
Voronoi Diagram



K-Means Clustering: Example

Example: $k = 3$ clusters, std. dev. $\sigma = 1.0$, $m = 1.5 \times 10^3$ samples.

K-Means Clustering



K-Means Clustering: Example

Example: $k = 3$ clusters, std. dev. $\sigma = 1.0$, $m = 1.5 \times 10^3$ samples.

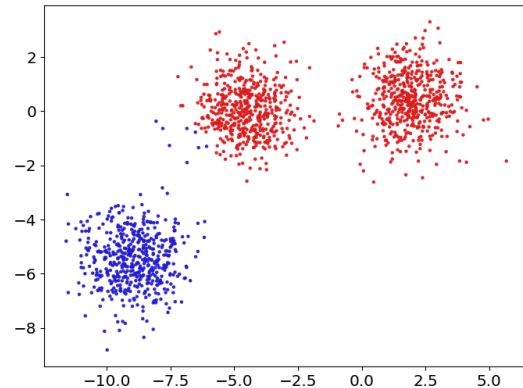
K-Means can be sensitive to hyperparameters and data distribution.

- Posing the wrong number k of clusters.
- Anisotropic cluster distribution.
- Unequal variance or size of clusters.

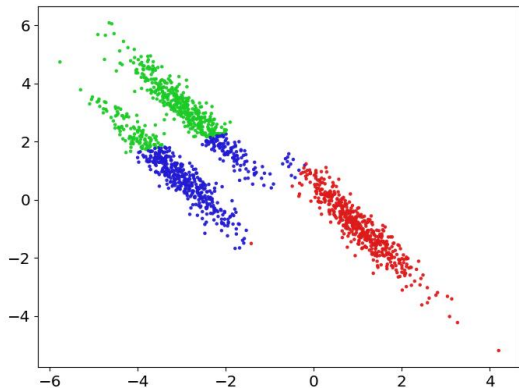
Cross validation can be used to help determine good hyperparameters.

Strong assumption that data is distributed in distinct clusters.

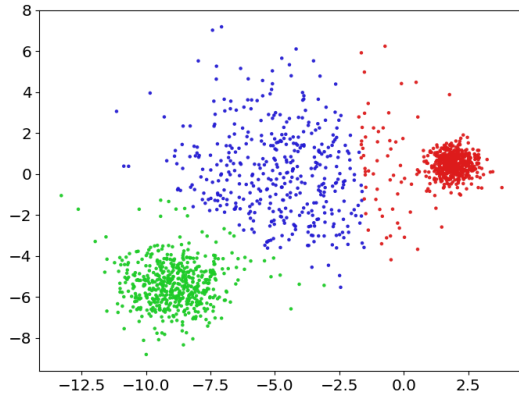
Wrong Number Clusters ($k=2$)



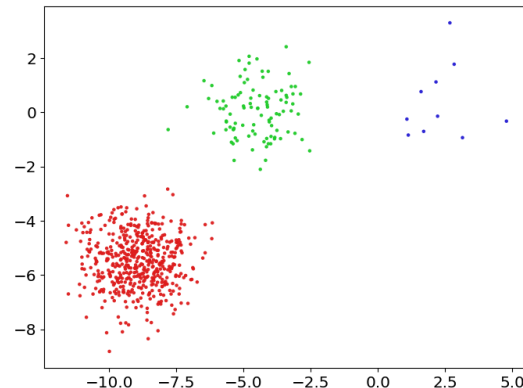
Anisotropic Clusters




Unequal Variance Clusters




Unequal Cluster Size





Singular Value Decomposition (SVD) and Principal Component Analysis (PCA)



Singular Value Decomposition (SVD)

Task: Given matrix C find the best approximating rank r matrix A_r .

Optimization Problem:

$$\begin{aligned} \min_A & \|C - A\|_2 \\ \text{subject } & \text{rank}(A) = r. \end{aligned}$$

Solution:

Singular Value Decomposition (SVD): $C = U\Lambda V^T \rightarrow A_r = U\Lambda_r V^T$

Features:

U, V are orthonormal matrices.

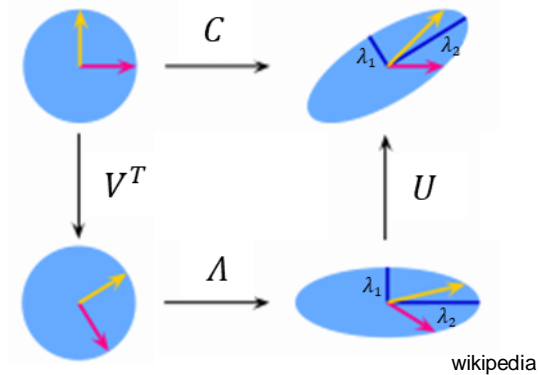
Λ is diagonal matrix of singular values $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$.

Rank r matrix $A_r = U\Lambda_r V^T$ where $\Lambda_r = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_r, 0, \dots, 0)$ matrix of largest r singular values.

Best approximation in sense $\|C - A_r\|_2 \leq \|C - A\|_2$ for any A with $\text{rank}(A) = r$.

Many useful applications.

Singular Value Decomposition



Example: Image Compression using SVD

Task: Given matrix C find the best approximating rank r matrix A ?

Optimization Problem:

$$\min_A \|C - A\|_2$$

subject $rank(A) = r$.

Image Compression:

View image as matrix C of column vectors \mathbf{x}_i .

Singular Value Decomposition: $C = UAV^T$.

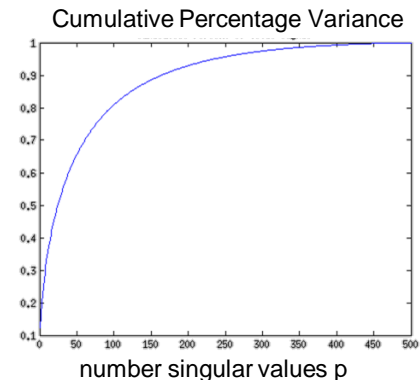
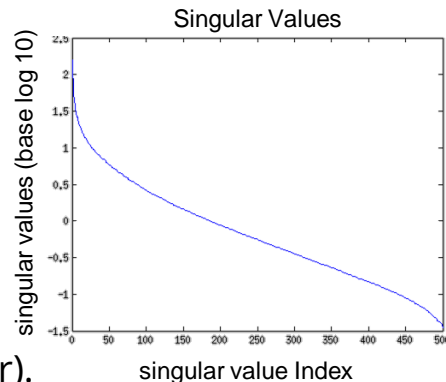
Compression \rightarrow find p -vectors that span columns of C .

Keep only first p singular vectors of U .

Many other compression techniques for images (better).



image 800x500 pixels



Example: Image Compression using SVD

Task: Given matrix C find the best approximating rank r matrix A ?

Optimization Problem:

$$\min_A \|C - A\|_2$$

subject $rank(A) = r$.

Image Compression:

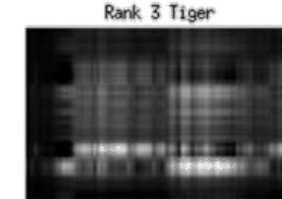
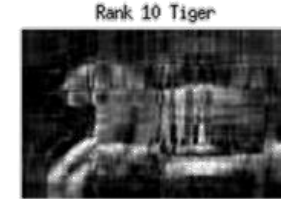
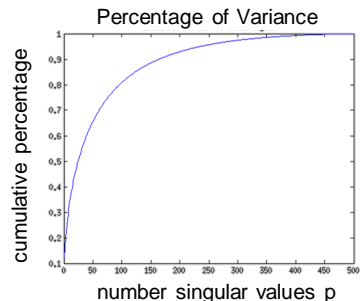
View image as matrix C of column vectors \mathbf{x}_i .

Singular Value Decomposition: $C = UAV^T$.

Compression \rightarrow find p -vectors that span columns of C .

Keep only first p singular vectors of U .

Many other compression techniques for images (better).



Example: Image Compression using SVD

Task: Given matrix C find the best approximating rank r matrix A ?

Optimization Problem:

$$\min_A \|C - A\|_2$$

subject $rank(A) = r$.

Image Compression:

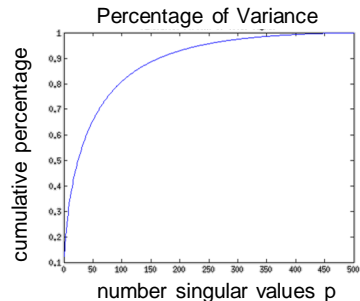
View image as matrix C of column vectors \mathbf{x}_i .

Singular Value Decomposition: $C = UAV^T$.

Compression \rightarrow find p -vectors that span columns of C .

Keep only first p singular vectors of U .

Many other compression techniques for images (better).



Rank Full



Rank 50



Rank 10



Rank 50 ~ 16% original size



Principal Component Analysis (PCA)



Principal Component Analysis (PCA)

Task: Given data set $S = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$ with $\mathbf{x}_i \in \mathbb{R}^N$ and $\sum_{i=1}^m \mathbf{x}_i = \mathbf{0}$ find the best approximating hyperplane $h = \text{span}\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_p\}$.

Optimization Problem:

$$\min_{W, Z} \sum_{i=1}^m \|\mathbf{x}_i - W\mathbf{z}_i\|_2^2$$

subject $W^T W = I_{p \times p}$, $W \in \mathbb{R}^{N \times p}$.

Finds optimal orthonormal basis \mathbf{w}_i best accounting for **covariance** $C = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^T \sim \sum_{a=1}^p \lambda_a \mathbf{w}_a \mathbf{w}_a^T$.

Closely related to the **Singular Value Decomposition (SVD)**: $C = U \Lambda V^T$.

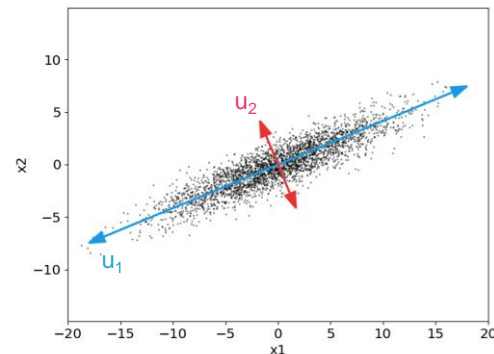
Solution: $W = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_p]$ first p singular vectors of U and $\mathbf{z}_i = \wp_W \mathbf{x}_i := W^T \mathbf{x}_i$, $\tilde{\mathbf{x}}_i = W \mathbf{z}_i = W W^T \mathbf{x}_i$.

Selects **p directions** of data $\{\mathbf{x}_i\}$ with **greatest variation**.

Projects data \mathbf{x}_i to the **nearest hyperplane point** to yield \mathbf{z}_i .

Low-dimensional linear structure then expect **good fit** for $p \ll N$.

Principal Component Analysis



Principal Component Analysis (PCA)

Task: Given data set $S = \{x_1, x_2, \dots, x_m\}$ with $x_i \in \mathbb{R}^N$ and $\sum_{i=1}^m x_i = 0$ find the best approximating hyperplane $h = \text{span}\{w_1, w_2, \dots, w_p\}$.

Optimization Problem:

$$\min_{W, Z} \sum_{i=1}^m \|x_i - Wz_i\|_2^2$$

subject $W^T W = I_{p \times p}$, $W \in \mathbb{R}^{N \times p}$.

Finds optimal orthonormal basis w_i best accounting for **covariance** $C = \frac{1}{m} \sum_{i=1}^m x_i x_i^T \sim \sum_{a=1}^p \lambda_a w_a w_a^T$.

Closely related to the Singular Value Decomposition (SVD): $C = U \Lambda V^T$

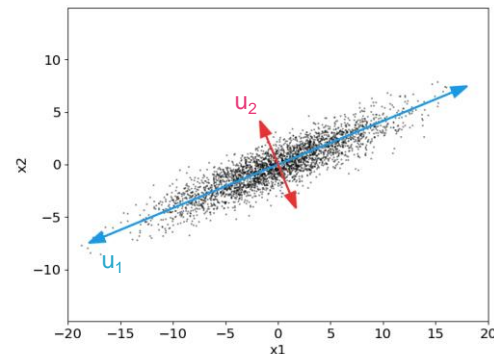
Solution: $W = [u_1, u_2, \dots, u_p]$ first p singular vectors of U and $z_i = \mathcal{P}_W x_i := W W^T x_i$.

Example: $m = 3 \times 10^3$ data points, exact: $u_1 = [\cos(\pi/8), \sin(\pi/8)]$, $u_2 = [-\sin(\pi/8), \cos(\pi/8)]$, $\Lambda = \text{diag}(36.0, 1.0)$.

Results: $\Lambda = \text{diag}(37.68, 1.01)$, $w_1 = [-0.926, -0.378]$, $w_2 = [-0.378, 0.926]$,
 $\cos(\pi/8) = 0.924$, $\sin(\pi/8) = 0.383 \rightarrow u_1 = [0.924, -0.383]$, $u_2 = [-0.383, 0.924]$.

Assumptions: Variation is indicative of importance of a feature, data distributed within a linear subspace.

Principal Component Analysis



Kernel Principal Component Analysis (KPCA):

Task: Given data set $S = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$ with $\mathbf{x}_i \in \mathbb{R}^N$, Φ feature map of K . Find in feature space $\mathcal{H}_{\text{RKHS}}$ the best approximating hyperplane $h = \text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_p\}$ of $\tilde{S} = \{\Phi(\mathbf{x}_1), \Phi(\mathbf{x}_2), \dots, \Phi(\mathbf{x}_m)\}$.

Optimization Problem:

$$\min_{W, Z} \sum_{i=1}^m \|\Phi(\mathbf{x}_i) - Wz_i\|_2^2$$

subject $W^T W = I_{p \times p}$, $W \in (\mathcal{H}_{\text{RKHS}})^p$.

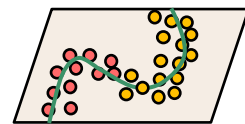
Reformulate: $\mathbf{q}_i = \Phi(\mathbf{x}_i)$, $C = \frac{1}{m} \sum_{i=1}^m \mathbf{q}_i \mathbf{q}_i^T$. Solution using SVD, $C = U \Lambda U^T$.

$$\lambda_a \mathbf{u}_a = C \mathbf{u}_a = \frac{1}{m} \sum_{i=1}^m \mathbf{q}_i \mathbf{q}_i^T \mathbf{u}_a = \frac{1}{m} \sum_{i=1}^m (\mathbf{q}_i^T \mathbf{u}_a) \mathbf{q}_i \Rightarrow \mathbf{u}_a = \sum_i \frac{\mathbf{q}_i^T \mathbf{u}_a}{m \lambda_a} \mathbf{q}_i = \sum_i \alpha_j^a \mathbf{q}_i, \mathbf{u}_a \in \text{span}\{\mathbf{q}_i\}.$$

Project the equations using

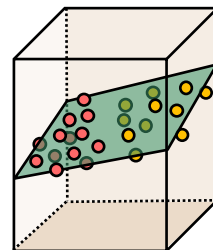
$$\begin{aligned} \mathbf{q}_i^T C \mathbf{u}_a &= \lambda_a \mathbf{q}_i^T \mathbf{u}_a, \Rightarrow \mathbf{q}_i^T \frac{1}{m} \sum_k \mathbf{q}_k \mathbf{q}_k^T \sum_j \alpha_j^a \mathbf{q}_j = \lambda_a \mathbf{q}_i^T \sum_j \alpha_j^a \mathbf{q}_j \\ &\Rightarrow \frac{1}{m} \sum_{j,k} \alpha_j^a (\mathbf{q}_i^T \mathbf{q}_k) (\mathbf{q}_k^T \mathbf{q}_j) = \lambda_a \sum_j \alpha_j^a (\mathbf{q}_i^T \mathbf{q}_j). \end{aligned}$$

Input Feature Space



$\Phi(\mathbf{X})$ ↓

New Feature Space



Kernel Principal Component Analysis (KPCA):

Optimization Problem:

$$\min_{W, Z} \sum_{i=1}^m \|\Phi(x_i) - Wz_i\|_2^2$$

$$\text{subject } W^T W = I_{p \times p}, W \in (\mathcal{H}_{\text{RKHS}})^p.$$

Reformulate: $q_i = \Phi(x_i)$, $C = \frac{1}{m} \sum_{i=1}^m q_i q_i^T$. Solution using SVD, $C = U \Lambda U^T$.

$$\lambda_a u_a = C u_a = \frac{1}{m} \sum_{i=1}^m q_i q_i^T u_a = \frac{1}{m} \sum_{i=1}^m (q_i^T u_a) q_i \Rightarrow u_a = \sum_i \frac{q_i^T u_a}{m \lambda_a} q_i = \sum_i \alpha_i^a q_i, u_a \in \text{span}\{q_i\}.$$

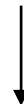
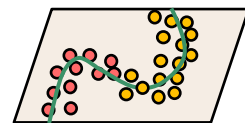
Project the equations using

$$\begin{aligned} q_i^T C u_a &= \lambda_a q_i^T u_a, \Rightarrow q_i^T \frac{1}{m} \sum_k q_k q_k^T \sum_j \alpha_j^a q_j = \lambda_a q_i^T \sum_j \alpha_j^a q_j \\ &\Rightarrow \frac{1}{m} \sum_{j,k} \alpha_j^a (q_i^T q_k) (q_k^T q_j) = \lambda_a \sum_j \alpha_j^a (q_i^T q_j). \end{aligned}$$

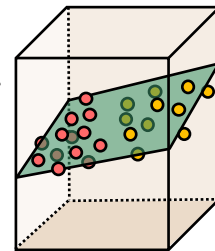
We now use $K(x_i, x_j) = \Phi^T(x_i) \Phi(x_j) = q_i^T q_j$, $\Rightarrow K^2 \alpha^a = m \lambda_a K \alpha^a$, $\Rightarrow K \alpha^a = \tilde{\lambda}_a \alpha^a$, $\tilde{\lambda}_a = m \lambda_a$.

Gives: $y_a = u_a^T q = \sum_i \alpha_i^a q_i^T q = \sum_i \alpha_i^a K(x_i, x)$.

Input Feature Space



New Feature Space



Kernel Principal Component Analysis (KPCA):

Task: Given data set $S = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$ with $\mathbf{x}_i \in \mathbb{R}^N$, Φ feature map of K . Find in feature space $\mathcal{H}_{\text{RKHS}}$ the best approximating hyperplane $h = \text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_p\}$ of $\tilde{S} = \{\Phi(\mathbf{x}_1), \Phi(\mathbf{x}_2), \dots, \Phi(\mathbf{x}_m)\}$.

KPCA Reformulation with Centering:

$$K^c \boldsymbol{\alpha}^a = \tilde{\lambda}_a \boldsymbol{\alpha}^a \text{ where } K^c(\mathbf{x}_i, \mathbf{x}) = K(\mathbf{x}_i, \mathbf{x}) - \boldsymbol{\kappa}^T(\mathbf{x}_i)\mathbf{1} - \boldsymbol{\kappa}^T(\mathbf{x})\mathbf{1} + \mathbf{1}^T \mathbf{k} \mathbf{1},$$

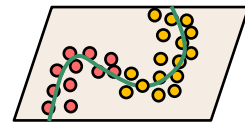
Notation: $[K^c]_{ij} = K^c(\mathbf{x}_i, \mathbf{x}_j)$ with $[\boldsymbol{\kappa}(\mathbf{x})]_l = \frac{1}{m} K(\mathbf{x}_l, \mathbf{x})$, $[\mathbf{k}]_{pl} = \frac{1}{m^2} K_{pl}$ and $[\mathbf{1}]_l = 1$.

Principal components: $\mathbf{u}_a = \sum_{i=1}^m \alpha_i^a \Phi(\mathbf{x}_i)$ with $\mathbf{u}_a^T \Phi(\mathbf{x}) = \sum_{i=1}^m \alpha_i^a K^c(\mathbf{x}_i, \mathbf{x})$,
 $K_{ij}^c = K_{ij} - \boldsymbol{\kappa}_i^T \mathbf{1} - \boldsymbol{\kappa}_j^T \mathbf{1} + \mathbf{1}^T \mathbf{k} \mathbf{1}$

Derivation: $K_{ij}^c = \left(\Phi(\mathbf{x}_i) - \frac{1}{m} \sum_{k=1}^m \Phi(\mathbf{x}_k) \right)^T \left(\Phi(\mathbf{x}_j) - \frac{1}{m} \sum_{l=1}^m \Phi(\mathbf{x}_l) \right)$

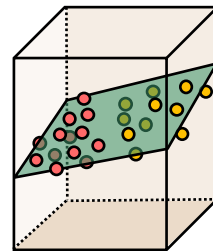
$$K_{ij}^c = \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_j) - \frac{1}{m} \sum_{k=1}^m \Phi(\mathbf{x}_k)^T \Phi(\mathbf{x}_j) - \frac{1}{m} \sum_{\ell=1}^m \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_\ell) + \frac{1}{m^2} \sum_{k,\ell} \Phi(\mathbf{x}_k)^T \Phi(\mathbf{x}_\ell).$$

Input Feature Space



$\Phi(\mathbf{X})$ ↓

New Feature Space



Kernel Principal Component Analysis (KPCA): Example

Task: Given data set $S = \{x_1, x_2, \dots, x_m\}$ with $x_i \in \mathbb{R}^N$, Φ feature map of K . Find in feature space $\mathcal{H}_{\text{RKHS}}$ the best approximating hyperplane $h = \text{span}\{v_1, v_2, \dots, v_p\}$ of $\tilde{S} = \{\Phi(x_1), \Phi(x_2), \dots, \Phi(x_m)\}$.

Example: Data-set $S = \{\text{two semi-circles centered around } (0,0) \text{ and } (1,0)\}$, $m = 10^3$ samples, noise = 0.075, kernel= Radial Basis Function, $\gamma = 15$.

Nearly one dimensional structure (non-linear).

Radial basis function kernel corresponds to **similarity by proximity**.

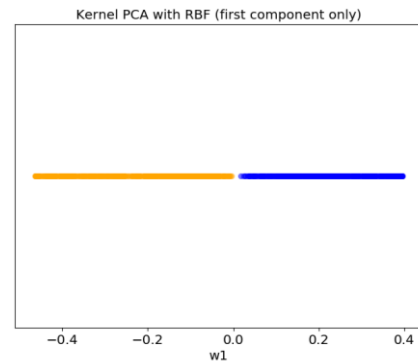
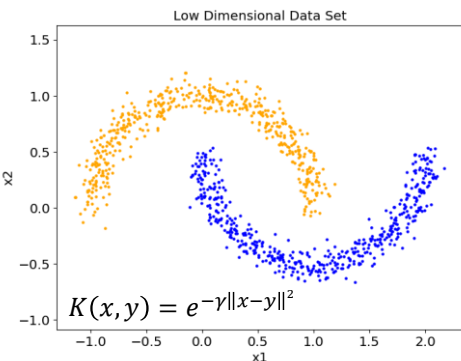
Results:

First component $w_1 := u_1$ captures enough information to **separate data set**.

Kernel methods can be sensitive to tuning of **hyperparameters**.

Cross-validation to get good hyperparameters but need metric (two common)

- if unsupervised learning can use reconstruction errors.
- if used for supervised learning (classifier) can use validation training errors.

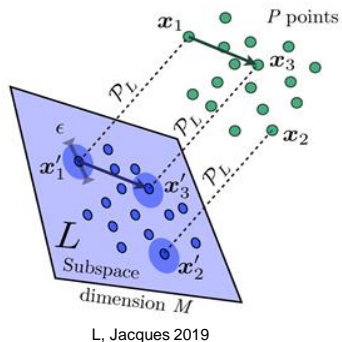




Random Projection



Random Projections



Random Projection

Consider the transform $x = W\tilde{x}$ with random $W \sim \mu(\mathbb{R}^{n \times d})$ fixed for all $\tilde{x} \in \mathbb{R}^d$.

We define the **distortion** r of a transformed vector x as

$$r(x) := \left| \frac{\|Wx\|}{\|x\|} - 1 \right|$$

Lemma: Distortion of a Gaussian Random Projection

For a fixed $x \in \mathbb{R}^d$, consider a random matrix $W \in \mathbb{R}^{n \times d}$ with each component $W_{ij} \sim \eta(0, 1)$ an independent Gaussian random variable. For every $\epsilon \in (0, 3)$, we have

$$\Pr \left\{ \left| \frac{\|n^{-1/2}Wx\|^2}{\|x\|^2} - 1 \right| > \epsilon \right\} \leq 2 \exp(-\epsilon^2 n/6).$$

Random Projections

Lemma: Distortion of a Gaussian Random Projection

For a fixed $\mathbf{x} \in \mathbb{R}^d$, consider a random matrix $W \in \mathbb{R}^{n \times d}$ with each component $W_{ij} \sim \eta(0, 1)$ an independent Gaussian random variable. For every $\epsilon \in (0, 3)$, we have

$$\Pr \left\{ \left| \frac{\|n^{-1/2} W \mathbf{x}\|^2}{\|\mathbf{x}\|^2} - 1 \right| > \epsilon \right\} \leq 2 \exp(-\epsilon^2 n/6).$$

Proof: Let $\mathbf{z} = \mathbf{x}/\|\mathbf{x}\|$, then we can express the probability as

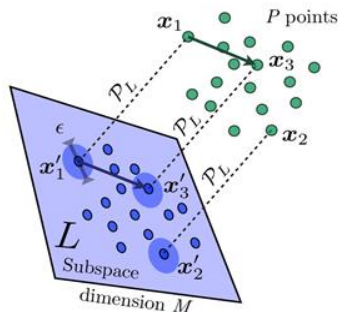
$$\Pr \left\{ (1 - \epsilon) n \leq \|W \mathbf{z}\|^2 \leq (1 + \epsilon) n \right\} \geq 1 - 2 \exp(-\epsilon^2 n/6).$$

For $w_i = W_{i,*}$ the i^{th} row of W , the product $\langle w_i, \mathbf{z} \rangle$ is a Gaussian random variable with mean 0 and variance $\sum_j z_j^2 = 1$. As a consequence, $\|W \mathbf{z}\|^2 = \sum_{i=1}^n (\langle w_i, \mathbf{z} \rangle)^2$ is a sum of squared Gaussians.

This has the well-known χ_n^2 distribution, which satisfies $\Pr \{ \chi_n^2 \leq (1 - \epsilon)n \} \leq \exp(-\epsilon^2 n/6)$ and $\Pr \{ \chi_n^2 \geq (1 + \epsilon)n \} \leq \exp(-\epsilon^2 n/6)$.

This yields the result. ■

Random Projections



L. Jacques 2019

Lemma: Johnson-Lindenstrauss

Consider m data points $x_i \in \mathbb{R}^d$ mapped to $y_i = Tx_i \in \mathbb{R}^n$. For any $0 < \epsilon < 1/2$, $m > 4$, and $n = 20 \log(m)/\epsilon^2$, consider the random projection $T : \mathbb{R}^d \rightarrow \mathbb{R}^n$ given by $T = n^{-1/2}W$ with $n < d$ and W having components $W_{ij} \sim \eta(0, 1)$ i.i.d. Gaussians with mean 0 and variance 1. For any two data points x_i and x_j we have

$$(1 - \epsilon) \|x_i - x_j\|^2 \leq \|Tx_i - Tx_j\|^2 \leq (1 + \epsilon) \|x_i - x_j\|^2$$

holds with probability $1 - \delta$ where $\delta = 2 \exp(-n\epsilon^2/40)$.

This shows **low distortion** T can be achieved provided the dimension $n \sim \log(m)/\epsilon^2$.

Random projections allow for **reducing the dimension** from d to n while **nearly preserving the pair-wise distances** between points.



Manifold Learning

Manifold Learning

Task: Given data set $S = \{x_1, x_2, \dots, x_m\}$ with $x_i \in \mathbb{R}^N$, Find a p -dimensional manifold that approximates the data set.

Kernel PCA is one form of manifold learning.

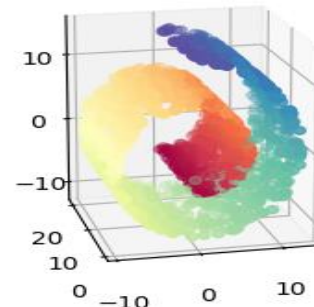
Many other (related) ways to try to approximate data by manifold.

- **Locally Linear Embedding (LLE)**
- **ISOMAP**
- **Spectral Embedding (SE)**
- **t-distributed Stochastic Neighbor Embedding (t-SNE)**

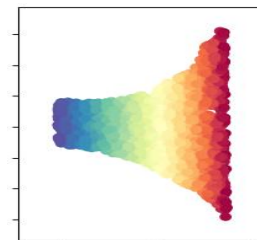
Can be viewed in some cases as a form of KPCA (i.e. SE, ISOMAP).

Treats data as having relevant features close to some smooth low-dimensional manifold for inference, data exploration, visualization.

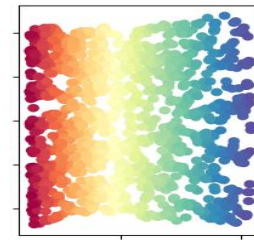
Swiss Roll Dataset



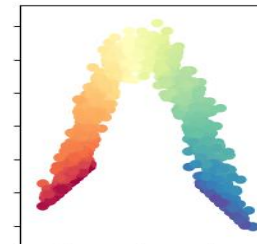
LLE



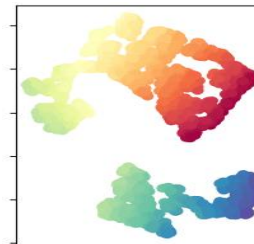
ISOMAP



SE



t-SNE



Manifold Learning: t-SNE

Method:

Consider for all distinct pairs of points \mathbf{x}_i and \mathbf{x}_j the Gaussian probability distribution

$$p_{ij} = \frac{\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / 2\sigma^2)}{\sum_{k \neq l} \exp(-\|\mathbf{x}_k - \mathbf{x}_l\|^2 / 2\sigma^2)}$$

Want mapping $x \in \mathbb{R}^N \rightarrow y \in \mathbb{R}^p$ into a smaller dimensional $p \ll N$.

Consider probability t-distribution for neighbors under mapping

$$q_{ij} = \frac{(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}}{\sum_{k \neq l} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1}}$$

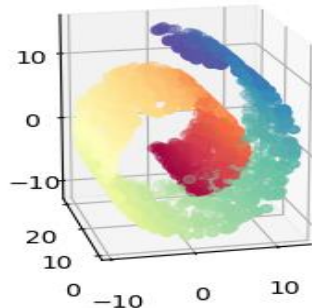
Optimization Problem:

$$\min_{\mathbf{y}} KL(P||Q) = \sum_{i,j} p_{ij} \log\left(\frac{p_{ij}}{q_{ij}}\right)$$

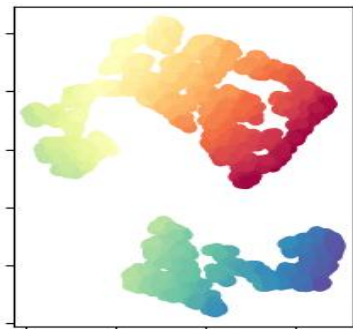
Finds points \mathbf{y}_i minimizing Kullback–Leibler (KL) Divergence of pair distributions.

t-SNE used primarily for visualization (good at showing crowding / clustering).

Swiss Roll Dataset



t-SNE



Manifold Learning: Locally Linear Embedding (LLE)

Method:

Want mapping $x \in \mathbb{R}^N \rightarrow y \in \mathbb{R}^p$ into a smaller dimensional $p \ll N$.

Find for a point x_i the k -nearest neighbors x_j .

Optimization Problem (step I):

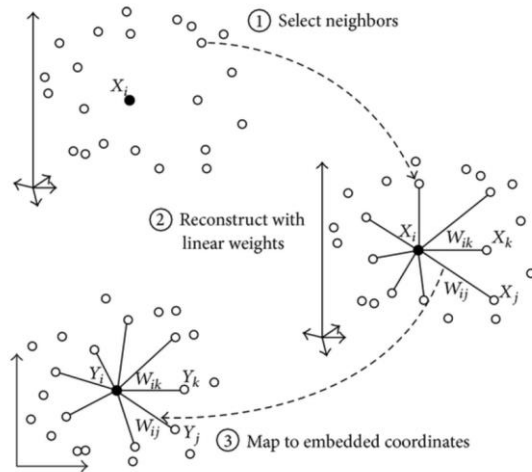
$$\min_W \sum_i \left\| x_i - \sum_{j \in \mathcal{N}(x_i)} W_{ij} x_j \right\|_2^2$$

subject $W \cdot \mathbf{1} = \mathbf{1}$

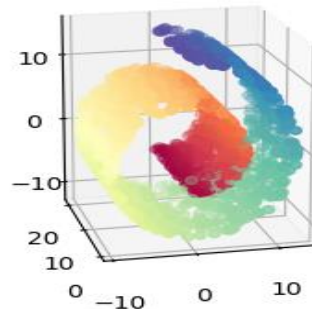
Optimization Problem (step II):

$$\min_y \sum_i \left\| y_i - \sum_{j \in \mathcal{N}(x_i)} W_{ij} y_j \right\|_2^2$$

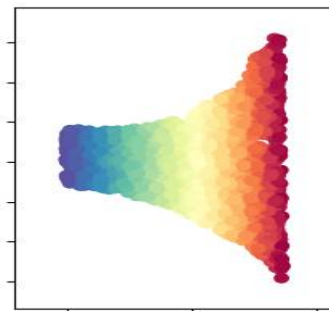
Finds points y_i so that embedding has similar linear weighting to reconstruct y_i from nearby points y_j .



Swiss Roll Dataset



LLE



Manifold Learning: Isomap

Method:

Consider for \mathbf{x}_i the neighborhood of points \mathbf{x}_j within distance ε .

Form graph \mathcal{G} with edges \mathbf{x}_j in ε -neighborhood of \mathbf{x}_i .

Compute distance matrix Δ_{ij} using shortest path in graph between points $\mathbf{x}_i \rightarrow \mathbf{x}_j$. The $\Delta_{ij} \sim$ squared geodesic distance.

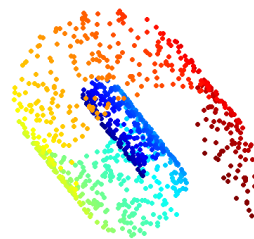
Want mapping $x \in \mathbb{R}^N \rightarrow y \in \mathbb{R}^p$ into a smaller dimensional $p \ll N$.

Optimization Problem:

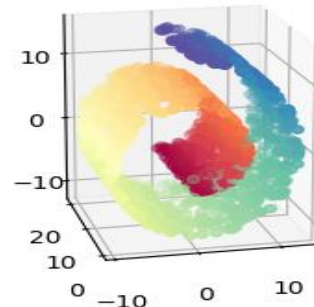
$$\min_{\mathbf{y}} \sum_{i,j} \left(\|\mathbf{y}_i - \mathbf{y}_j\|_2^2 - \Delta_{ij} \right)^2$$

Finds points \mathbf{y}_i so that embedding distance between pairs is close to Δ_{ij} .

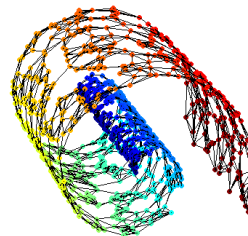
Related to KCPA with kernel $\mathbf{K}_{\text{Iso}} = -\frac{1}{2}\mathbf{H}\Delta\mathbf{H}$ and $\mathbf{H} = \mathbf{I}_m - \frac{1}{m}\mathbf{1}\mathbf{1}^\top$.



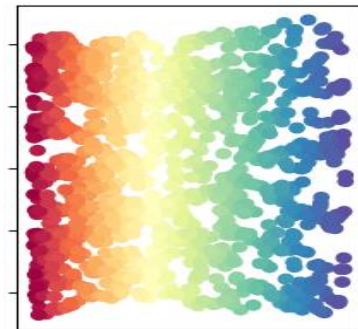
Swiss Roll Dataset



Construct Graph of
k-Nearest Neighbors



ISOMAP



Manifold Learning: Spectral Embedding

Method:

Consider for \mathbf{x}_i the neighborhood of points \mathbf{x}_j within distance ε .

Form graph \mathcal{G} with edges \mathbf{x}_j in ε -neighborhood of \mathbf{x}_i .

Compute weight matrix $W_{ij} = e^{-\|\mathbf{x}_i - \mathbf{x}_j\|_2^2 / 2\sigma^2}$

Want mapping $x \in \mathbb{R}^N \rightarrow y \in \mathbb{R}^p$ into a smaller dimensional $p \ll N$.

Optimization Problem:

$$\min_{\mathbf{y}} \sum_{i,j} W_{ij} \|\mathbf{y}_i - \mathbf{y}_j\|_2^2$$

subject $\mathbf{y}^T D \mathbf{y} = 1$ with $D = \text{diag}(W \cdot \mathbf{1})$.

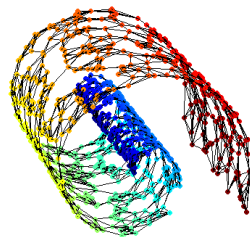
Finds points \mathbf{y}_i so that embedding distance between pairs minimized. Note the penalty W_{ij} is stronger for closer pairs of points and weaker for further pairs of points.

Related to KCPA with kernel $K_L = L^\dagger$ where graph Laplacian

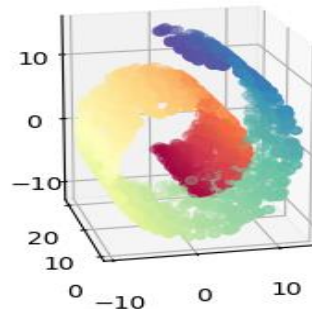
$L = D - W$ and $D = \text{diag}(W \cdot \mathbf{1})$. Solution: $\mathbf{Y} = U_{L,p}$ from SVD of $L = U\Lambda U^T$.

Approximately preserves diffusion commute times on manifold from \mathbf{x}_i to \mathbf{x}_j .

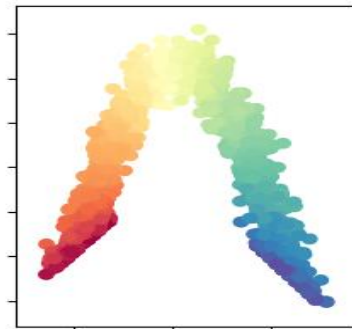
Construct Graph of ε -Neighborhoods




Swiss Roll Dataset



SE






Examples

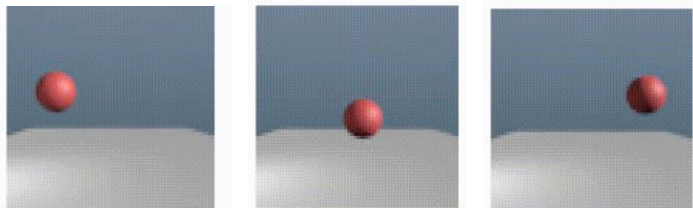
Unsupervised Learning

Manifold Learning / Clustering



Pendulum Dynamics from Video: Spectral Embedding

Input: Image frames from video.

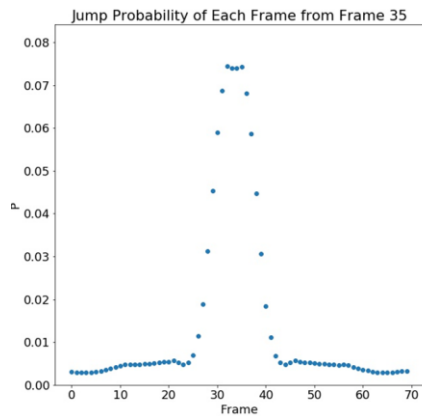


Spectral Embedding:

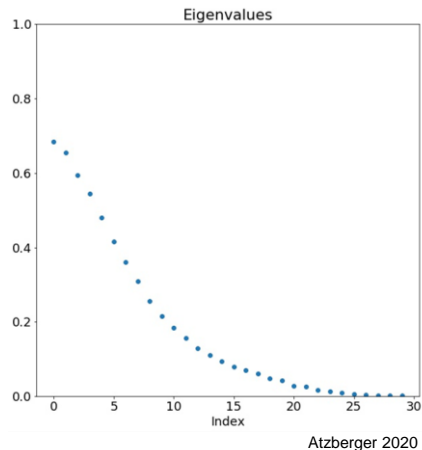
$$\min_{\mathbf{y}} \sum_{i,j} W_{ij} \|\mathbf{y}_i - \mathbf{y}_j\|_2^2$$

subject $\mathbf{y}^T D \mathbf{y} = 1$ with $D = \text{diag}(W \cdot \mathbf{1})$.

**Kernel for frames
(similarity)**

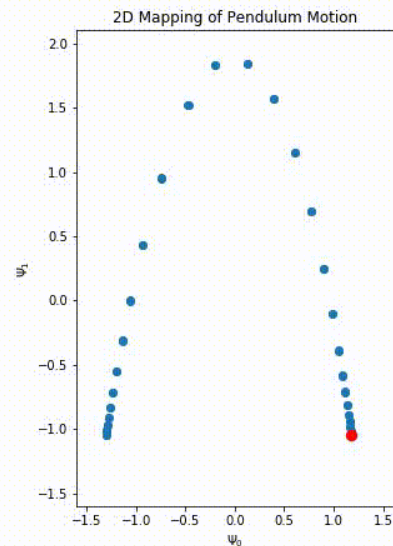


**Eigenvalues
Diffusion Process**

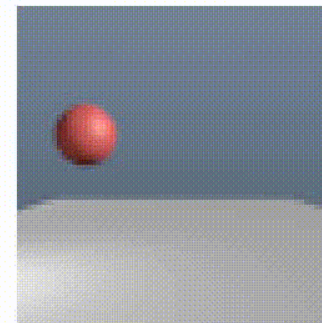


Atzberger 2020

Results



Frame 0



Atzberger 2020

Manifold Learning and Clustering

MNIST Digits Database:

Widely used benchmark for classification methods.

Task: Classify image to digit 0 - 9.

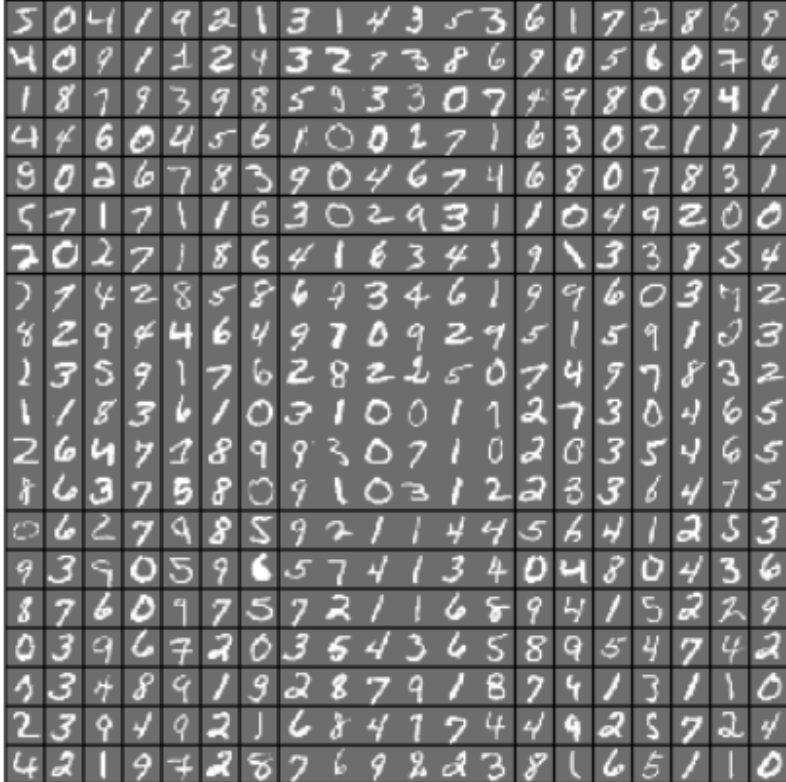


Data: train + test ~ 80K+ images 28x28.

Considered as points in 784-dimensional space.

Dimension reduction methods used to visualize, better understand structure, perform tasks.

t-SNE vs Local Linear Embedding (LLE).



Manifold Learning and Clustering: LLE vs t-SNE

MNIST Digits Database:

Widely used benchmark for classification methods.

Task: Classify image to digit 0 - 9.

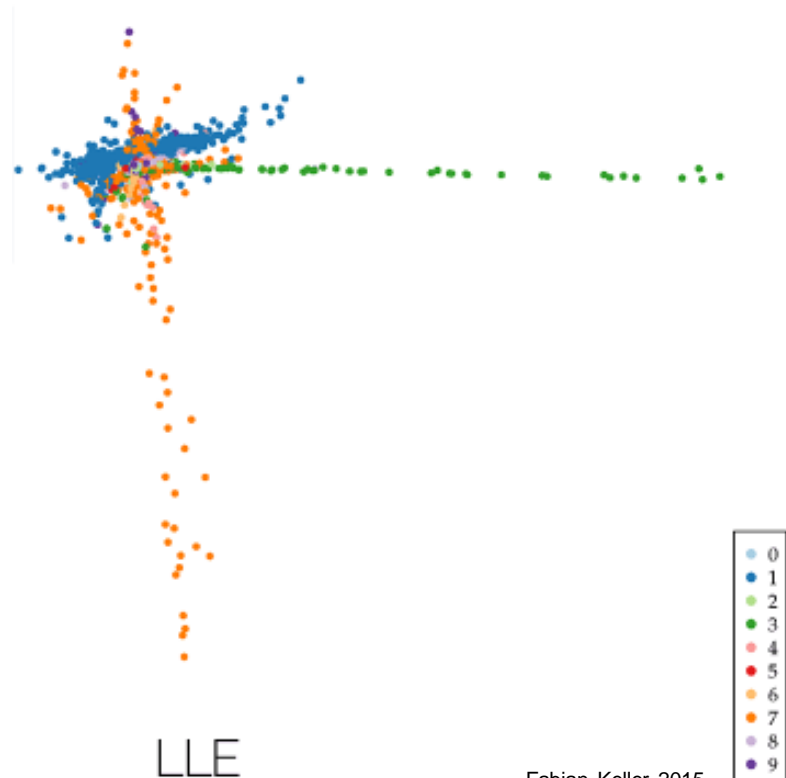


Data: train + test ~ 80K+ images 28x28.

Considered as points in 784-dimensional space.

Dimension reduction methods used to visualize, better understand structure, perform tasks.

t-SNE vs Local Linear Embedding (LLE).



Fabian Keller 2015

Manifold Learning and Clustering: LLE vs t-SNE

MNIST Digits Database:

Widely used benchmark for classification methods.

Task: Classify image to digit 0 - 9.

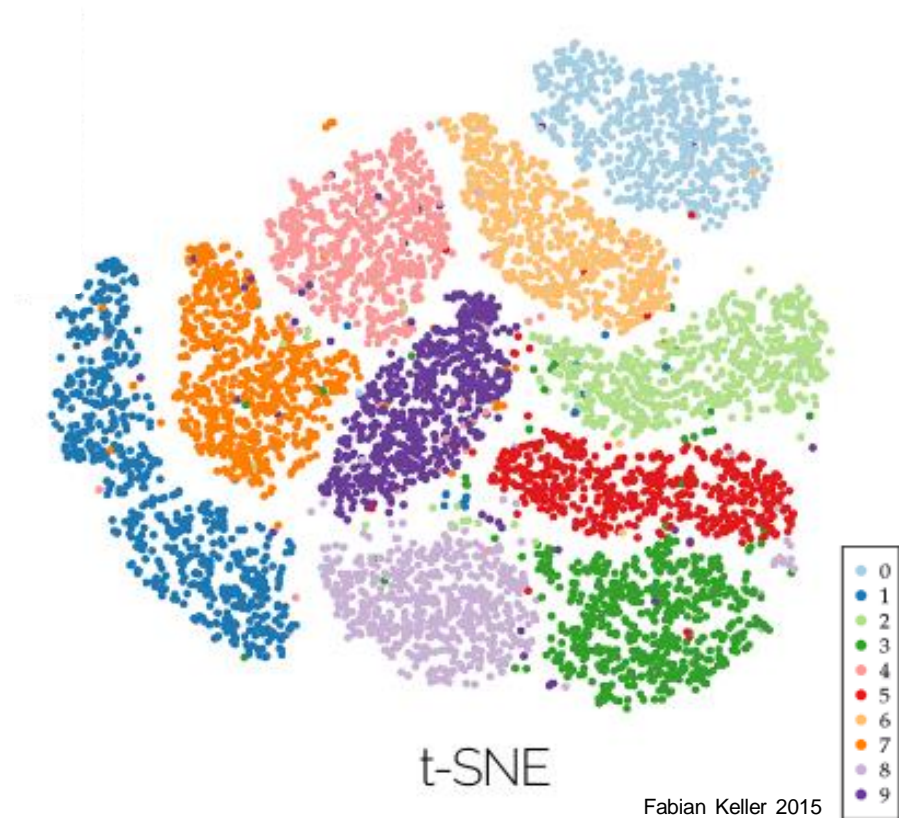


Data: train + test ~ 80K+ images 28x28.

Considered as points in 784-dimensional space.

Dimension reduction methods used to visualize, better understand structure, perform tasks.

t-SNE vs Local Linear Embedding (LLE).



Manifold Learning: LLE

Face Database:

Local Linear Embedding (LLE).

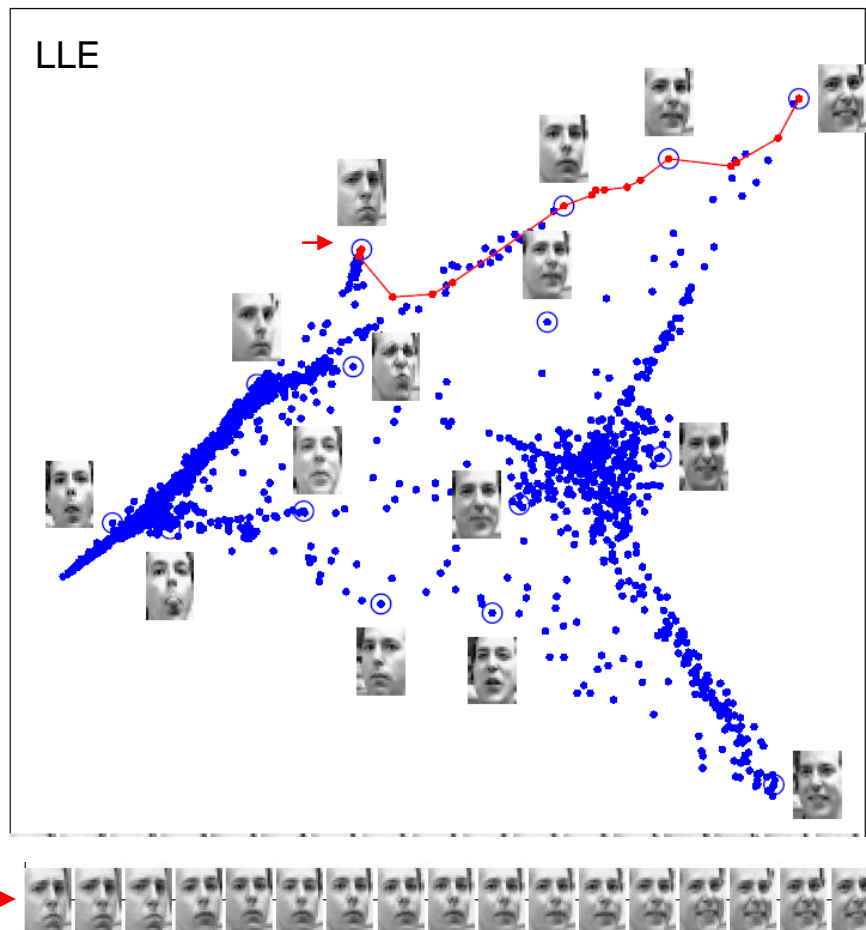
Task: Organize data set of faces by “similarity”

Data: 1965 images at 20×28 .

Considered as points in 560-dimensional space.

Dimension reduction method useful
to visualize, better understand structure, classify,
perform tasks.

Red path shows progression through
embedding space.



Hastie 2017, Chen and Buja 2008.



Summary



Summary: Dimension Reduction

Motivations for Unsupervised Learning:

Abstractly trying to learn characteristics of $\mathcal{D} \sim \mathcal{X}$.

Find Structure and patterns in data $S = \{x_1, x_2, \dots, x_m\}$.

A few features are often extremal in determining key properties.

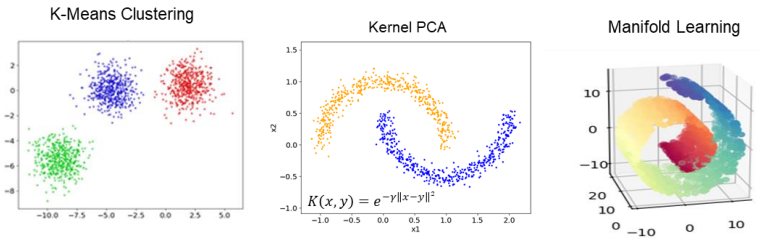
Transforms input data allowing for incorporating **prior knowledge**.

Filtering to capture essential aspects of data, improves generalization.

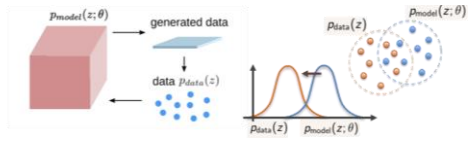
Methods (depends on the task):

- **Clustering Methods (K-means, Spectral Graph)**
- **Principal Component Analysis (PCA / KPCA)**
- **Manifold Learning (Isomap, LLE, Spectral)**
- **Generative Adversarial Networks (GANs).**
- (many more methods)

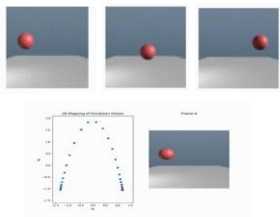
Provides ways to learn structure from data, incorporate prior knowledge, improve generalization, improve computational efficiency.



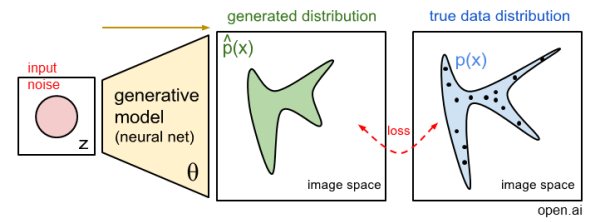
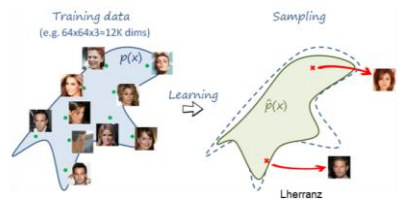
Generative Models



Spectral Embedding



GANs



open.ai

