

The background of the slide is a dense field of 3D-rendered numbers in various shades of blue. The numbers are scattered across the frame, creating a sense of depth and data. Some numbers are larger and more prominent than others, while many are smaller and recede into the background. The overall effect is a vibrant, digital landscape of numerical information.

Big Data Course

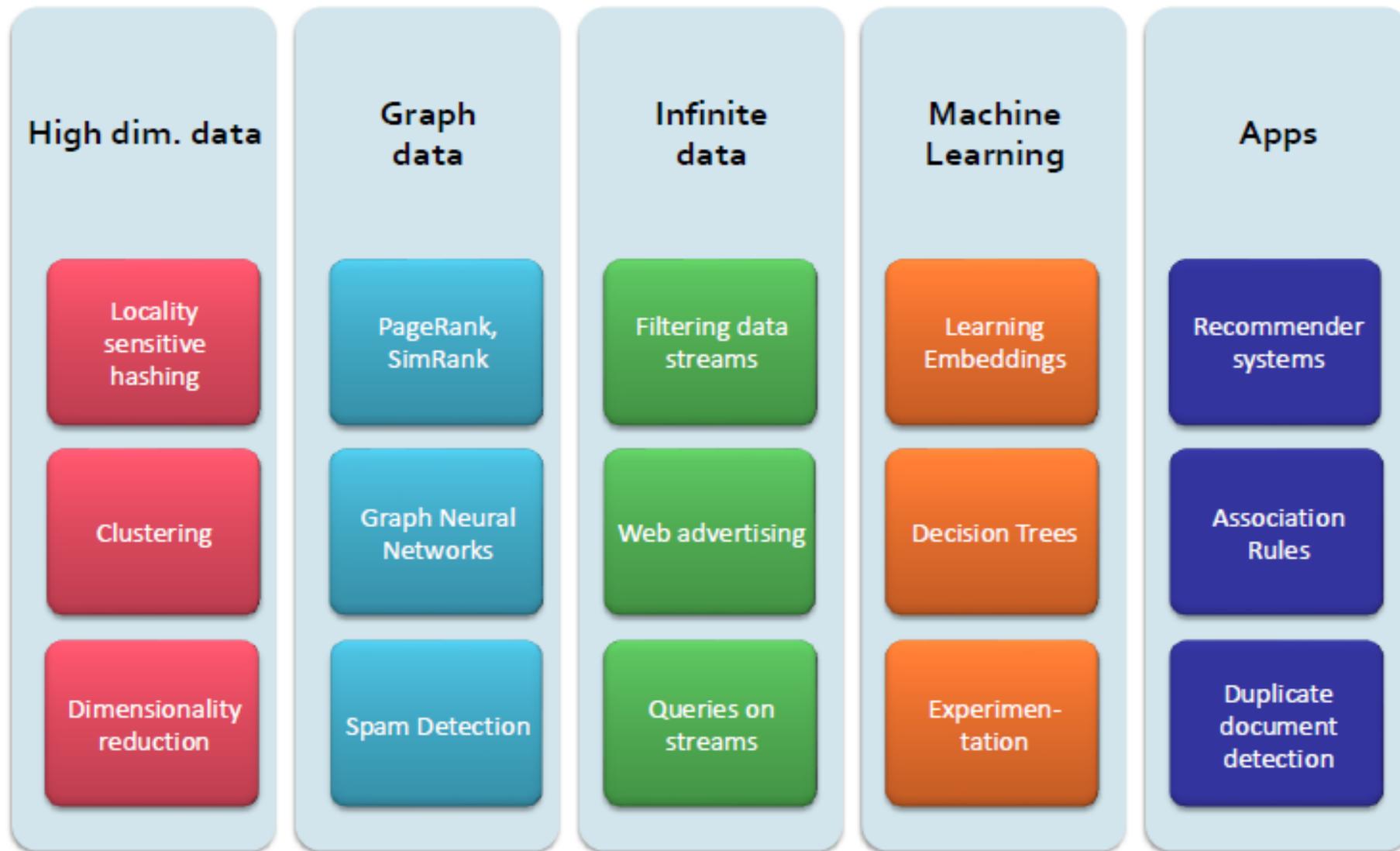
Chap 2- Dimension Reduction

Electrical Engineering department of
Amirkabir University of technology

Dr. Mohammadreza Pourfard

August 2025

How the Class Fits Together



تحم بالایی داده مادر حوزه پزشکی

و لزوم کاهش بعد

DNA

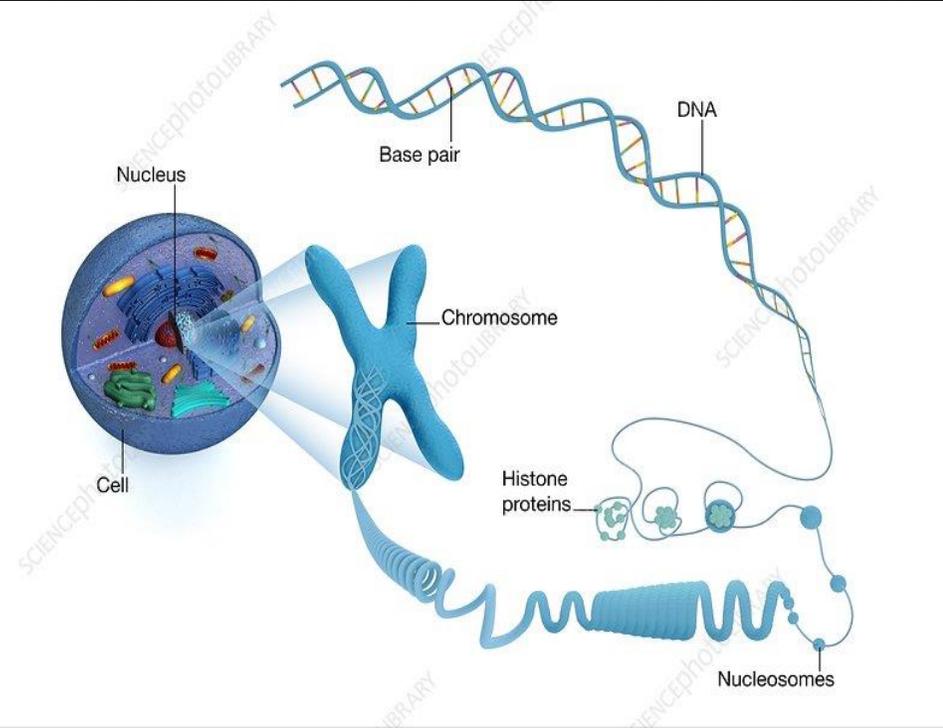
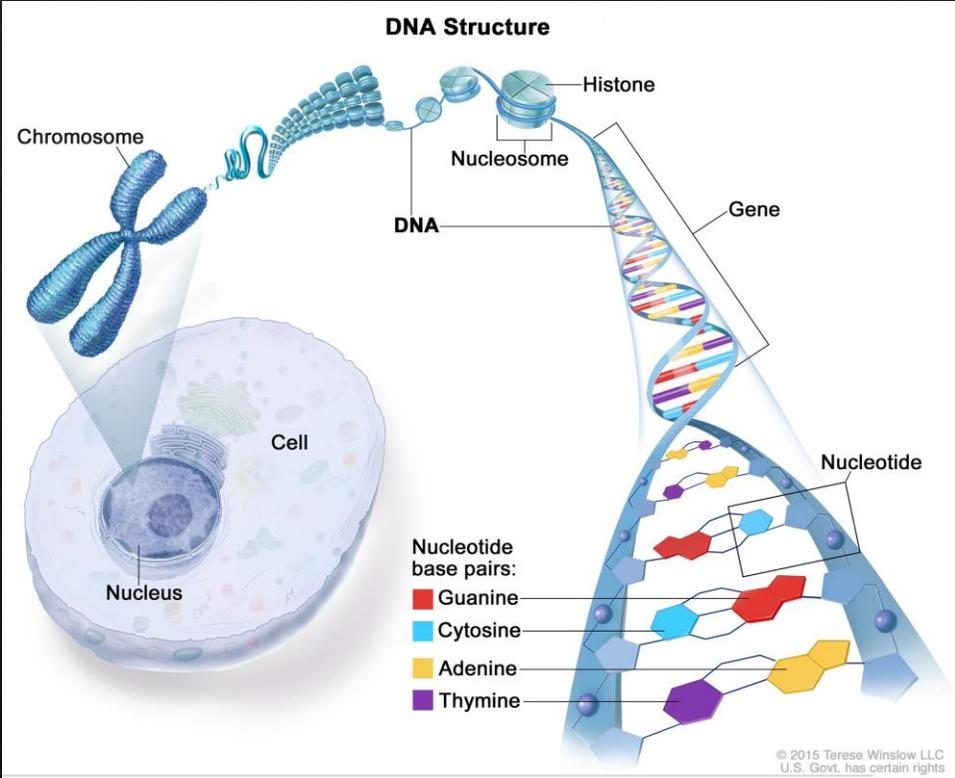
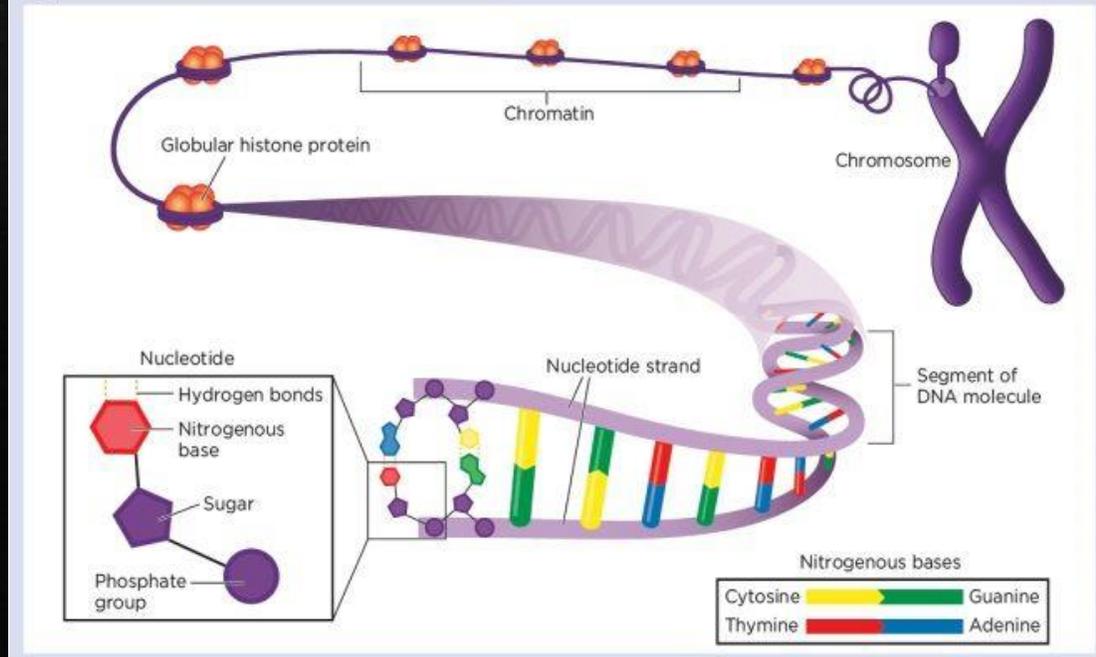
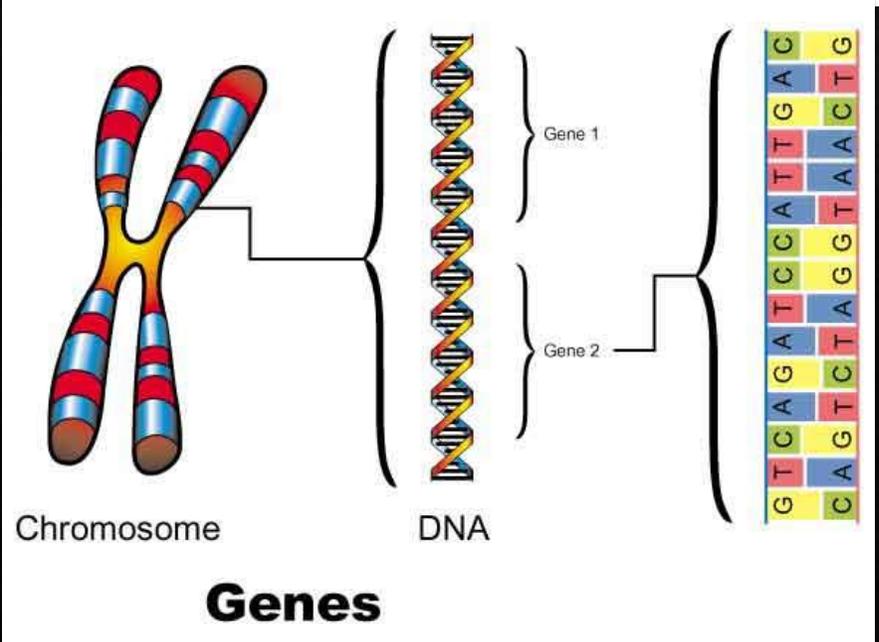
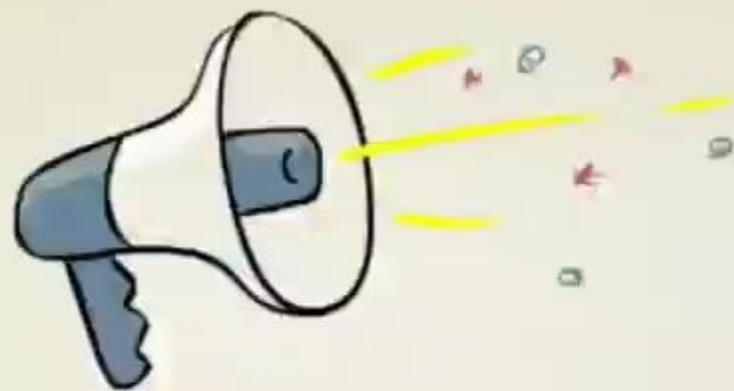


Fig 1. Structure of DNA and chromosomes

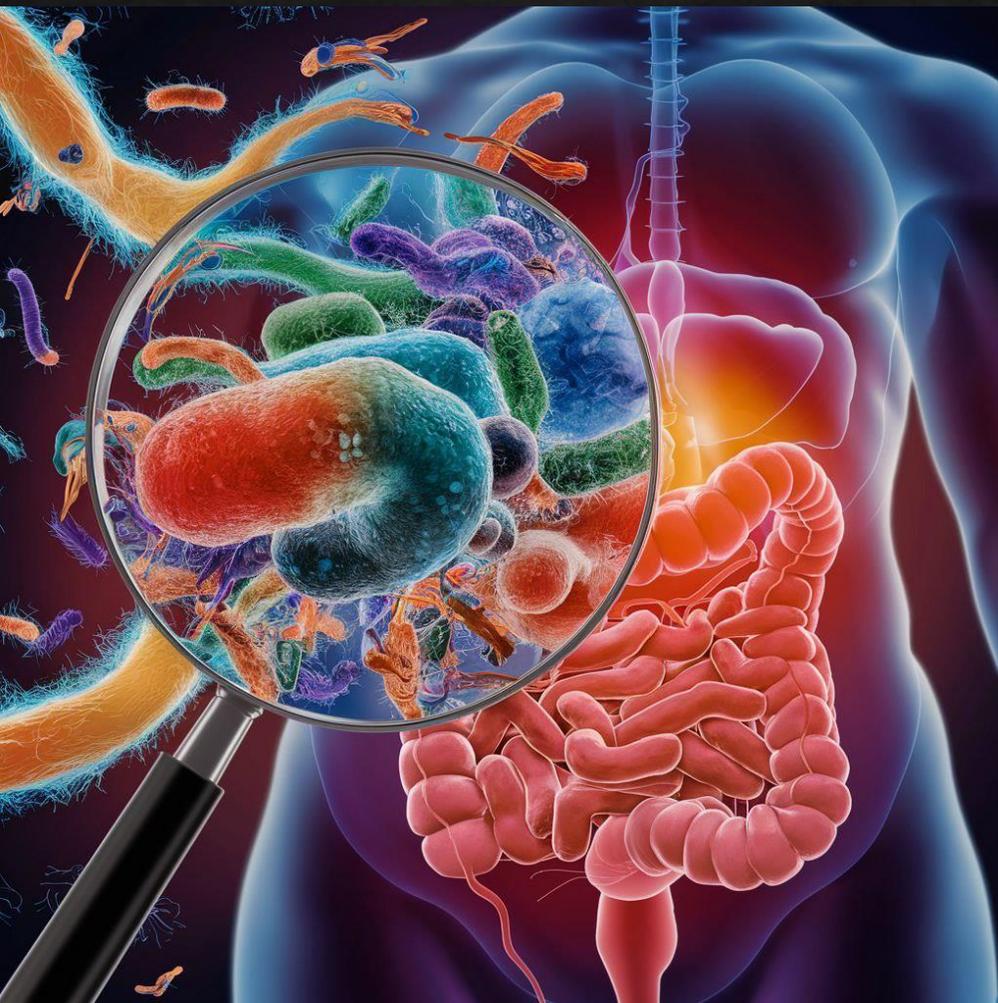


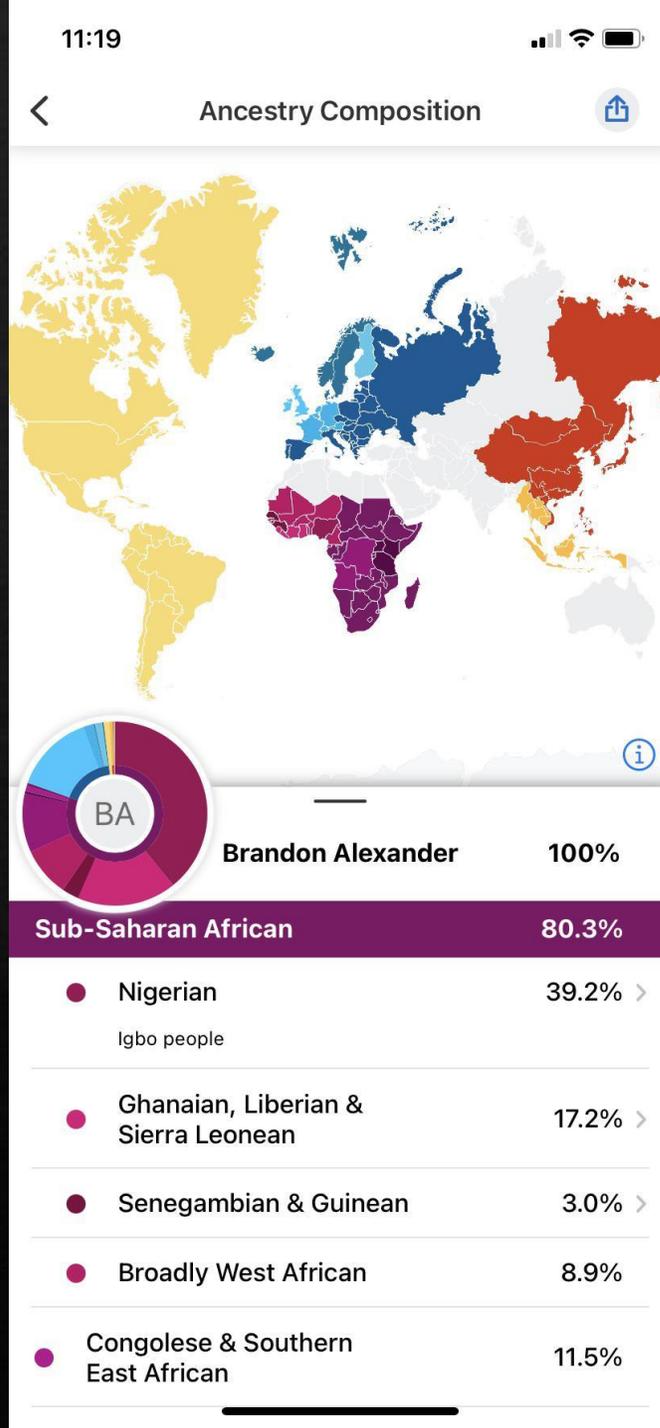
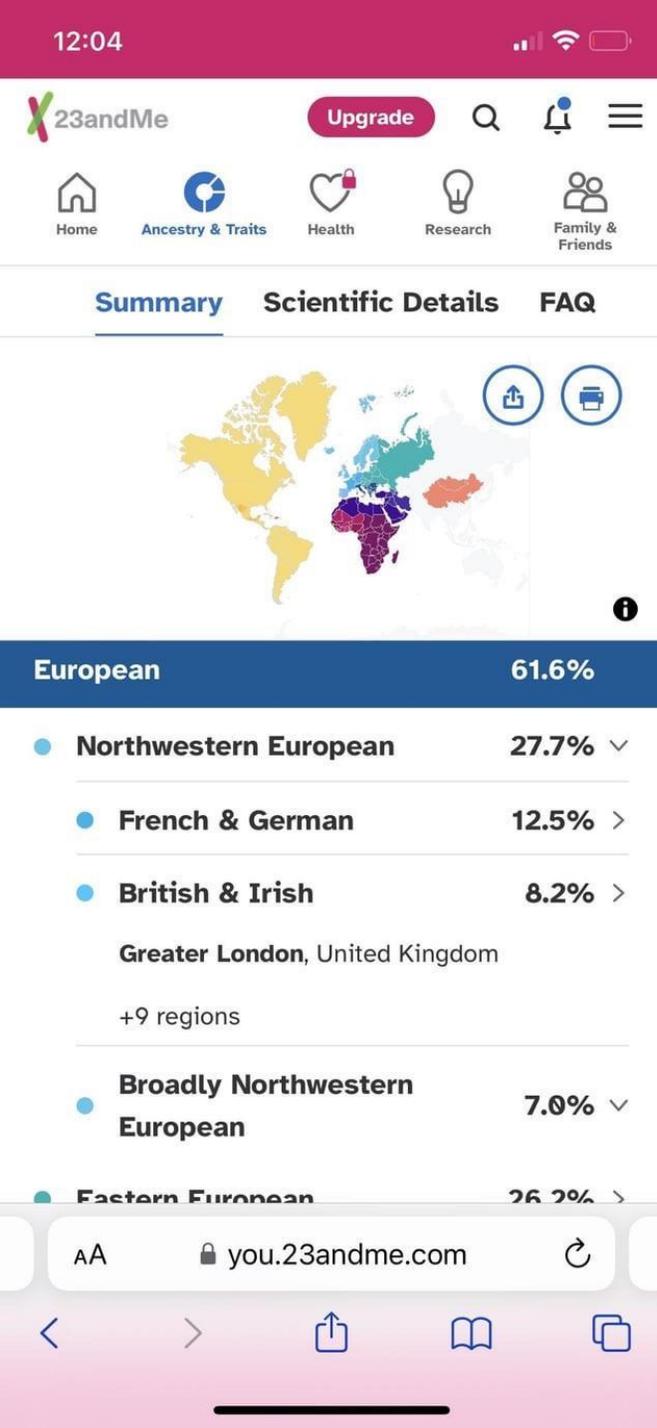




Stated Clearly

حوزه Gut Microbiota بر اساس داده های ژنتیک



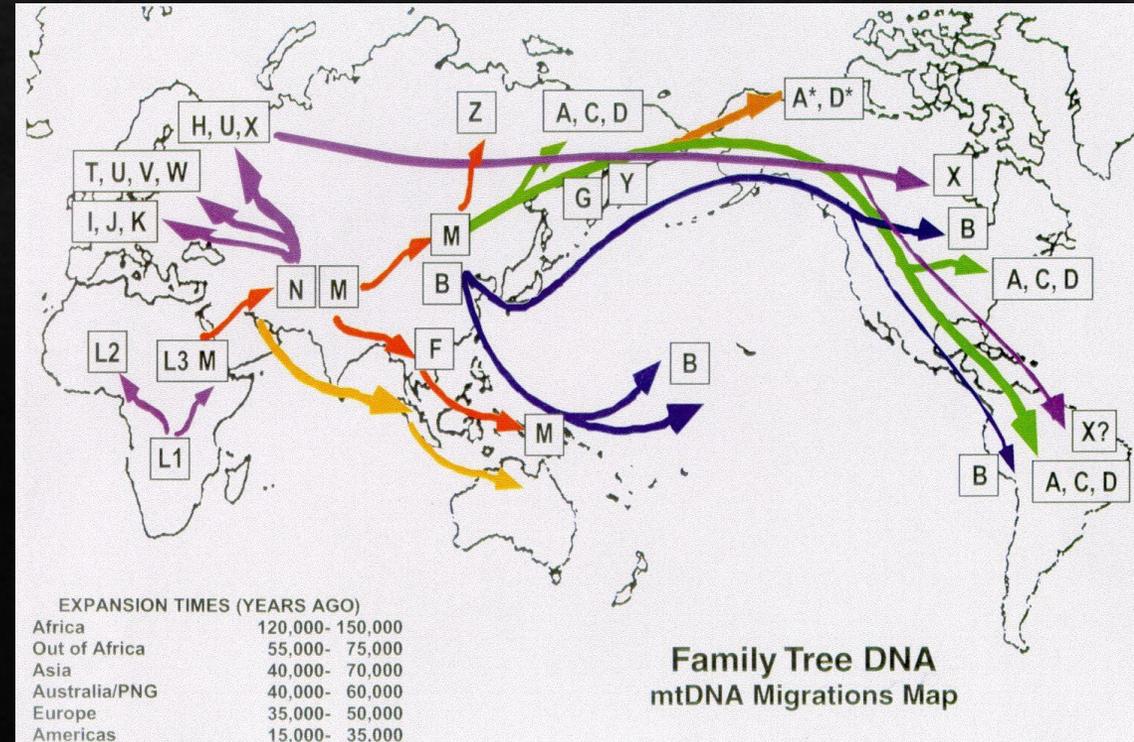
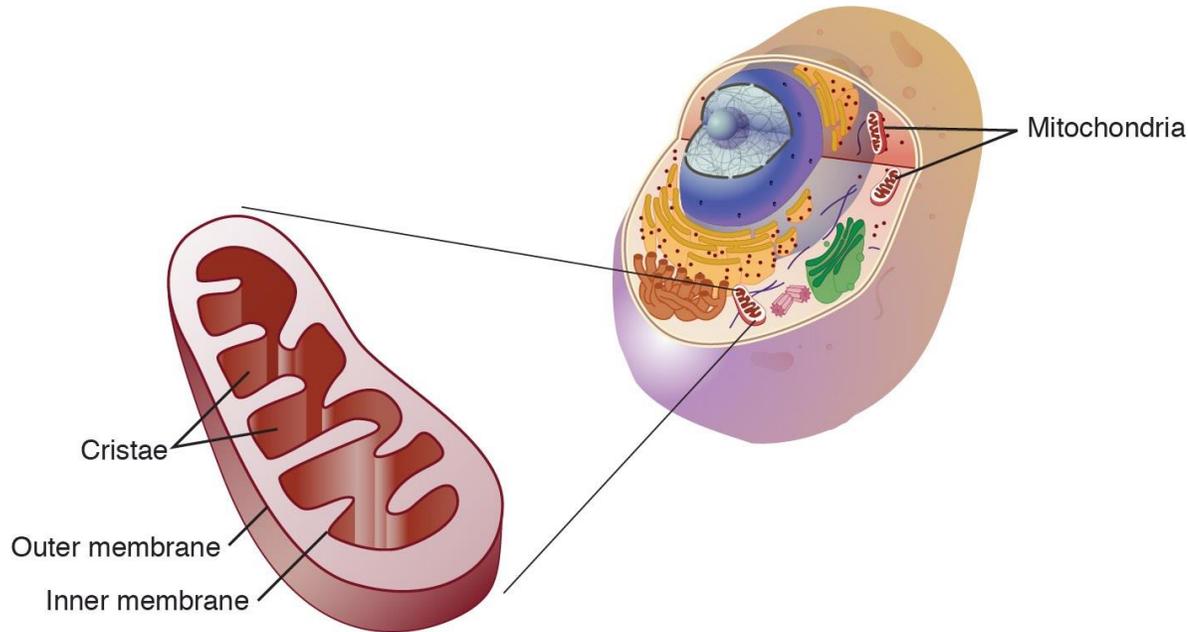


تحليل جمعيت اجدادی با داده ژنتیک



تشخیص مسیر مهاجرت مادری

◊ تعیین مسیر مهاجرت اجداد به کمک داده ژنتیک میتوکندری



Analysis

Integration of multi-omics data

Phenomics

Metabolomics

Proteomics

Transcriptomics

Epigenomics

Genomics

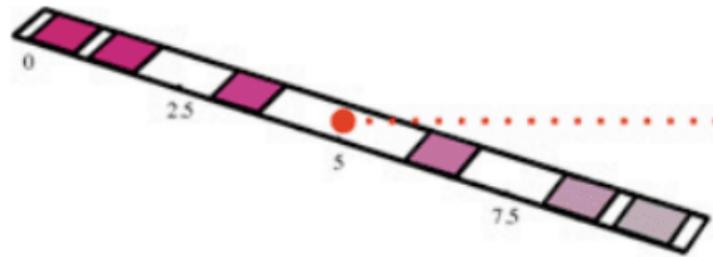
Age
/environment

Organization
(organelle, cell, organ, organism)

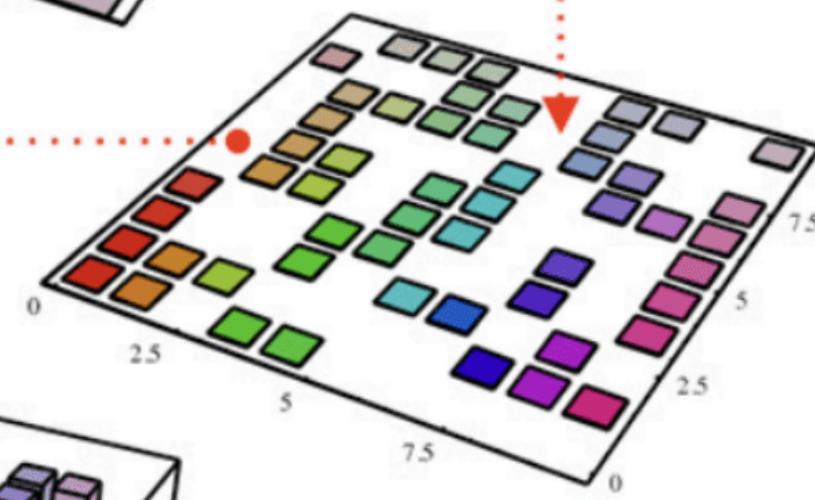
سطوح تشخیص
بیماری (شروع
از پیچیده ترین
سطح)

ضرورت کاهش بعد

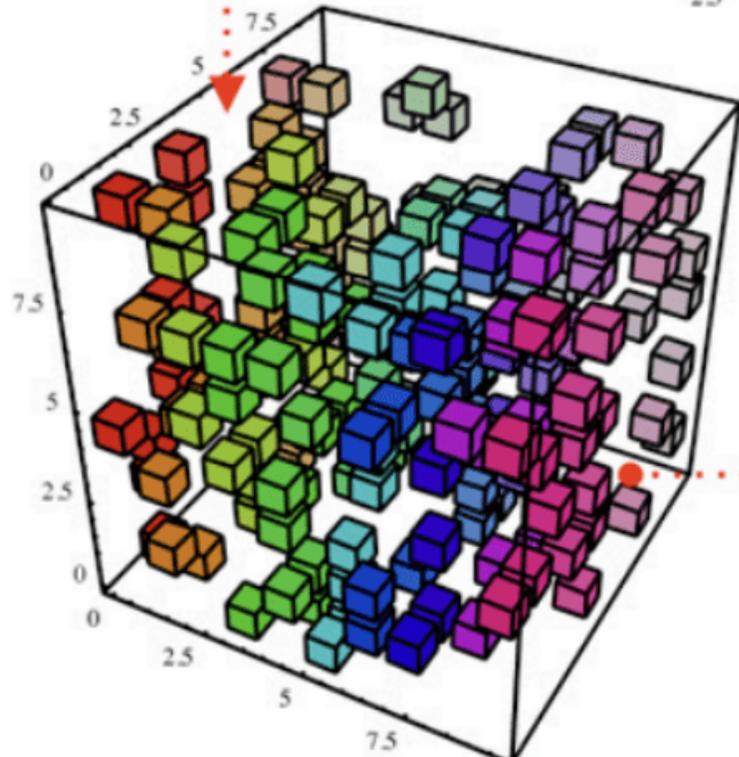
Dimension reduction



1 dimension:
10 positions

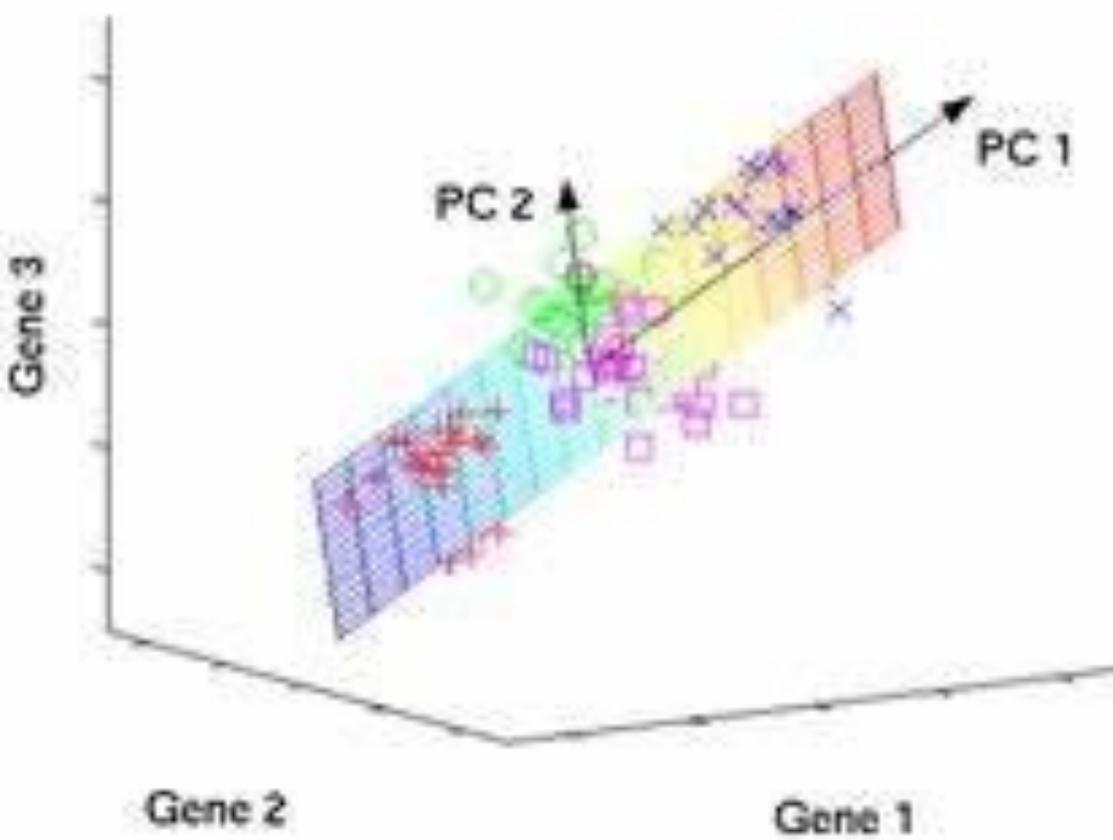


2 dimensions:
100 positions



3 dimensions:
1000 positions!

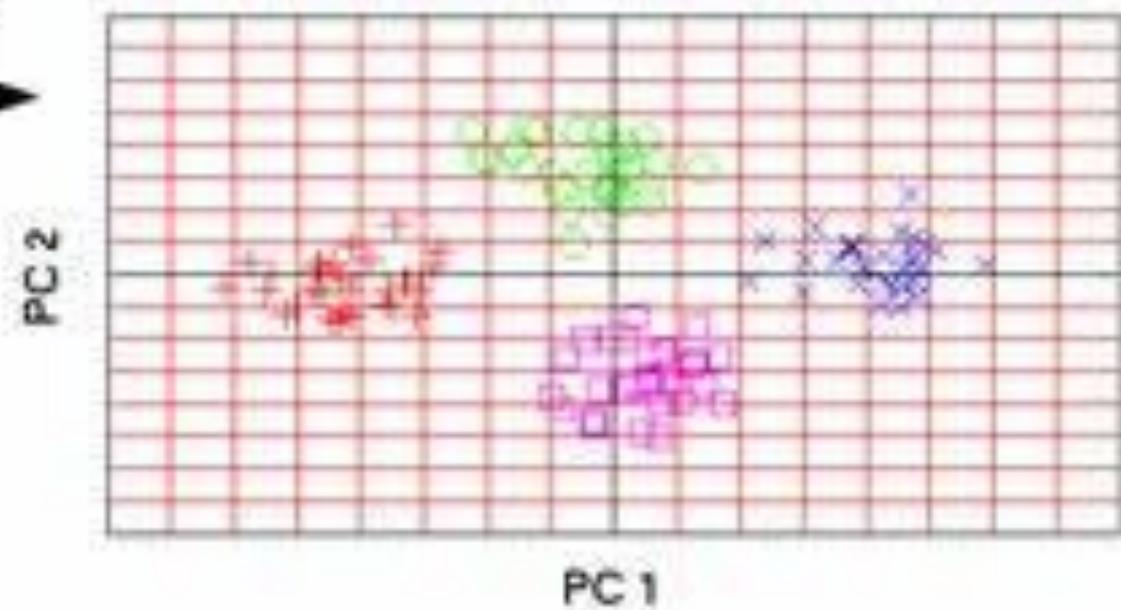
original data space



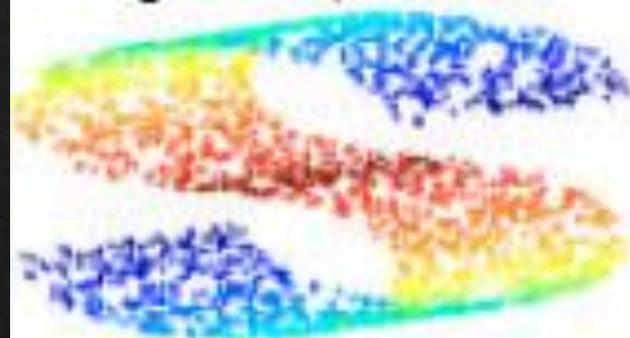
PCA



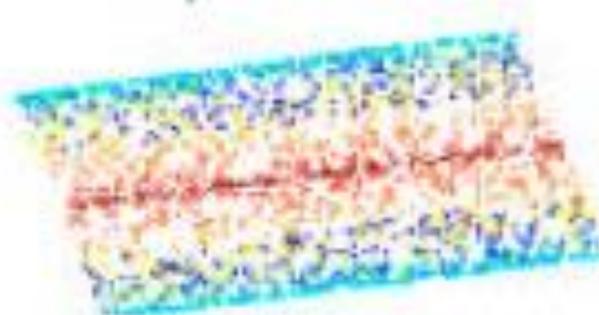
component space



Original data, $N = 2000$



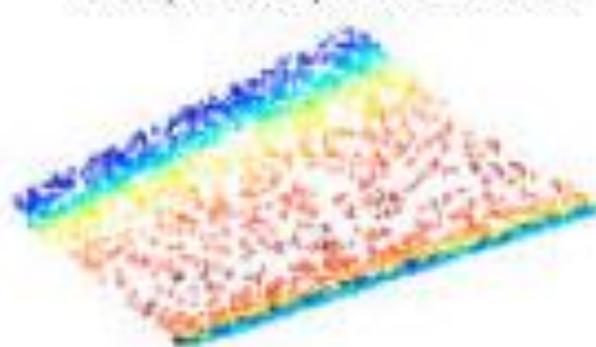
MDS, $\text{corr} = 0.7911$



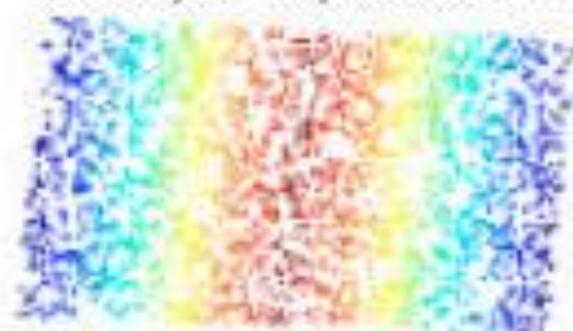
Isomap, $k = 12$, $\text{corr} = 0.9995$



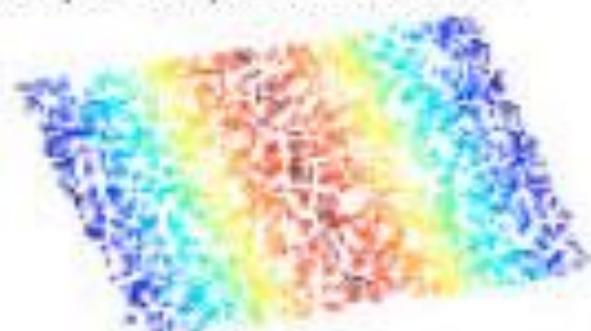
LLE, $k = 12$, $\text{corr} = 0.8261$



HLLE, $k = 12$, $\text{corr} = 0.9069$



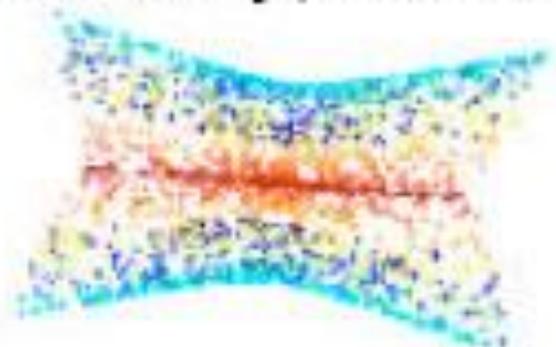
LTSA, $k = 12$, $\text{corr} = 0.9068$



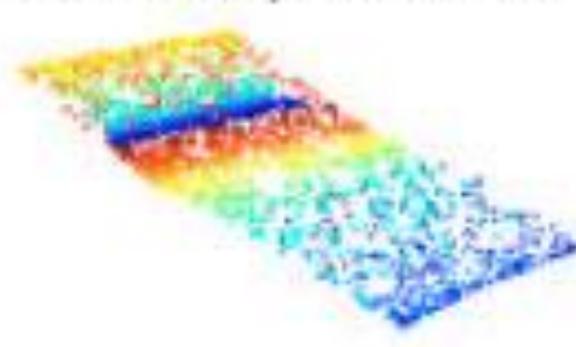
Kernel PCA, poly, $\text{corr} = 0.4487$



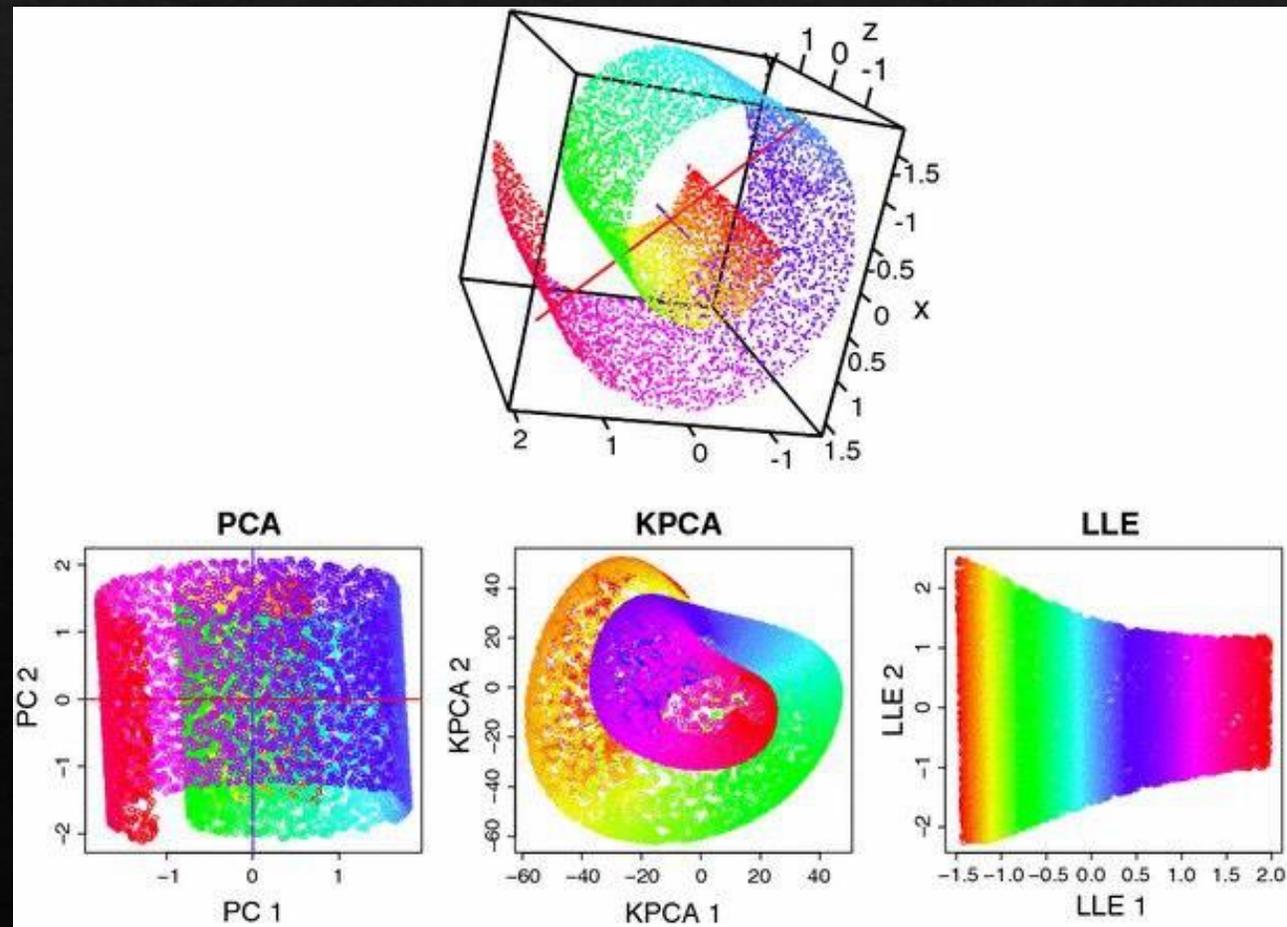
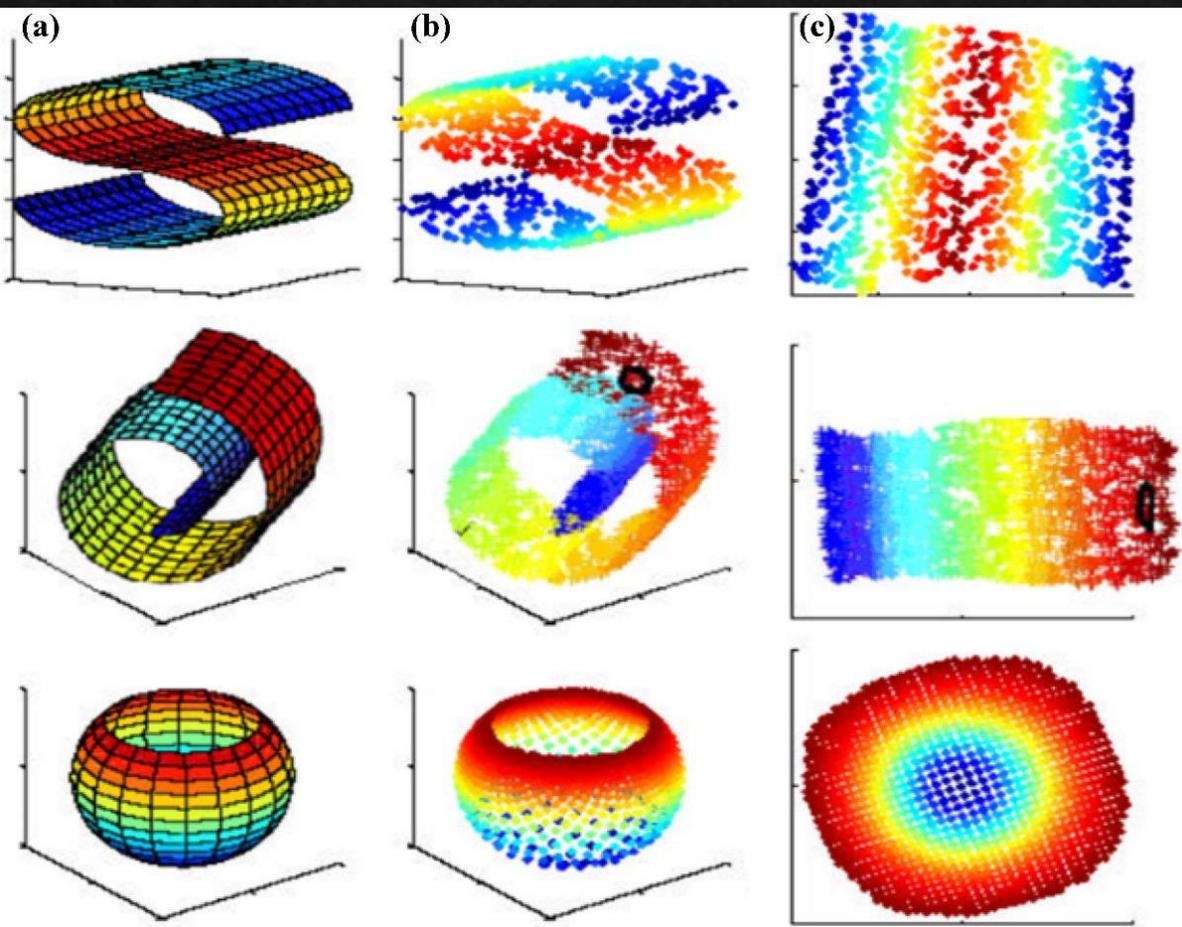
Diffusion maps, $\text{corr} = 0.6995$

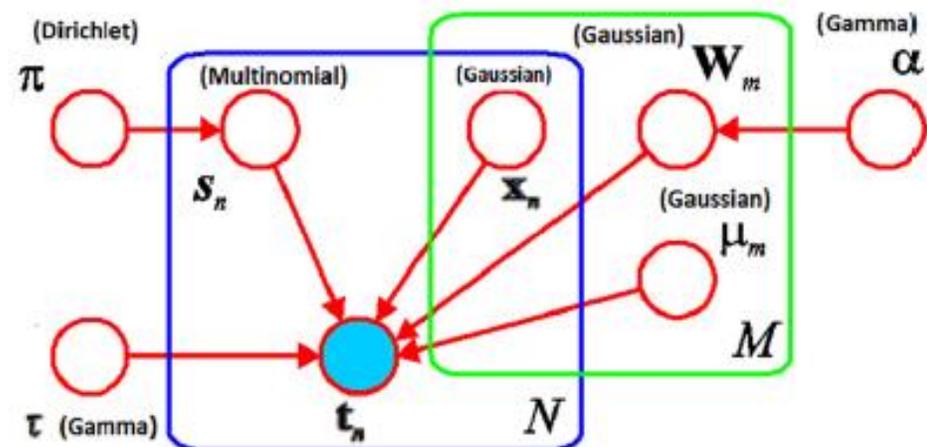
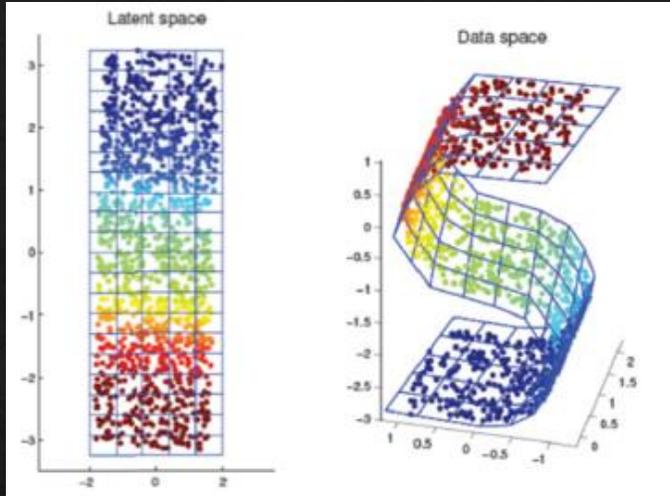
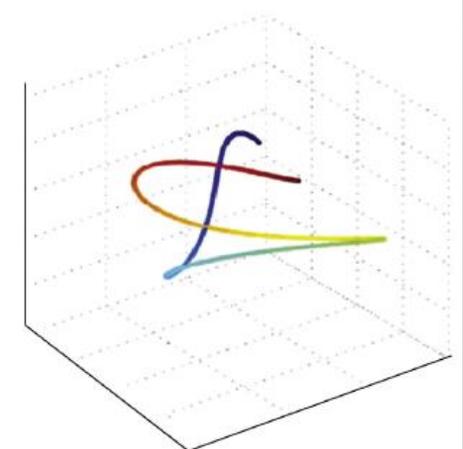


Auto encoder RBM, $\text{corr} = 0.6441$

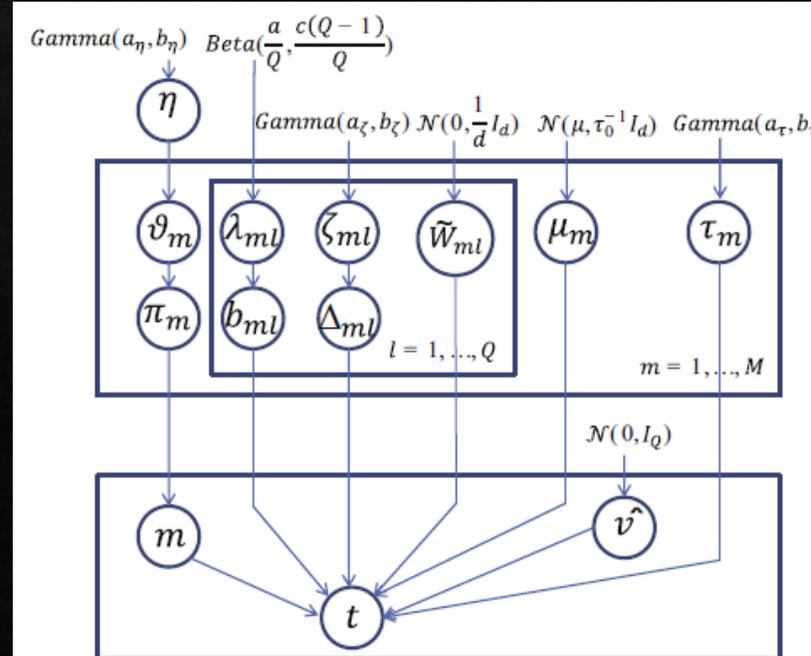
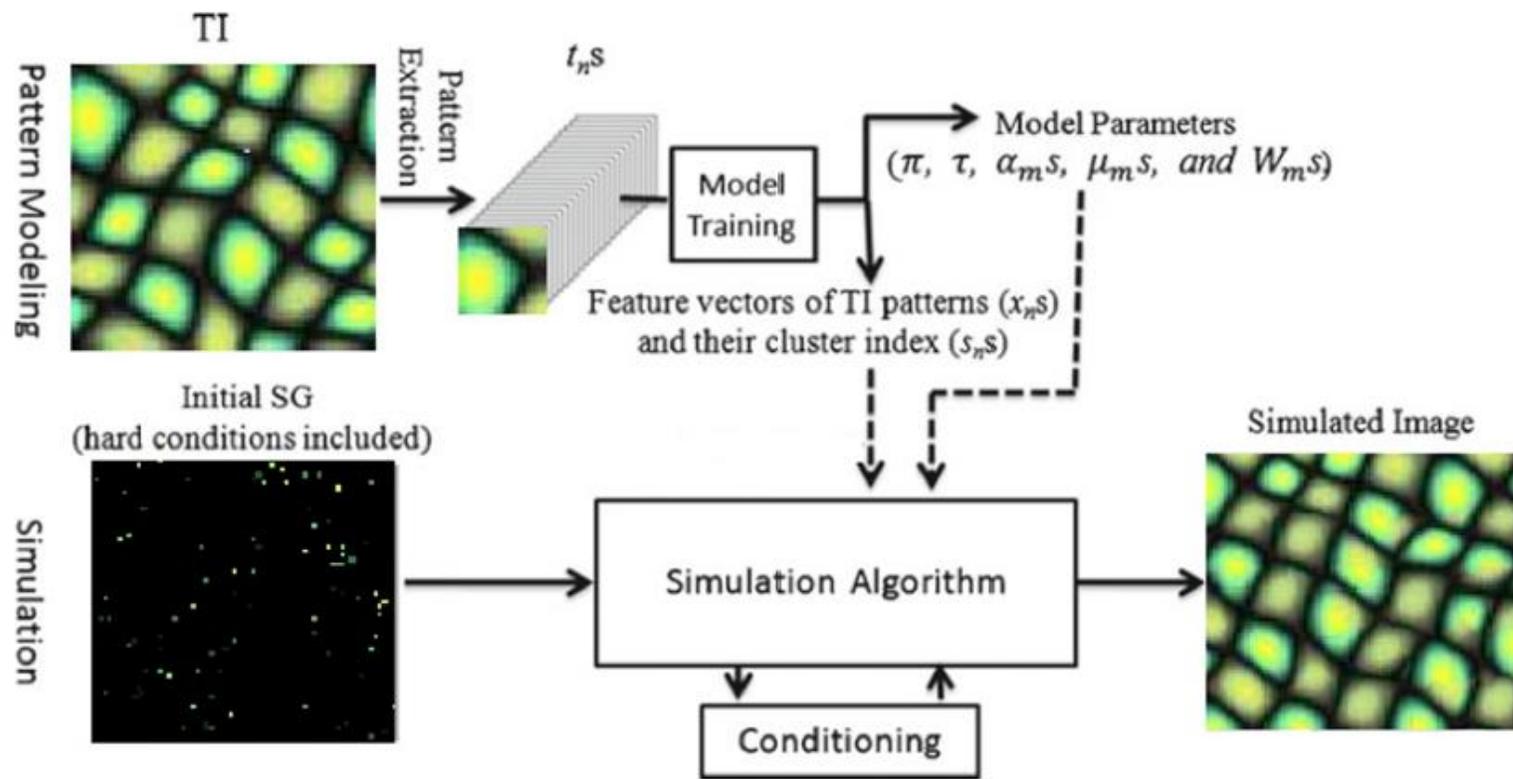


Advanced Dimension reduction techniques to Intrinsic dimension





Bayesian MPPCA (ζ)



Data Reduction

Dimensionality Reduction

Wavelet Transform

Principal Component Analysis

Attribute Subset Selection

Numerosity Reduction

Nonparametric

Data Cube Aggregation

Clustering

Sampling

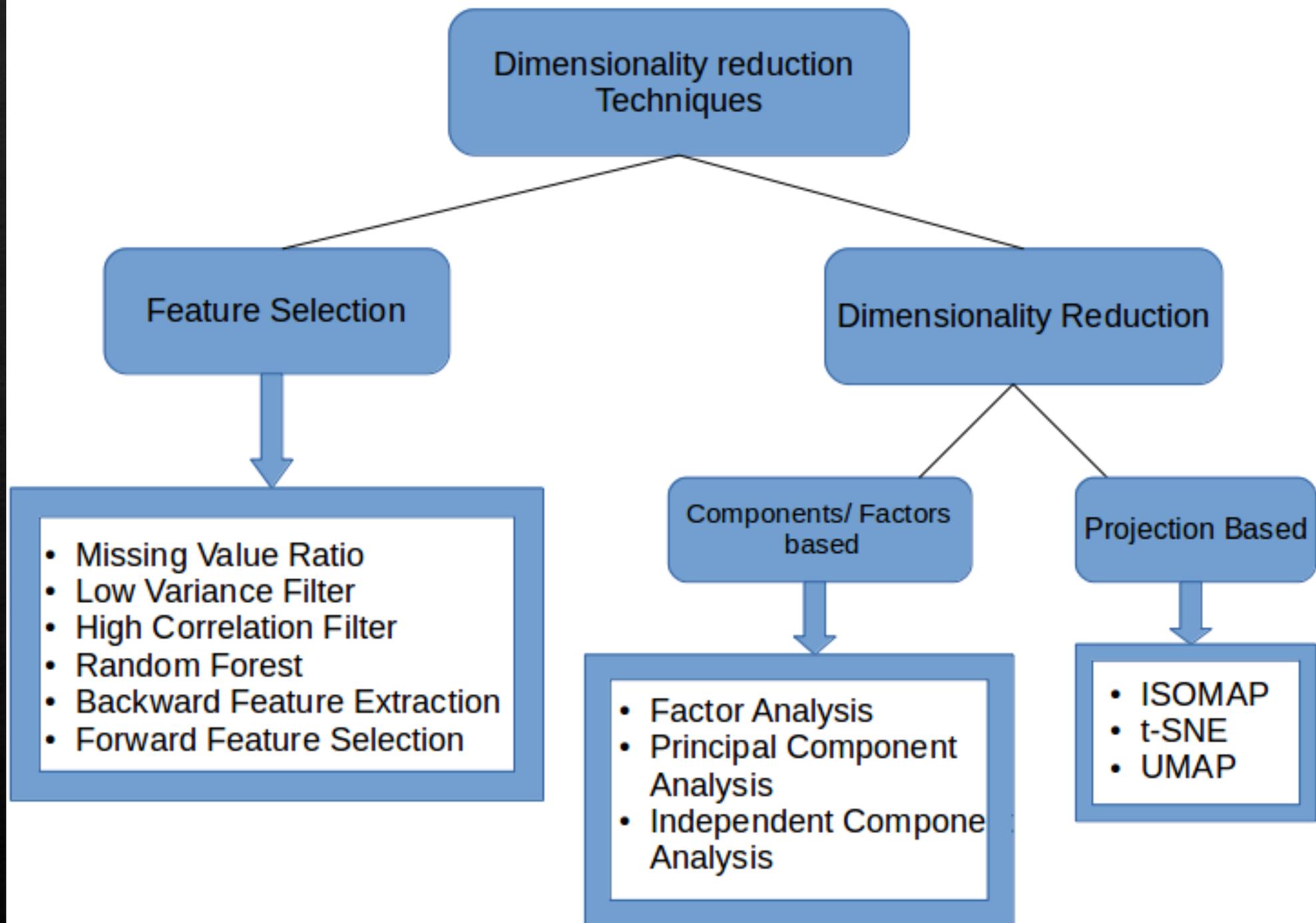
Histogram

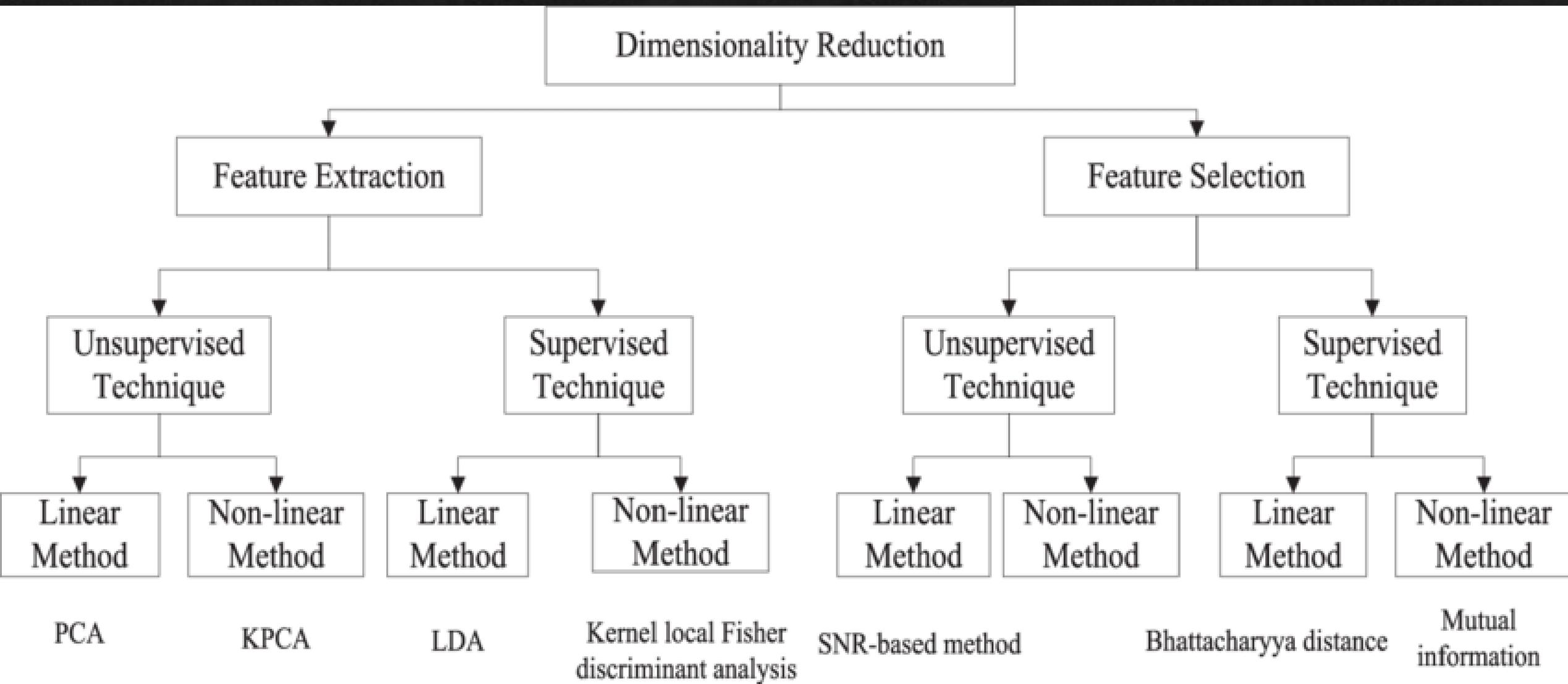
Data Compression

Parametric

Regression

Graphical Models



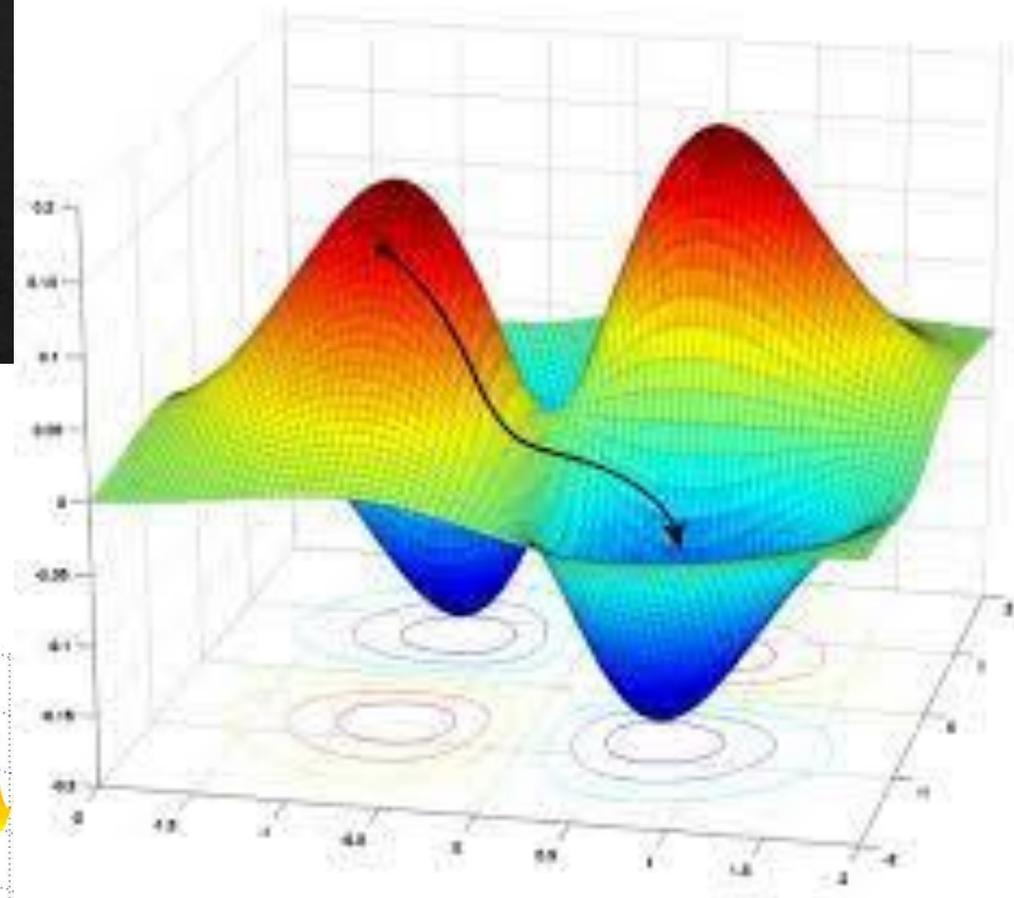
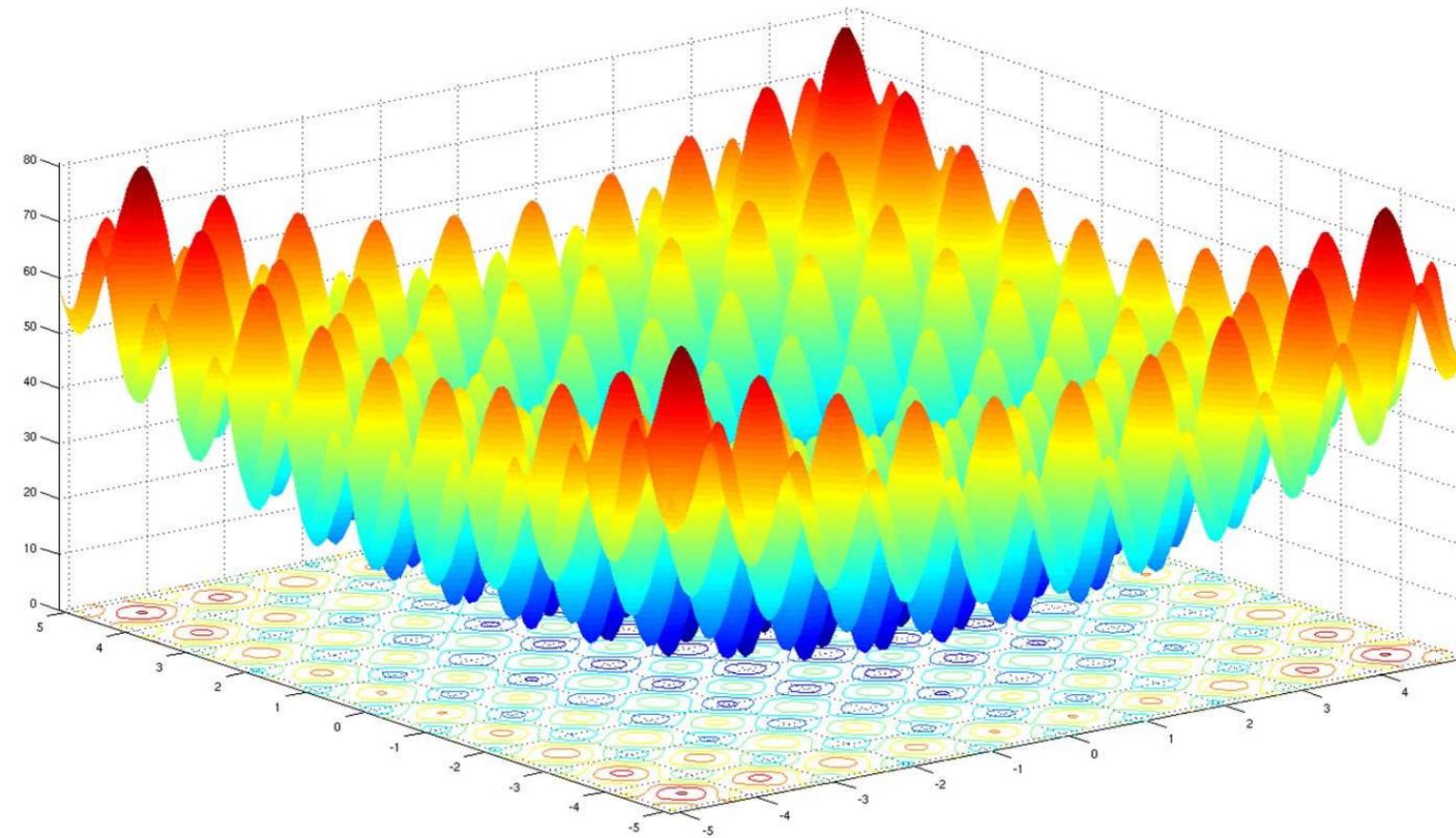


چاپگاه روشهای فرااستگاری در

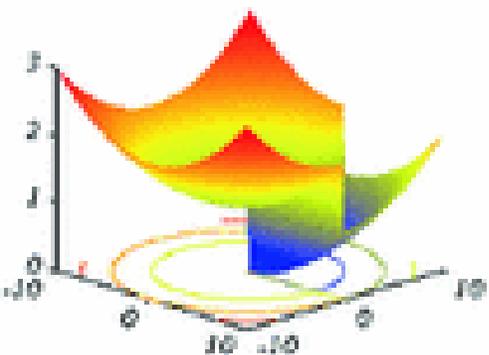
Optimization

Optimization

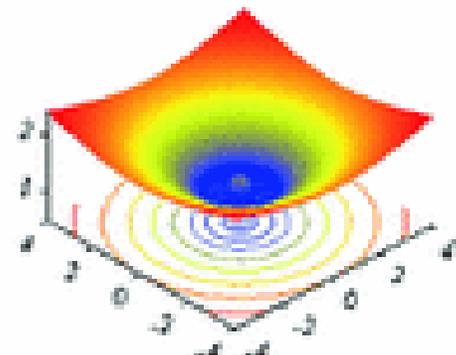
Rastrigin function



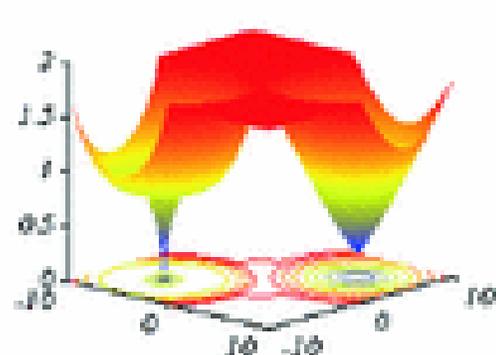
TP1



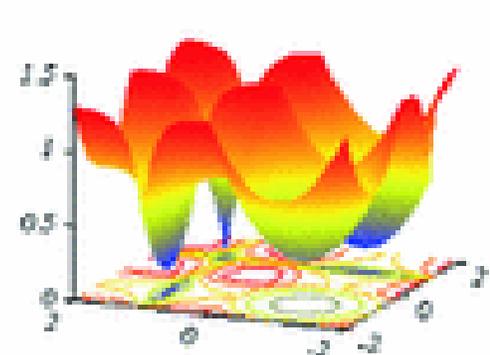
TP2



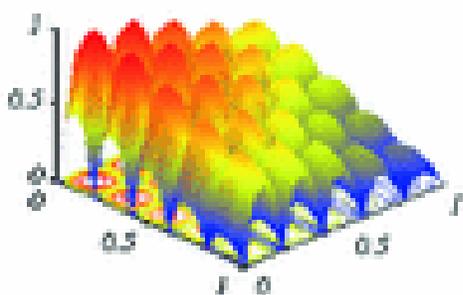
TP3



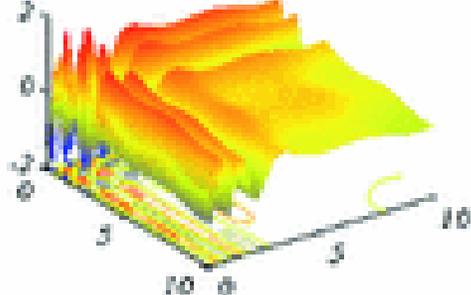
TP4



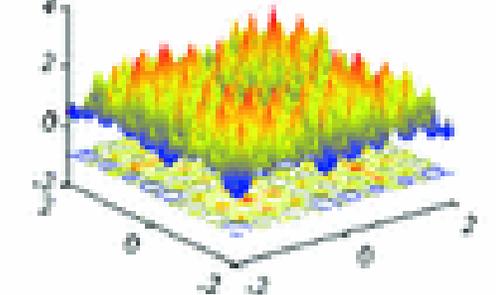
TP5



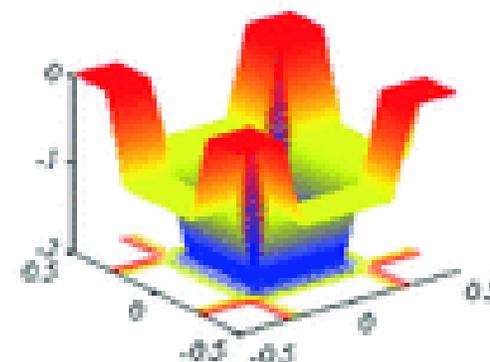
TP6



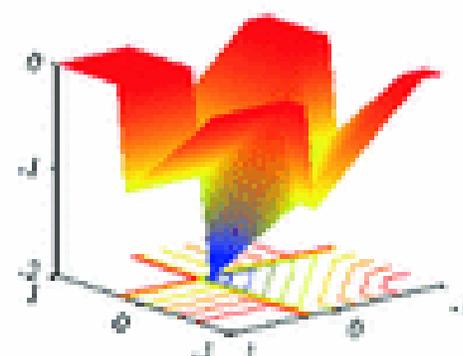
TP7



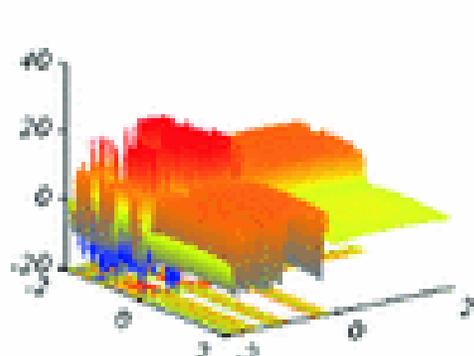
TP8



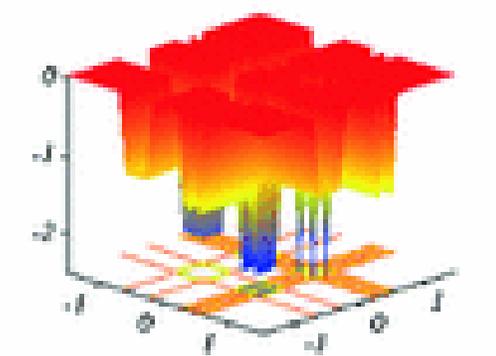
TP9



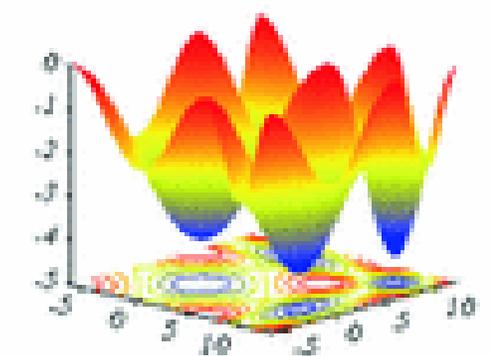
TP10



TP11



TP12



Output

Mapping from features

Hand-designed program

Hand-designed features

Features

More abstract features

Simple features

Deep Learning

Representation Learning

Machine Learning

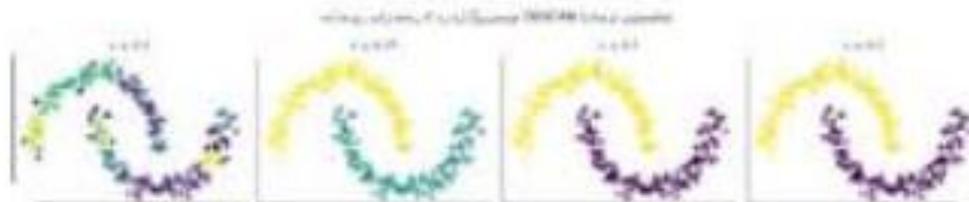
Artificial Intelligence

Input

Curse of Dimensionality

• تأثیر در الگوریتم‌های مختلف

نوع الگوریتم	اثر نفرین بُعدی
K-Means, KNN	فاصله‌ها معنی‌دار نیستند
Decision Trees	شاخه‌های بیش‌ازحد زیاد → Overfitting
DBSCAN	چگالی واقعی پیدا نمی‌شود
SVM	نیاز به کرنل مناسب یا کاهش بُعد
Neural Networks	نیاز به داده‌ی زیاد و زمان آموزش طولانی



در این تصویر می بینیم که مقدار ϵ (epsilon) تأثیر مستقیمی بر نحوه‌ی تشکیل خوشه‌ها در DBSCAN دارد:

- $\epsilon = 0.1$: خیلی کوچک است \rightarrow خوشه‌ها شکسته و نقاط زیادی به عنوان نویز (رنگ خاکستری) در نظر گرفته می‌شوند.
- $\epsilon = 0.15$: شروع به اتصال نقاط هم‌جوار می‌کند، اما هنوز تعدادی نویز وجود دارد.
- $\epsilon = 0.2$: بهترین حالت! دو خوشه‌ی نعل‌اسبی درست شناسایی شده‌اند.
- $\epsilon = 0.3$: بیش از حد بزرگ \rightarrow دو خوشه به هم چسبیده‌اند و به صورت یک خوشه‌ی بزرگ دیده می‌شوند.

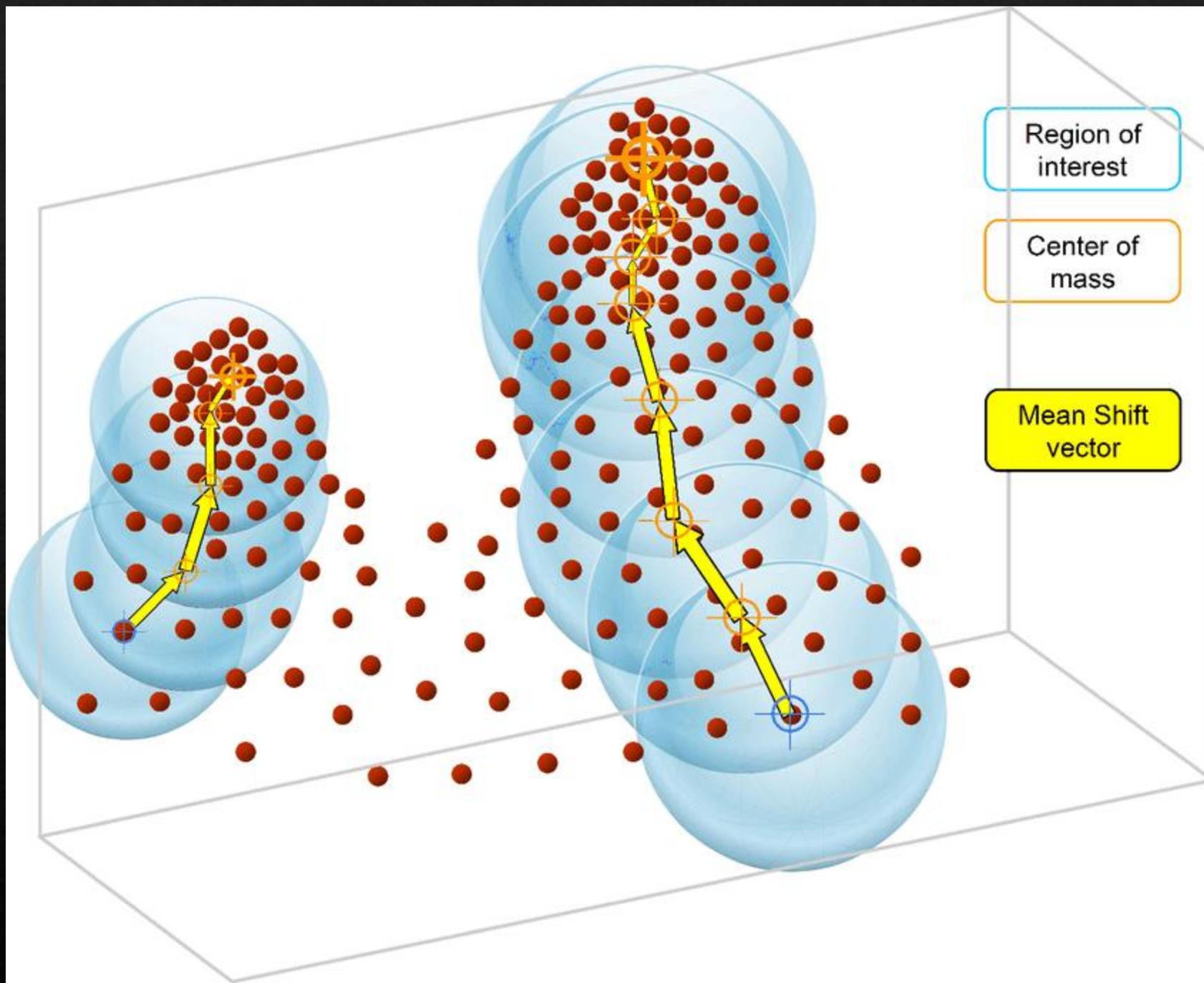
تفاوت CLARANS با CLARA و K-Medoids

ویژگی	K-Medoids	CLARA
داده‌ی ورودی	کل داده‌ها	نمونه‌گیری از داده‌ها
سرعت	کند	سریع‌تر
دقت	بالا	بستگی به نمونه دارد
روش جستجو	کامل (exhaustive)	روی نمونه‌ها

خلاصه مقایسه

نوع الگوریتم	نیاز به تعیین k	برای شکل‌های غیرکروی
K-Means	بله <input checked="" type="checkbox"/>	خیر <input type="checkbox"/>
Hierarchical	اختیاری <input type="checkbox"/>	گاهی <input checked="" type="checkbox"/>
DBSCAN	خیر <input type="checkbox"/>	بله <input checked="" type="checkbox"/>
GMM	بله <input checked="" type="checkbox"/>	بله <input checked="" type="checkbox"/>
Spectral	بله <input checked="" type="checkbox"/>	بله <input checked="" type="checkbox"/>

Clustering methods
Mean shift



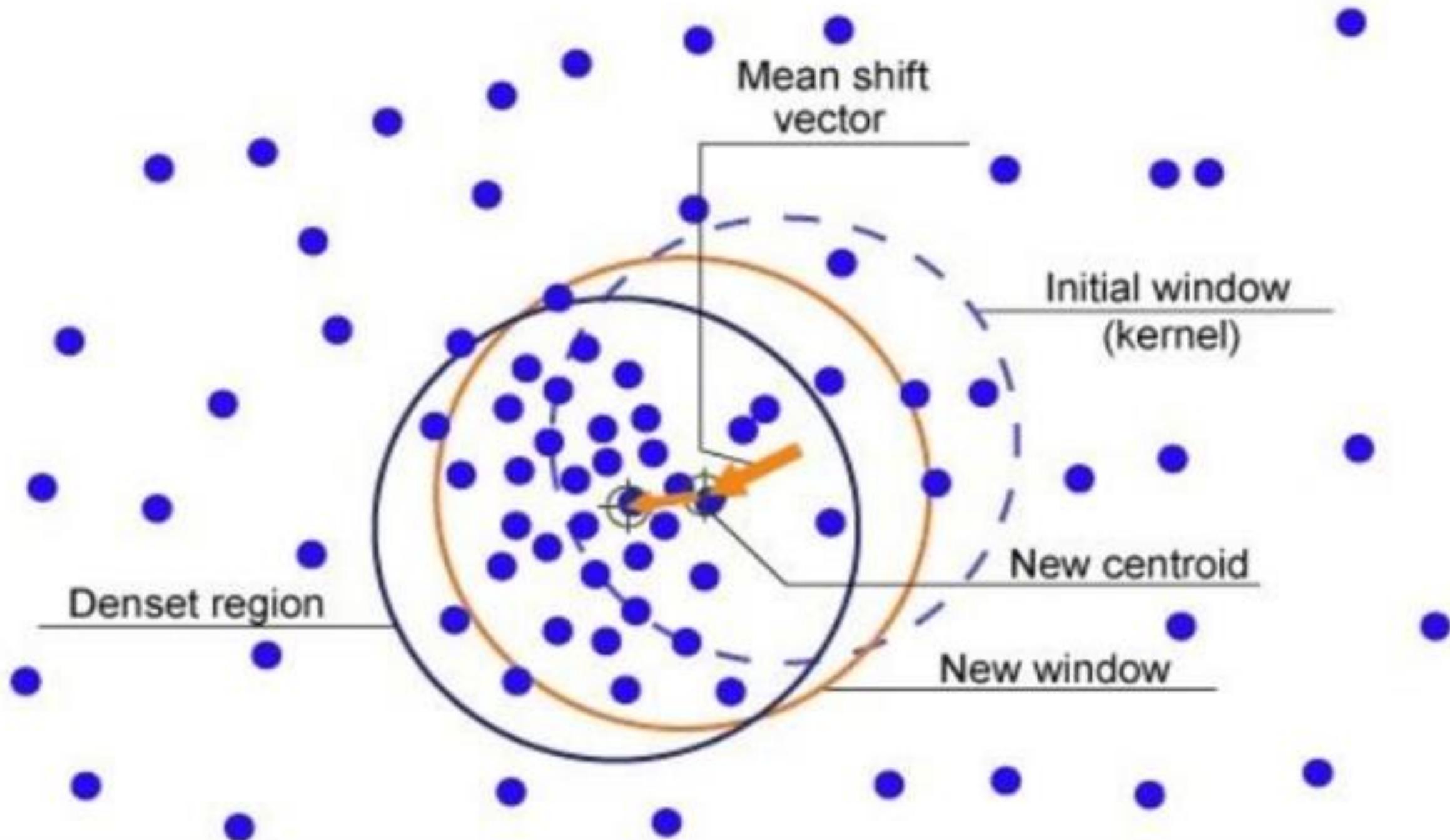
Mean shift vector

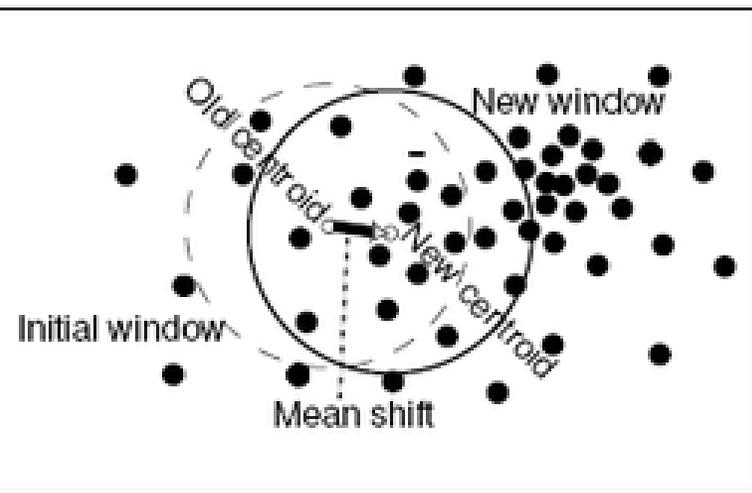
Initial window (kernel)

New centroid

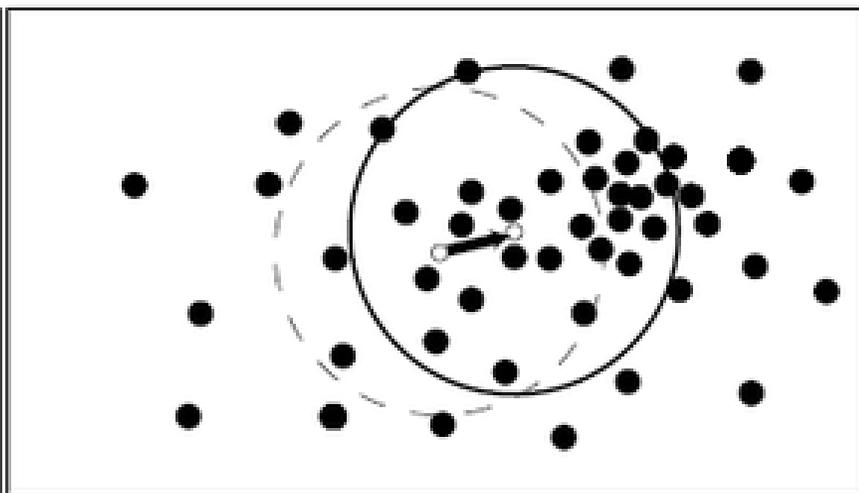
New window

Denset region

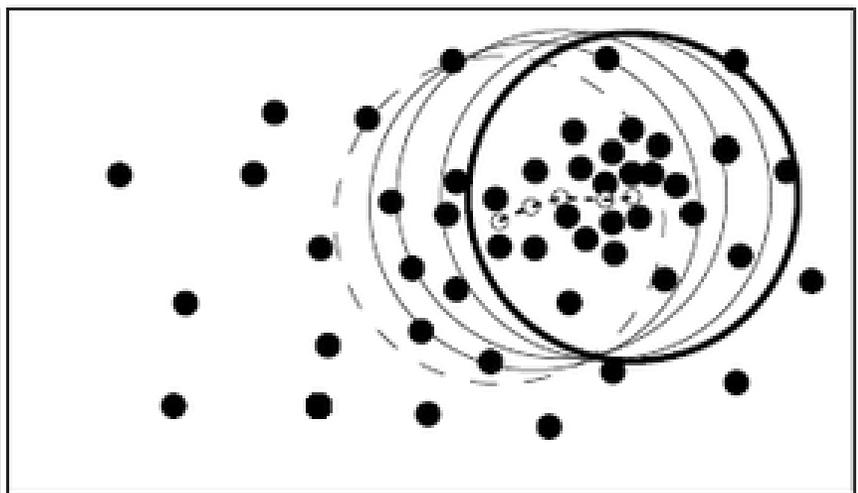




1st Iteration

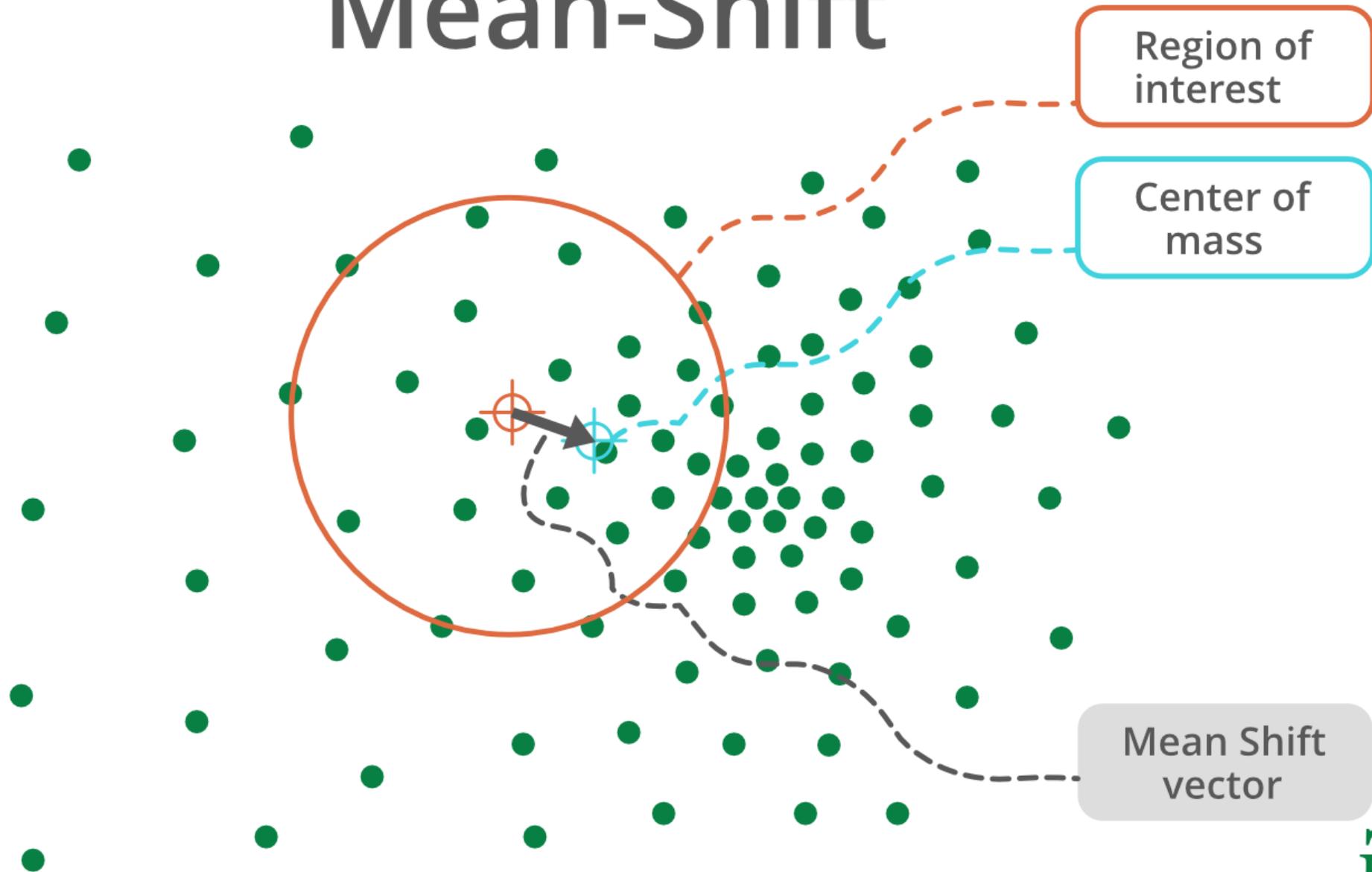


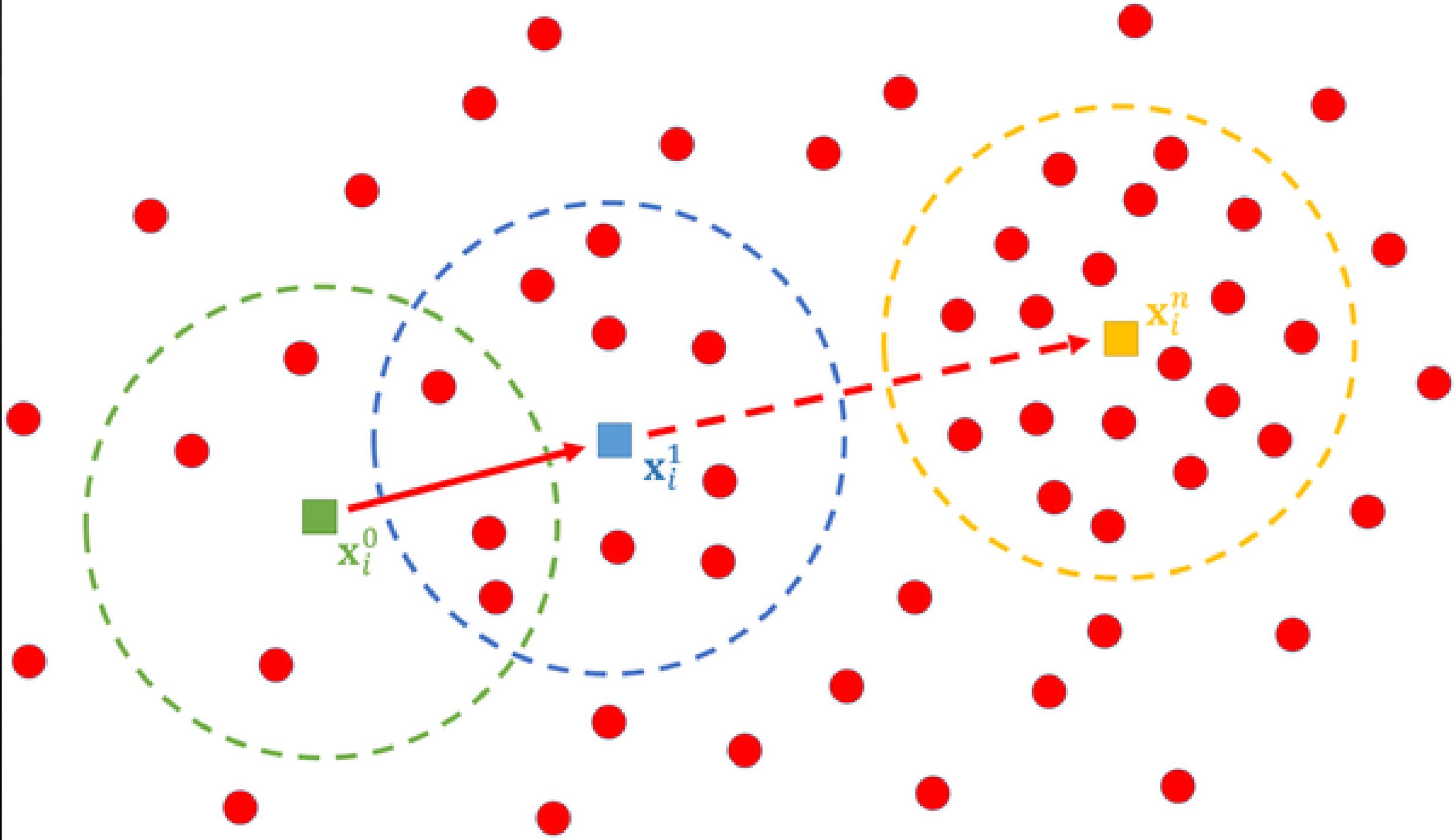
2nd Iteration



Nth iteration

Mean-Shift





A. Gaussian mean-shift (MS) algorithm

```
for  $n \in \{1, \dots, N\}$   
   $\mathbf{x} \leftarrow \mathbf{x}_n$   
  repeat  
     $\forall n: p(n|\mathbf{x}) \leftarrow \frac{\exp\left(-\frac{1}{2}\|(\mathbf{x}-\mathbf{x}_n)/\sigma\|^2\right)}{\sum_{n'=1}^N \exp\left(-\frac{1}{2}\|(\mathbf{x}-\mathbf{x}_{n'})/\sigma\|^2\right)}$   
     $\mathbf{x} \leftarrow \sum_{n=1}^N p(n|\mathbf{x})\mathbf{x}_n$   
  until stop  
   $\mathbf{z}_n \leftarrow \mathbf{x}$   
end  
connected-components( $\{\mathbf{z}_n\}_{n=1}^N, \epsilon$ )
```

C. Gaussian MS algorithm in matrix form

```
 $\mathbf{Z} = \mathbf{X}$   
repeat  
   $\mathbf{W} = \left(\exp\left(-\frac{1}{2}\|(\mathbf{z}_m - \mathbf{x}_n)/\sigma\|^2\right)\right)_{nm}$   
   $\mathbf{D} = \text{diag}\left(\sum_{n=1}^N w_{nm}\right)$   
   $\mathbf{Q} = \mathbf{W}\mathbf{D}^{-1}$   
   $\mathbf{Z} = \mathbf{X}\mathbf{Q}$   
until stop  
connected-components( $\{\mathbf{z}_n\}_{n=1}^N, \epsilon$ )
```

B. Gaussian blurring mean-shift (BMS) algorithm

```
repeat  
  for  $m \in \{1, \dots, N\}$   
     $\forall n: p(n|\mathbf{x}) \leftarrow \frac{\exp\left(-\frac{1}{2}\|(\mathbf{x}_m - \mathbf{x}_n)/\sigma\|^2\right)}{\sum_{n'=1}^N \exp\left(-\frac{1}{2}\|(\mathbf{x}_m - \mathbf{x}_{n'})/\sigma\|^2\right)}$   
     $\mathbf{y}_m \leftarrow \sum_{n=1}^N p(n|\mathbf{x}_m)\mathbf{x}_n$   
  end  
   $\forall m: \mathbf{x}_m \leftarrow \mathbf{y}_m$   
until stop  
connected-components( $\{\mathbf{x}_n\}_{n=1}^N, \epsilon$ )
```

D. Gaussian BMS algorithm in matrix form

```
repeat  
   $\mathbf{W} = \left(\exp\left(-\frac{1}{2}\|(\mathbf{x}_m - \mathbf{x}_n)/\sigma\|^2\right)\right)_{nm}$   
   $\mathbf{D} = \text{diag}\left(\sum_{n=1}^N w_{nm}\right)$   
   $\mathbf{P} = \mathbf{W}\mathbf{D}^{-1}$   
   $\mathbf{X} = \mathbf{X}\mathbf{P}$   
until stop  
connected-components( $\{\mathbf{x}_n\}_{n=1}^N, \epsilon$ )
```

Algorithm 1 Structure of the Mean Shift algorithm.

Data: $\mathbf{X} = \{\mathbf{x}_i\}$ - input image, $\mathbf{Y} = \{\mathbf{y}_i\}$ - output image, $\mathbf{x}_i^{(t)} = (\boldsymbol{\xi}_i^{(t)}, \boldsymbol{\eta}_i^{(t)})$ - pixel at position $\boldsymbol{\xi}_i^{(t)}$ with color vector $\boldsymbol{\eta}_i^{(t)}$,
 $\mathcal{B}_i^{(t)}$ - processing block centered at $\boldsymbol{\xi}_i^{(t)}$, $\varepsilon = 10^{-3}$ - stopping criterion

- 1: **for** each image pixel $\mathbf{x}_i \in \mathbf{X}$ **do**
- 2: $t = 1$, $\mathbf{x}_i^{(1)} = \mathbf{x}_i$ (initialization)
- 3: **do**
- 4: $t++$
- 5: calculate $\mathbf{x}_i^{(t)} = (\boldsymbol{\xi}_i^{(t)}, \boldsymbol{\eta}_i^{(t)})$ according to (13)
- 6: move center of $\mathcal{B}_i^{(t+1)}$ to $\lfloor \boldsymbol{\xi}_i^{(t)} \rfloor$ (rounding to the nearest integer)
- 7: **while** $\|\mathbf{x}_i^{(t)} - \mathbf{x}_i^{(t-1)}\| \geq \varepsilon$
- 8: $\mathbf{y}_i = (\boldsymbol{\xi}_i^{(1)}, \boldsymbol{\eta}_i^{(t)})$
- 9: **end for**

Mean-shift offset

Spatial kernel

Color weight

Local pixel offset

$$\Delta \mathbf{x} = \frac{\sum_i \mathbf{K}(\mathbf{x}_i - \mathbf{x}) w(\mathbf{x}_i) (\mathbf{x}_i - \mathbf{x})}{\sum_i \mathbf{K}(\mathbf{x}_i - \mathbf{x}) w(\mathbf{x}_i)}$$

Normalization factor

The Mean-Shift

$$\nabla P(x) = \frac{1}{N} 2c \sum_n g_n \left(\frac{\sum_n x_n g_n}{\sum_n g_n} - x \right)$$

Mean

Shift