

# Introduction to Clustering and its Applications

Dongfeng Li

Instructor: Yikun Zhang

# Index

- Introduction
- Materials and Methods
- Code Implementation
- Results
- Discussion
- References

# Introduction

## Clustering

- In many real world contexts, there aren't clearly defined labels so we won't be able to do classification
- We will need to come up with methods that uncover structure from the (unlabeled) input data  $X$ .
- **Clustering** is an automatic process of trying to find related groups within the given dataset.

*Input:  $x_1, x_2, \dots, x_n$*

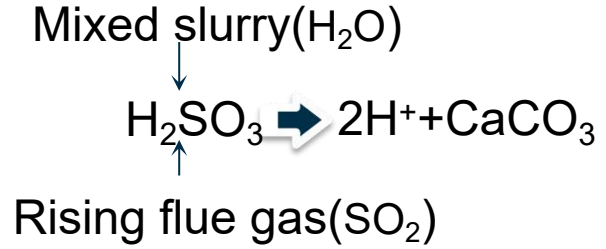


*Output:  $z_1, z_2, \dots, z_n$*





# Materials and Methods Data



## Huangtai power plant desulfurization process

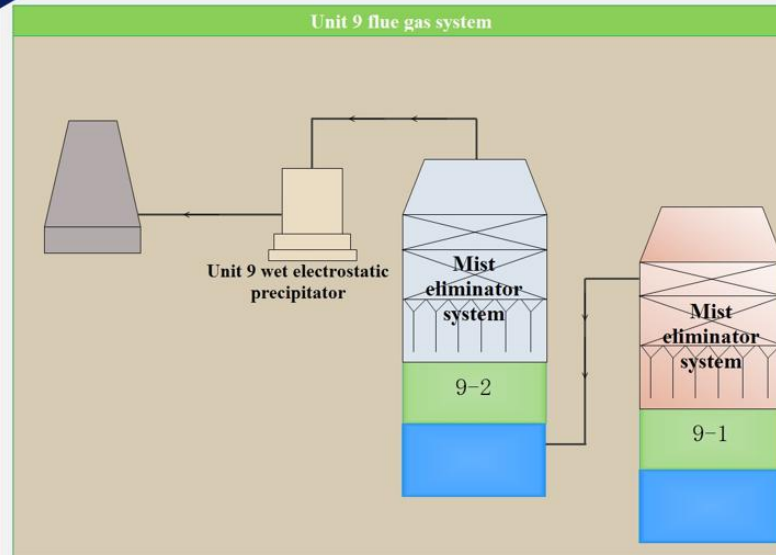


Figure 6. The framework of system

Research object is Unit 9 of Huangtai with double absorption towers.

The system includes:  
SO<sub>2</sub> absorption sys,  
flue-gas sys,  
absorbent preparation sys,  
gypsum treatment sys,  
processed water sys,  
accident slurry sys, etc

# Methods: Variable Selection

Variable selection by variance

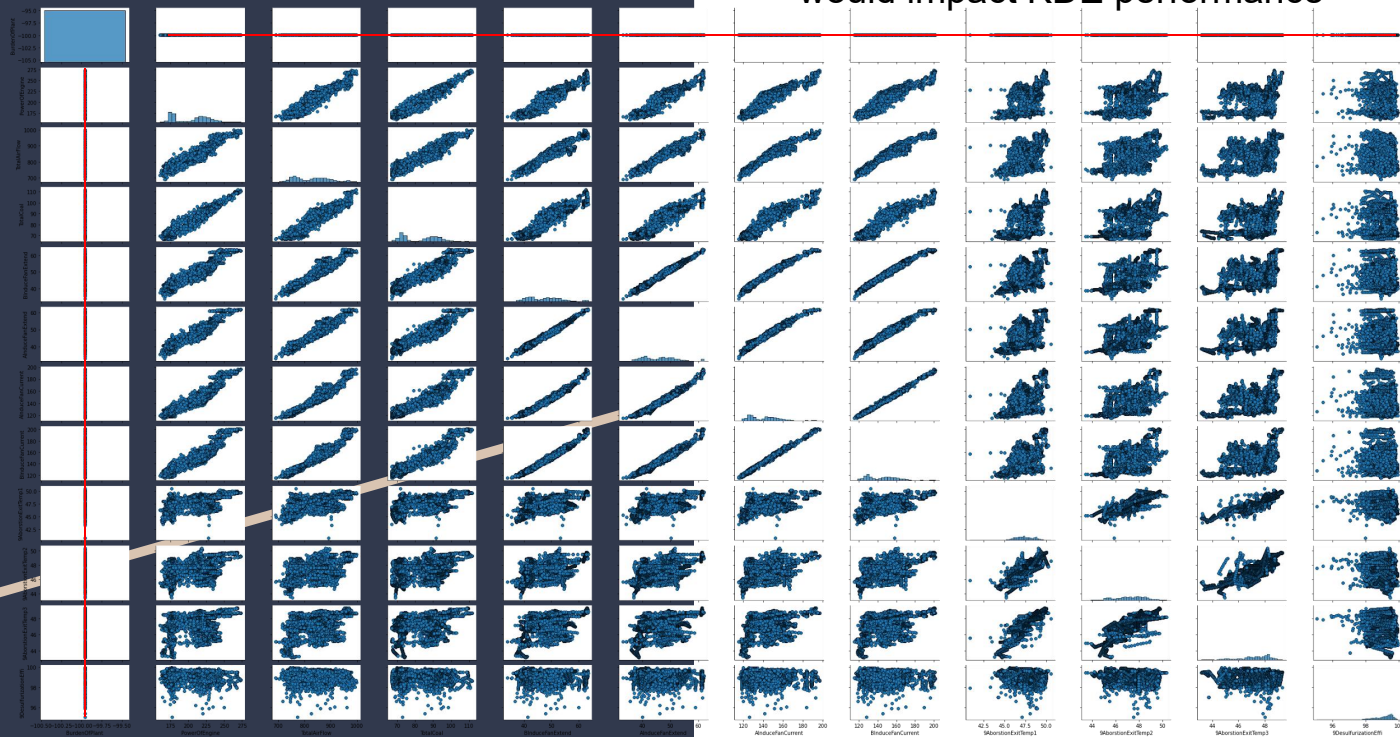
Why?

- Too many data points, some of which do not provide useful information.
- High dimensions bring curse of dimensionality.
- Perform clustering on scattered data instead of too much concentrated data.

# Methods

pre-selection

does not give useful information and would impact KDE performance



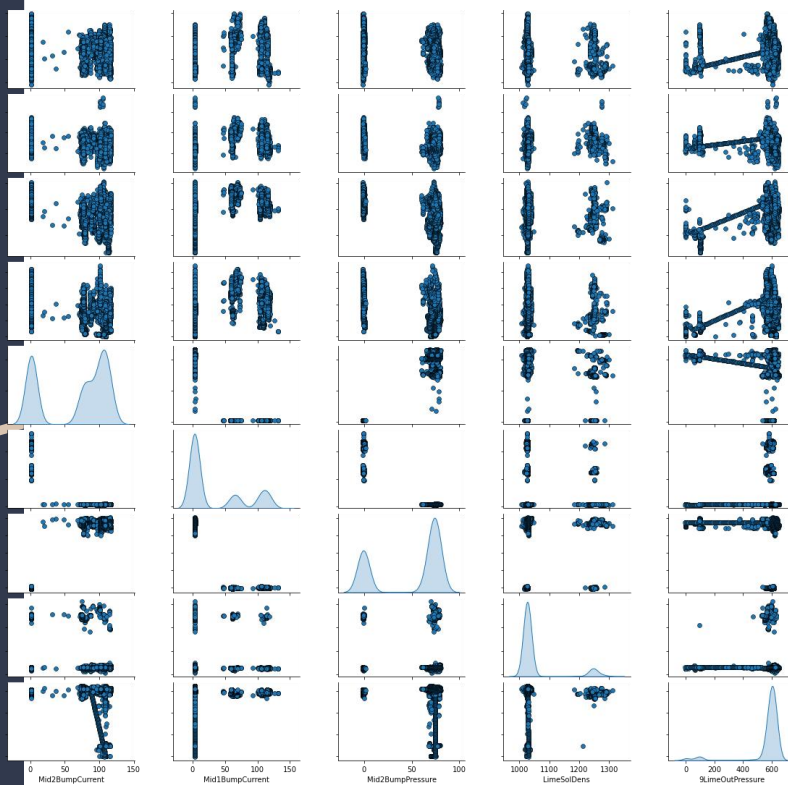


# Methods



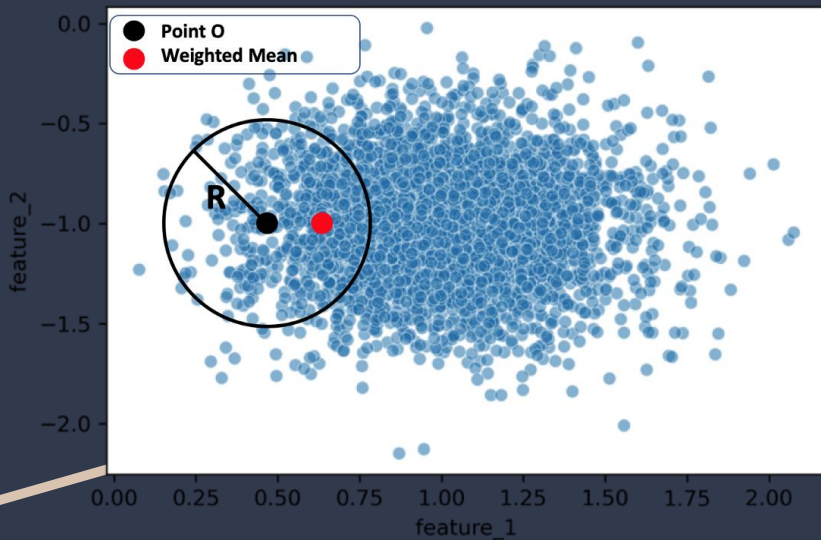
Test Plots

## post-selection





# Methods: Mean-shift Clustering



Density-based clustering

Weighted mean:

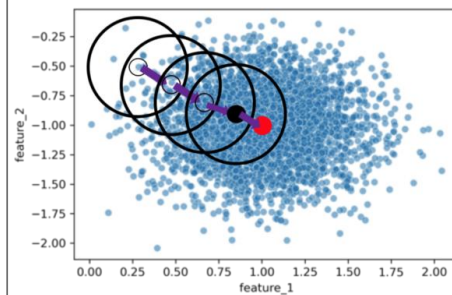
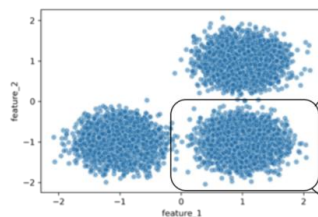
$$M_W = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$$

Using Gaussian weight function:

$$w(d) = e^{-\frac{d}{2\sigma^2}}$$

$d$  is distance from the current point toward any other points.

$\sigma$  is the parameter to adjust for how fast the weight decreases with the increasing of distance



# Methods: Mean-shift Clustering

KDE is kind of similar to PDF

Then by gradient descent we could  
get cluster updating

Density gradient estimator:

$$\hat{\nabla} f_{h,K}(\mathbf{x}) \equiv \nabla \hat{f}_{h,K}(\mathbf{x}) = \frac{2c_{k,d}}{nh^{d+2}} \sum_{i=1}^n (\mathbf{x} - \mathbf{x}_i) k' \left( \left\| \frac{\mathbf{x} - \mathbf{x}_i}{h} \right\|^2 \right).$$

- Kernel density estimator(KDE):

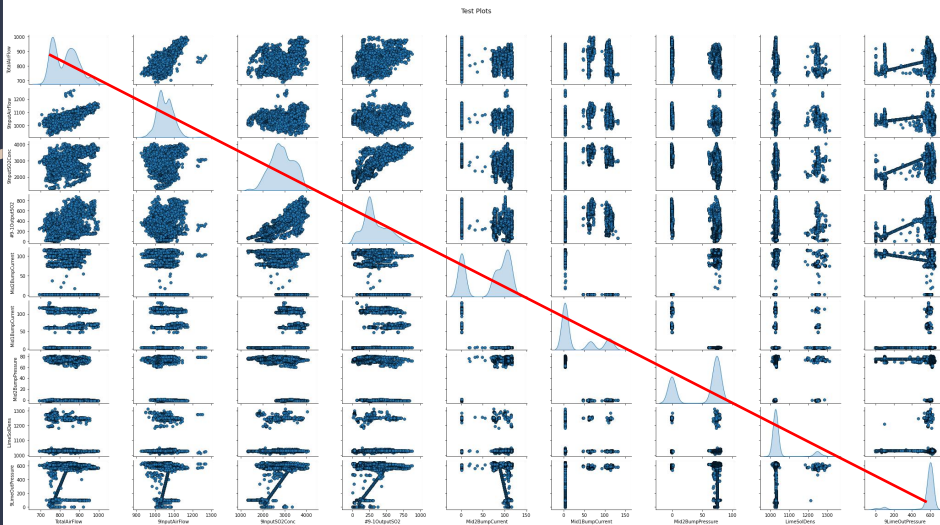
$$\hat{f}(\mathbf{x}) = \frac{1}{nh^d} \sum_{i=1}^n K \left( \frac{\mathbf{x} - \mathbf{x}_i}{h} \right).$$

$n$  is number of data,

$h$  is bandwidth,

$d$  is dimension of data,

$K$  is the kernel function(we use Gaussian kernel).

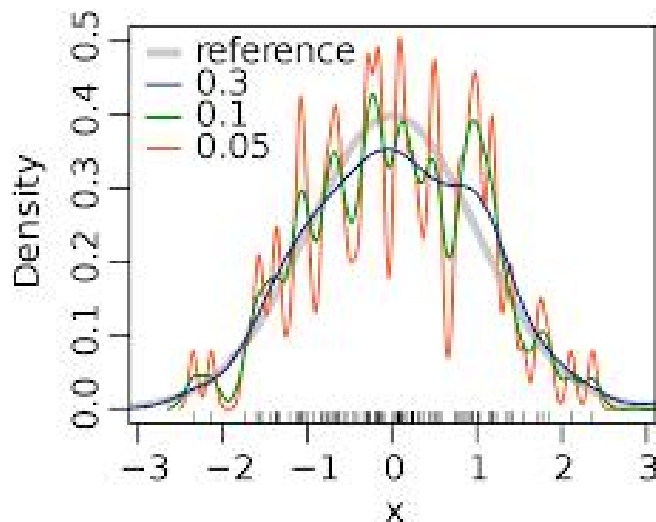


# Methods: Mean-shift Clustering

- Selection of bandwidth

Applying Silverman's rule of thumb to self-select bandwidth since we are using Gaussian Kernel

$$h = \left(\frac{4}{d+2}\right)^{\frac{1}{d+4}} \cdot n^{\frac{-1}{d+4}} \cdot \sum_0^n \frac{\sigma}{n}$$



Under-smooth  
vs  
Over-smooth

Examples of different bandwidths on KDE

# Methods: K-means Clustering

From STAT 416 I learned about K-means clustering:

Define the Score for assigning a point to a cluster is

$$\text{Score}(x_i, \mu_j) = \text{dist}(x_i, \mu_j)$$

Lower score => Better Clustering.

## **k-means Algorithm at a glance**

Step 0: Initialize cluster centers

Repeat until convergence:

Step 1: Assign each example to its closest cluster centroid

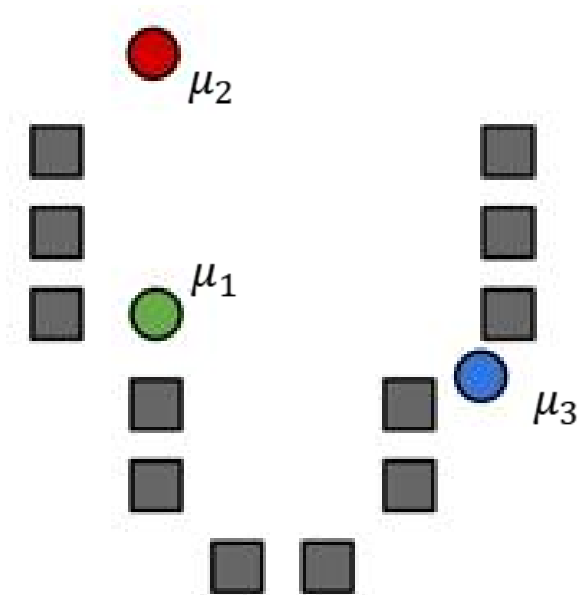
Step 2: Update the centroids to be the average of all the points

assigned to that cluster

# Step 0:

Start by choosing the initial cluster centroids

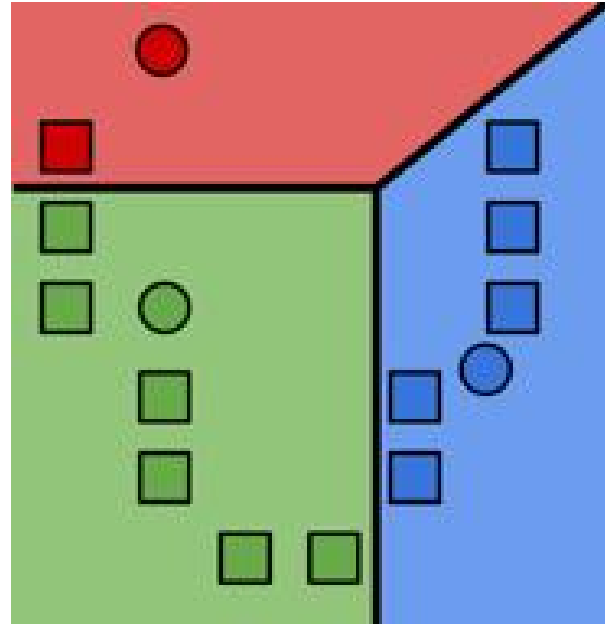
- A common default choice is to choose centroids at random
- Will see later that there are smarter ways of initializing



# Step 1:

Assign each example to its closest cluster centroid

●  $z_i \leftarrow \operatorname{argmin}_{j \in [k]} \|\mu_j - x_i\|^2$

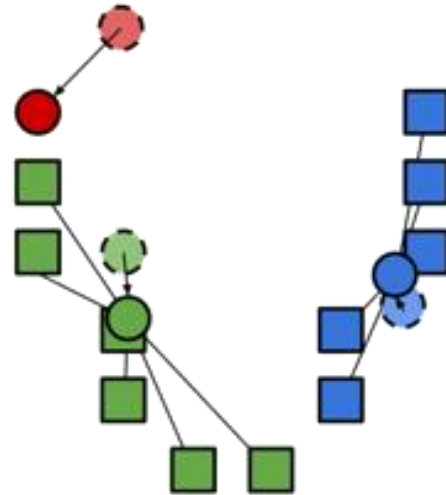


## Step 2:

- Update the centroids to be the mean of all the points assigned to that cluster.

$$\mu_j \leftarrow \frac{1}{n_j} \sum_{i:z_i=j} x_i$$

Computes center of mass for cluster!

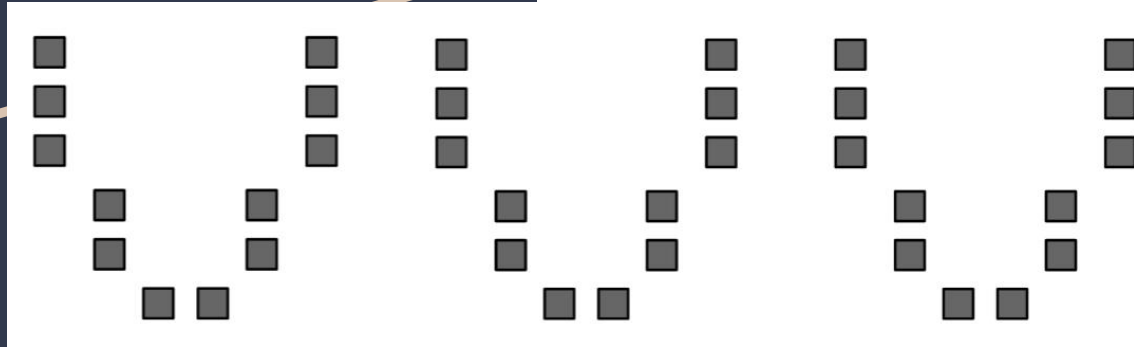




# K-means++: Smart Initialization

● **means++** does a slightly smarter random initialization

1. Choose first cluster  $\mu_1$  from the data uniformly at random
2. For the current set of centroids (starting with just  $\mu_1$ ), compute the distance between each datapoint and its closest centroid
3. Choose a new centroid from the remaining data points with probability of  $x_i$  being chosen proportional to  $d(x_i)^2$
4. Repeat 2 and 3 until we have selected  $k$  centroids

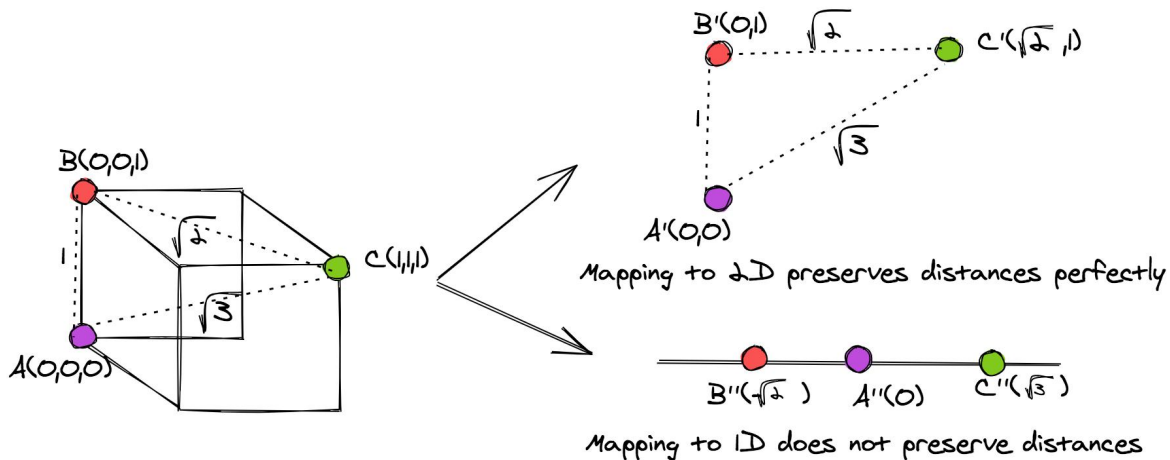


# Methods: Multidimensional Scaling(MDS)

Multidimensional Scaling(MDS) gives a deduction in dimensionality so that we could have a better visualization while still keep the distance or pattern somehow.

We use Euclidean distance and metric MDS which preserve the pairwise distance here.

$$d(p, q) = \sqrt{(q_1 - p_1)^2 + \dots + (q_n - p_n)^2}$$



A possible mapping of points from 3D to 2D and 1D. The mapping is not optimized.

# Code Implementation

```
def MS(mesh_0, data, h=None, eps=1e-7, max_iter=1000, wt=None):
    """
    Mean Shift Algorithm with the Gaussian kernel.

    Parameters
    -----
    mesh_0: a (m,d)-array
        The coordinates of m initial points in the d-dim Euclidean space.

    data: a (n,d)-array
        The coordinates of n data sample points in the d-dim Euclidean space.

    h: float
        The bandwidth parameter. (Default: h=None. Then the Silverman's
        rule of thumb is applied. See Chen et al.(2016) for details.)

    eps: float
        The precision parameter.

    max_iter: int
        The maximum number of iterations for the SCMS algorithm on each
        initial point.

    wt: (n,)-array
        The weights of kernel density contributions for n random sample
        points. (Default: wt=None, that is, each data point has an equal
        weight "1/n".)

    Return
    -----
    MS_path (m,d,T)-array
        The entire iterative MS sequence for each initial point.
    """
    n = data.shape[0] # Number of data points
    d = data.shape[1] # Dimension of the data

    if h is None:
        # Apply Silverman's rule of thumb to select the bandwidth parameter
        # (Only works for Gaussian kernel)
        h = (4*(d+2))**(1/(d+4))*n**(1/(d+4))*np.mean(np.std(data, axis=0))
        print("The current bandwidth is " + str(h) + ".\n")

    MS_path = np.zeros((mesh_0.shape[0], d, max_iter))
    # Create a vector indicating the convergent status of every mesh point
    conv_sign = np.zeros((mesh_0.shape[0],))
    if wt is None:
        wt = np.ones((n,))
    else:
        wt = n*wt
    MS_path[:, :, 0] = mesh_0
    for t in range(1, max_iter):
        if all(conv_sign == 1):
            print("The MS algorithm converges in " + str(t-1) + " steps!")
            break
        for i in range(mesh_0.shape[0]):
            if conv_sign[i] == 0:
                x_pt = MS_path[i, :, t-1]
                ker_w = wt*np.exp(-np.sum(((x_pt-data)/h)**2, axis=1)/2)
                # Mean shift update
                x_new = np.sum(data*ker_w.reshape(n,1), axis=0) / np.sum(ker_w)
                if LA.norm(x_pt - x_new) < eps:
                    conv_sign[i] = 1
                    MS_path[i, :, t] = x_new
            else:
                MS_path[i, :, t] = MS_path[i, :, t-1]

    if t >= max_iter-1:
        print("The MS algorithm reaches the maximum number of iterations, \
        +str(max_iter)+, and has not yet converged.")
    return MS_path[:, :, t]
```

- Variable selection:  
sklearn.feature\_selection.VarianceThreshold

```
In [7]: def smart_initialize(data, k, seed=None):
        """
        Use k-means++ to initialize a good set of centroids
        """
        if seed is not None: # useful for obtaining consistent results
            np.random.seed(seed)

        centroids = np.zeros((k, data.shape[1]))

        # Randomly choose the first centroid.
        # Since we have no prior knowledge, choose uniformly at random
        idx = np.random.randint(data.shape[0])
        centroids[0] = data[idx, :].toarray()

        # Compute distances from the first centroid chosen to all the other data points
        distances = pairwise_distances(data, centroids[0:1], metric='euclidean').flatten()

        for i in range(1, k):
            # Choose the next centroid randomly, so that the probability for each data point to be chosen
            # is directly proportional to its squared distance from the nearest centroid.
            # Roughly speaking, a new centroid should be as far as from other centroids as possible.
            idx = np.random.choice(data.shape[0], 1, p=distances/sum(distances))
            centroids[i] = data[idx, :].toarray()

            # Now compute distances from the centroids to all data points
            distances = np.min(pairwise_distances(data, centroids[0:i+1], metric='euclidean'), axis=1)

        return centroids
```

# Code Implementation: Plotting

```
import plotly.express as px

# Create a scatter plot
fig = px.scatter(cated, x='x', y='y', opacity=1, color="dataset")

# Change chart background color
fig.update_layout(dict(plot_bgcolor = 'white'))

# Update axes lines
fig.update_xaxes(showgrid=True, gridwidth=1, gridcolor='lightgrey',
                 zeroline=True, zerolinewidth=1, zerolinecolor='lightgrey',
                 showline=True, linewidth=1, linecolor='black')

fig.update_yaxes(showgrid=True, gridwidth=1, gridcolor='lightgrey',
                 zeroline=True, zerolinewidth=1, zerolinecolor='lightgrey',
                 showline=True, linewidth=1, linecolor='black')

# Set figure title
fig.update_layout(title_text="MDS Transformation")

# Update marker size
fig.update_traces(marker=dict(size=5,
                              line=dict(color='black', width=0.2)))

fig.show()
```

```
In [59]: from mpl_toolkits.mplot3d import Axes3D
import matplotlib as mpl
import matplotlib.pyplot as plt
import seaborn as sns

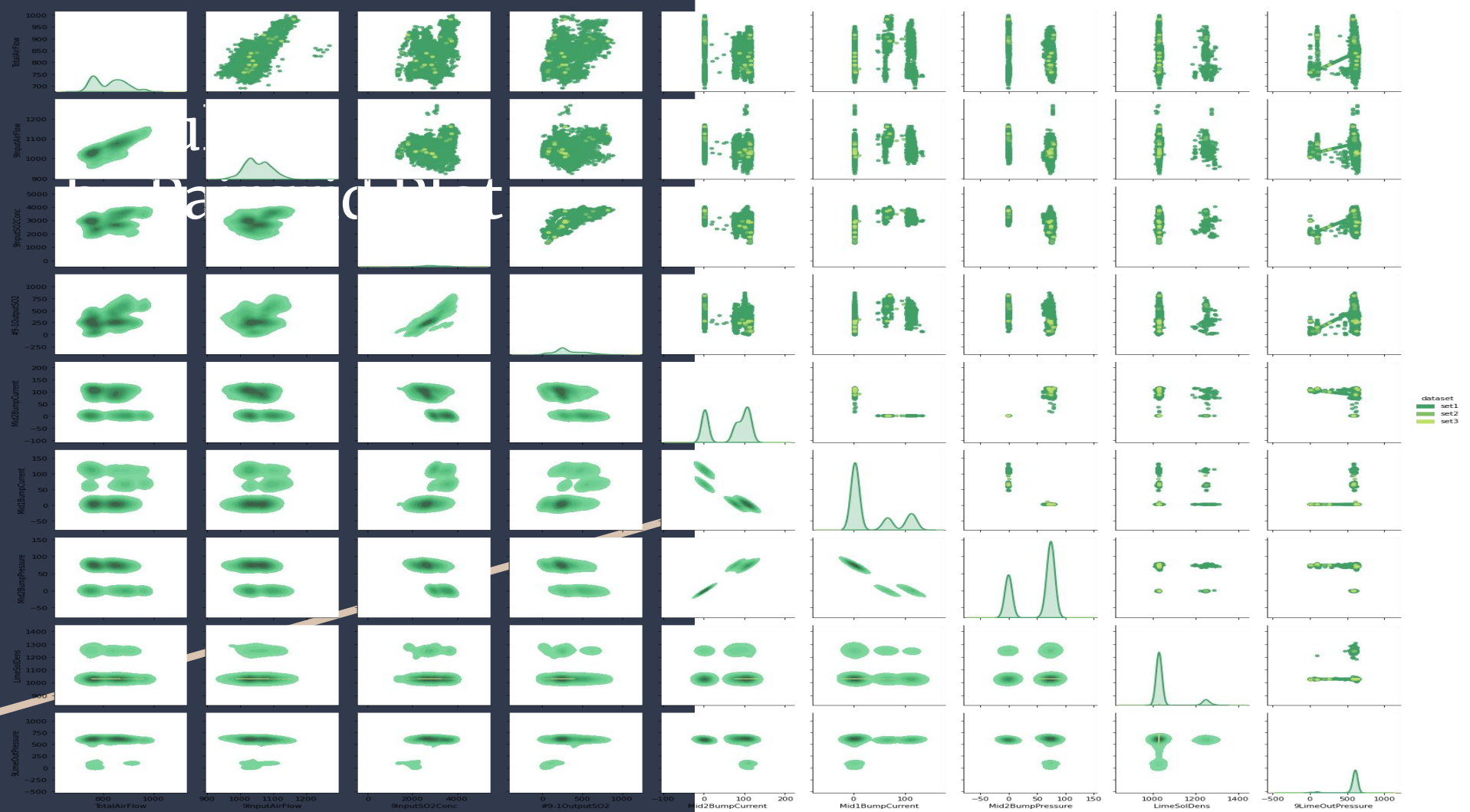
dff = pd.DataFrame(transformed, columns = labels)
dfm = pd.DataFrame(oo[0], columns = labels)
dfk = pd.DataFrame(mm, columns = labels)

concatenated = pd.concat([dff.assign(dataset='set1'), dfm.assign(dataset='set2'), dfk.assign(dataset='set3')], ignore_index = True)

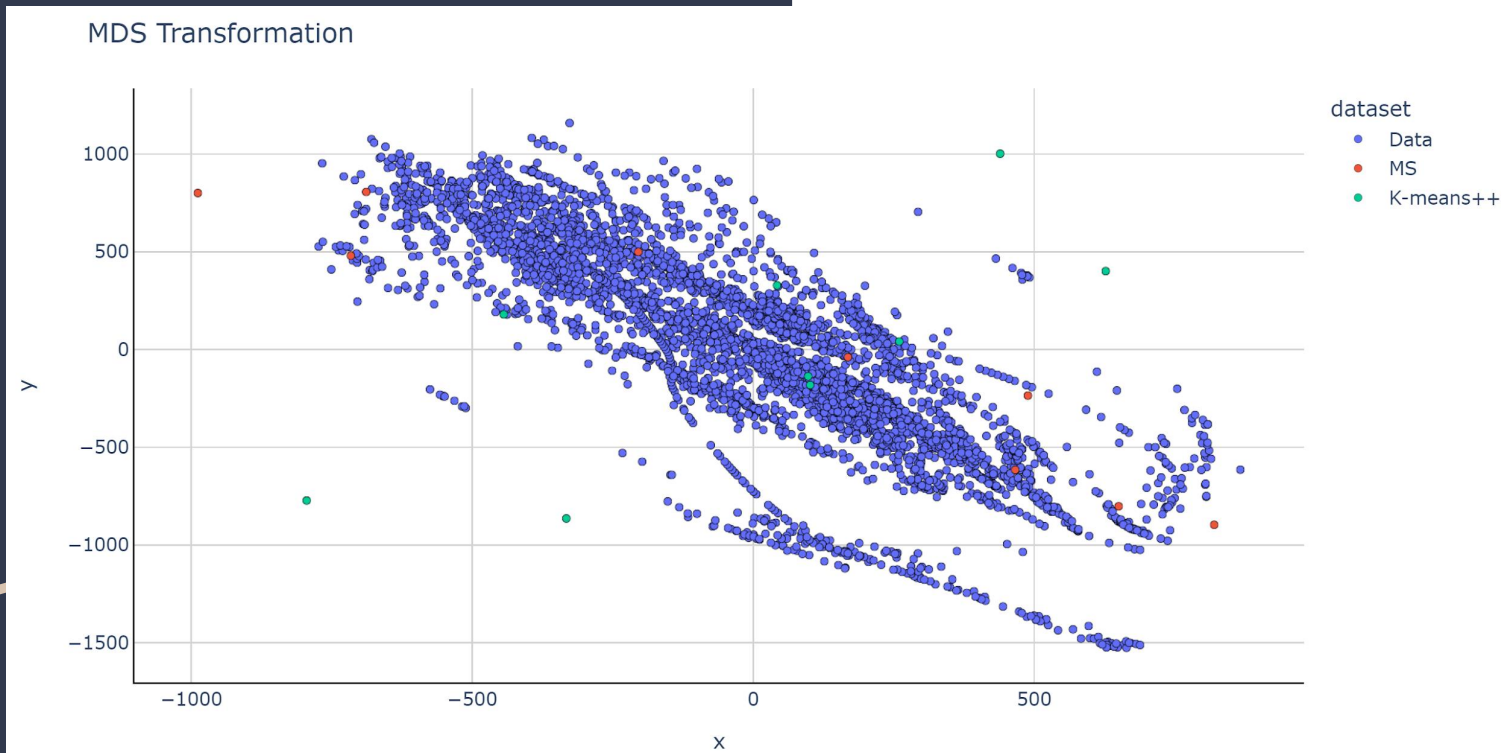
'''pp = sns.pairplot(concatenated, height=1.8, aspect=1.8,
                    plot_kws=dict(edgecolor="k", linewidth=0.5),
                    diag_kind="kde", diag_kws=dict(shade=True),
                    hue="dataset")'''

grid = sns.PairGrid(data=concatenated, hue='dataset', palette = 'summer')
grid = grid.map_upper(plt.scatter, s=20, alpha=0.8)
grid = grid.map_diag(sns.kdeplot, fill=True, lw=2) #Gaussian Kernel
grid = grid.map_lower(sns.kdeplot, fill=True)
grid.add_legend()
#grid = grid.map_lower(sns.kdeplot, Levels=4, color=".2")
#pp.map_lower(sns.kdeplot, Levels=4, color=".2")

'''fig = pp.fig
fig.subplots_adjust(top=0.93, wspace=0.3)
t = fig.suptitle('Test Plots', fontsize=14)'''
```



# Results: by MDS



# Discussion on Performance

- Time
- Who's better?



# Performance: Time Complexity

Time of MS: 240.09044551849365

Time of K-means: 0.018936634063720703

Mean-shift algorithm:

$$\hat{f}(\mathbf{x}) = \frac{1}{nh^d} \sum_{i=1}^n K\left(\frac{\mathbf{x} - \mathbf{x}_i}{h}\right).$$

Time complexity of  $O(tn^2d)$

For high dimension it is time-consuming.

K-means algorithm:

$$\arg \min_{\mathbf{S}} \sum_{i=1}^{\kappa} \sum_{\mathbf{x}_j \in S_i} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2$$

is a NP-hard problem that has complexity  $O(tnd)$

where t is iteration times, k is number of centroids, n is number of data and d is dimension.

# Performance: Who's better?

- Silhouette Score
  - Separation between different clusters.

$$s = \frac{b - a}{\max(a, b)}$$

Not for DB

- Rand Index/Adjusted Rand Index
  - Calculate difference pairs between ground-truth clusters and model clusters.

$$ARI = \frac{RI - \text{Expected RI}}{\text{Max}(RI) - \text{Expected RI}}$$

No true label

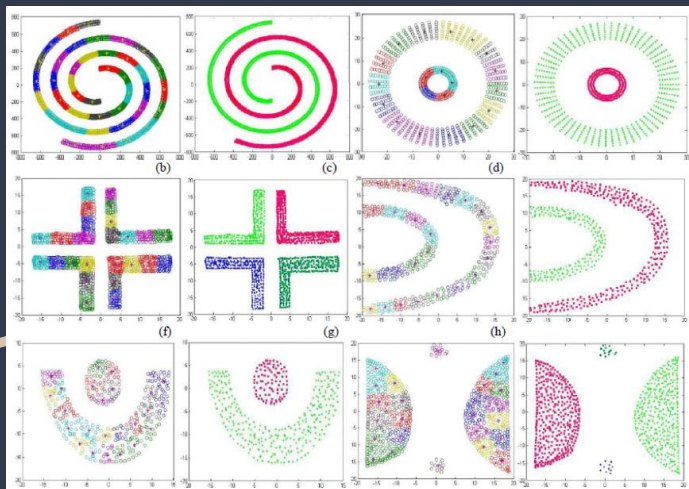
- Calinski-Harabasz Index(Variance Ratio Criterion)

$$s = \frac{\text{tr}(P_k^2) - k}{\text{tr}(W_k) - k}$$

DB has always higher score

# Performance: Who's better?

BUT Generally



- Pros of mean shift (compare to k-means):
  - Mean Shift is quite better at clustering as compared to K Means, mainly due to the fact that we don't need to specify the value of 'K', i.e. the number of clusters.
  - Output of mean shift is not strongly dependent on initialization
  - The algorithm only takes one input, the bandwidth of the window.
  - Has smaller limitation to shape of the data.
- Cons of mean shift:
  - Mean Shift performs a lot of steps, so it can be computationally expensive, with a time complexity of  $O(n^2)$  while k-means is  $O(n)$
  - The selection of the bandwidth itself can be non-trivial, although we introduce the silverman's rule.
  - The performance on high dimension would not be good since for MISE we have:

$$\int \mathbb{E} \left( (\hat{p}_{n,\text{opt}}(x) - p(x))^2 \right) dx = \inf_{h>0} \int \mathbb{E} \left( (\hat{p}_n(x) - p(x))^2 \right) dx = O \left( n^{-\frac{2}{d+4}} \right)$$



# Further Discussion

- Silverman's rule performs well on normal density. Otherwise might be bad-performance.
- Running time for Mean-shift is too long. Introducing parallel programming to be faster.
- Do not have a clear physical explanation of found: what do clusters of data patterns do?
- Performance measurement is not solved. Maybe could try Davies-Bouldin Index for non-convex clusters.
- Do not have ground-truth labels on data set. Introducing labels would be better for quantitative measurement.
- More

# Index

Questions?

- Introduction
- Materials and Methods(data and methods)
- Code Implementation
- Results(graphs by pairgrid and MDS)
- Discussion(Performance)
- References

# References

Sheth, Jash. "Mean Shift Clustering Algorithm." OpenGenus IQ: Computing Expertise & Legacy, OpenGenus IQ: Computing Expertise & Legacy, 6 Apr. 2019, <https://iq.opengenus.org/mean-shift-clustering-algorithm/>.

6.869 Advances in Computer Vision: Learning and Interfaces, <https://courses.csail.mit.edu/6.869/>.

Yufeng. "Understanding Mean Shift Clustering and Implementation with Python." Medium, Towards Data Science, 22 Feb. 2022, <https://towardsdatascience.com/understanding-mean-shift-clustering-and-implementation-with-python-6d5809a2ac40>.

Zuccarelli, Eugenio. "Performance Metrics in Machine Learning-Part 3: Clustering." Medium, Towards Data Science, 31 Jan. 2021, <https://towardsdatascience.com/performance-metrics-in-machine-learning-part-3-clustering-d69550662dc6>.

Landup, David. "Guide to Multidimensional Scaling in Python with Scikit-Learn." Stack Abuse, Stack Abuse, 21 July 2022, <https://stackabuse.com/guide-to-multidimensional-scaling-in-python-with-scikit-learn/>.

Chen, Yen-Chi. "A Tutorial on Kernel Density Estimation and Recent Advances." ArXiv.org, 12 Sept. 2017, <https://arxiv.org/abs/1704.03924>.