

## Fast, Robust and Scalable Clustering Algorithms

#### with Applications in Computer Vision

Vahan Petrosyan

Department of Automatic Control, KTH Royal Institute of Technology

Supervisor: Alexandre Proutiere, Co-supervisor: Mikael Johansson







- Introduction to Clustering Algorithms
- Motivation
- Cluster Validation
- Cluster Initialization
- Density Transformation
- Clustering Pixels in the Images
  - Application: Image Annotation





- Introduction to Clustering Algorithms
- Motivation
- Cluster Validation
- Cluster Initialization
- Density Transformation
- Clustering Pixels in the Images
  - Application: Image Annotation



- Clustering is the task of grouping a set of objects such that more "similar" object will appear in the same group
- > Applications:
  - Document Classification
  - Insurance Fraud Detection
  - Gene Expression
  - Image Segmentation
  - etc.





## **Clustering Algorithms**

- ➤ K-means
- Spectral Clustering
- Level Set Clustering
- Hierarchical Clustering
- Density-Based Clustering
- Mode-Based Clustering
- Expectation Maximization
- ➢ etc.



## **Clustering Algorithms**

- ➤ K-means
- Spectral Clustering
- Level Set Clustering
- Hierarchical Clustering
- Density-Based Clustering
- Mode-Based Clustering
- Expectation Maximization
- ➢ etc.



• Objective: 
$$J(C) = \sum_{j=1}^{k} \sum_{x_i \in c_j} ||x_i - \mu_j||^2$$

#### > Algorithm

#### **K-means**

- 1. Select an initial partition with k clusters.
- 2. Repeat steps 3 and 4 until cluster membership stabilizes.
- 3. Generate a new partition by assigning each pattern to its closest cluster center.
- 4. Compute new cluster centers.



## K-means cont.

- Parameters
  - Number of Clusters, k
  - Cluster Initialization
  - Distance Metric
- Drawbacks

4. Sensitive to outliers



#### 3. Sensitive to initialization







2. k is not known





## K-means cont.

- Parameters
  - Number of Clusters, k
  - Cluster Initialization
  - Distance Metric
- Drawbacks











## **Spectral Clustering**

#### **Spectral Clustering**

- 1. Input dataset X and select the number of clusters k.
- 2. Construct S
- 3. Compute normalized Laplacian L
- 4. Compute the first k eigenvectors  $u_1, \dots, u_k$  of L
- 5. Let  $U \in \mathbb{R}^{n \times k}$  be the matrix containing the vectors  $u_1, \dots, u_k$  as columns
- 6. Form the matrix  $T \in \mathbb{R}^{n \times k}$  from U by normalizing the rows to norm 1
- 7. Use the matrix T as an input to k-means algorithm
- 8. Set output of k-means as the cluster assignment of input data X



where 
$$S_{i,j}= e^{rac{|x_i-x_j|}{\sigma^2}}$$

١/



## **Spectral Clustering**

#### **Spectral Clustering**

- 1. Input dataset X and select the number of clusters k.
- 2. Construct S
- 3. Compute normalized Laplacian L
- 4. Compute the first k eigenvectors  $u_1, \dots, u_k$  of L
- 5. Let  $U \in \mathbb{R}^{n \times k}$  be the matrix containing the vectors  $u_1, \dots, u_k$  as columns
- 6. Form the matrix  $T \in \mathbb{R}^{n \times k}$  from U by normalizing the rows to norm 1
- 7. Use the matrix T as an input to k-means algorithm
- 8. Set output of k-means as the cluster assignment of input data X



where 
$$S_{i,j} = e^{\frac{\left|x_i - x_j\right|^2}{\sigma^2}}$$



## Spectral Clustering cont.

- Parameters
  - Same as K-means
  - σ
- Drawbacks



## 4. Sensitive to Outliers

3. Sensitive to Initialization





## Spectral Clustering cont.

- Parameters
  - Same as K-means



Drawbacks



#### Tuned $\sigma$



#### Tuning $\sigma$ is a challenging problem



## Level Set Clustering

▷ Given:  $X = \{x_1, \cdots, x_n\}$  i.i.d. sample drawn from  $p(x|\theta)$ 

► Upper Level Set 
$$L_{\lambda}(p) = \{x \in \mathbb{R}^d : p(x|\theta) \ge \lambda\}.$$

- $\succ$  Connected components of  $L_{\lambda}(p)$  form a cluster at each level  $\lambda$
- > Challenge: The estimation of  $p(x|\theta)$  is not trivial



## Level Set Clustering cont.



Visualize Clusters with Cluster Trees





#### Problems:

1. Fails when clusters overlap





#### 2. Suffers from curse of dimensionality





All clustering algorithms come with their advantages and disadvantages.



## Outline

- Introduction to Clustering Algorithms
- Motivation
- Cluster Validation
- Cluster Initialization
- Density Transformation
- Clustering Pixels in the Images
  - Application: Image Annotation



## Motivation

User Dilemma:

- How do we define the pairwise similarity/distance?
- How many clusters are present in the data?
- Does the algorithm requires multiple runs?
- Does the data contain any outliers?
- Does it suffer from curse of dimensionality?
- Does the algorithm scale well with increase in datasize?
- Is it easy to tune the parameters?



## Outline

- Introduction to Clustering Algorithms
- Motivation
- Cluster Validation
- Cluster Initialization
- Density Transformation
- Clustering Pixels in the Images
  - Application: Image Annotation



## **Determining the number of clusters**

- Choose a clustering algorithm (usually K-means is used)
- > Choose a cluster validation metric (Gap, BIC, AIC, Silhouette or other indices)



## **Determining the number of clusters**

- Choose a clustering algorithm (usually K-means is used)
- > Choose a cluster validation metric (Gap, Jump, Silhouette or other indices
- > Run the given clustering algorithm for  $k = 1, 2, ..., k_{max}$  (often K-means)
- $\succ$  For each k run the algorithm multiple times for different initialization
- > Compute the validation index for each value of k and initialization
- Choose the number of clusters by maximizing the value of a given validation index.



## Determining the number of clusters

- Choose a clustering algorithm (usually K-means is used)
- > Choose a cluster validation metric (Gap, Jump, Silhouette or other indices
- > Run the given clustering algorithm for  $k = 1, 2, ..., k_{max}$  (often K-means)
- $\succ$  For each k run the algorithm multiple times of better initialization
- $\succ$  Computer the validation index for each value of k and initialization
- Choose the number of clusters by maximizing/minimizing the value of a given validation index.



#### Run the algorithm once

#### Start each point as a separate cluster

> Alternate between Spread Virus and Suppress Virus until convergence



- 1. Randomly sample a point from the smallest cluster without replacement
- 2. Change the label of sampled point with the label of one of its NN, randomly
- 3. Repeat 1 and 2, until all the points in the dataset are sampled
- > Interpretation: Spread the epidemic and kill the small and well-connected clusters





## **Suppress Virus**

- 1. perform one iteration of K-means
- > Interpretation: Stop the spreading epidemic based on the cluster center distance







## **Viral Clustering: Demo**

#### https://www.youtube.com/watch?v=h4xYrQZn5B4





## **Viral Clustering: Simulated Data**





## **Comparison: Number of Clusters**

	Exp	Т	Beta	Mixture	Gaussian
	$\mu_k \pm \sigma_k$	$\mu_k \pm \sigma_k$	$\mu_k\pm\sigma_k$	$\mu_k\pm\sigma_k$	$\mu_k\pm\sigma_k$
KL	$20.9 \pm 11$	$16 \pm 12.6$	$18.9 \pm 11.8$	$17.8 \pm 12.6$	$23.7 \pm 5.2$
СН	31.5±5.3	$25.5 \pm 4.5$	$24.9 \pm 3.8$	$26.2 \pm 6.3$	$24{\pm}2.5$
HR	$4.9 {\pm} 5.2$	$3.2 \pm 0.6$	$3.9{\pm}2.5$	$5.5 {\pm} 5.9$	$3.6 \pm 1.5$
SL	$19.2 \pm 7.4$	$3.6 \pm 2.3$	$5.7 {\pm} 4.5$	$9.9{\pm}5.8$	$2\pm0$
Gap	$2.1 \pm 1.5$	$1.6 {\pm} 0.9$	$1.9 \pm 1.1$	$1.9 \pm 1.2$	$2.4 \pm 1.7$
BIC	$23.3 \pm 1.5$	$22.4 \pm 5.7$	$21.3 \pm 2.3$	$20.6 \pm 2.6$	$1\pm0$
Jump	$26.4 \pm 5.6$	21.6±1.9	$20.3 \pm 4.2$	$23.7 \pm 6.6$	$39.2 \pm 0.9$
VC	$20\pm0$	$19.9 {\pm} 0.1$	19.7±0.6	$19.7 \pm 0.7$	$20\pm0$

Fig. Mean and Standard deviation of predicted number of clusters out of 50 experiment



## **Comparison: Accuracy and Initialization**



Accuracy of the various clustering algorithms on the artificial datasets (100 realizations for each dataset). Left column: Adjusted Rand Index (the higher, the better) and right column: Variation of Information Index (the lower, the better).





- We proposed a technique towards solving
   2 major drawbacks of K-means clustering
  - Number of Clusters, *k*
  - Cluster Initialization
  - Distance Metric
- Drawbacks

## 1. Finds only center based clusters











## What about the following 2 cases?

Center based clusters





#### Cannot be addressed with K-means algorithm, because of its objective



## **Viral Spectral Clustering**



#### More details are discuss in Chapter 4



## Outline

- Introduction to Clustering Algorithms
- Motivation
- Cluster Validation
- Cluster Initialization
- Density Transformation
- Clustering Pixels in the Images
  - Application: Image Annotation



➢ Given: 
$$X = \{x_1, \cdots, x_n\} \in \mathbb{R}^d$$
➢ The *k*-th nearest neighbor density estimate

The k-th nearest neighbor density estimate  $f_k(x)$  is defined

$$f_k(x) \triangleq \frac{k}{n \cdot v_d \cdot r_k(x)^d},$$



 $v_d$  is the volume of the unit sphere in  $\mathcal{R}^d$ , and  $r_k(x)$  is the *k*-th nearest neighbor distance of point *x* 

> The closest mode of point  $x_i$ 

$$\mathcal{M}(x_i) = \underset{x_j \in C_i}{\operatorname{argmax}} [f_k(x_j)] = \underset{x_j \in C_i}{\operatorname{argmin}} [r_k(x_j)]$$





Define Modal Density (MDR) Ratio

$$\theta_k(x) = \frac{r_k(x)}{r_k(\mathcal{M}(x))} = \left[\frac{f_k(\mathcal{M}(x))}{f_k(x)}\right]^{1/d}$$

> Transformed density  $\tilde{f}_k(x)$ 

$$\tilde{f}_k(x) = e^{-\gamma[\theta_k(x)-1]}$$





- > Transformed density  $\tilde{f}_k(x)$  preserves the local minima and maxima of the original density
- Fransformed density  $\tilde{f}_k(x) = 1$ , for every local maximum, i.e. when  $x = \mathcal{M}(x)$ .
- > Normalizes cluster modes to one, while increasing the gap between the clusters



Level Set Clustering requires distance matrix as an input

We modify the input of the distance matrix by the following

$$\tilde{d}_k(x_i, x_j) = \frac{d(x_i, x_j)}{\mathcal{K}_k(x_i, x_j)}$$

where

$$\mathcal{K}_k(x_i, x_j) = \left[\frac{\tilde{f}_k(x_i) + \tilde{f}_k(x_j)}{2}\right] \mathbb{I}\left[i \in \mathcal{N}_k(x_j) \mid \mid j \in \mathcal{N}_k(x_i)\right]$$

Instead we used the modified distance as an input of level set clustering









## Applications of $\tilde{f}_k(x)$ : T-sne Dim. Reduction

- T-sne is a dimensionality reduction technique that takes an input the similarity matrix
- > We modify the input of the similarity matrix by the following

$$\tilde{S}_k(x_i, x_j) = S(x_i, x_j) \cdot \mathcal{K}_k(x_i, x_j)$$

where

$$\mathcal{K}_k(x_i, x_j) = \left[\frac{\tilde{f}_k(x_i) + \tilde{f}_k(x_j)}{2}\right] \mathbb{I}\left[i \in \mathcal{N}_k(x_j) \mid \mid j \in \mathcal{N}_k(x_i)\right]$$

Instead we used the modified similarities as an input of T-sne





- The transformation increases the gap between the clusters
- Clusters are less uniform when transformation is applied





- The transformation increases the gap between the clusters
- Clusters are less uniform when transformation is applied



## Summary

- > Transformed Density,  $\tilde{f}_k(x)$
- > Defining similarity or distance using  $\tilde{f}_k(x)$

$$\begin{split} \tilde{S}_k(x_i, x_j) &= S(x_i, x_j) \cdot \mathcal{K}_k(x_i, x_j), \qquad \tilde{d}_k(x_i, x_j) = \frac{d(x_i, x_j)}{\mathcal{K}_k(x_i, x_j)} \\ \text{where,} \quad \left[ \mathcal{K}_k(x_i, x_j) = \left[ \frac{\tilde{f}_k(x_i) + \tilde{f}_k(x_j)}{2} \right] \mathbb{I}\left[ i \in \mathcal{N}_k(x_j) \mid \mid j \in \mathcal{N}_k(x_i) \right] \end{split}$$

This can be viewed as semi-supervised approach for defining distance or similarity i.e. better we estimate  $\tilde{f}_k(x)$ , better the similarity or distance metric would be.



## Outline

- Introduction to Clustering Algorithms
- Motivation
- Cluster Validation
- Cluster Initialization
- Density Transformation
- Clustering Pixels in the Images
  - Application: Image Annotation



## **Clustering Pixels in the Images**

Superpixel algorithms group neighboring pixels with similar color intensities into small patches such that object boundaries are well preserved.





## **Desirable properties of Superpixel Algorithms**

- Boundary adherence
- Compactness
- Smoothness

- Size homogeneity
- ➤ Runtime
- Number of superpixels



## **Desirable properties of Superpixel Algorithms**

- Boundary adherence
- Compactness
- Smoothness

- Size homogeneity
- ➤ Runtime
- Number of superpixels

Motivation: Create an algorithm that will satisfy all the requirements above



## Setup

- > Given
  - The image  $X = \{x_1, \dots, x_n\}$
  - Image contour  $C = \{c_1, \dots, c_n\}$
  - Weighted neighborhood graph of degree 4,  $G = \{\mathcal{V}, \mathcal{E}, W\}$

> Define:

- Contour pixel similarity: 
$$w_{l,l'} = e^{-\frac{c_l + c_{l'}}{\sigma_1}}$$
, if  $(l, l') \in \mathcal{E}$ 



## **Cut and Average Linkage**

> For the given superpixels  $S_i$  and  $S_j$  define

$$\operatorname{Cut}(S_i, S_j) = \sum_{\ell \in S_i} \sum_{\ell' \in S_j} w_{\ell, \ell'}$$

Define average linkage as

$$\operatorname{Link}(S_i, S_j) = e^{-\frac{||\mu_i - \mu_j||_2}{\sigma_2}}$$

where  $\mu_i$  and  $\mu_j$  are average color intensity of  $S_i$  and  $S_j$ 



## **Superpixel Similarity**

Define the Superpixel Similarity

$$A(S_i, S_j) = \frac{\operatorname{Cut}(S_i, S_j) \times \operatorname{Link}(S_i, S_j)}{|S_i| |S_j|}$$



## **Superpixel Similarity**

#### Define the Superpixel Similarity





## **Choice of the Superpixel Similarity**

Algorithm: greedy hierarchical merge, which promotes the characteristics we discussed earlier.

$$A(S_i, S_j) = \frac{\operatorname{Cut}(S_i, S_j) \times \operatorname{Link}(S_i, S_j)}{|S_i| |S_j|}$$





## **PALC Superpixels**

Algorithm 7: PALC Superpixels 1 **Input:** image  $\mathcal{X}$ , contour  $C, k, \sigma_1, \sigma_2, h$ ; 2  $C_G = C + \operatorname{Grid}(\mathcal{X}, k, \sigma_1);$ 3  $\mathcal{S} \leftarrow \text{Watershed}(C_G);$ 4  $K \leftarrow \max(\mathcal{S});$ 5 Compute  $A = (A(S_i, S_j))_{i,j \in \{1,...,K\}}$  from  $\mathcal{X}, C_G$  and  $\mathcal{S}$  using Eq. (6.1); 6 while K > k do Choose  $(i, j) \in \operatorname{argmax} A(S_i, S_j);$ 7 if  $\frac{|S_i| + |S_j|}{|S_{(1)}|} < h$  then 8  $| S_i \leftarrow S_i \cup S_j, \mathcal{S} \leftarrow \mathcal{S} \setminus \{S_j\}$ 9 else 10  $i' \in \operatorname{argmax} A(S_{(1)}, S_{i'});$ 11  $S_{i'} \leftarrow S_{(1)} \cup S_{i'}, \mathcal{S} \leftarrow \mathcal{S} \setminus \{S_{(1)}\}$ 12 end 13 Update A; 14  $K \leftarrow K - 1$ ; 15 16 end 17 **Output:** S;





## Effect of parameter *h*

Algorithm 7: PALC Superpixels

1 **Input:** image  $\mathcal{X}$ , contour  $C, k, \sigma_1, \sigma_2$  h;

- 2  $C_G = C + \operatorname{Grid}(\mathcal{X}, k, \sigma_1);$ 3  $\mathcal{S} \leftarrow \operatorname{Watershed}(C_G);$
- 4  $K \leftarrow \max(\mathcal{S});$

5 Compute  $A = (A(S_i, S_j))_{i,j \in \{1,...,K\}}$  from  $\mathcal{X}$ ,  $C_G$  and  $\mathcal{S}$  using Eq. (6.1); 6 while K > k do

7 | Choose 
$$(i, j) \in \operatorname{argmax} A(S_i, S_j);$$

8 **if** 
$$\frac{|S_i|+|S_j|}{|S_{(1)}|} < h$$
 then

9 
$$| S_i \leftarrow S_i \cup S_j, \mathcal{S} \leftarrow \mathcal{S} \setminus \{S_j\}$$

#### 10 else

$$\begin{array}{c|c|c|c|c|c|c|} \mathbf{u} & i' \in \operatorname{argmax} A(S_{(1)}, S_{i'}); \\ \mathbf{s}_{i'} \leftarrow S_{(1)} \cup S_{i'}, \mathcal{S} \leftarrow \mathcal{S} \setminus \{S_{(1)}\} \end{array}$$

14 Update A;

15 
$$K \leftarrow K - 1;$$

#### 16 end

17 **Output:** S;

#### **Original Image**



PALC with h=10



#### PALC with h=30



PALC with h=100







## **Experiments on DAVIS dataset**





## Visual Comparison

- Preserves object boundaries
- PALC Has smooth boundaries
  - Has fast runtime

- Promotes Similar Color
- Promotes Size Homogeneity
- Promotes Compactness





## **Applications: Image Annotation**





## **Applications: Image Annotation**

- > Deep learning based computer vision algorithms are data hungry
- > Data quality often is more important than state-of-the-art algorithm
- Hence, image annotation software became a hot research topic in the computer vision community



## **Applications: Image Annotation**

- Deep learning based computer vision algorithms are data hungry
- Data quality often is more important than state-of-the-art algorithm
- Hence, image annotation software became a hot research topic in the computer vision community

Image Annotation Software based on PALC Superpixels!



## **Advantages of PALC Annotation**

- PALC runs 10x faster than some of the fastest segmentation-based technologies used in annotation.
- PALC is extremely accurate and generates non-homogeneous regions. This allows to select both large and small objects with just one click.
- > PALC allows to **change** the number of segments **instantly**.
- PALC has a manual correction feature for the regions that need further editing.
- PALC also has a self-learning feature (i.e. the segmentation quality improves with more and more annotated data).





Chapter 3





#### Chapter 3,4 Initialization



Chapter 4 Finds only center

based clusters



Chapter 5

Dimensionality reduction



Chapter 5 Density estimation and transformation



Chapter 6 Clustering pixels in the images



Chapter 7

Applications





## **Future Work**

- **1. Clustering**: combine various components in Chapter 3-6 towards designing an algorithm with the following features:
  - Probabilistic outcome
  - Scale easily to millions of datapoints
  - Not sensitive to curse of dimensionality
  - Not sensitive to initialization
  - Not sensitive to outliers
  - Automatically adjust the number of clusters
  - Have easy interpretable and tunable parameters

Developed the code and already tested. Works like a charm ©



## Future Work cont.

Video Annotation: combine various components in Chapter 6-7 towards designing a pixel accurate video annotation software

# Thank you for your attention.