

CS109 – Data Science

Verena Kaynig-Fittkau

vkaynig@seas.harvard.edu

staff@cs109.org

AWS Clusters

- New and updated instructions for Spark 1.5 are on Piazza:

<https://piazza.com/class/icf0cypdc3243c?cid=1369>

Avoid Unnecessary Charges!

- Look at AWS console > Services > EMR
- There should be some terminated clusters there
- Check the region on the top right corner
- Make sure to change it to US East

<https://piazza.com/class/icf0cypdc3243c?cid=1256>

Region Setting in AWS

The screenshot shows the AWS CloudWatch Metrics console. At the top, there is a navigation bar with the user name "Verena Kaynig-Fittkau", the region "N. Virginia", and a "Support" link. Below the navigation bar, there is a sidebar on the left containing links for "Internet of Things", "AWS IoT BETA" (with a description "Connect Devices to the cloud"), "Mobile Services" (with links to "Mobile Hub BETA", "Cognito", "Device Farm", "Mobile Analytics", and "SNS"), and "Logs". On the right side, there is a main content area with a sidebar titled "Resources" and a "Create New" button. A dropdown menu is open over the "N. Virginia" link, listing various AWS regions:

- US East (N. Virginia) (selected)
- US West (Oregon)
- US West (N. California)
- EU (Ireland)
- EU (Frankfurt)
- Asia Pacific (Singapore)
- Asia Pacific (Tokyo)
- Asia Pacific (Sydney)
- South America (São Paulo)

At the bottom of the dropdown menu, there is a link "learn more about AWS".

Announcements

- Final project
 - Team assignments have been posted to piazza
 - Make sure you are in a 3-4 person team
 - Try and date on the piazza thread
 - If you have problems write to staff@cs109.org
 - Project proposals are due on Thursday
<https://piazza.com/class/icf0cypdc3243c?cid=1317>

Final Project Proposal

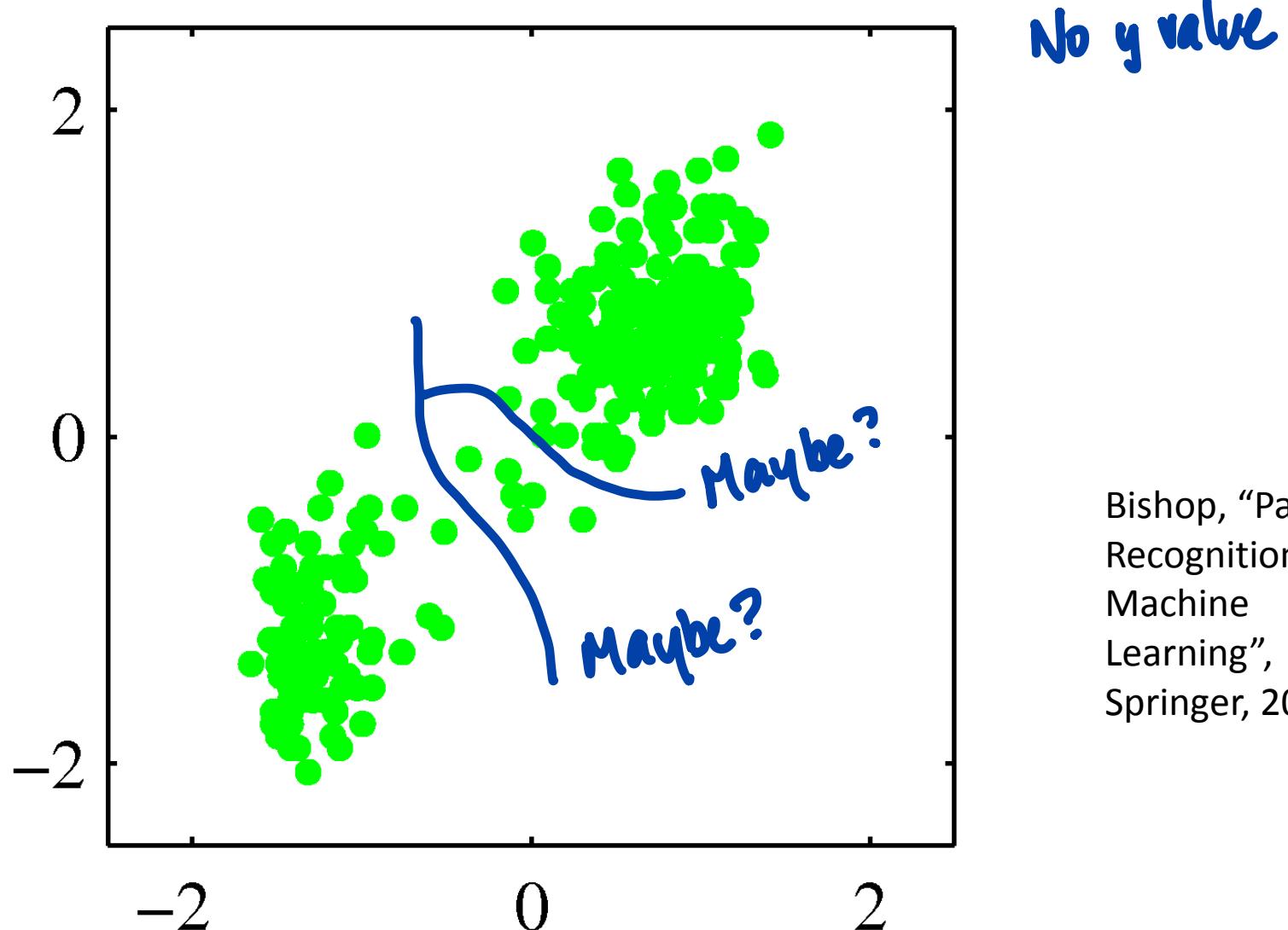
- Submit just **one form per team**.
- Do it as **early as possible!**
- No project approval until you meet your TF

<https://piazza.com/class/icf0cypdc3243c?cid=1317>

Supervised vs. Unsupervised

- We mainly talked about supervised learning so far
- Joe already moved to unsupervised with LDA
- In these settings we have **no labels** in our training data.

Unsupervised Setting



Unsupervised Learning

- Find patterns in unlabeled data
- Sometimes used for a supervised setting in which labels are hard to get
- Can identify new patterns that you were not aware of.

→ find existing
and new
patterns

Clustering Applications

- Google image search categories
- Author Clustering:
<http://academic.research.microsoft.com/VisualExplorer#1048044>
- Opening a new location for a hospital, police station, etc. → optimizing:
Can often be unsupervised or supervised.
- Outlier detection

Above to general list of pictures

Unsupervised Learning

- K-means
- Mean-shift
- Hierarchical Clustering
- Rand index, stability } *is it a good solution*

K-means – Algorithm

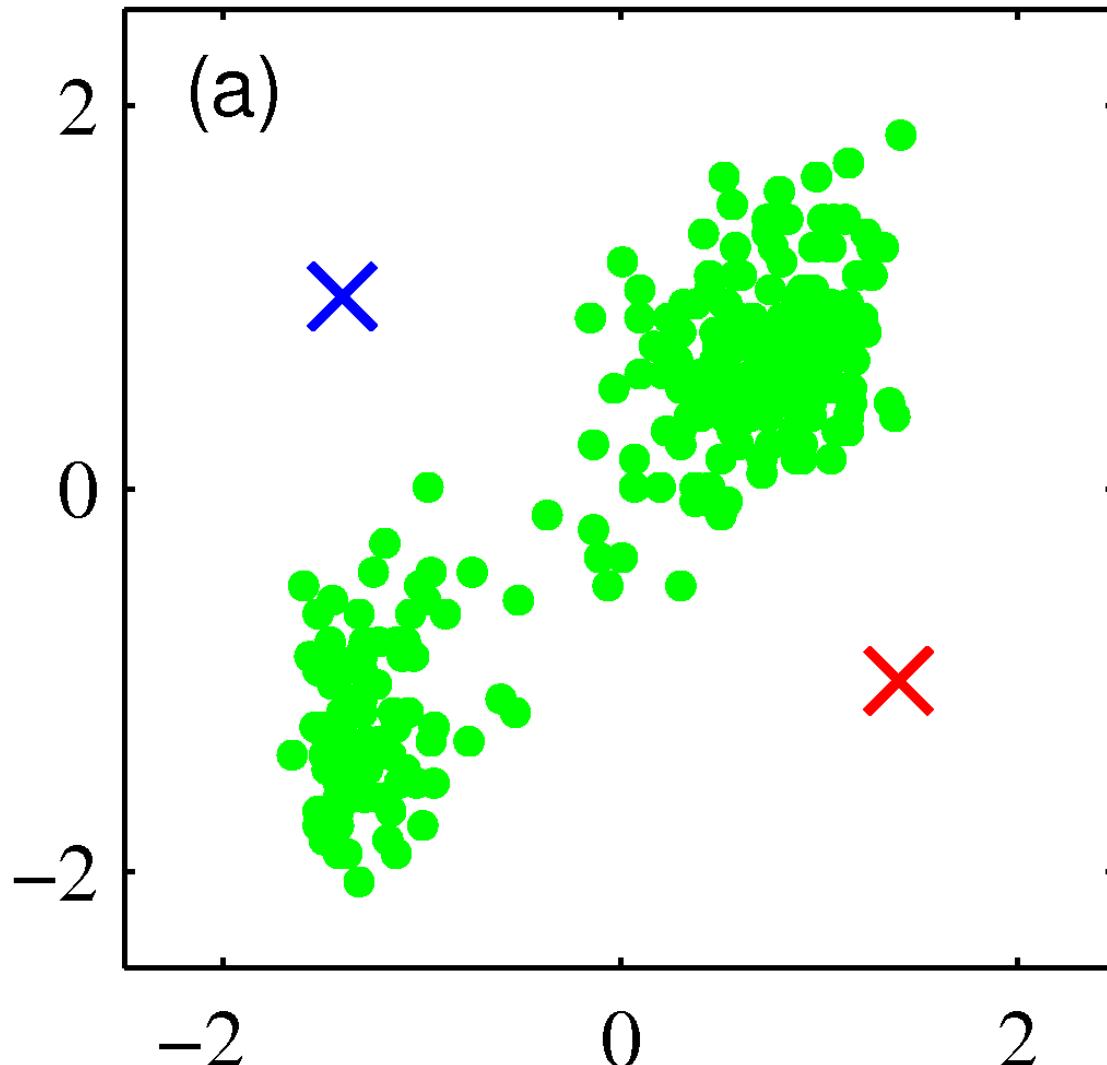
(popular 1st choice)

- Initialization:

→ Not related to kNN

- choose k random positions
- assign cluster centers $\mu^{(j)}$ to these positions

K-means



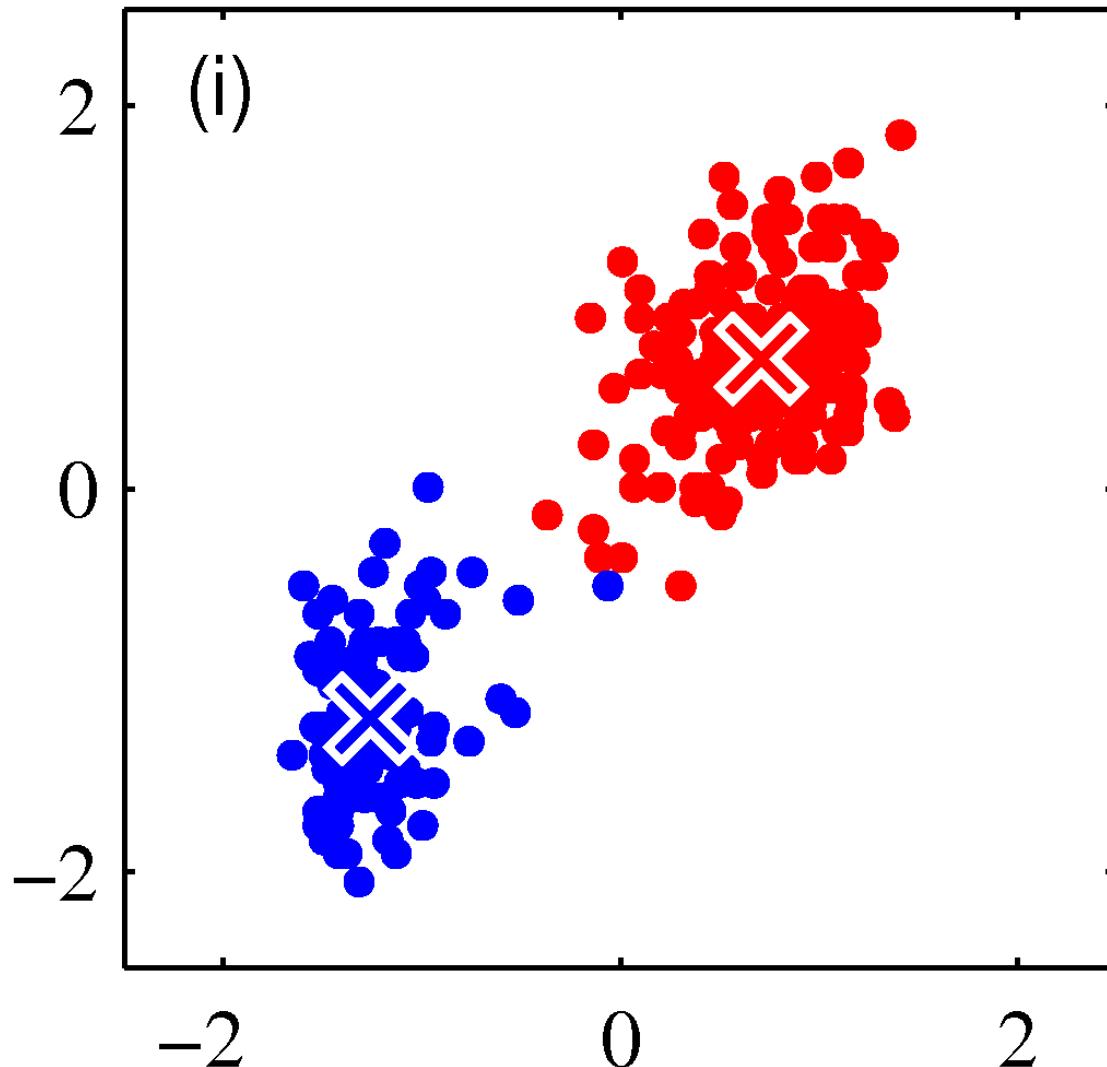
Bishop, "Pattern
Recognition and
Machine
Learning",
Springer, 2006

K-means

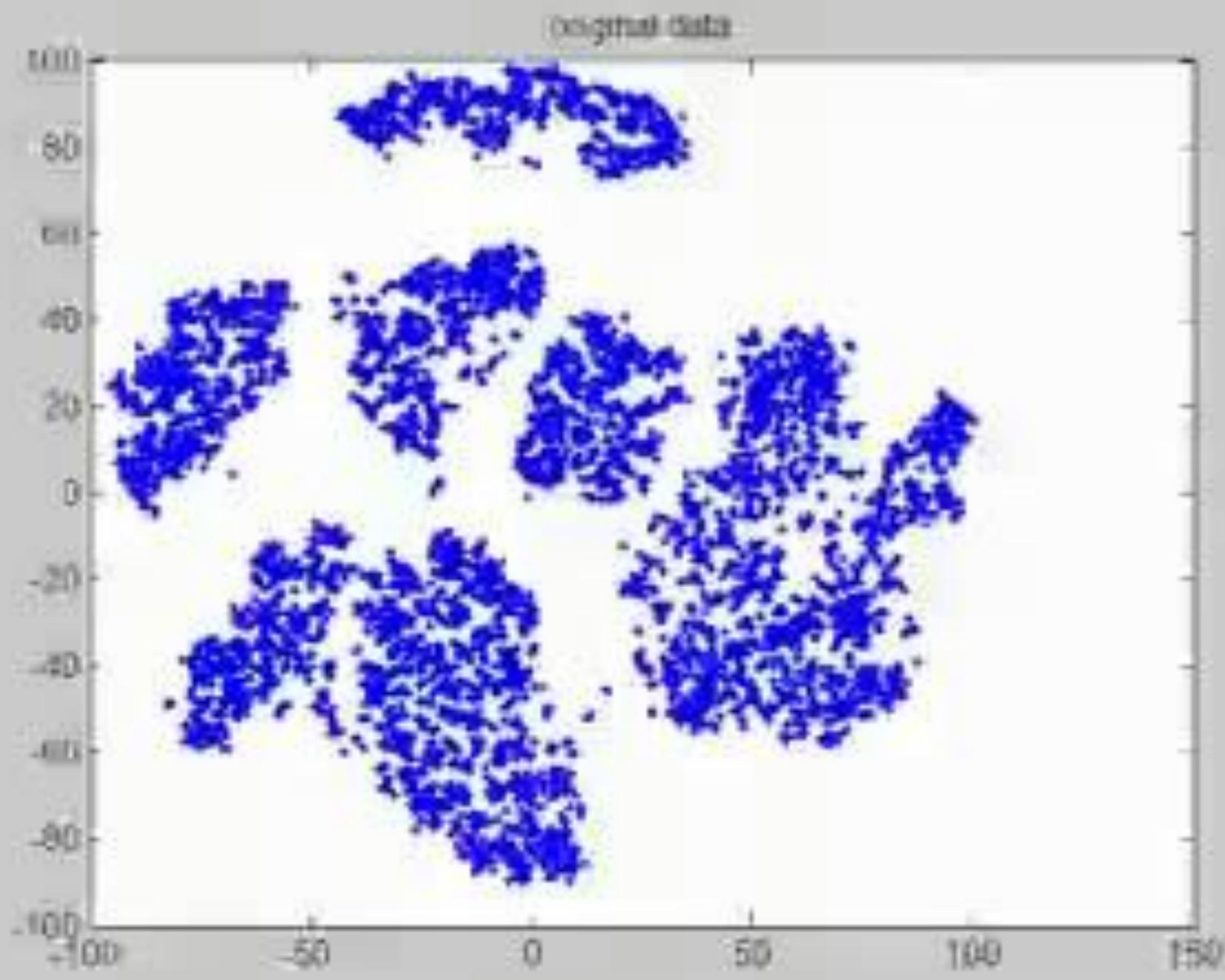
- Until Convergence:
 - Compute distances $\|x^{(i)} - \mu^{(j)}\|$
 - Assign points to nearest cluster center
 - Update Cluster centers:

$$\mu^{(j)} = \frac{1}{N_j} \sum_{x_i \in C_j} x_i$$

K-means



Bishop, "Pattern
Recognition and
Machine
Learning",
Springer, 2006



K-means Example



R

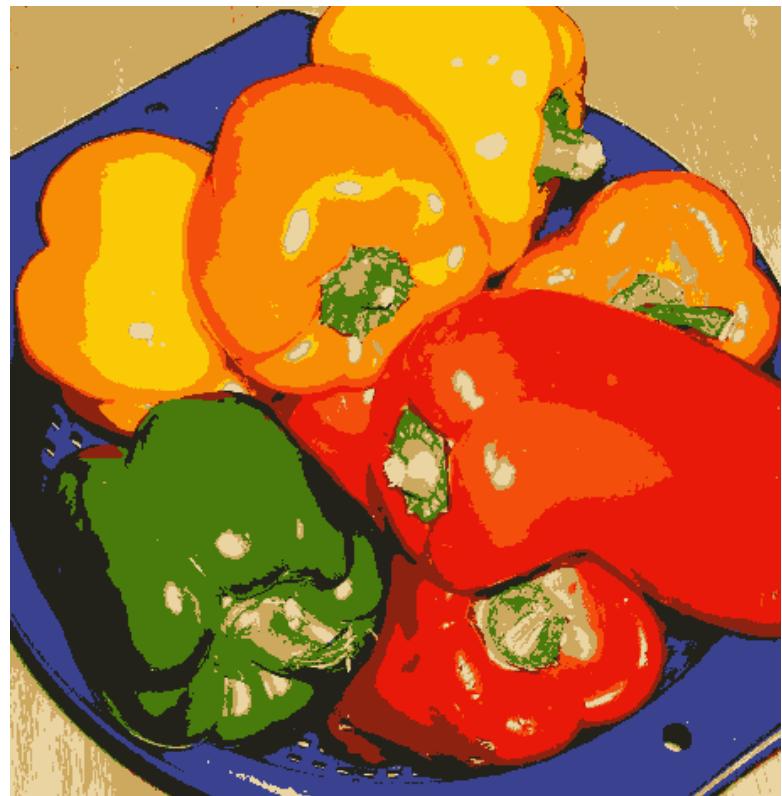


G



B

K-means Example



K-means Example

(goal is 'image compression')



Same
 $k = 10$

K-means Summary

- Guaranteed to converge
- Result depends on initialization (which might be random).
- Number of clusters is important
↳ You have to specify. Binary is straightforward.
- Sensitive to outliers
 - Use median instead of mean for updates

Initialization Methods

- Random Positions
- Random data points as Centers
- Random Cluster assignment to data points
- Start several times *-> see if one answer comes up multiple times.*

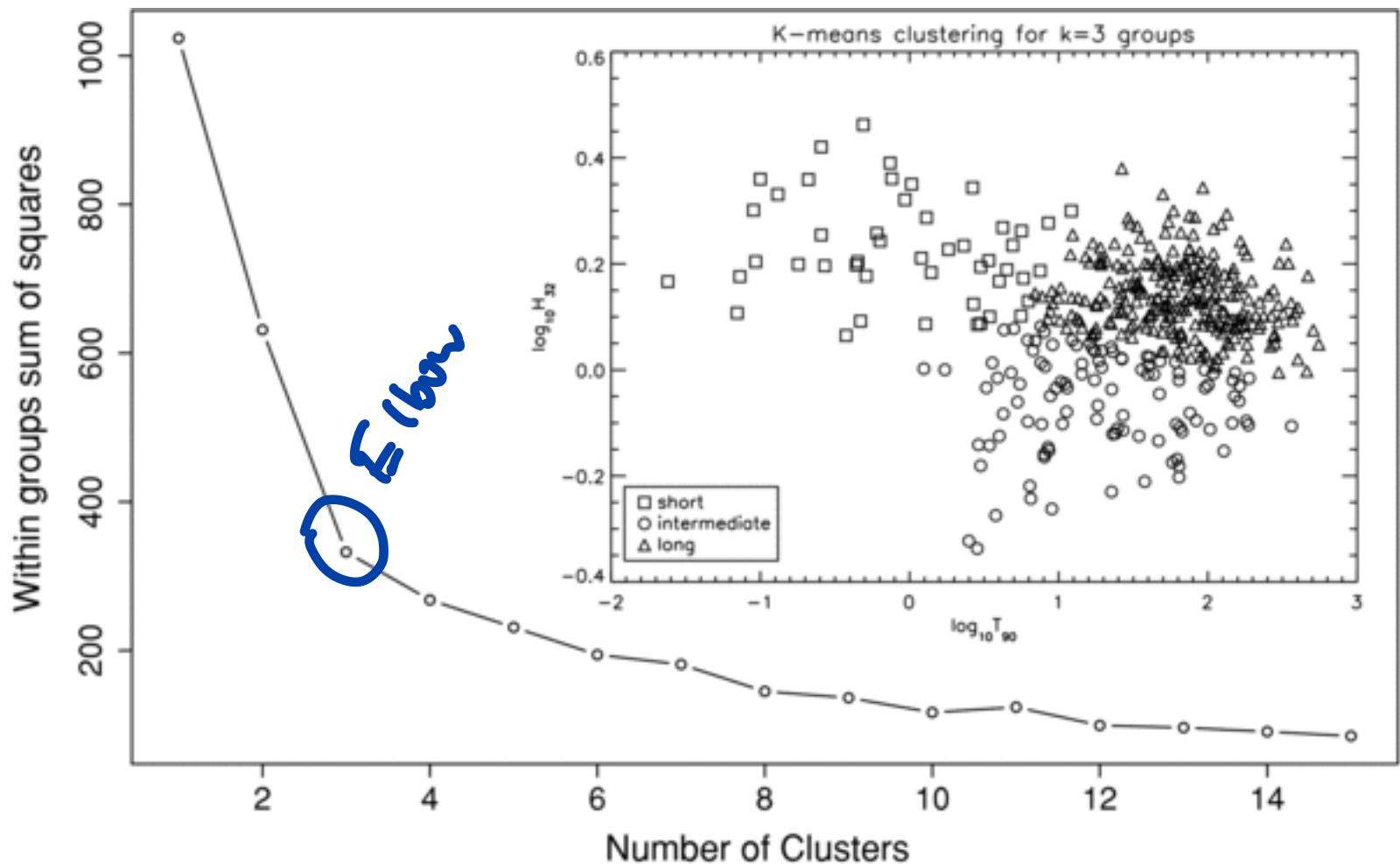
How to find K

- Extreme cases:
 - $K=1 \rightarrow$ likely won't explain.
 - $K=N$
- Choose K such that increasing it does not model the data much better.

\hookrightarrow will start to split clusters if K is too high

\hookrightarrow also is overfitting

“Knee” or “Elbow” method



Cross Validation

- Use this if you want to apply your clustering solution to new unseen data
- Partition data into n folds
- Cluster on n-1 folds
- Compute sum of squared distances to centroids for validation set

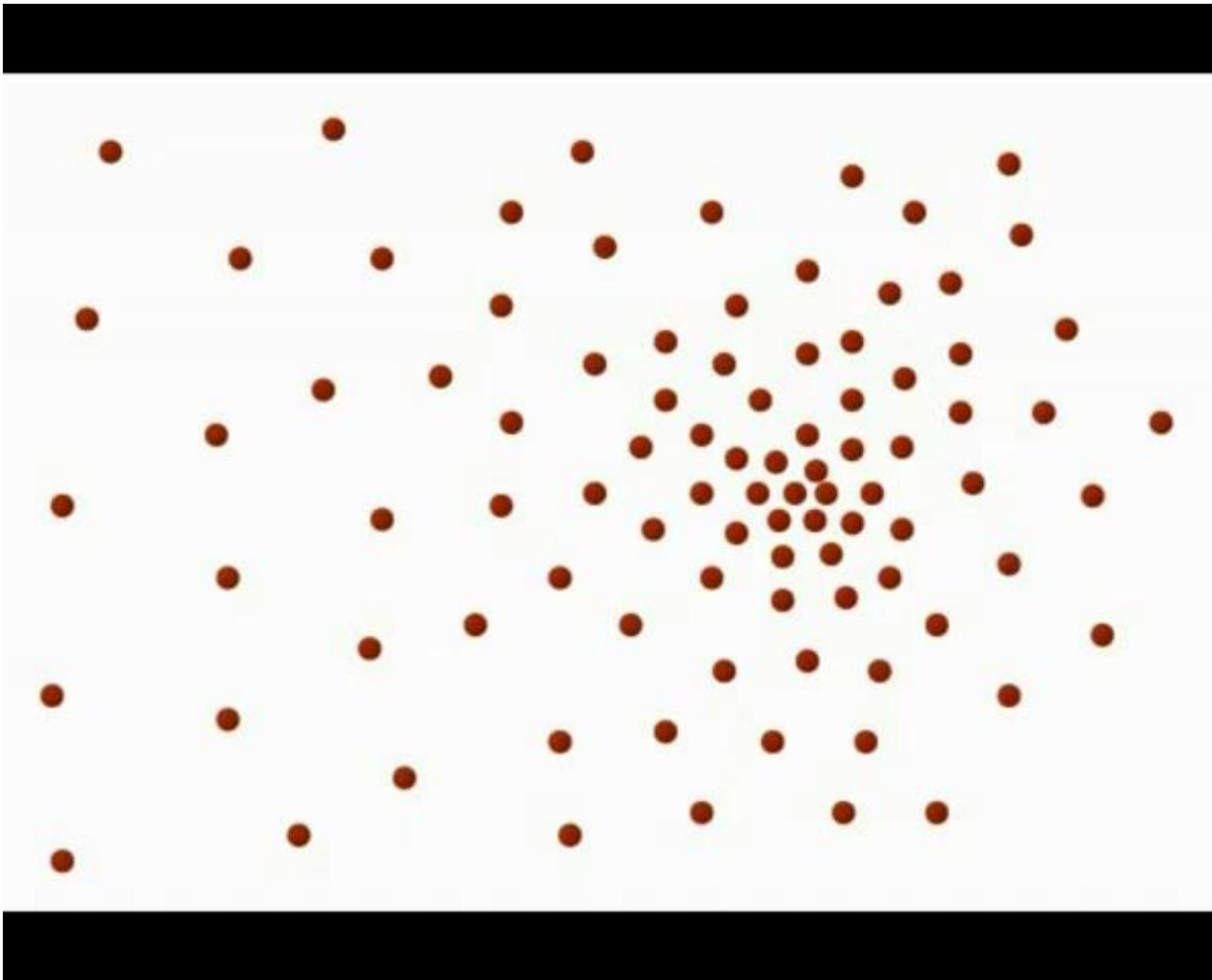
Getting Rid of K

- Having to specify K is annoying
- Can we do without?

Mean Shift

1. Put a window around each point
2. Compute mean of points in the frame.
3. Shift the window to the mean
4. Repeat until convergence

Mean Shift

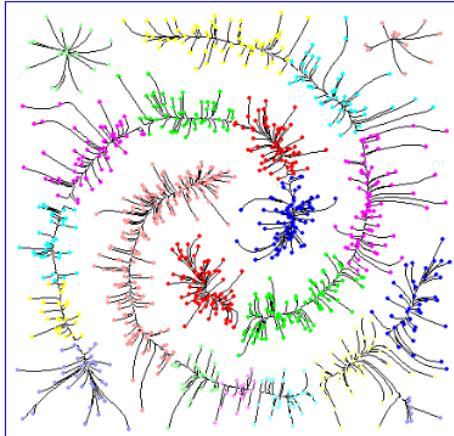


<http://www.youtube.com/watch?v=kmQAsotT9s>

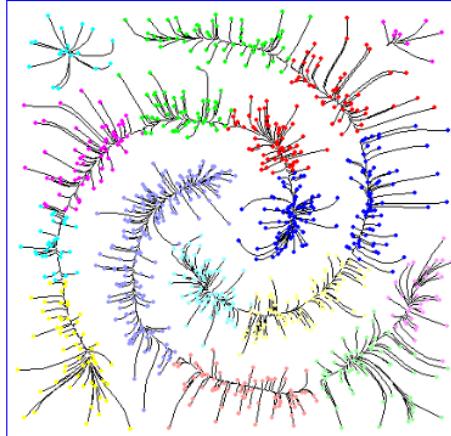
Mean Shift

$k=19$

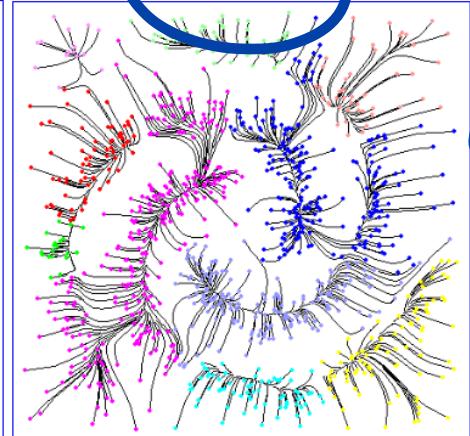
$\sigma = 0.6$



$\sigma = 0.7$

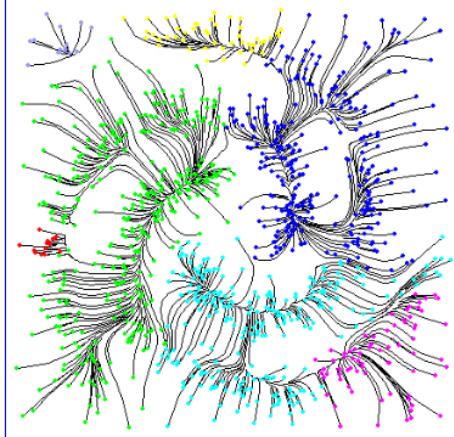


$\sigma = 0.8$

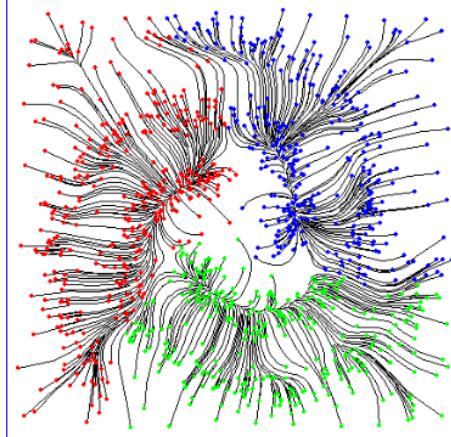


— window size;
of clusters
not specified.

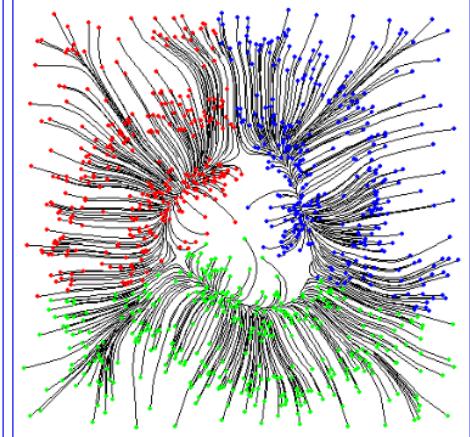
$\sigma = 0.83$



$\sigma = 0.9$



$\sigma = 1.0$



$k=3$

Mean Shift Summary

- Does not need to know number of clusters
- Can handle arbitrary shaped clusters
- Robust to initialization *Removes choosing a random K value*
- Needs bandwidth parameter (window size)
- Computationally expensive
↳ Blc done fr every single data point. *In KMeans.*
- Very good article:

<http://saravananthirumuruganathan.wordpress.com/2010/04/01/introduction-to-mean-shift-algorithm/>

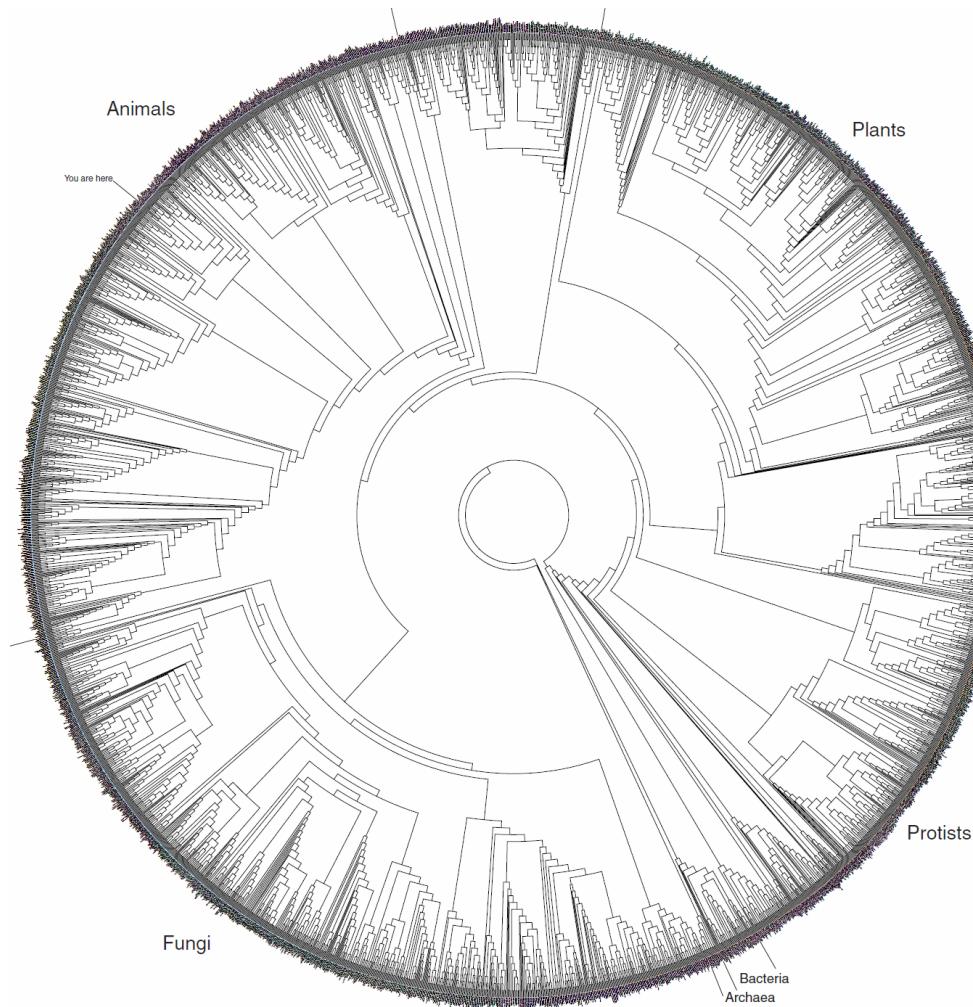
Multi-feature object trajectory clustering for video analysis

Nadeem Anjum Andrea Cavallaro

Parameters parameters

- For K means we need K and result depends on initialization
- For mean shift we need the window size and a lot of computation
- Hierarchical Clustering keeps a history of all possible cluster assignments

Tree of Life



<http://www.zo.utexas.edu/faculty/antisense/DownloadfilesToL.html>

Hierarchical Clustering

Each point is its own cluster to start.



Hierarchical Clustering



Hierarchical Clustering



Hierarchical Clustering

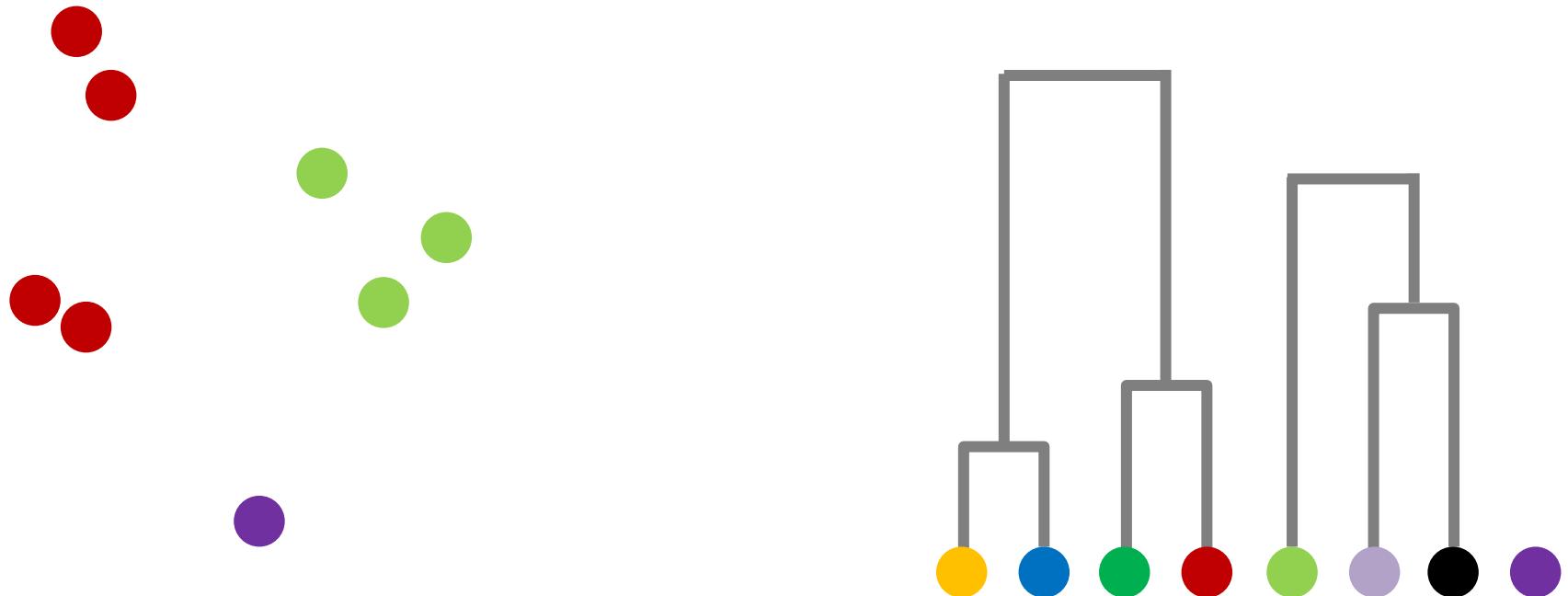


Hierarchical Clustering

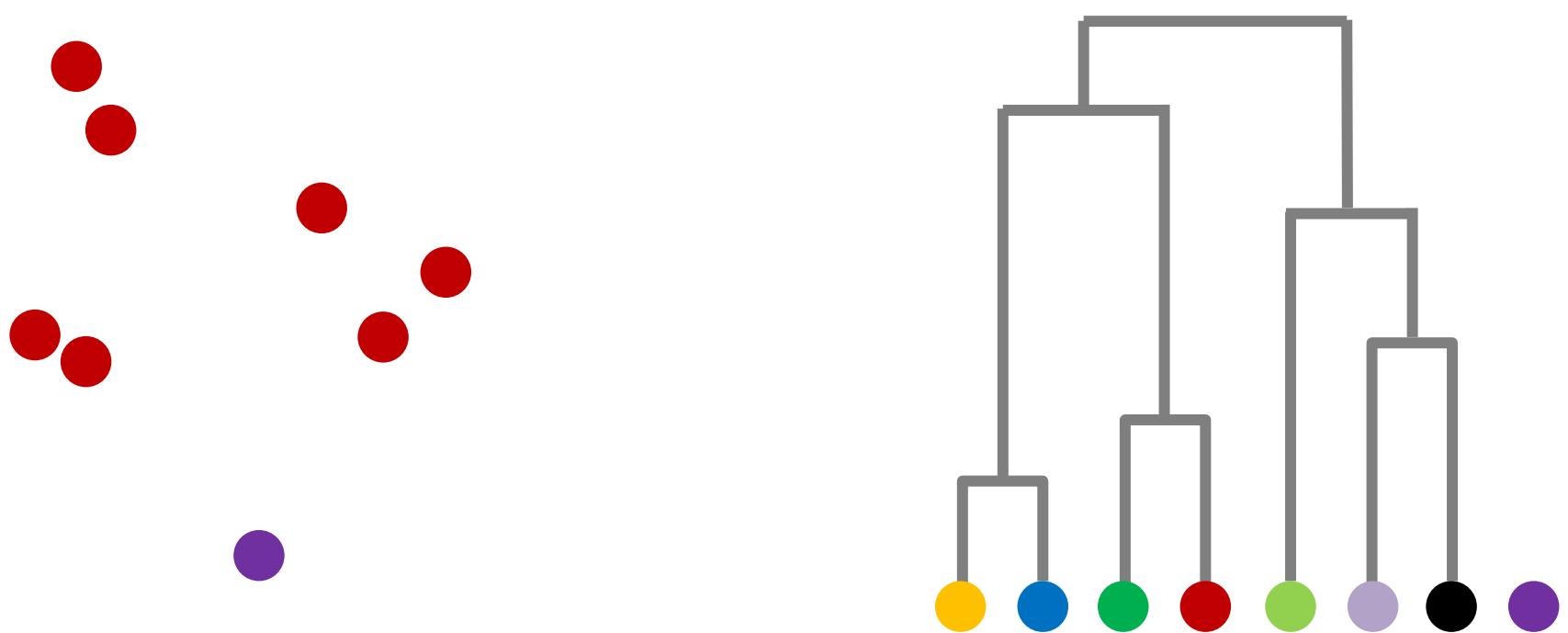
Growing cluster



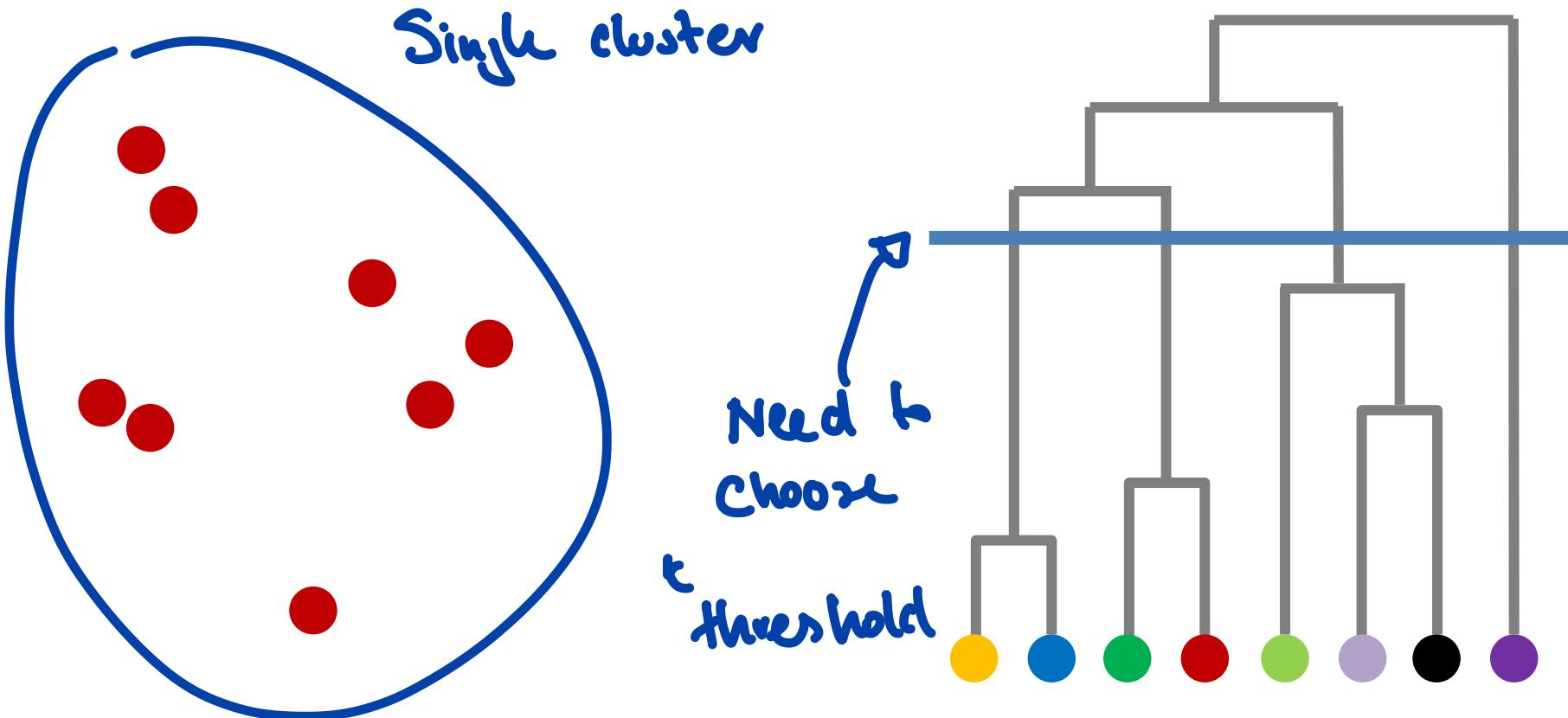
Hierarchical Clustering



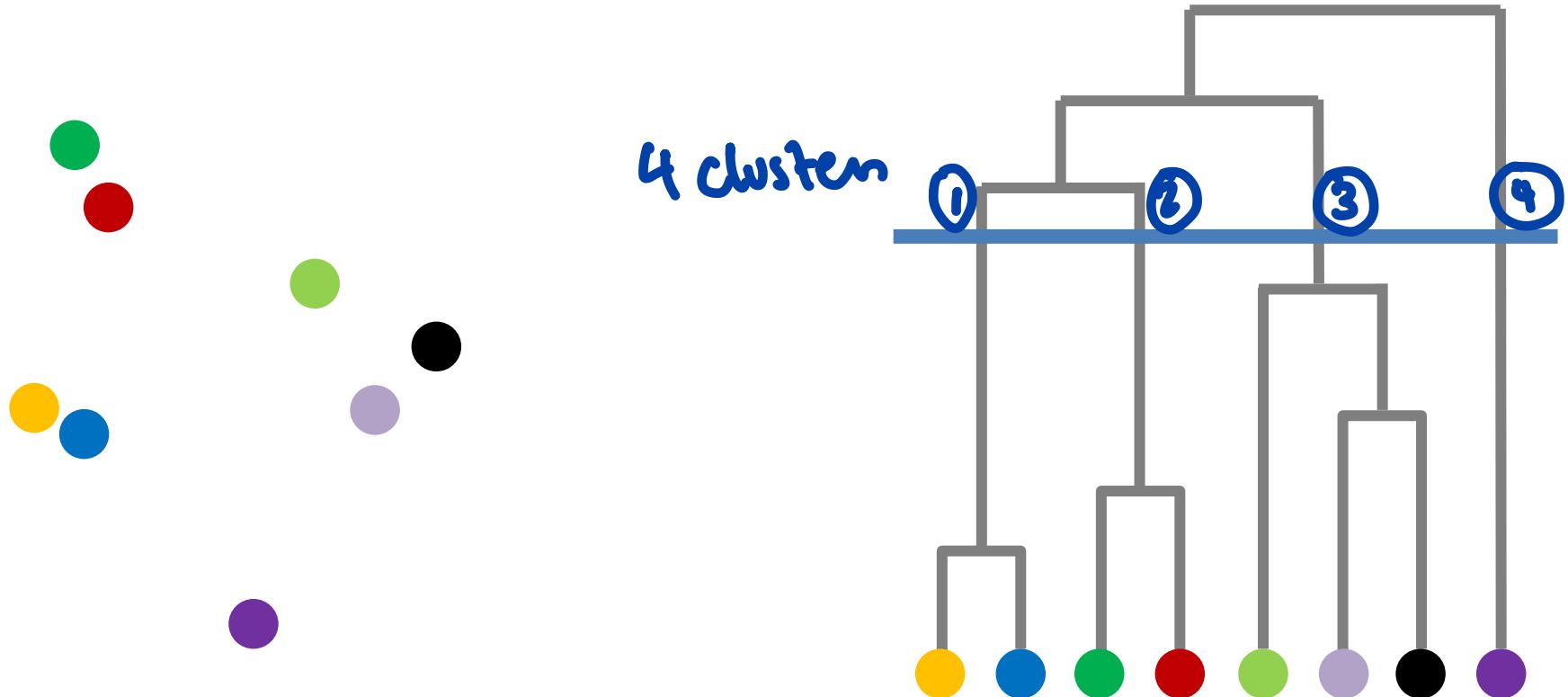
Hierarchical Clustering



Hierarchical Clustering



Hierarchical Clustering



Hierarchical Clustering

- Produces complete structure

Can be useful to have no threshold.

- No predefined number of clusters

↳ like a tree of life.

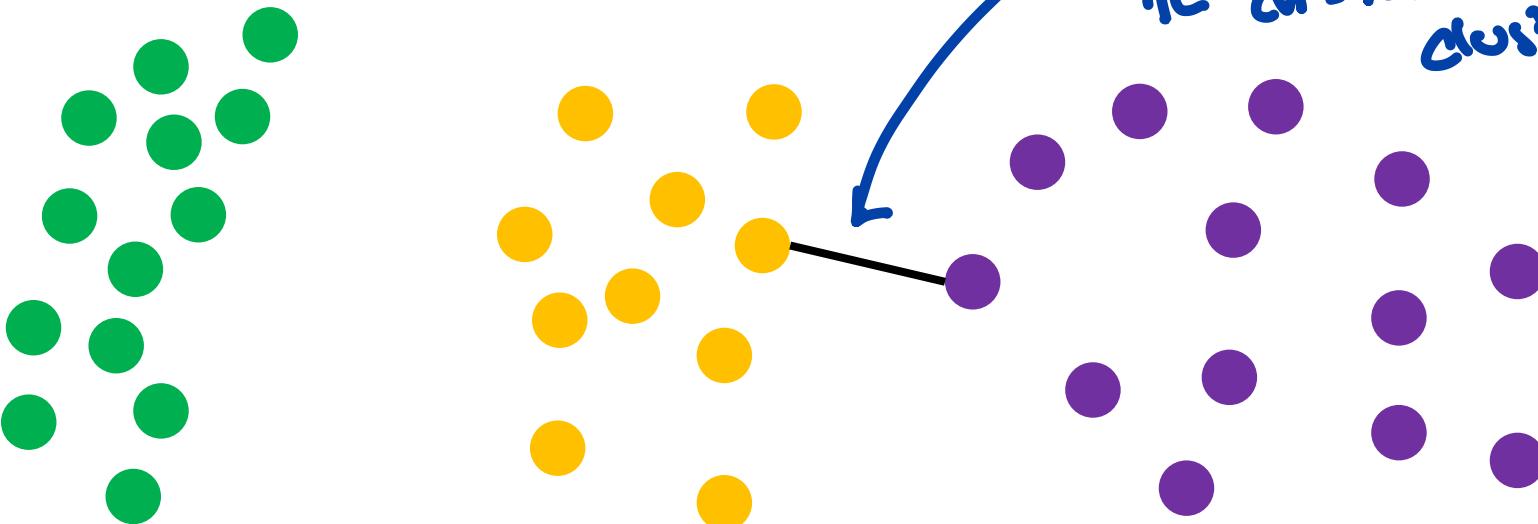
- Similarity between clusters:

- single-linkage: $\min\{d(x,y) : x \in \mathcal{A}, y \in \mathcal{B}\}$ (As seen w/ dots)

- complete-linkage: $\max\{d(x,y) : x \in \mathcal{A}, y \in \mathcal{B}\}$

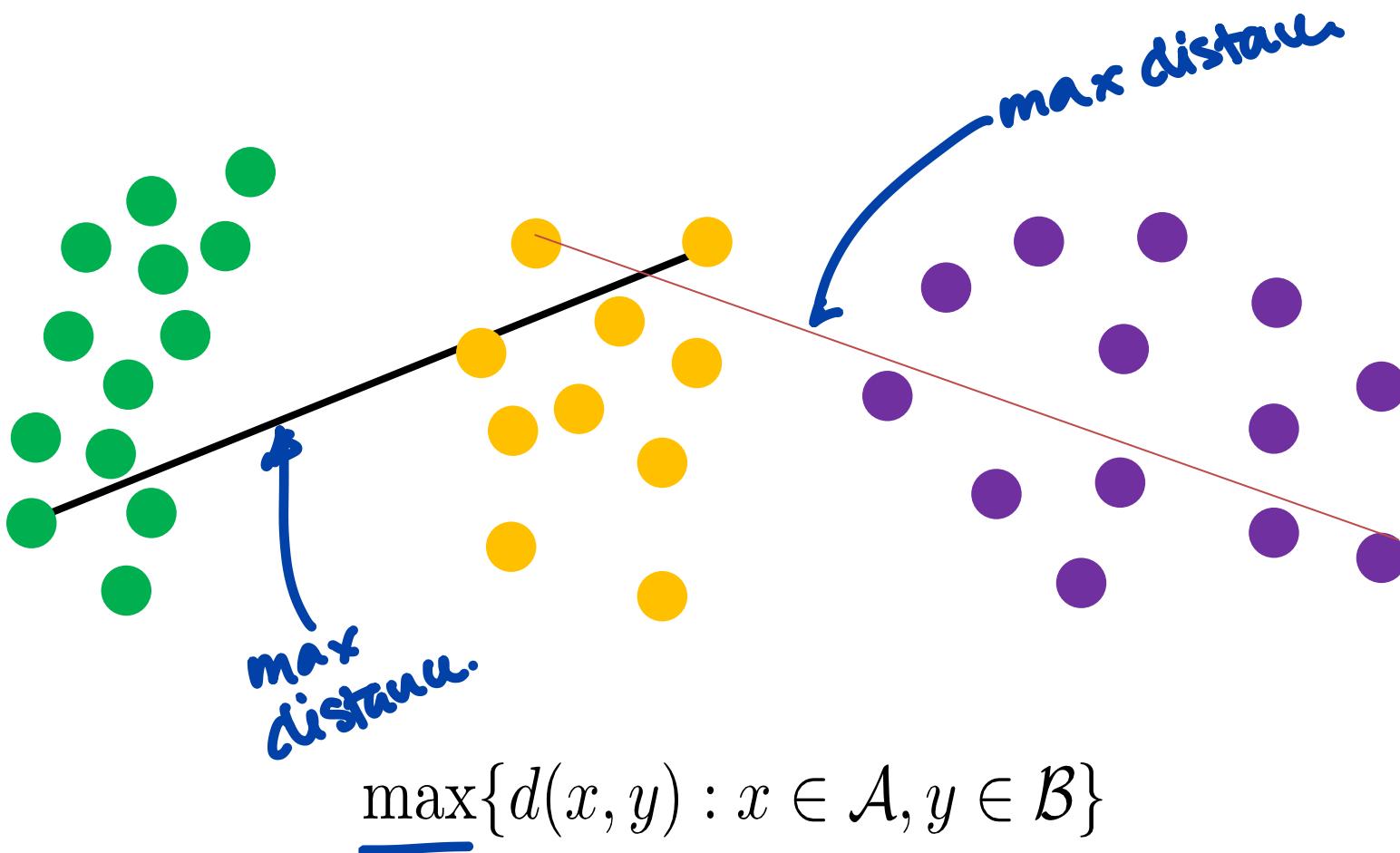
- average linkage: $\frac{1}{|\mathcal{A}| \cdot |\mathcal{B}|} \sum_{x \in \mathcal{A}} \sum_{y \in \mathcal{B}} d(x, y)$

Single Linkage



$$\min\{d(x,y) : x \in \mathcal{A}, y \in \mathcal{B}\}$$

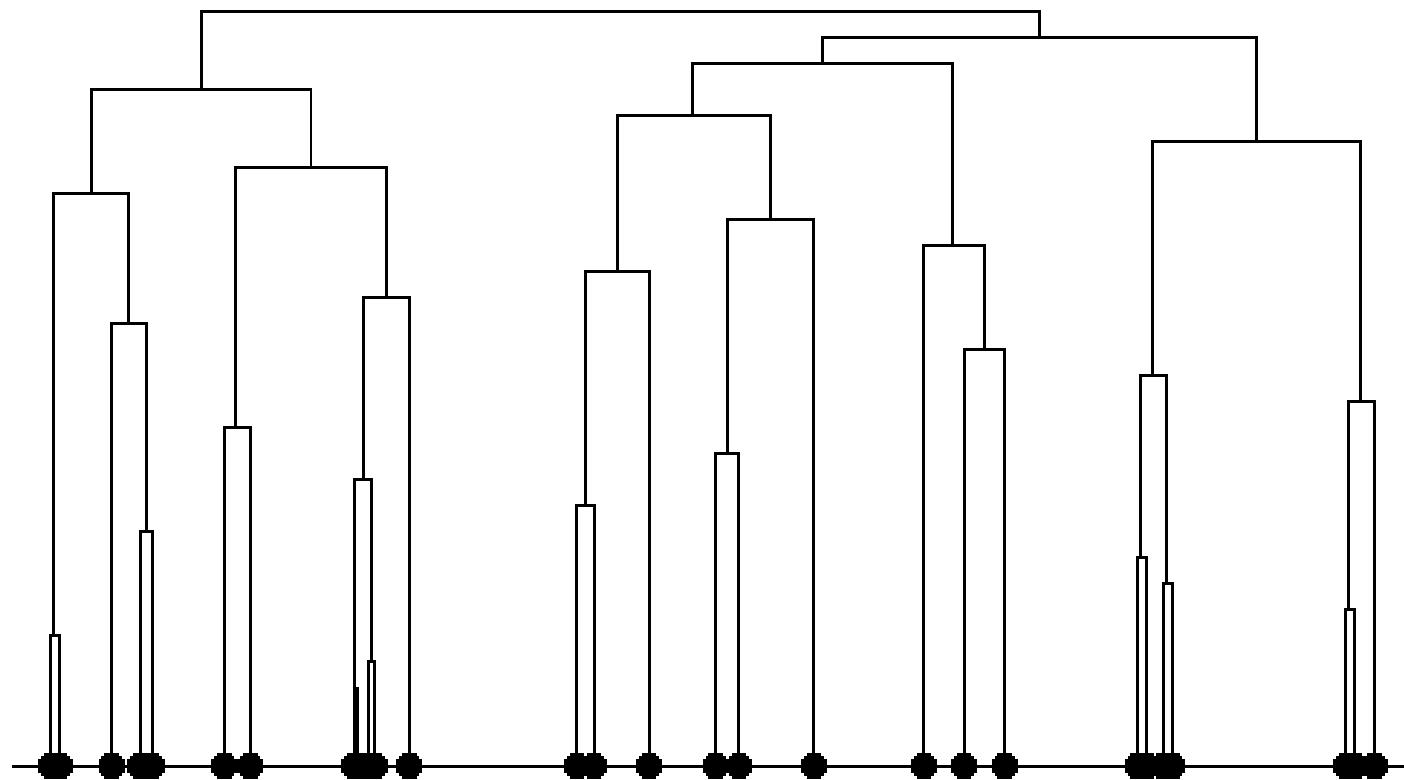
Complete Linkage



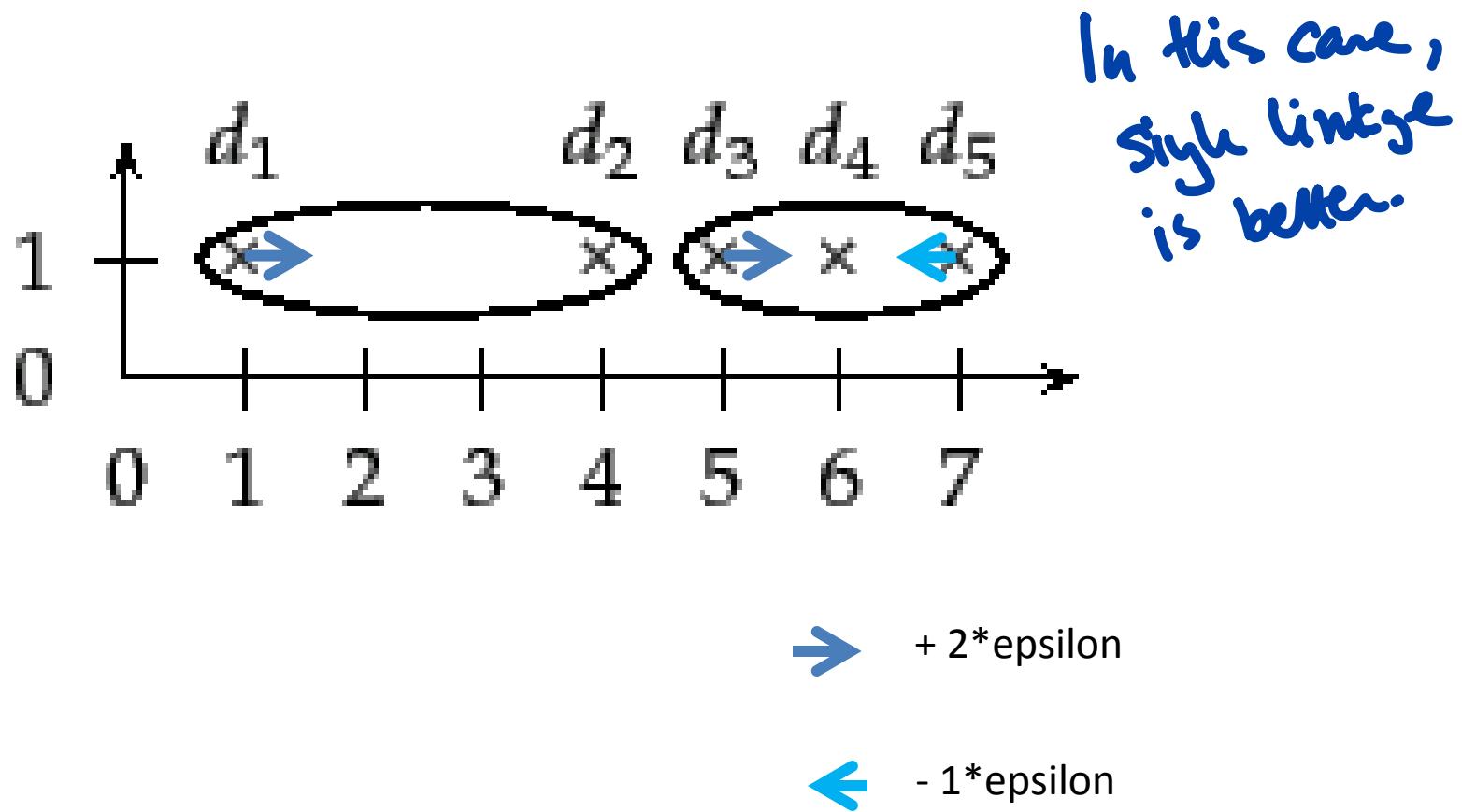
Linkage Matters

- Single linkage: tendency to form long chains
- Complete linkage: Sensitive to outliers (*Addressee chaining*)
- Average-link: Trying to compromise between the two

Chaining Phenomenon



Outlier Sensitivity



Efficient Hierarchical Graph-Based Video Segmentation

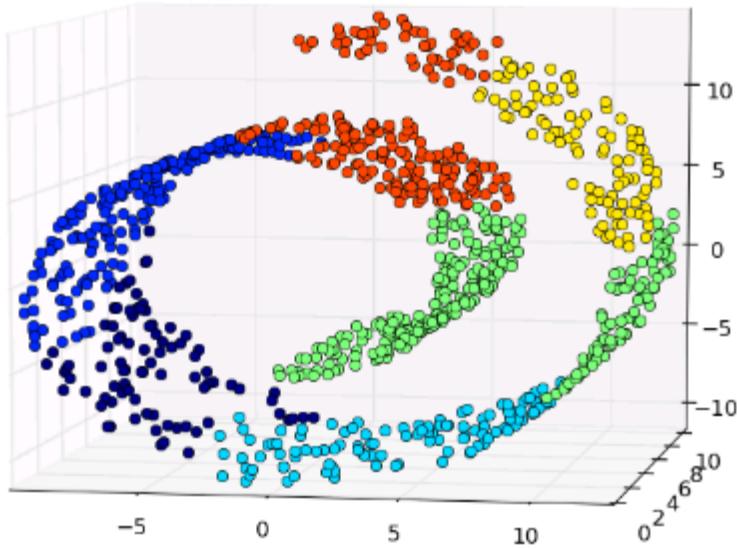
Matthias Grundmann^{1,2}, Vivek Kwatra²,
Mei Han² and Irfan Essa¹

¹Georgia Tech ²Google Research

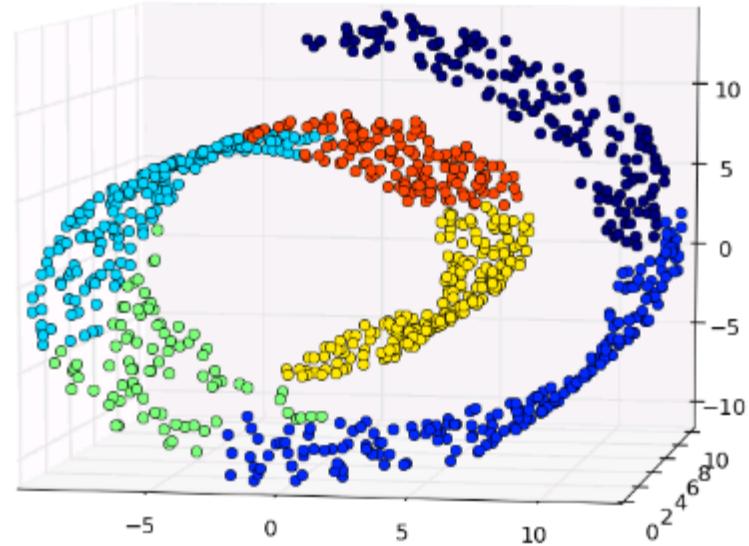
IEEE CVPR, San Francisco, USA, June 2010



Swiss Role Problem



without connectivity
constraints



with connectivity
constraints

only adjacent clusters can be merged together

Evaluation Criteria

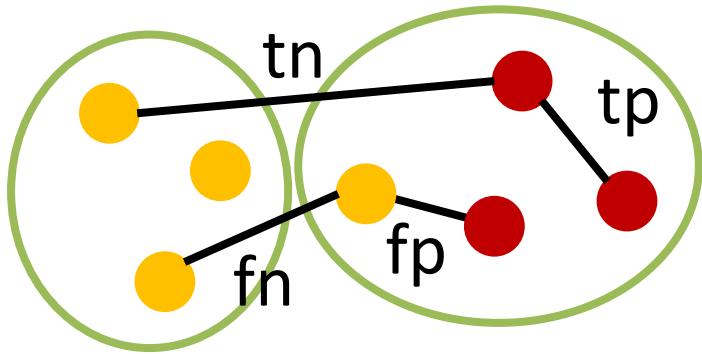
*Very hard
for
unsupervised
in general.*

- Based on expert knowledge
- Debatable for real data
- Hidden Unknown structures could be present
- Do we even want to just reproduce known structure?

Rand Index

- Percentage of correct classifications
- Compare pairs of elements:

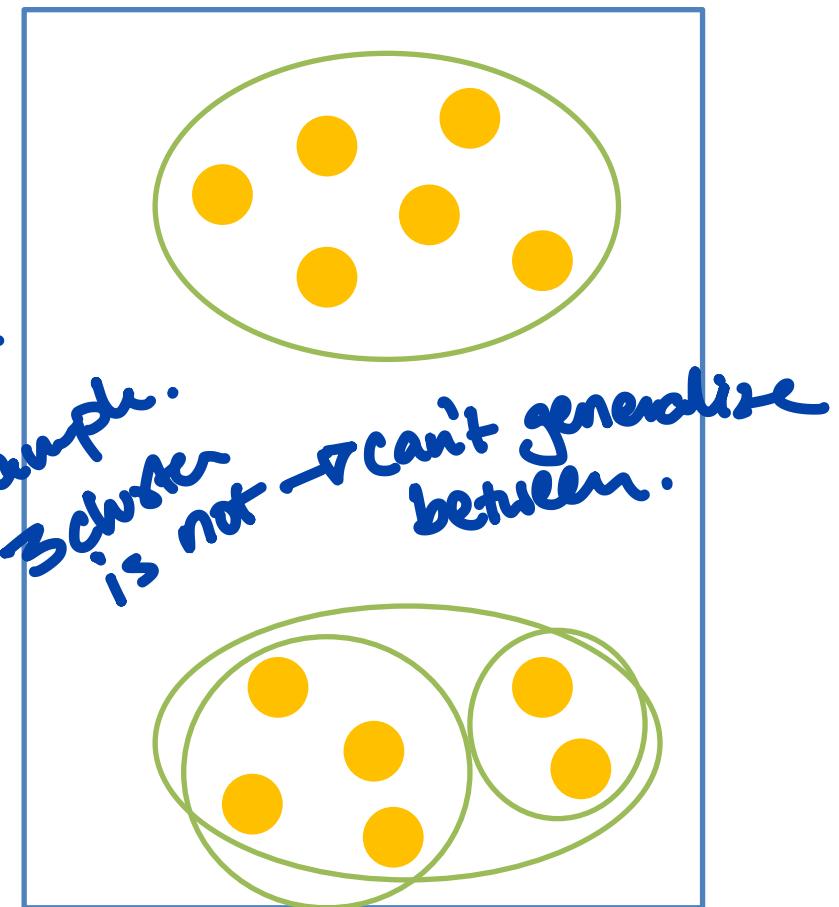
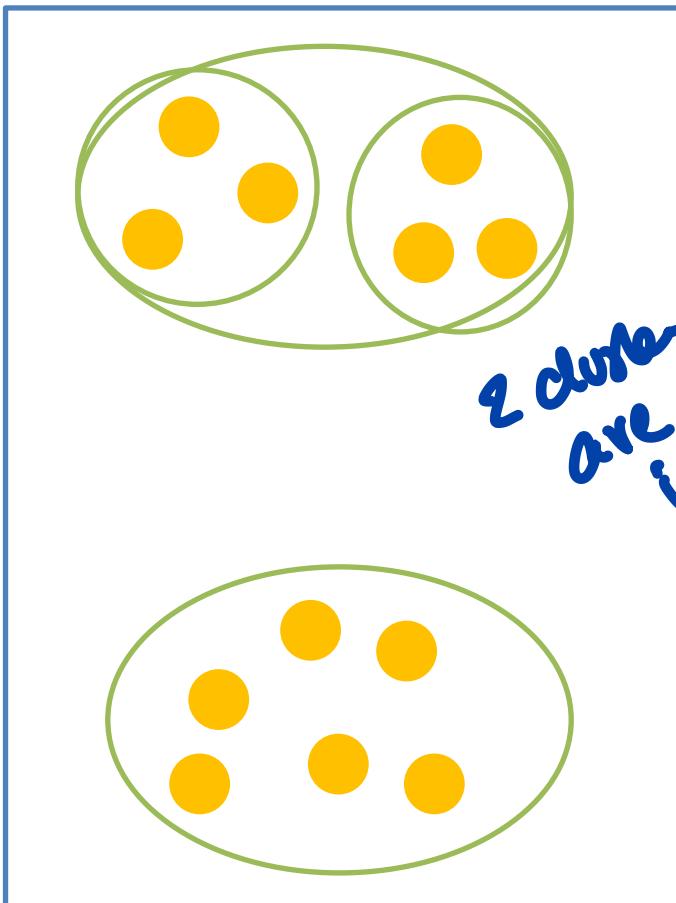
$$R = \frac{tp+tn}{tp+tn+fp+fn}$$



- Fp and fn are equally weighted

Stability

A good clustering solution should generalize



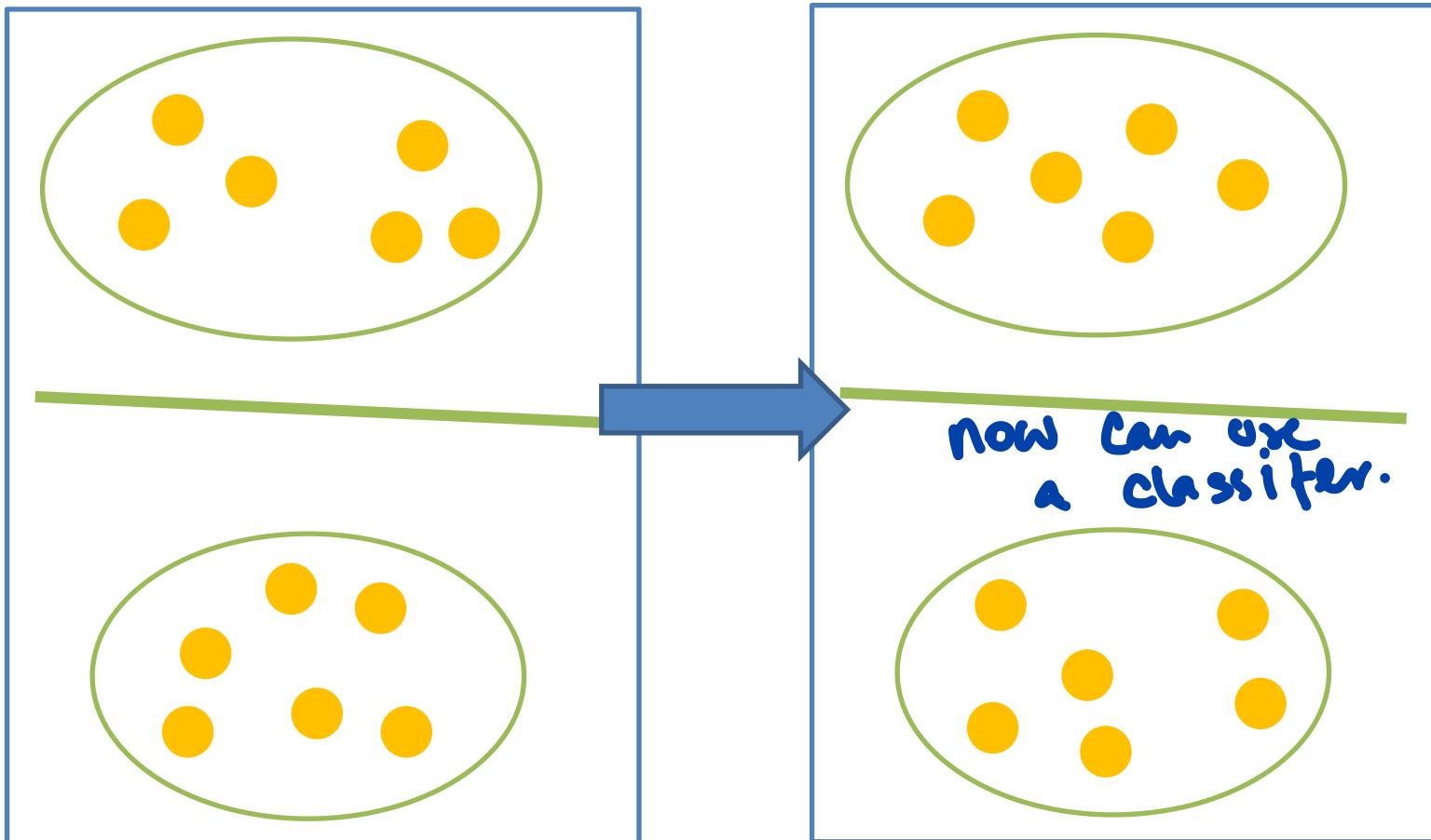
Stability

- What is the right number of clusters?
- What makes a good clustering solution?
- Clustering should generalize!

Can turn into
a supervised set.
if stable, the
error should
be small.

Use clustering to
turn an unsupervised
to supervised.

Stability



Summary

- We have covered a lot today
- Clustering
 - K-means
 - Mean-shift
 - Hierarchical clustering
- Evaluation criteria
 - Rand index
 - Stability