

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 "Regions" section. The "Regions" section lists various AWS regions, each with a small icon and a description. A dropdown menu is open over the "US East (N. Virginia)" entry, showing a list of other regions: "US West (Oregon)", "US West (N. California)", "EU (Ireland)", "EU (Frankfurt)", "Asia Pacific (Singapore)", "Asia Pacific (Tokyo)", "Asia Pacific (Sydney)", and "South America (São Paulo)". The "Create New Region" button is visible at the bottom of the dropdown menu.

Verena Kaynig-Fittkau | N. Virginia | Support

Resources

Regions

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

Create New Region

Addition

Getting Started

Read ou

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>

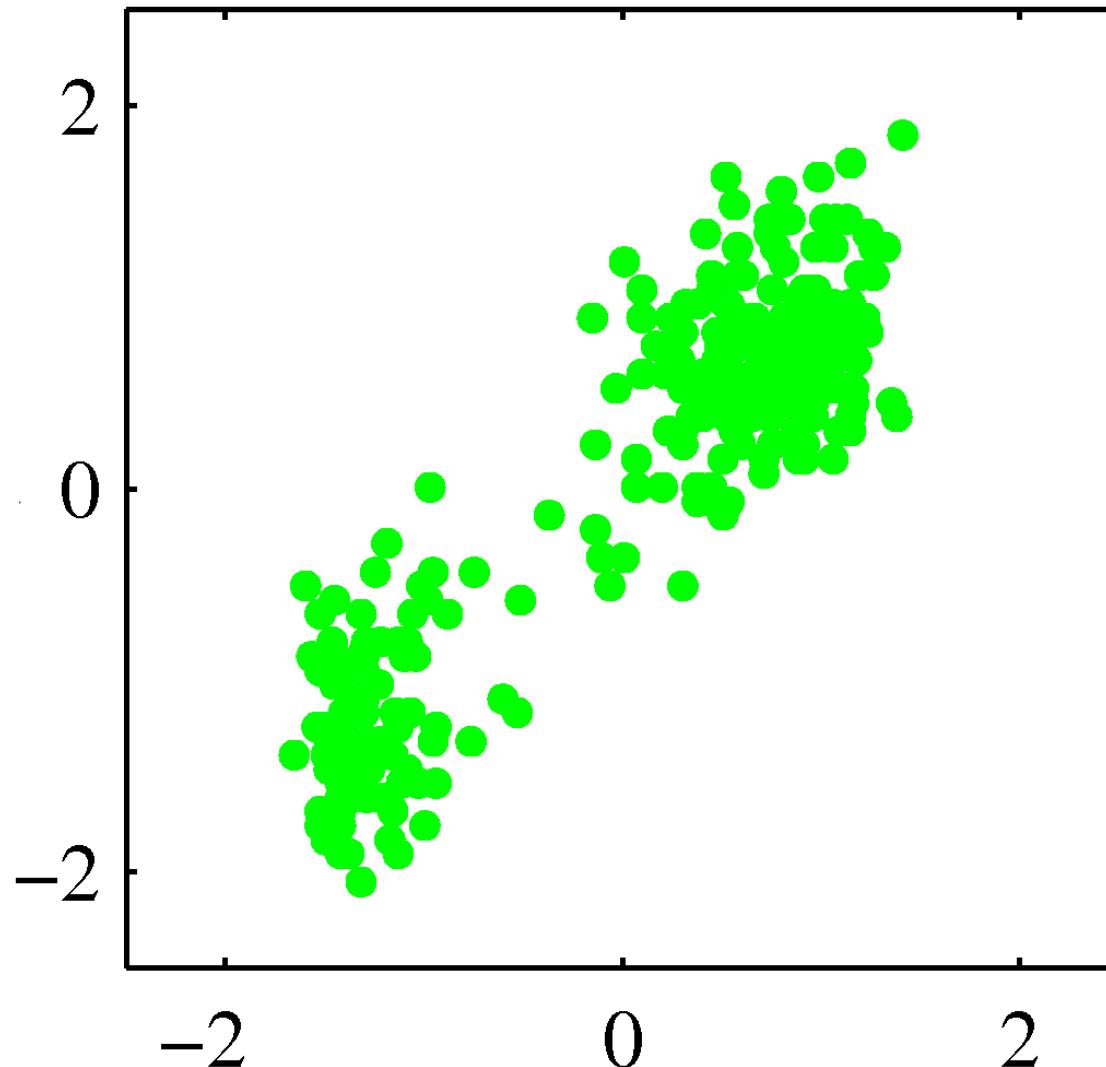
Final Project Proposal

- 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>

Where before we had the y , or 'labels', now we don't have any of that and the task becomes much more difficult, because we cannot use the y /label to guide us defining the hyperplane.

Unsupervised Setting



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

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.

Clustering Applications

In clustering, you're trying to find a pattern that you don't already know ahead of time.

- Google image search categories
- Author Clustering:
[http://academic.research.microsoft.com/Visu
alExplorer#1048044](http://academic.research.microsoft.com/VisualExplorer#1048044)
- Opening a new location for a hospital, police station, etc.
- Outlier detection

In this scenario, some institutions throw out nearly all of their information and only keep the outlier data or only the significant events data.

Unsupervised Learning

- K-means
- Mean-shift
- Hierarchical Clustering
- Rand index, stability

Because we don't have y labels, this is how to evaluate how well the above methods performed.

K-means – Algorithm

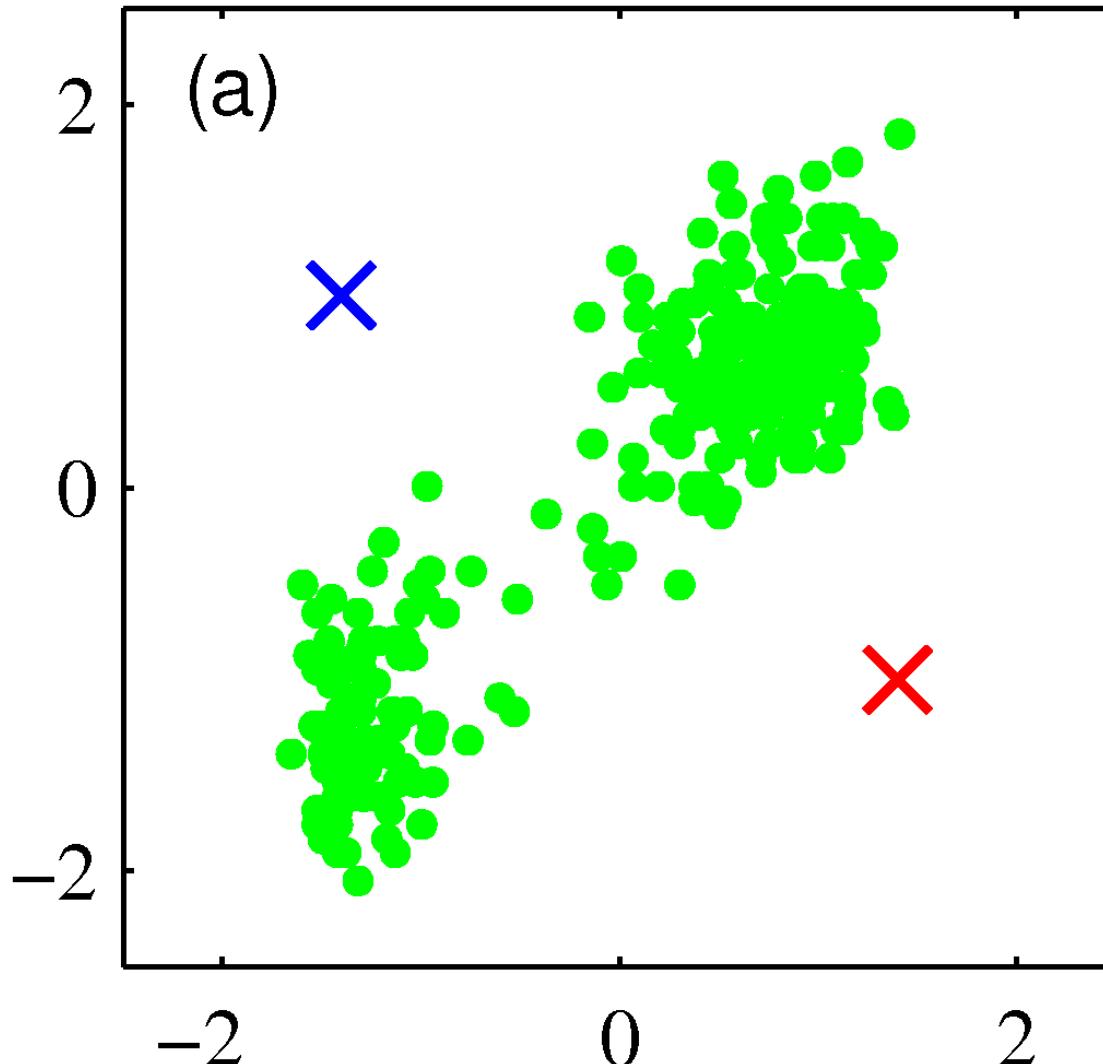
Where before k = number of neighbors, here it's the number of random positions

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

K-means

Here, you randomly choose two points, and say this is now the center of the cluster.

We initialize this algorithm by choosing two arbitrary centers.



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

K-means

- Until Convergence:

- Compute distances $\|x^{(i)} - \mu^{(j)}\|$

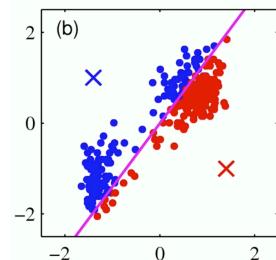
Now we compute the distances of all the data points we have, to these clusters and assign the points to the nearest cluster center.

- Assign points to nearest cluster center

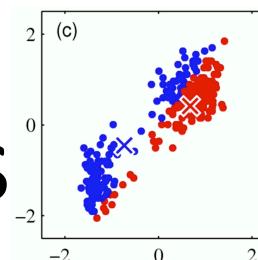
- Update Cluster centers:

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

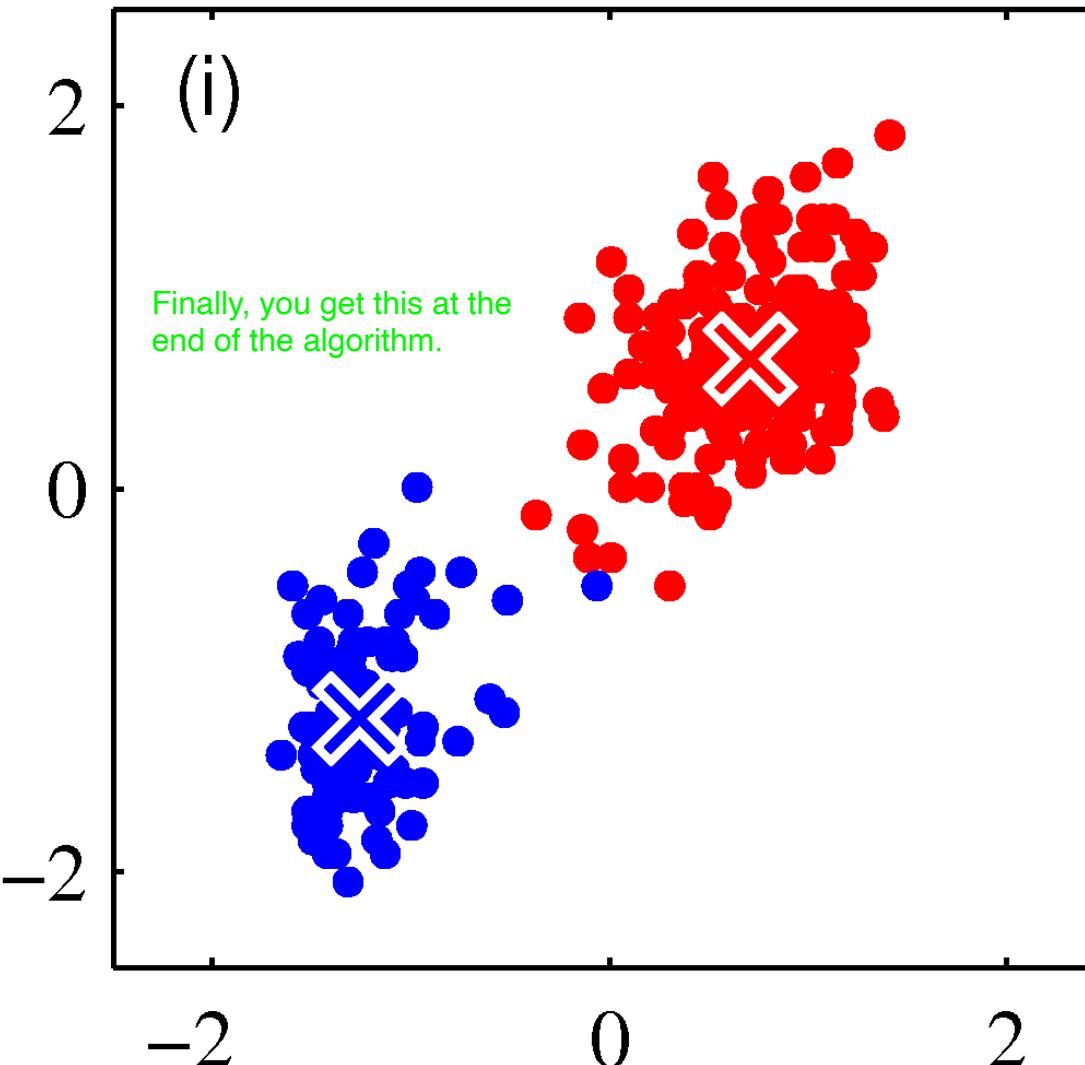
K-means



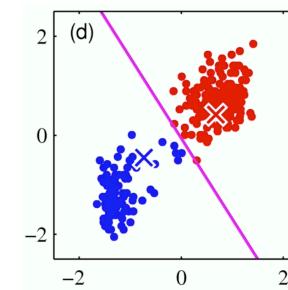
Before updating the cluster centers, the data points nearest to each of the randomly-assigned cluster centers get classified as such.



Here, the centers of the clusters get assigned to be in the middle of all the data points that were assigned to it.



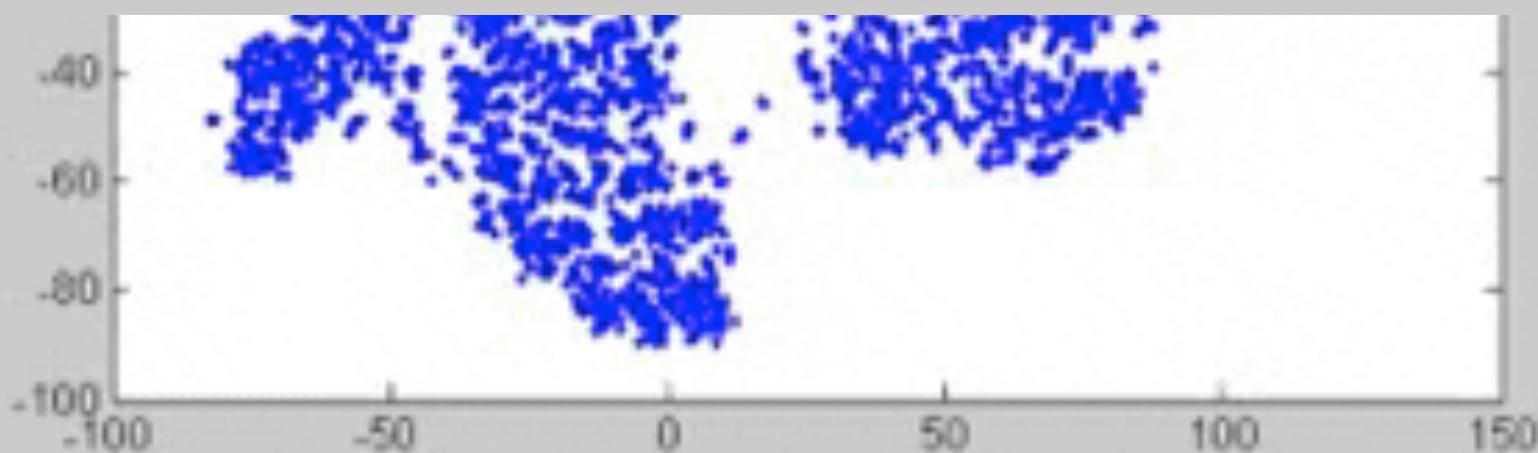
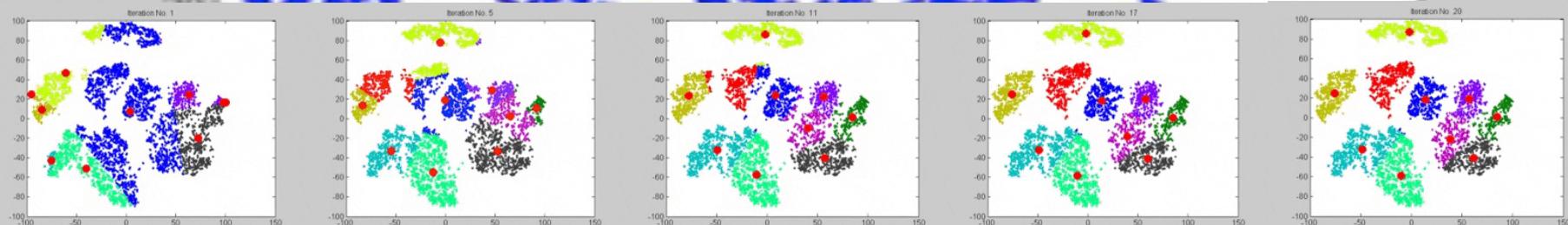
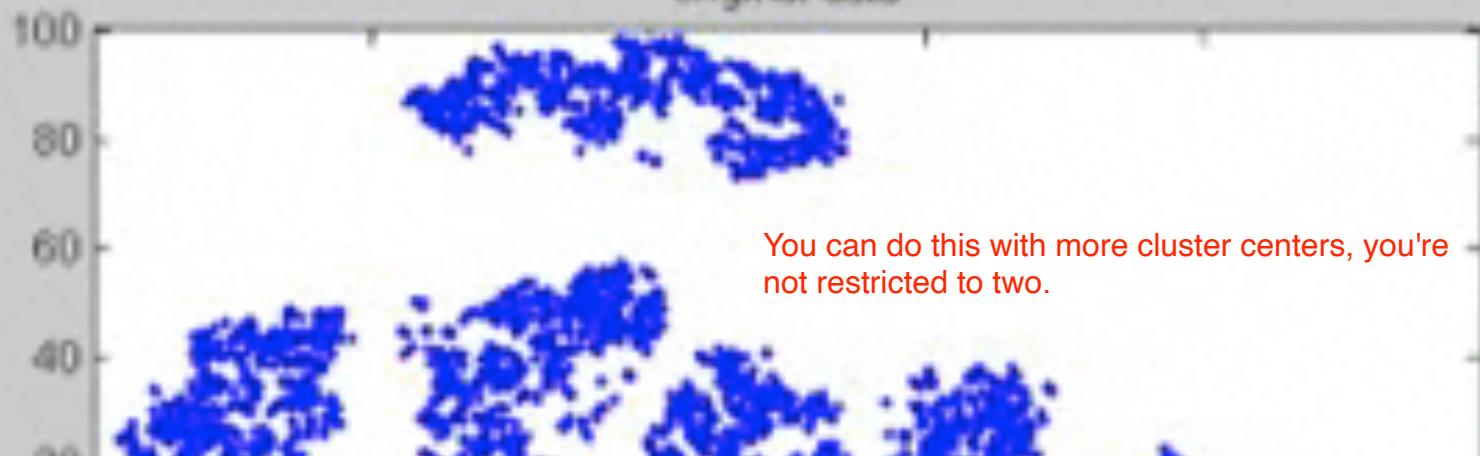
Finally, you get this at the end of the algorithm.

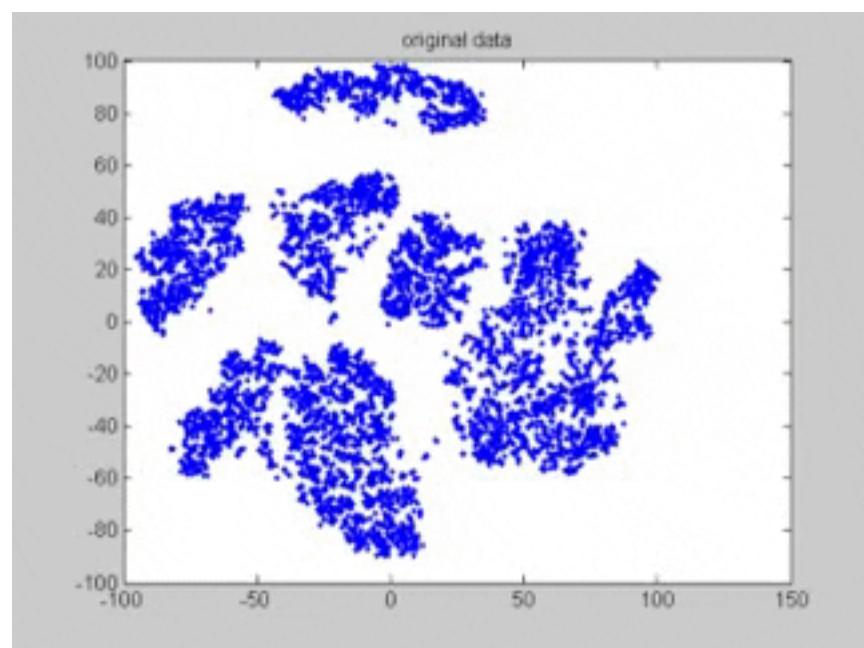


The whole thing starts again and the distances are calculated to determine the boundary

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

original data





K-means Example



R



G



B

K-means Example



K-means Example



Despite both images using $k = 10$ (random) positions, they appear different.



The difference is due to the initializing with a random position for the cluster centers.

K-means Summary

- Guaranteed to converge
- Result depends on initialization
- Number of clusters is important
- Sensitive to outliers
 - Use median instead of mean for updates

How do you know when you've converged? Set a value for Epsilon, and once your cluster centers have moved less than epsilon, then you're done.

This makes it hard to decide which pattern is the one you ought to be choosing/looking for.

You can mediate some of that sensitivity to outliers by using median instead of the mean during the update phase from the centers of your clusters.

Initialization Methods

- Random Positions
- Random data points as Centers
- Random Cluster assignment to data points
- Start several times

Instead of going into the feature space and picking two positions, this says the points have to belong to a cluster so it makes sense to pick random data points as your initial cluster centers.

First do the random cluster assignment, then do the update step, and see where the centers end up.

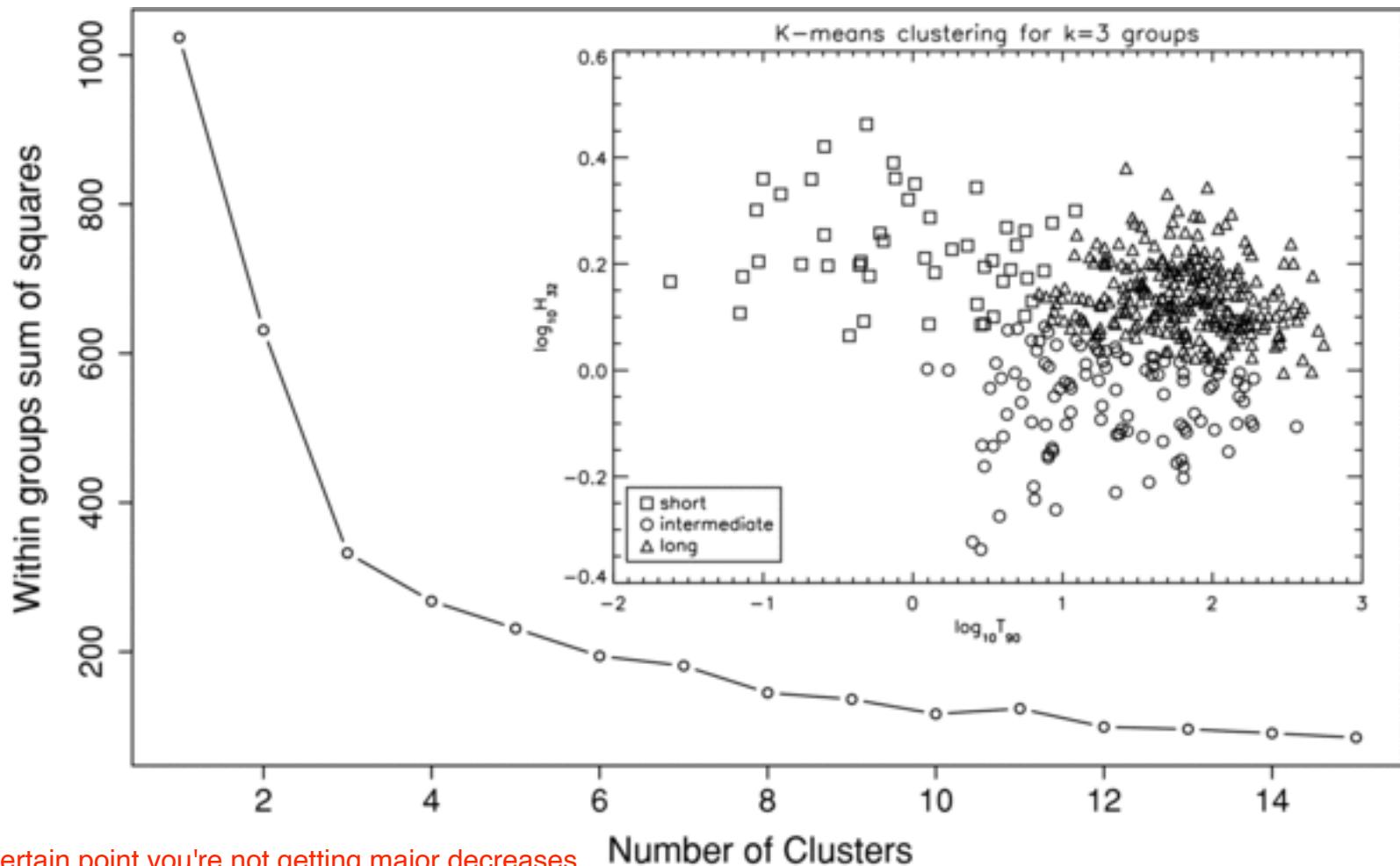
With clustering, there is this idea of stability. If you do 100 runs of k=10 and you get a solution that pops up 90 times and another that comes up 10 times, you'd want to go with the 90x solution since it's more strongly held to by the data.

How to find K

- Extreme cases:
 - $K=1$
 - $K=N$
- Choose K such that increasing it does not model the data much better.

“Knee” or “Elbow” method

Here, within groups sum of squares measures the distance from the data points to their cluster group's center.



So, at a certain point you're not getting major decreases
in sum of squares for the number of K/random cluster centers.

Cross Validation

If you want to be able to generalize to new data, then cross validation will help you pick the best k/number of cluster centers to avoid overfitting.

- 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

Instead of using the same data to compute the sum of squares, you have training data that you use to solve your kmeans and then you have the validation data coming in where you compute all the distances to the groups that you just identified.

Cross-validation gives you confidence about generalization or about the new data coming in or that you haven't seen yet.

Getting Rid of K

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

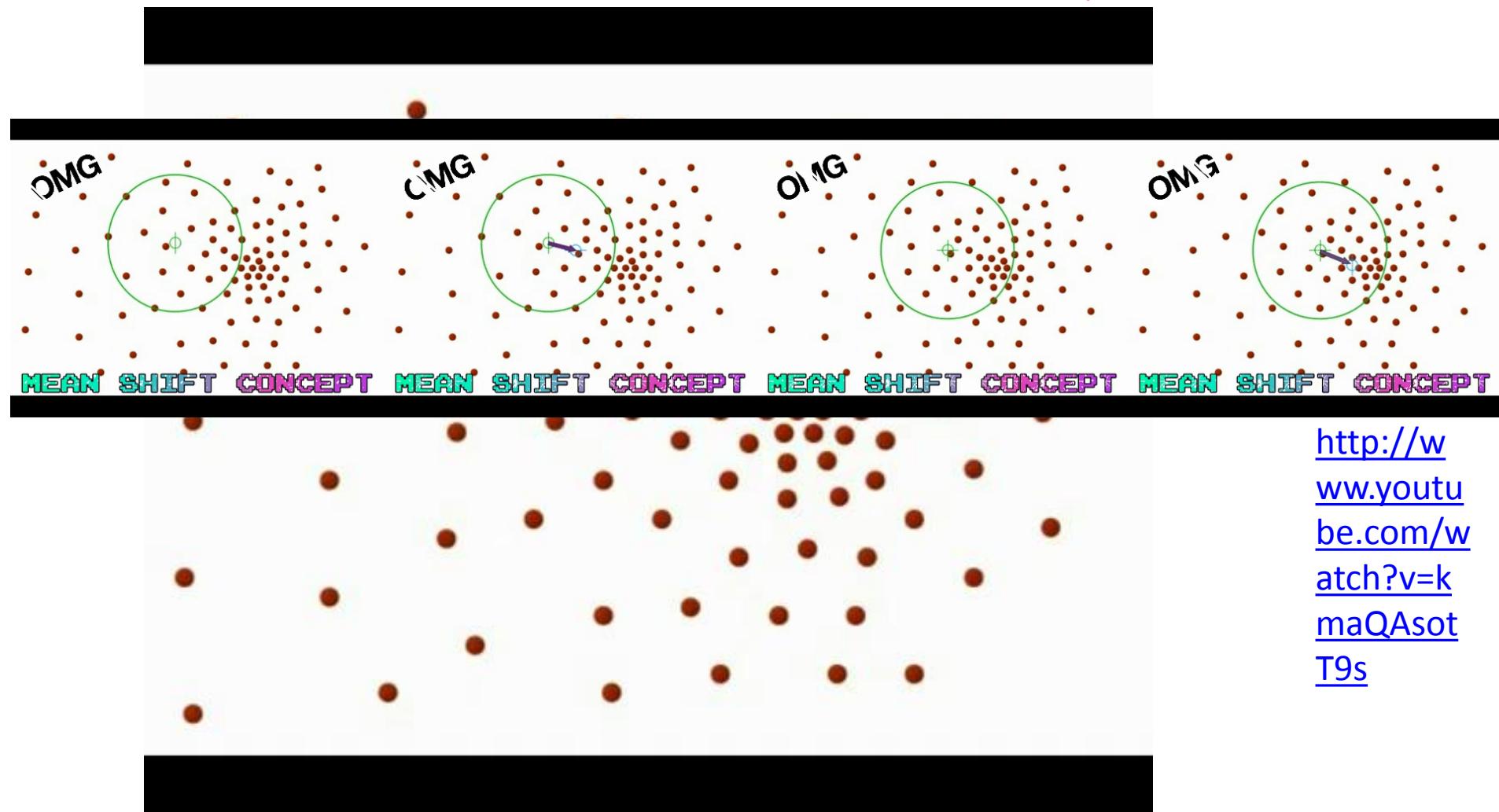
For each point in the dataset, you specify the window around it. And that window is the same size for all points. Now this point is essentially becomes a cluster center and you look for all the points that are in a specific distance to that point, inside this window.

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

You end up doing this for each individual data point, and shift, shift, shift.



<http://www.youtube.com/watch?v=kmaQAsotT9s>

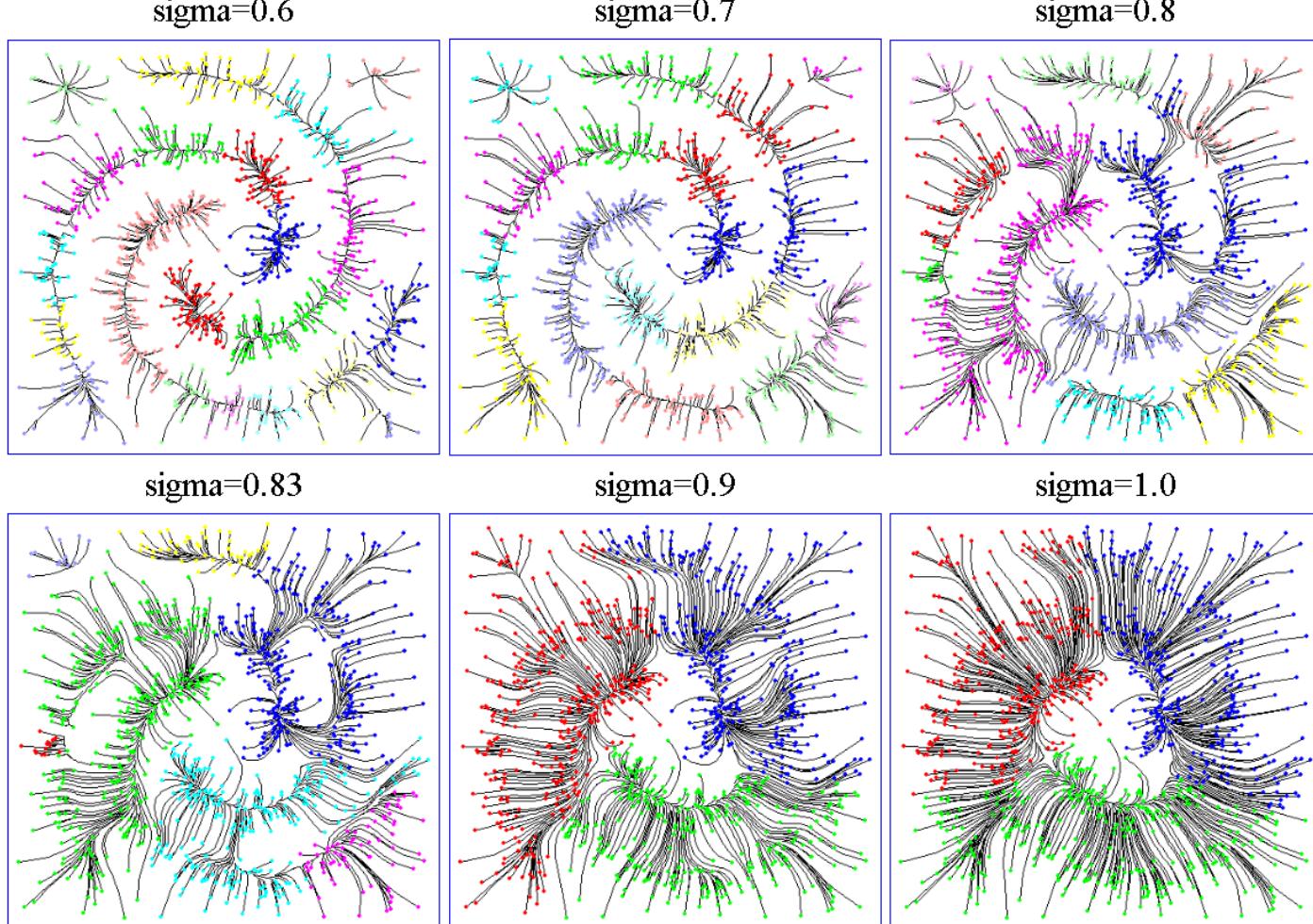
When you have a gradient of points, where they're more clustered/denser, the window will shift towards that density of points. All the points will end up with the same center in the end because of this, and that center is your cluster.

This is a demonstration for different window sizes and the lines here show the path that the window took in each iteration.

Here the number of clusters hasn't been defined, just the size of the window.

Mean Shift

All the points where the window shifted to the window shifted to the same place are assumed to be part of the same cluster.



The size of the window has a lot to do with the number of clusters that're being found.

Mean Shift Summary

- Does not need to know number of clusters
- Can handle arbitrary shaped clusters
- Robust to initialization
- Needs bandwidth parameter (window size)
- Computationally expensive

The reason why is because you have to do it for each and every single data point.

- Very good article:

Calculating the mean for each of the data points within the window is computationally expensive.

<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

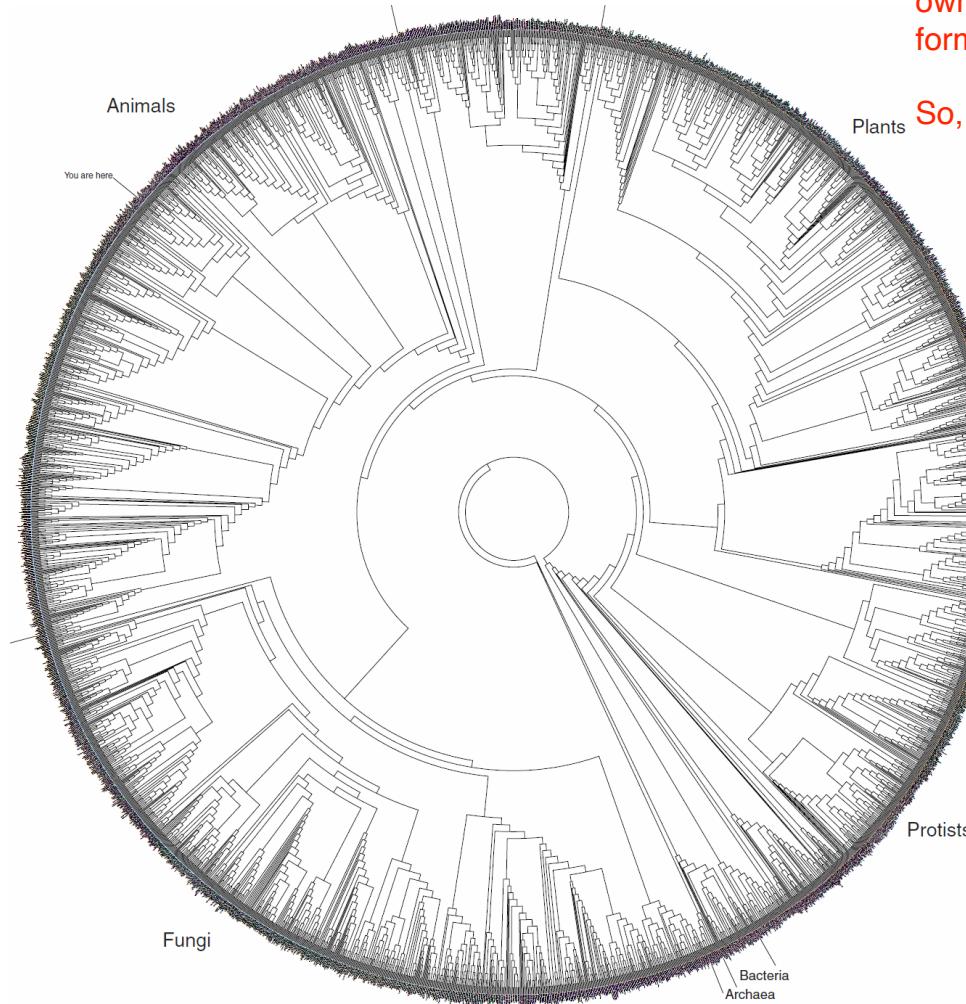
Again, there is no 'free lunch' paramaters, meaning one of them has to be set for the model to function.

- 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

This approach uses neither k nor windows to function.

This is a listing of organisms and how similar they are in terms of evolvement.

Tree of Life

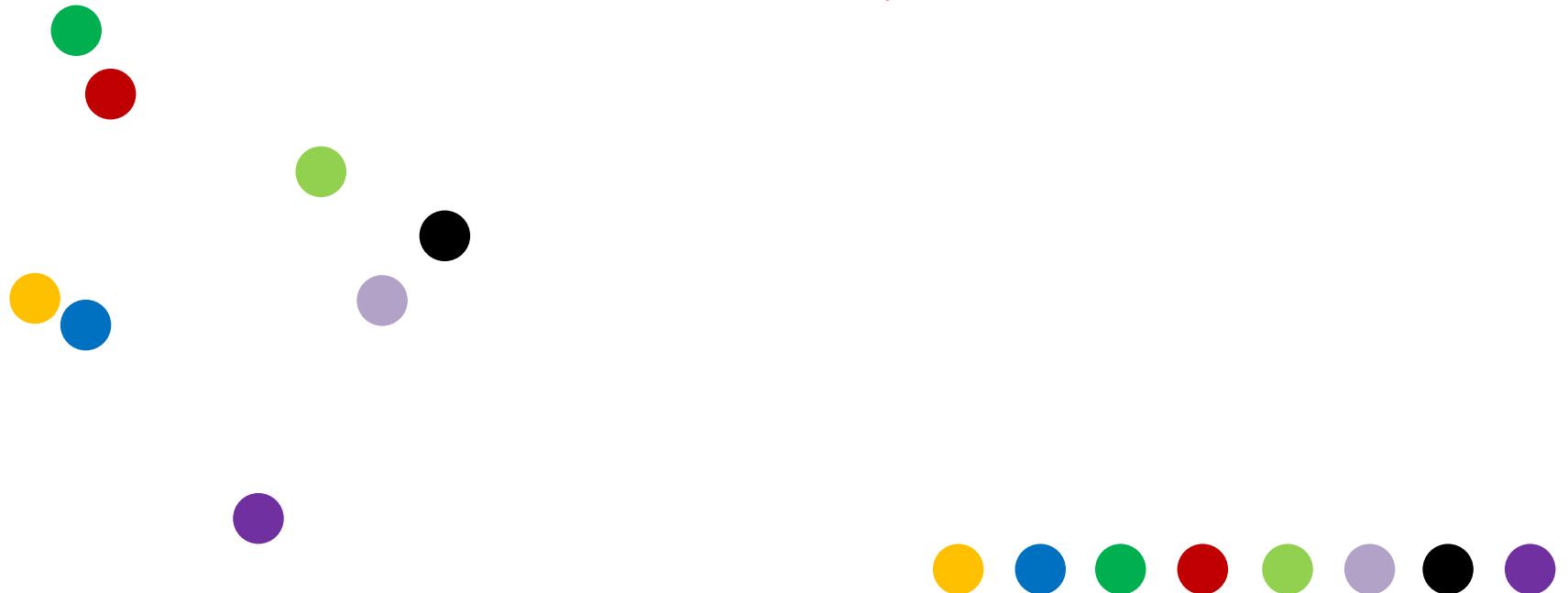


On the lowest level, every organism is its own cluster. Then the groupings start to form larger and larger clusters of points.

So, plants are in one cluster then fungi, etc.

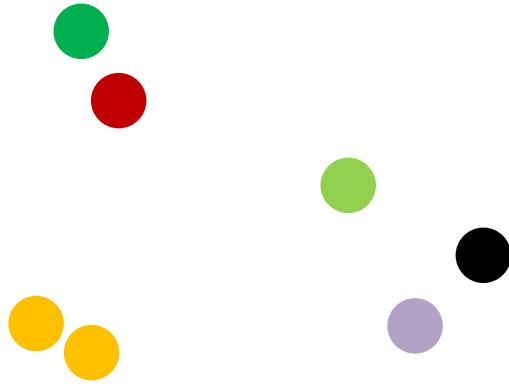
Hierarchical Clustering

Each point is its own cluster.



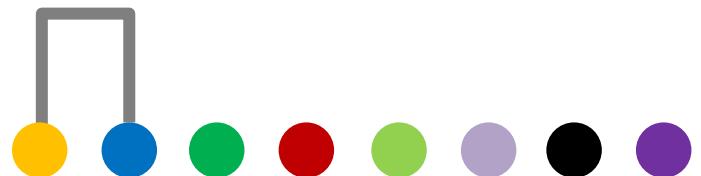
Hierarchical Clustering

The distance between the two closest points is calculated.

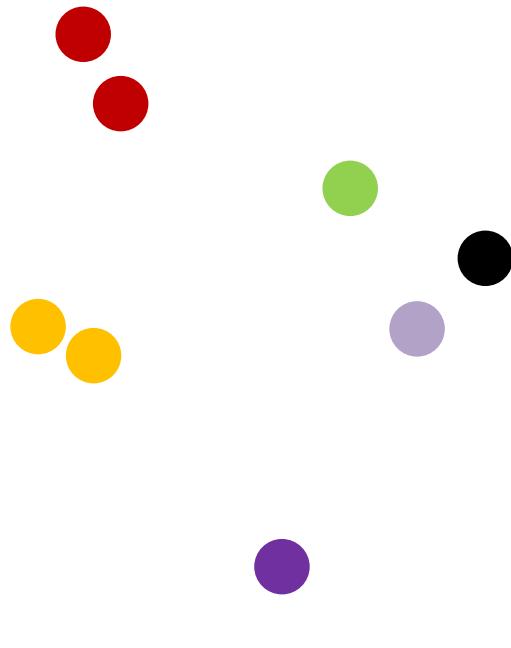


At that point, only those two points are grouped together,
and are marked as such on the dendrogram.

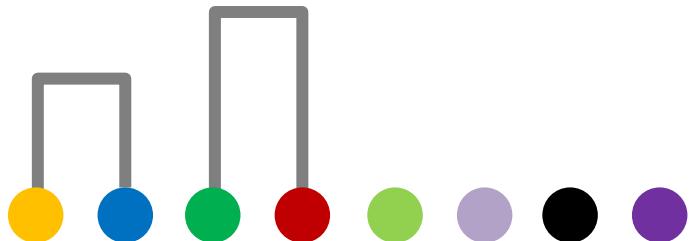
These can now be considered to be in the same cluster.



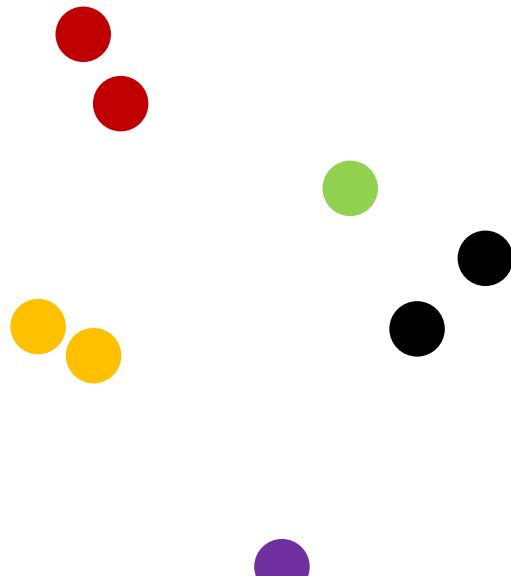
Hierarchical Clustering



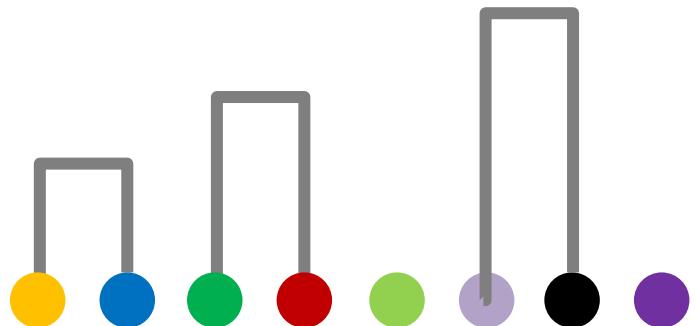
So then the next closest data points have their distance calculated and then assigned to the same cluster.



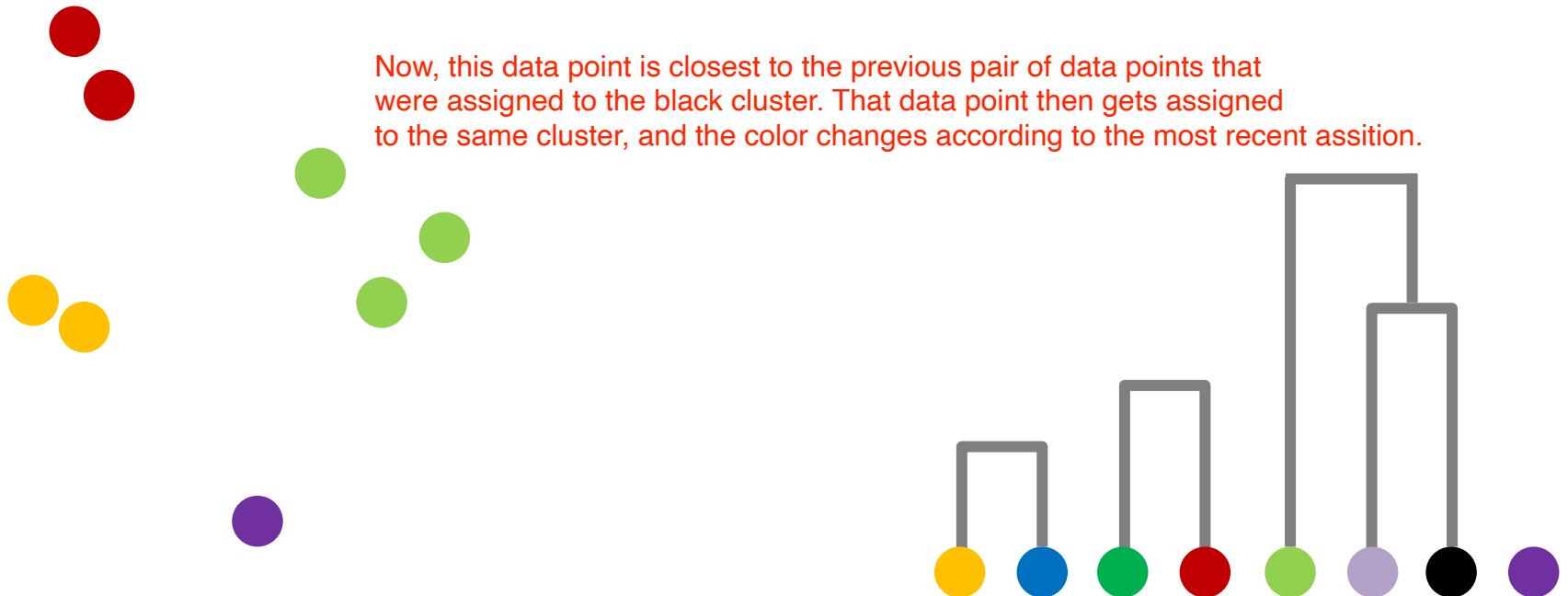
Hierarchical Clustering



And the next pair that're closest together are assigned to the same cluster.

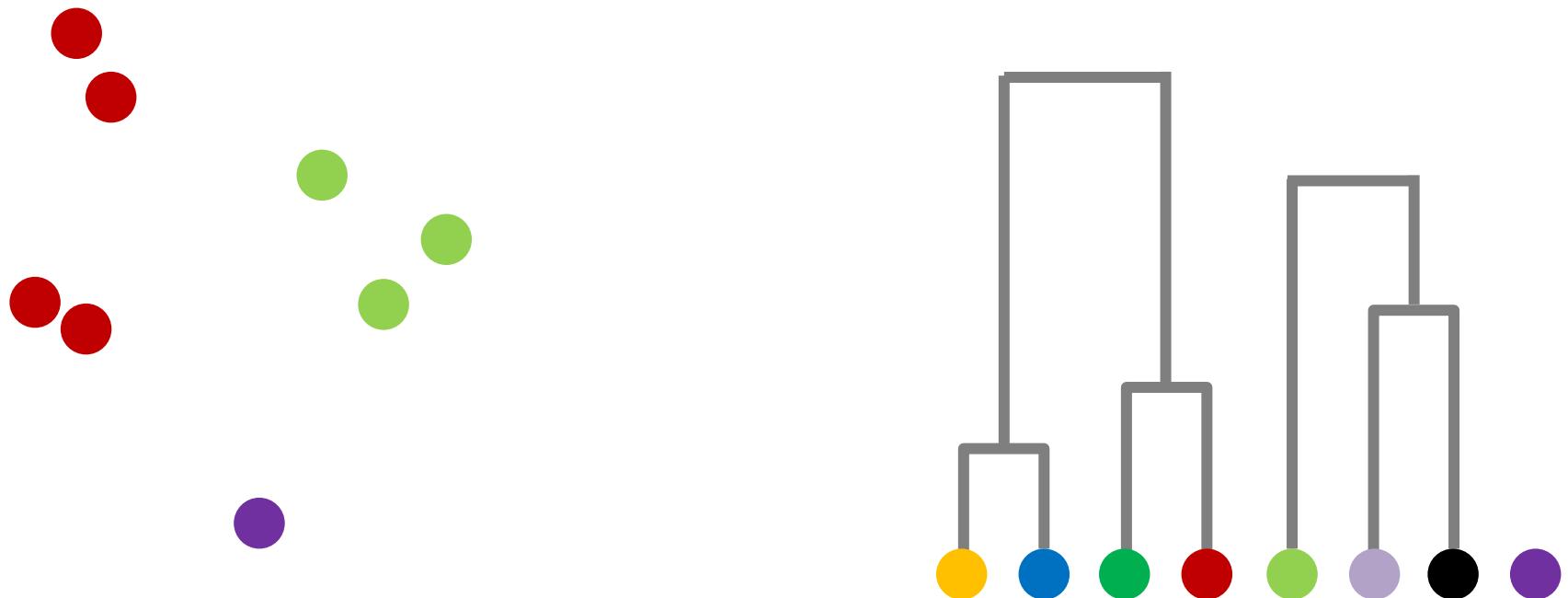


Hierarchical Clustering



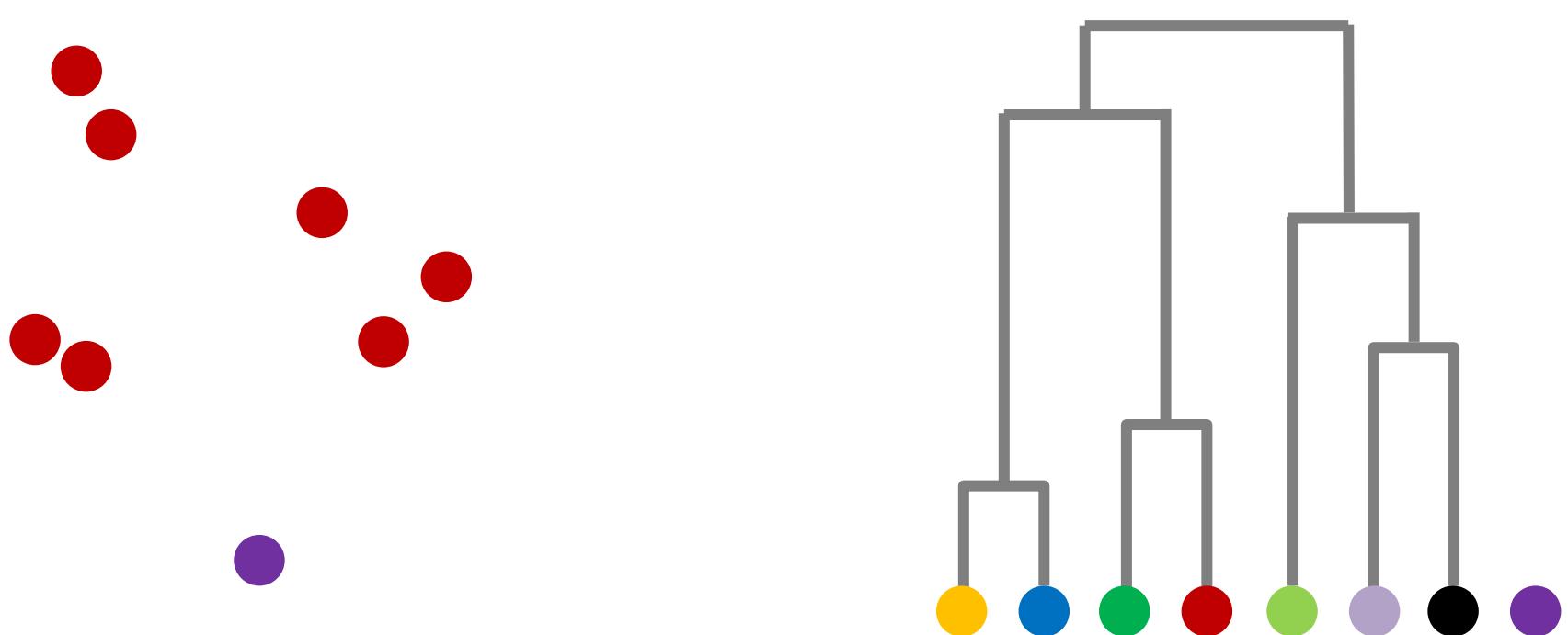
Hierarchical Clustering

And these two (precious) clusters are closest to one another and are grouped into a cluster themselves.



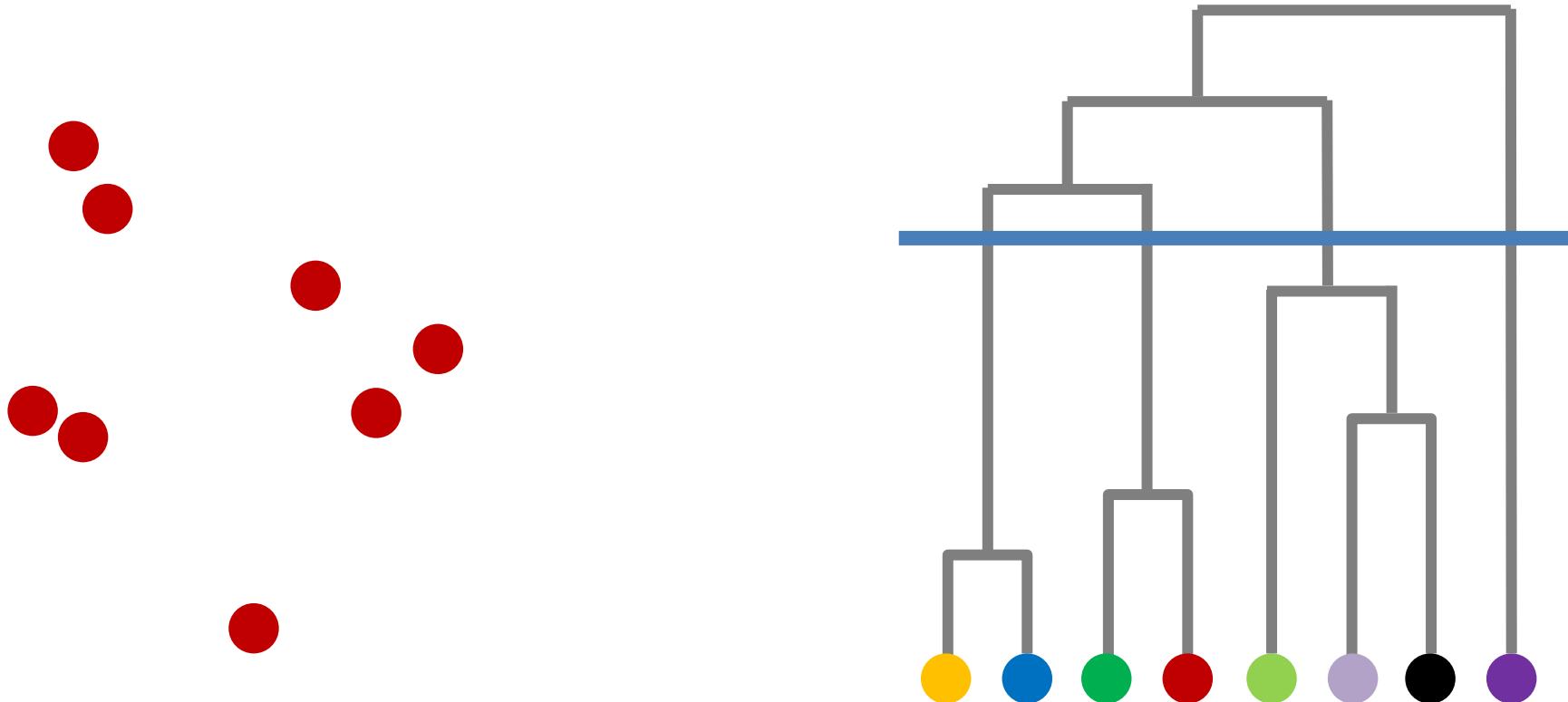
Hierarchical Clustering

So, this is a whole solution where you're going from each data point being its own cluster to having one large cluster of all the data points. That is hierarchical clustering.



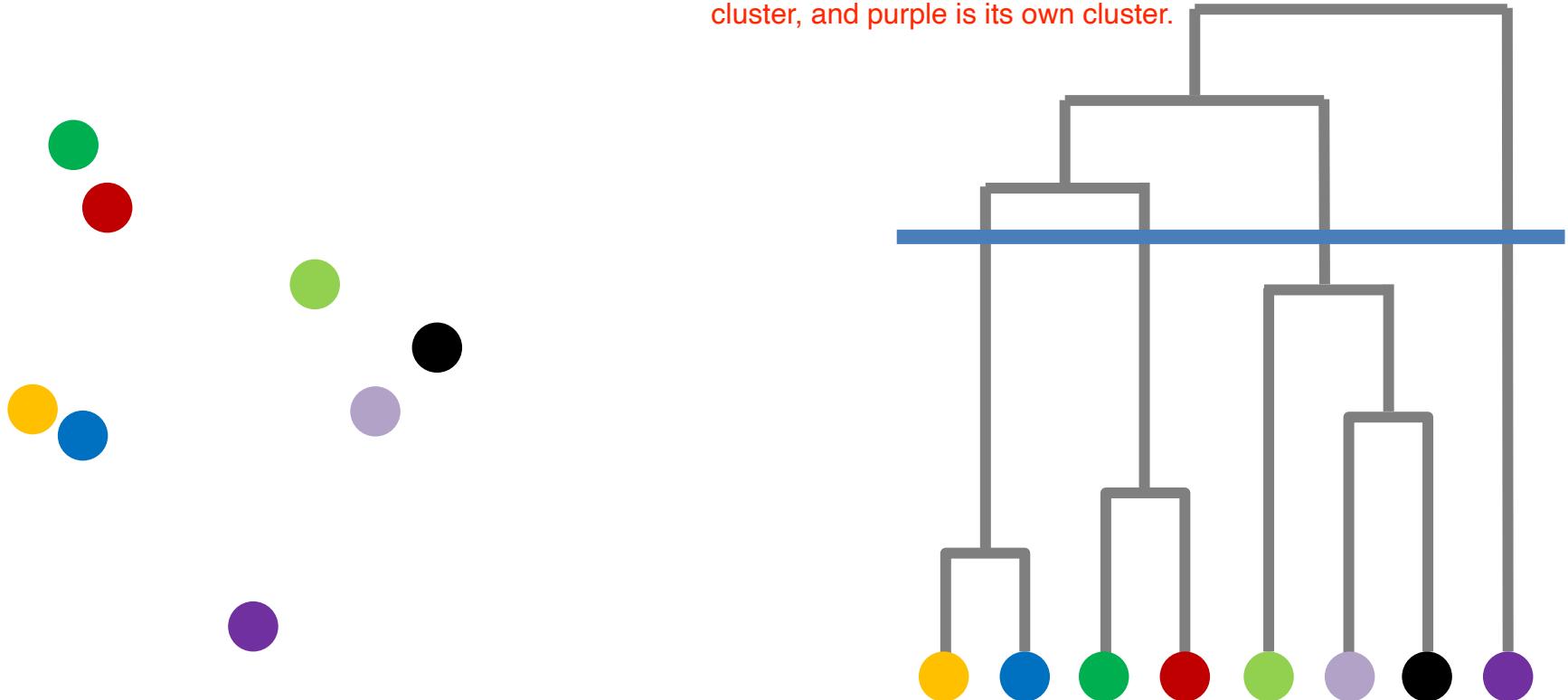
Hierarchical Clustering

The trick is that you need to ask the user to choose a threshold, and then the coloring that corresponds to that threshold.



Hierarchical Clustering

At this threshold, yellow + blue are in one cluster, green + red = one cluster, lime + gray + black = 1 cluster, and purple is its own cluster.

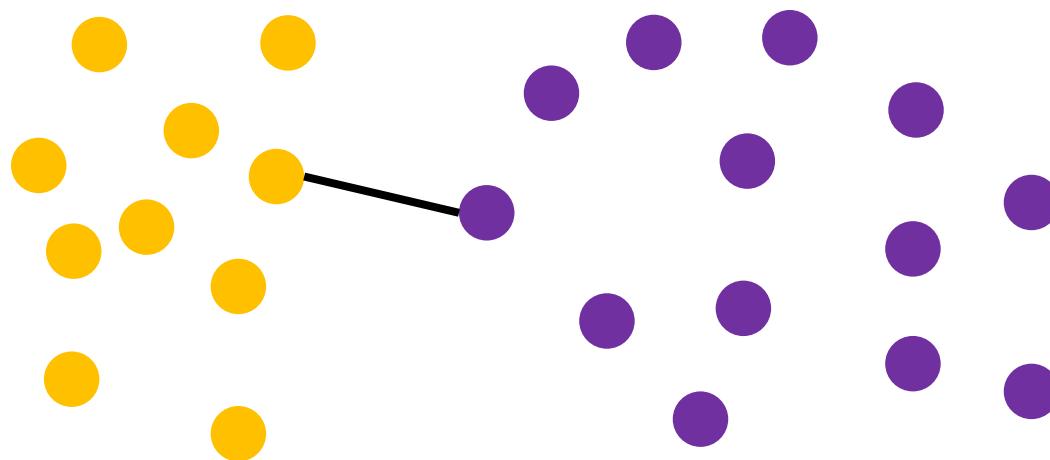
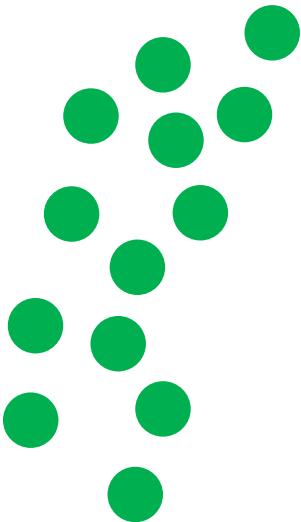


Hierarchical Clustering

This can be beneficial to see how all the data points relate to one another, like in the tree of life example above.

- Produces complete structure
- No predefined number of clusters
- Similarity between clusters:
 - single-linkage: $\min\{d(x,y) : x \in \mathcal{A}, y \in \mathcal{B}\}$
This is the case when you're looking at the points closest to each other.
 - 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



Here, we're looking at the minimum distance between clusters.
In this example, these two data points are closest together.

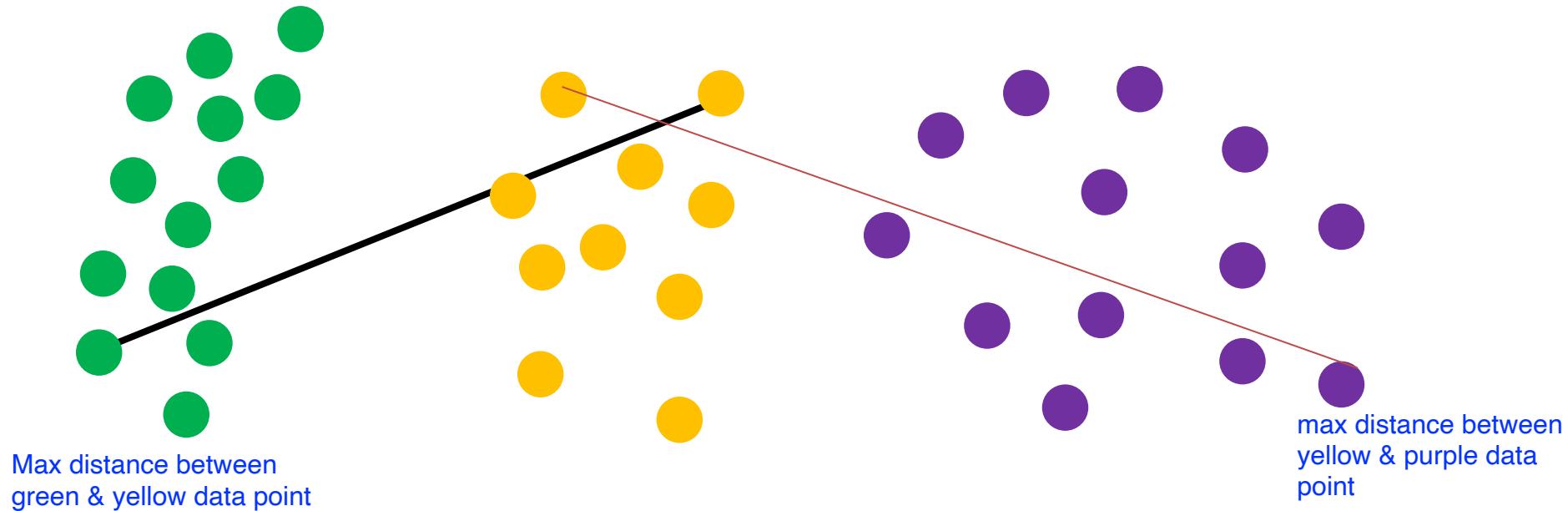
The distance between the closest data points determines the distance between the entire cluster from one another. This means you get one number between yellow and purple and one number for the distance between green and yellow.

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

And then you look at which of those three numbers is the smallest, and that's the next grouping or clustering that will happen.

Complete Linkage

So this ends up being the opposite of the single linkage. Here, you look at the maximum distance between points in this cluster in comparison to the other clusters.



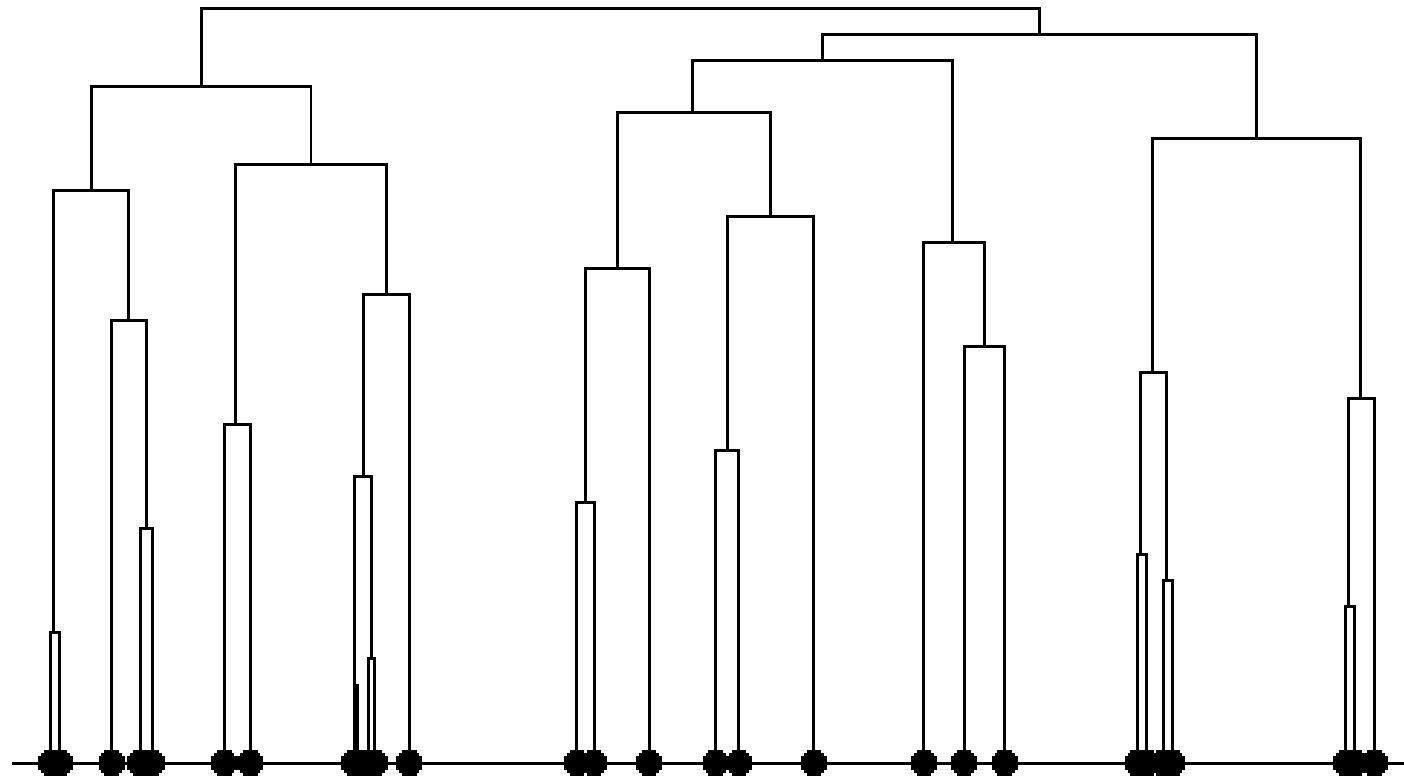
$$\max\{d(x, y) : x \in \mathcal{A}, y \in \mathcal{B}\}$$

And then you choose the smallest of the maximum distances to base your next grouping/clustering on.

Linkage Matters

- Single linkage: tendency to form long chains
- Complete linkage: Sensitive to outliers
- Average-link: Trying to compromise between the two

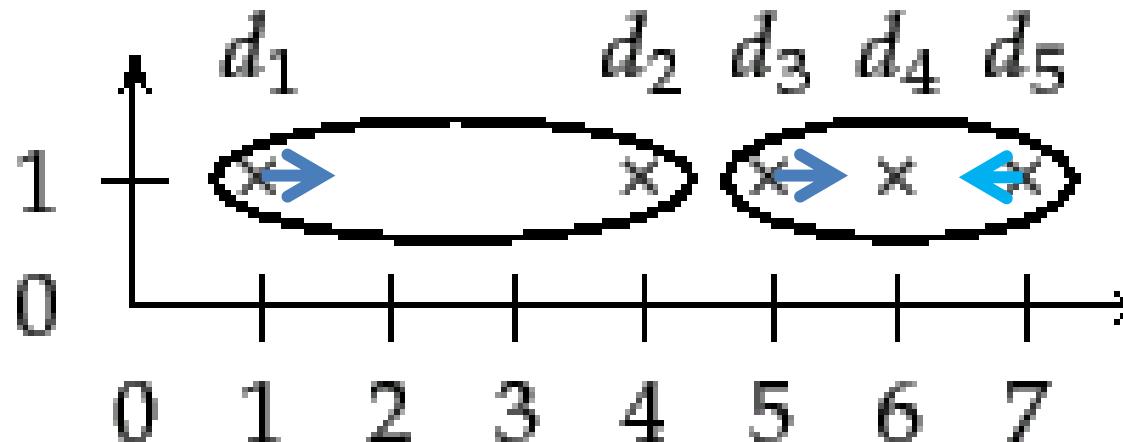
Chaining Phenomenon



Outlier Sensitivity

Complete linkage has a problem with outliers given that it's looking at maximum values, but it has the benefit of giving you more balanced clusters over the hierarchy.

Single linkage is more robust towards outliers, but produces chaining clusters.



→ + 2*epsilon

← - 1*epsilon

Efficient Hierarchical Graph-Based Video Segmentation

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

¹Georgia Tech ²Google Research

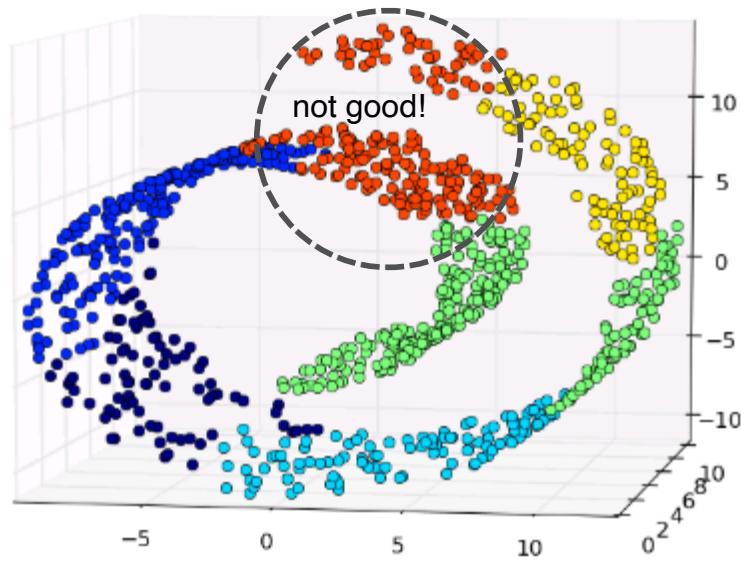
IEEE CVPR, San Francisco, USA, June 2010



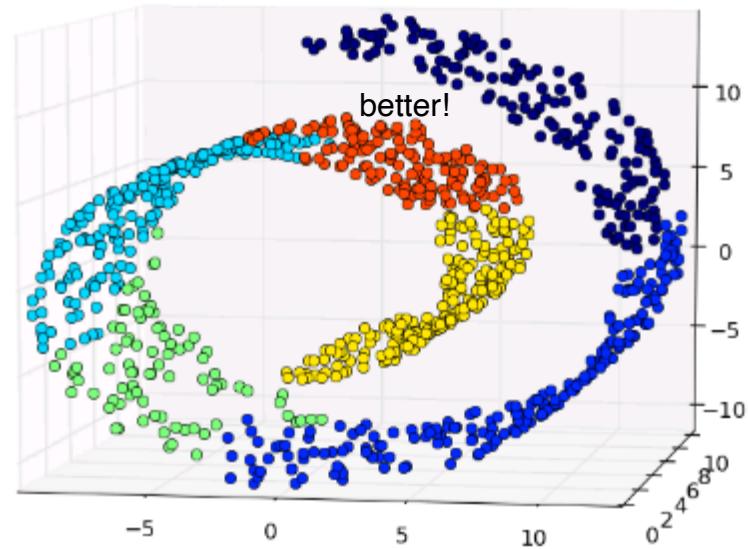
k-means looks for blobs of data and spirals are not blobs so it presented a challenge to clustering via k-means.

Swiss Role Problem

scikit-learn allows you to specify a connectivity between points and to have a constraint that points that're grouped, need to be adjacent. That can solve for you the 'swiss roll problem'.



without connectivity
constraints



with connectivity
constraints

only adjacent clusters can be merged together

The idea is to unravel the spiral and find which points are adjacent. With connectivity contraints you get the clusters along the spirals. YOu will have to specify the connectivity along the points.

Evaluation Criteria

This is hard, because you're attempting to answer the question of whether the pattern you've demonstrated is sensible. How do you know that your clustering solution actually makes sense?

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

How would you evaluate for novelty and that that novelty makes sense?

Rand Index

tp = two points in the same cluster than belong in the same cluster

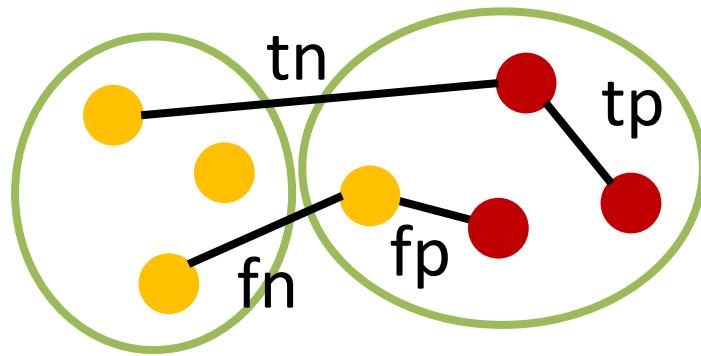
tn = two points in separate clusters than belong in separate clusters

fp = two points in the same cluster that should've been in separate clusters.

fn = two points in separate clusters that should've been in the same cluster.

- Percentage of correct classifications
- Compare pairs of elements:

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

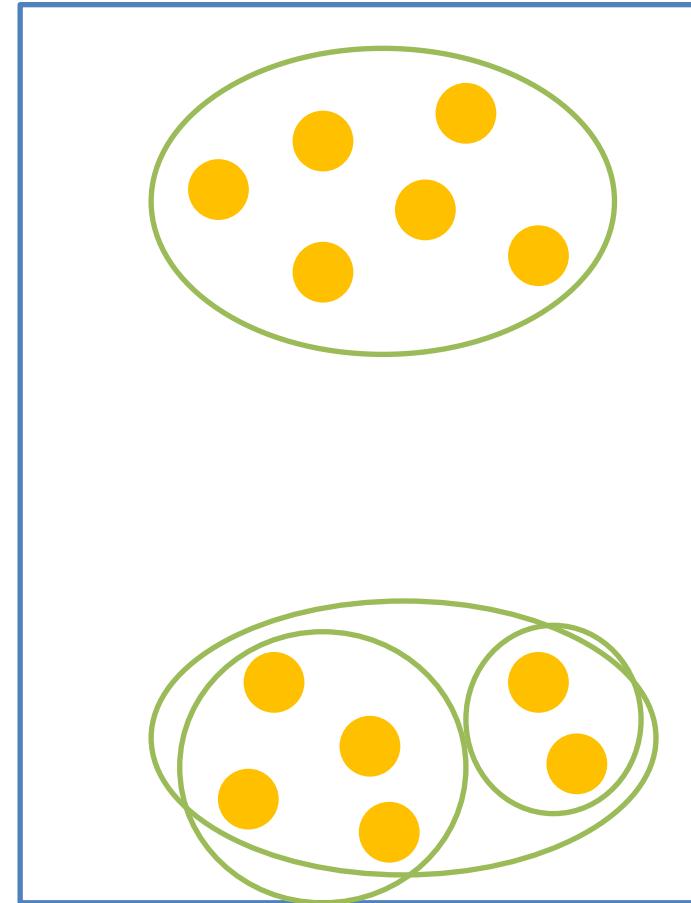
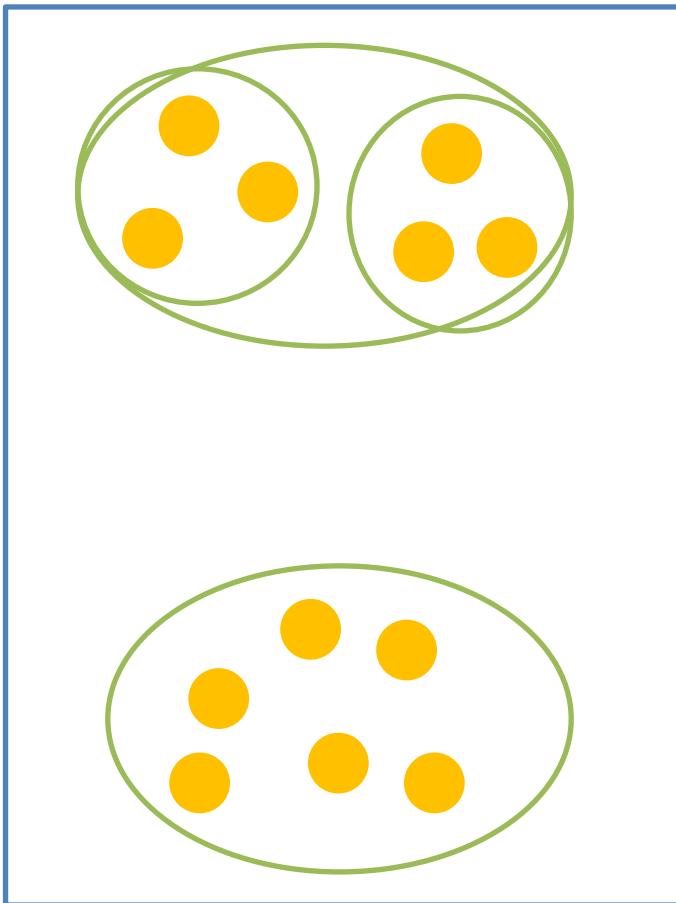


- Fp and fn are equally weighted

Basically, what percentage did you get rightly clustered together.

A good clustering solution generalizes to unseen data. But how do you know that your clustering will generalize well?

Stability



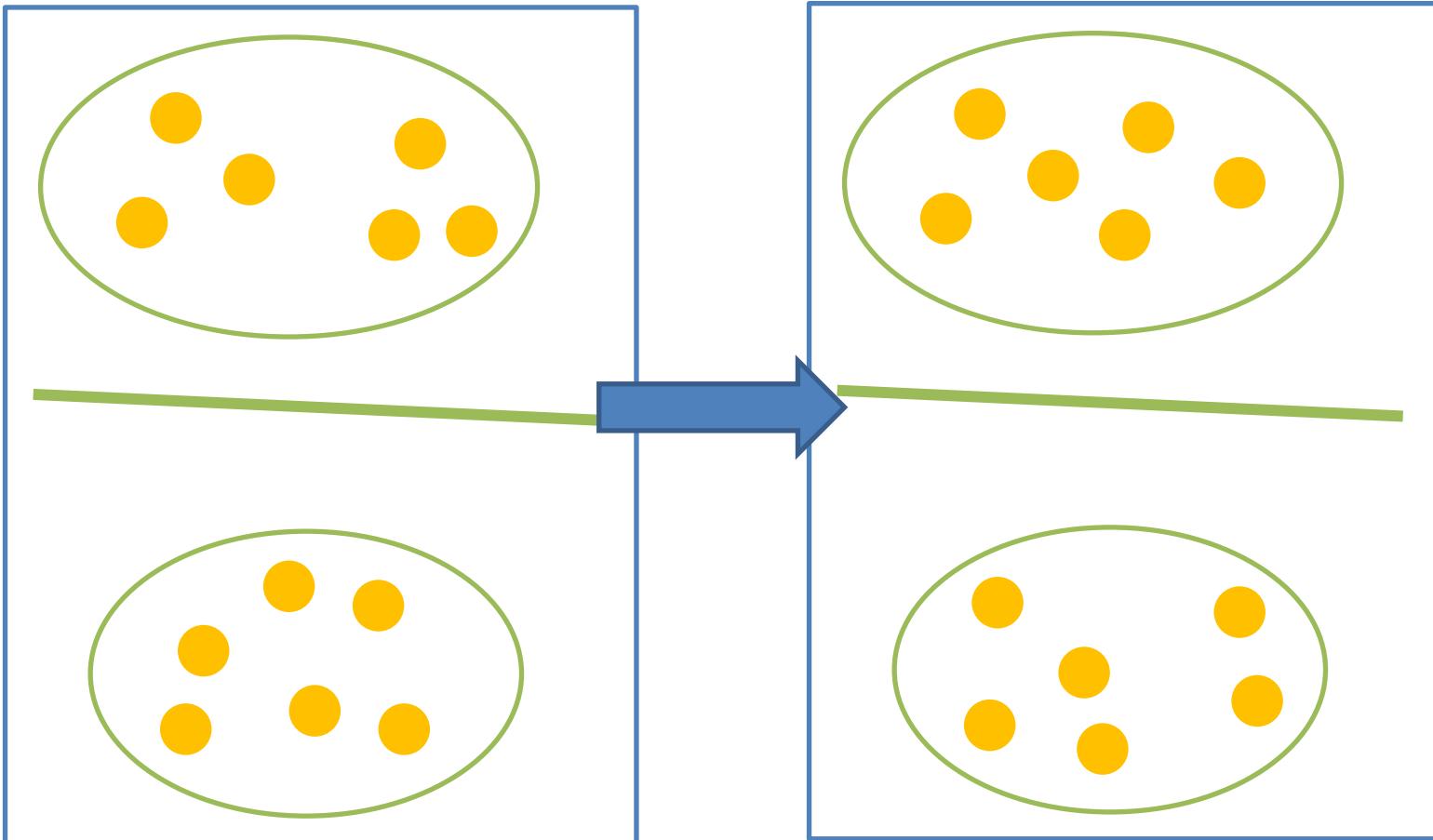
The idea is to split your data into a second set. Whatever clustering you choose in the first set, should generalize well and explain the second dataset, because it's coming from the same source. That's the notion that a clustering should be stable, so that your k cluster centers persist across the new/unseen data.

Stability

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

When you discover clusters, you're classifying data points, essentially creating labels/y.
This allows you to define a hyperplane/decision boundary and train a classifier... random forest support vector machine, within there.

Stability



So if the clustering from the initial data set was a good idea and is comparable to that clustering in the second dataset. Use that classifier to compute the test error on the second set. If the clustering is stable, the error should be small.

Use the clustering to transform the unsupervised learning problem into a supervised one. Then see if the labels assigned by the clustering solution were useful to generalize to unseen data.

Summary

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