Clustering

Matt Dickenson



Agenda

- Introduction
- Why Clustering Works
- Popular Approaches
 - DBSCAN
 - K-Means
 - Gaussian EM
- When to Use Clustering
- When Clustering Fails



Image: https://xkcd.com/388/

Overview



Overview



Overview

- Unsupervised approach
 - No labeled data
- Group similar points together
 - How to define similarity?



Why Clustering Works

• Often datasets consist of underlying groups



Why Clustering Works

- Often datasets consist of underlying groups
- Usually simple to define similarity as a distance



Popular Clustering Algorithms







DBSCAN neighborhood-based

K-Means centroid-based Gaussian E-M distribution-based

Images: http://shapeofdata.wordpress.com



Image: http://shapeofdata.wordpress.com

"Density-based spatial clustering of applications with noise"



Image: http://shapeofdata.wordpress.com

- 2 parameters
 - \circ ϵ controls the size of a point's "neighborhood"
 - How many neighbors a point must have to be considered "core"
- Result
 - List of core points
 - Assigned points & outliers



ε = 4, minPts = 30



Iterate through points, count neighbors within ε distance



 ε = 4, minPts = 30

If the point has |neighbors| > minPts, consider it "core"



 ε = 4, minPts = 30

If the point has |neighbors| > minPts, consider it "core"



 ε = 4, minPts = 30

Check neighbors' neighbors for density



 ε = 4, minPts = 30

Assign core point's neighbors to its cluster



 ε = 4, minPts = 30

Iterate through points until you find next core point



 ε = 4, minPts = 30



Check neighbors' neighbors for density

 ε = 4, minPts = 30

Assign unlabeled neighbors to new core point



 ε = 4, minPts = 30

Continue until all points have been checked



 ε = 4, minPts = 30

Result: 3 clusters



 ε = 4, minPts = 30

Result: 2 clusters



 ε = 5, minPts = 30

Result: 9 clusters



ε = 2, minPts = 10

• Pros

- Don't have to specify number of clusters
- Fairly robust to outliers
- Handles non-spherical clusters
- Cons
 - Non-deterministic
 - For neighbors of two core points, cluster assignment depends on order that data is processed
 - Non-probabilistic
 - Returns cluster assignments without uncertainty/probability
 - Tuning parameters can be difficult



Image: http://shapeofdata.wordpress.com

- 1 parameter, *k*, controls number of clusters
- Distance function
 - Typically Euclidean distance in feature space
- Result: list of centroids and point labels
- Caveat: there are many ways to implement K-Means





Pick *k* random data points as initial centroids



Compute the distance of each point from each centroid



Compute the distance of each point from each centroid



Assign each point to the closest cluster



Update centroids to the mean of new clusters



Update centroids to the mean of new clusters










• Pros

- Quick to implement (and usually quick to run)
- Simple, intuitive algorithm
- Cons
 - Can be difficult to choose *k* in practice
 - Clusters will all have the same radius
 - Sensitive to outliers

How do I choose a good value of k?

- Visually inspect the data
 - This is impractical for high-dimensional data



How do I choose a good value of k?

- Visually inspect the data
 - This is impractical for high-dimensional data
- Try several different values
 - Which value of k minimizes the total error?
 - Increasing k will always decrease the total error
 - Pick the value that gives the greatest gain ("elbow method")





How do I choose a good value of k?

- Visually inspect the data
 - This is impractical for high-dimensional data
- Try several different values
 - Which value of *k* minimizes the total error?
 - Increasing *k* will always decrease the total error
 - Pick the value that gives the greatest gain ("elbow method")
- Other metrics
 - Silhouette score (more detail next week)





Expectation-Maximization



Image: http://shapeofdata.wordpress.com

Expectation-Maximization

- You can use EM with a variety of distributions
- We will use it with a mixture of Gaussians (GMM)
- 2 parameters per cluster
 - µ represents cluster center
 - σ or represents cluster shape & scale (standard deviation)
- Result
 - Parameters for final clusters
 - Labels for each point ("soft" assignments)



Make initial guesses of parameter values



Make initial guesses of parameter values



Assign each point to its most likely cluster (E-step)



Update cluster parameters (M-step)



Repeat M-Step



Repeat M-Step



iteration 3

Repeat M-Step



Repeat M-Step



iteration 5

Repeat M-Step



iteration 6

Continue until clusters converge



Expectation-Maximization

- Pros
 - Each cluster is a statistical distribution
 - "Soft" assignment with probabilistic uncertainty
 - Robust to poor initial clusters
 - Clusters can be different sizes
 - Very flexible approximation
- Cons
 - Must choose number of clusters in advance
 - Distance function is less flexible
 - Somewhat sensitive to outliers





Truth

DBSCAN





Truth

DBSCAN





Truth

K-Means





Truth

K-Means



Truth

Gaussian EM



Truth

Gaussian EM

Choosing a clustering algorithm

- DBSCAN
 - \circ You have a good sense of what an appropriate ϵ threshold is
 - You are concerned about outliers
- K-Means
 - You want to get started very quickly
 - You have time to try different values of *k*
- Expectation-Maximization
 - You want to convey uncertainty about your predictions
 - You know the number of clusters or have time to try several values

Other Clustering Approaches



Hierarchical

http://www.sthda.com/sthda/RDoc/figure/clustering/cluster-analysi s-in-r-hierarchical-clustering2-1.png

https://hemberg-lab.github.io/scRNA.seq.course/figures/graph_net work.jpg

Graph-Based

Neural Model

http://www.pitt.edu/~is2470pb/Spring05/FinalProjects/Group1a/tut orial/kohonen1.gif

Other Clustering Approaches



Image: http://scikit-learn.org/stable/modules/clustering.html

When Clustering Can Fail



	non-spherical	different densities	high-dimensional data
DBSCAN	\checkmark	×	×
K-Means	×	OK	×
Gaussian EM	×	 ✓ 	OK

http://hameddaily.blogspot.com/2015/03/when-not-to-use -gaussian-mixtures-model.html https://stats.stackexchange.com/questions/133656/how-to-understand-the -drawbacks-of-k-means https://upload.wikimedia.org/wikipedia/commons/thumb/a/a2/Schle gel_wireframe_8-cell.png/400px-Schlegel_wireframe_8-cell.png

When to Use Clustering

- When you...
 - suspect your data consists of underlying groups
 - cannot collect labels for classification
 - can define "similar" points via a distance metric

Resources

- Geometric view of clustering: https://shapeofdata.wordpress.com/category/clustering/
- K-Means
 - <u>https://www.datascience.com/blog/k-means-clustering</u>
 - <u>https://www.naftaliharris.com/blog/visualizing-k-means-clustering/</u>
 - Elbow method: <u>https://dataskeptic.com/blog/episodes/2016/the-elbow-method</u>
- DBSCAN
 - <u>https://www.naftaliharris.com/blog/visualizing-dbscan-clustering/</u>
 - http://lineardigressions.com/episodes/2017/11/19/dbscan
- EM
 - http://hameddaily.blogspot.com/2015/03/when-not-to-use-gaussian-mixtures-model.html
 - https://pythonmachinelearning.pro/clustering-with-gaussian-mixture-models/
- Other clustering approaches
 - http://lineardigressions.com/episodes/2016/2/29/t-sne-reduce-your-dimensions-keep-your-clusters
- TF Example: <u>http://playground.tensorflow.org</u>
- Image compression: <u>http://scikit-learn.org/stable/auto_examples/cluster/plot_color_quantization.html</u>

Thanks!



Image: https://toptal.io/

Lab Session
Recap

- Clustering is an unsupervised approach
- DBSCAN
 - Intuitive parameters
 - Difficult to control number of clusters
- K-Means
 - Easy to control number of clusters
 - Sensitive to outliers
- Gaussian EM
 - Probabilistic cluster assignment

Evaluating Clustering Models

- How can we find the "best" parameters?
 - Grid Search
- How can we evaluate cluster quality without labeled data?
 - Silhouette Score
- Code Lab

Grid Search



Image: http://jmlr.csail.mit.edu/papers/volume13/bergstra12a/bergstra12a.pdf

Grid Search



K-Means

DBSCAN

Grid Search

- Pros
 - Very simple
 - Works well with a small number parameters
- Cons
 - Assumes all parameters are equally important
 - Not always true!
 - When this assumption is violated, consider random search

There are many other advanced parameter search methods, too

- Revisiting the question "how do I choose k?"
- General scoring method for clustering algorithms
 - Works for models other than K-Means too



- *a(i)*: average dissimilarity of *i* to all other objects in cluster A
- *d(i, X)*: average dissimilarity between *i* and objects in cluster X
- $b(i) = \min_{X \neq A} d(i, X)$

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

- *a*: mean distance between a sample & all other points in the **same cluster**
- *b*: mean distance between a sample & all other points in the **next nearest cluster**

$$s(i) = \begin{cases} 1 - a(i)/b(i), & \text{if } a(i) < b(i) \\ 0, & \text{if } a(i) = b(i) \\ b(i)/a(i) - 1, & \text{if } a(i) > b(i) \end{cases}$$

 $-1 \leq s(i) \leq 1$

• Only defined for n>1 and k>1

$$s(i) = \begin{cases} 1 - a(i)/b(i), & \text{if } a(i) < b(i) \\ 0, & \text{if } a(i) = b(i) \\ b(i)/a(i) - 1, & \text{if } a(i) > b(i) \end{cases}$$

- Close to 1: $a \ll b$
 - Dissimilarity *within* clusters ≪ smallest dissimilarity *between* clusters
 - Objects are well-clustered
- Close to 0: *a* ≅ *b*
 - This object could belong to either of two clusters without affecting quality
- Close to -1: *b* ≪ *a*
 - \circ $\;$ This object is very different from the others in its cluster
 - \circ $\$ Likely assigned to the wrong cluster





К = З









• Pros

- Simple to evaluate and compare
- Allows comparison of different clustering methods
- Tends to correlate with qualitative judgment of "good" clusters
- Cons
 - Can be sensitive to data scaling
 - Scale your features to equal "importance"
 - Average silhouette score can mask differences between clusters
 - Biased toward many small clusters

Many other cluster quality metrics exist, too

Resources

- Parameter search
 - <u>https://medium.com/rants-on-machine-learning/smarter-parameter-sweeps-or-why-grid-search-is-pl</u> <u>ain-stupid-c17d97a0e881</u>
 - <u>https://blog.sigopt.com/posts/evaluating-hyperparameter-optimization-strategies</u>
- Silhouette score
 - <u>http://scikit-learn.org/stable/auto_examples/cluster/plot_kmeans_silhouette_analysis.html#sphx-gl</u> <u>r-auto-examples-cluster-plot-kmeans-silhouette-analysis-py</u>
 - <u>http://scikit-learn.org/stable/modules/clustering.html#silhouette-coefficient</u>
 - <u>https://kapilddatascience.wordpress.com/2015/12/08/k-mean-clustering-using-silhouette-analysis-</u> with-example-part-3/
 - <u>https://en.wikipedia.org/wiki/Determining_the_number_of_clusters_in_a_data_set</u>