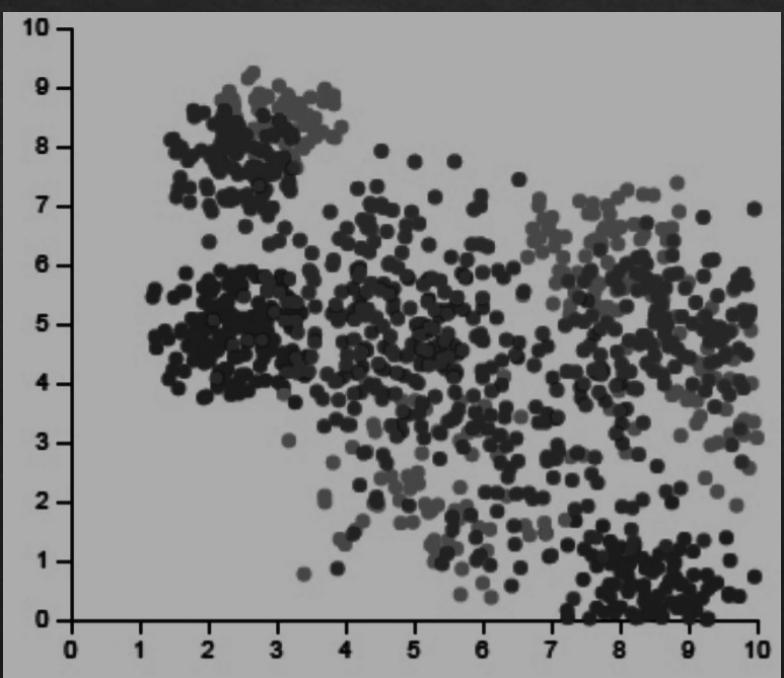
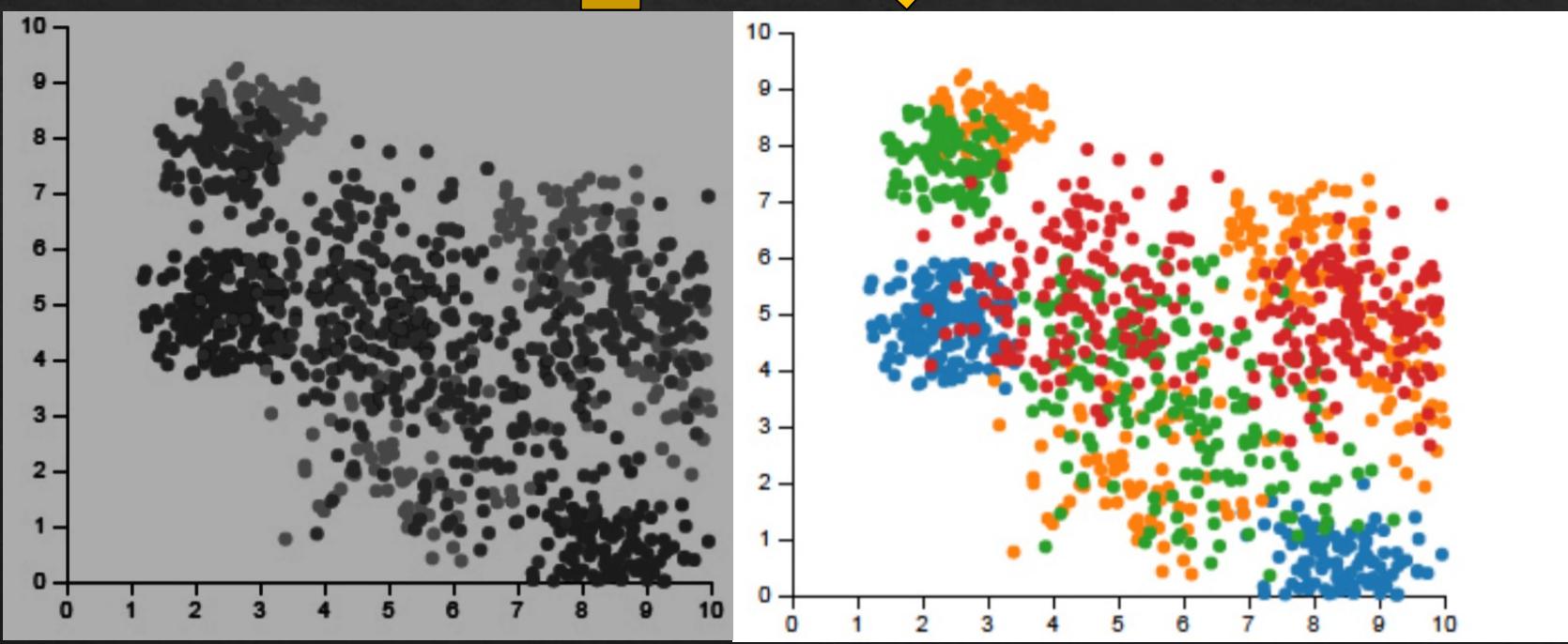
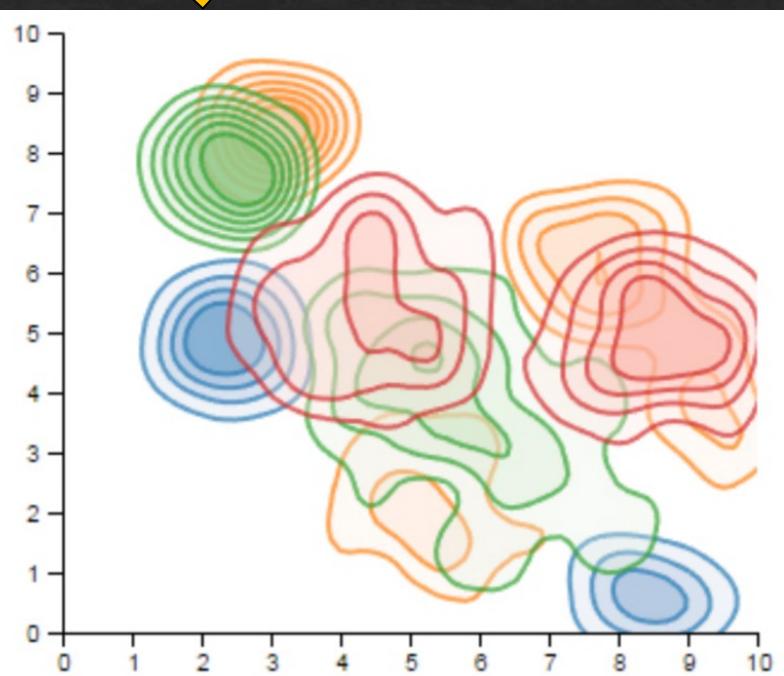
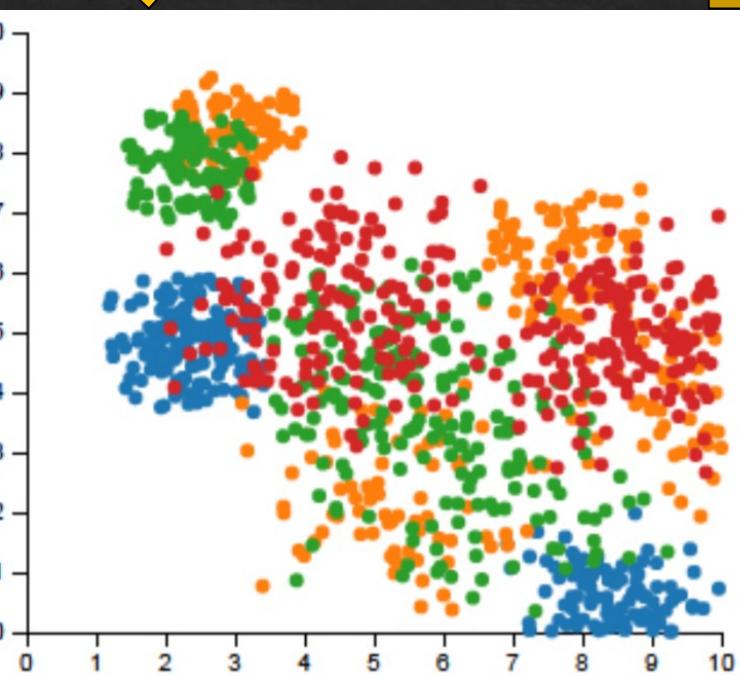
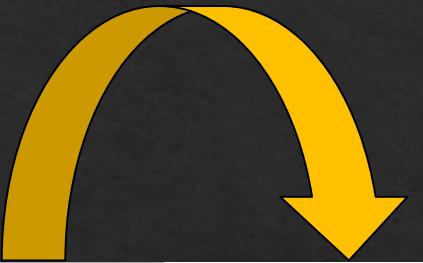
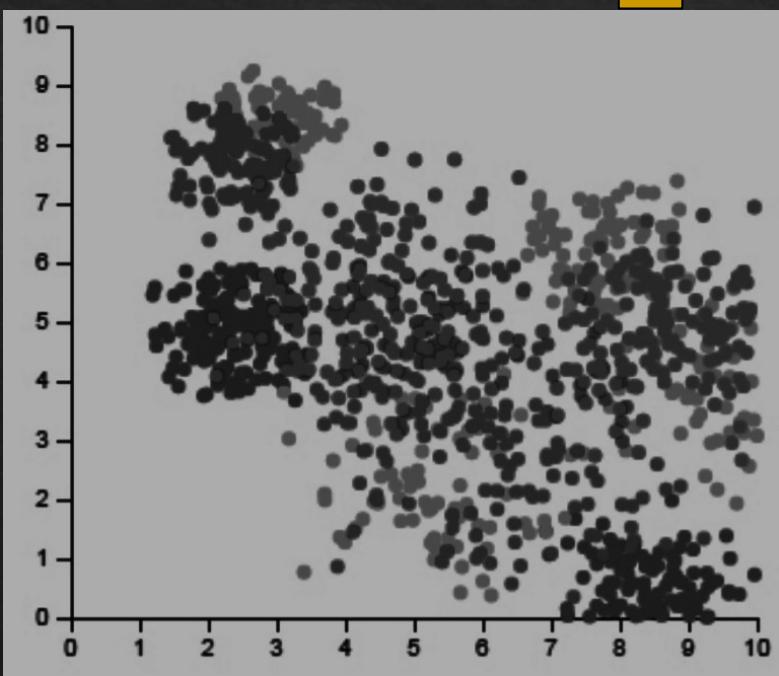


Cluster before You Color

Debajyoti Mondal
University of Saskatchewan







<https://scikit-learn.org/stable/modules/clustering.html>

<https://github.com/d3/d3-contour>

K-means Clustering

K-Means Clustering

Input: A point set S and an integer k

Output: A set of k points: c_1, c_2, \dots, c_k , where the points of S which are closer to c_i form the i th cluster.

K-Means Clustering

Input: A point set S and an integer k

Output: A set of k points: c_1, c_2, \dots, c_k , where the points of S which are closer to c_i form the i th cluster.

Why not just output a set of k random points?

K-Means Clustering

Input: A point set S and an integer k

Output: A set of k points: $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k$, where the points of S which are closer to \mathbf{c}_i form the i th cluster.

Why not just output a set of k random points?

Minimize the mean distance between data points and their cluster centroid, i.e., within-cluster sum of squared distances.
$$\sum_{j=1}^k \sum_{x \in C_j} \|x - \mathbf{c}_j\|^2$$

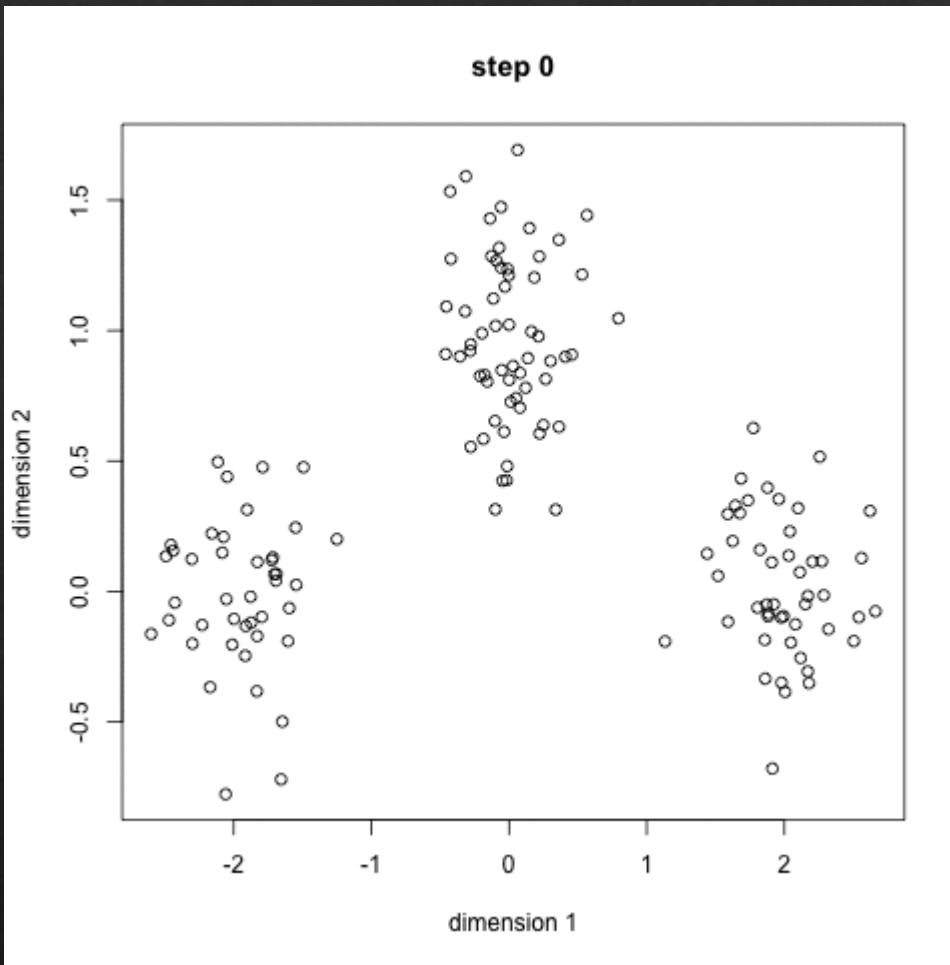
K-Means Clustering

- To begin, we first select a number of classes/groups to use and randomly initialize their respective center points.
- Classify each data point: classify the point to be in the group whose center is closest to it.
- Based on these classified points, we recompute the group center by taking the mean of all the points in the group.
- Repeat these steps for a set number of iterations or until the group centers do not change much between iterations.

K-Means Clustering

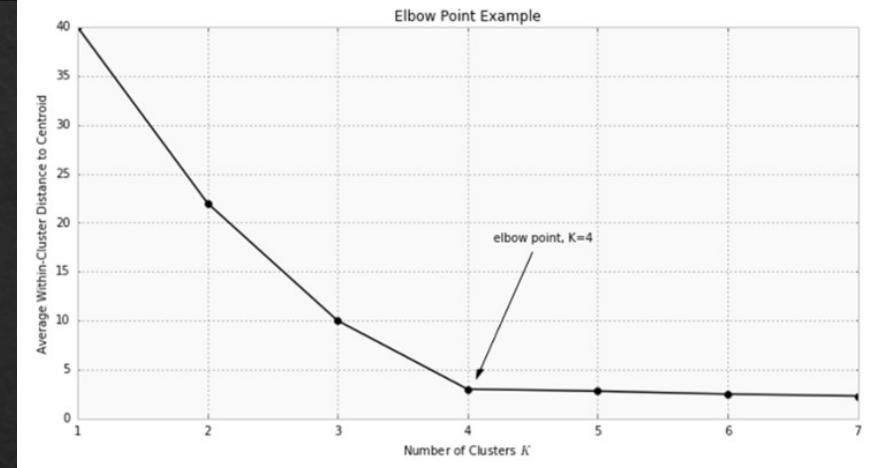
- 
- To begin, we first select a number of classes/groups to use and randomly initialize their respective center points.
 - Classify each data point: classify the point to be in the group whose center is closest to it.
 - Based on these classified points, we recompute the group center by taking the mean of all the points in the group.
 - Repeat these steps for a set number of iterations or until the group centers do not change much between iterations.

K-Means Clustering



Choosing K

- One common technique is to use the mean distance between data points and their cluster centroid.
- Increasing K decreases this mean distance metric
- We plot mean distance as a function of K and the elbow point, where the rate of decrease sharply shifts, can be used to roughly determine K.

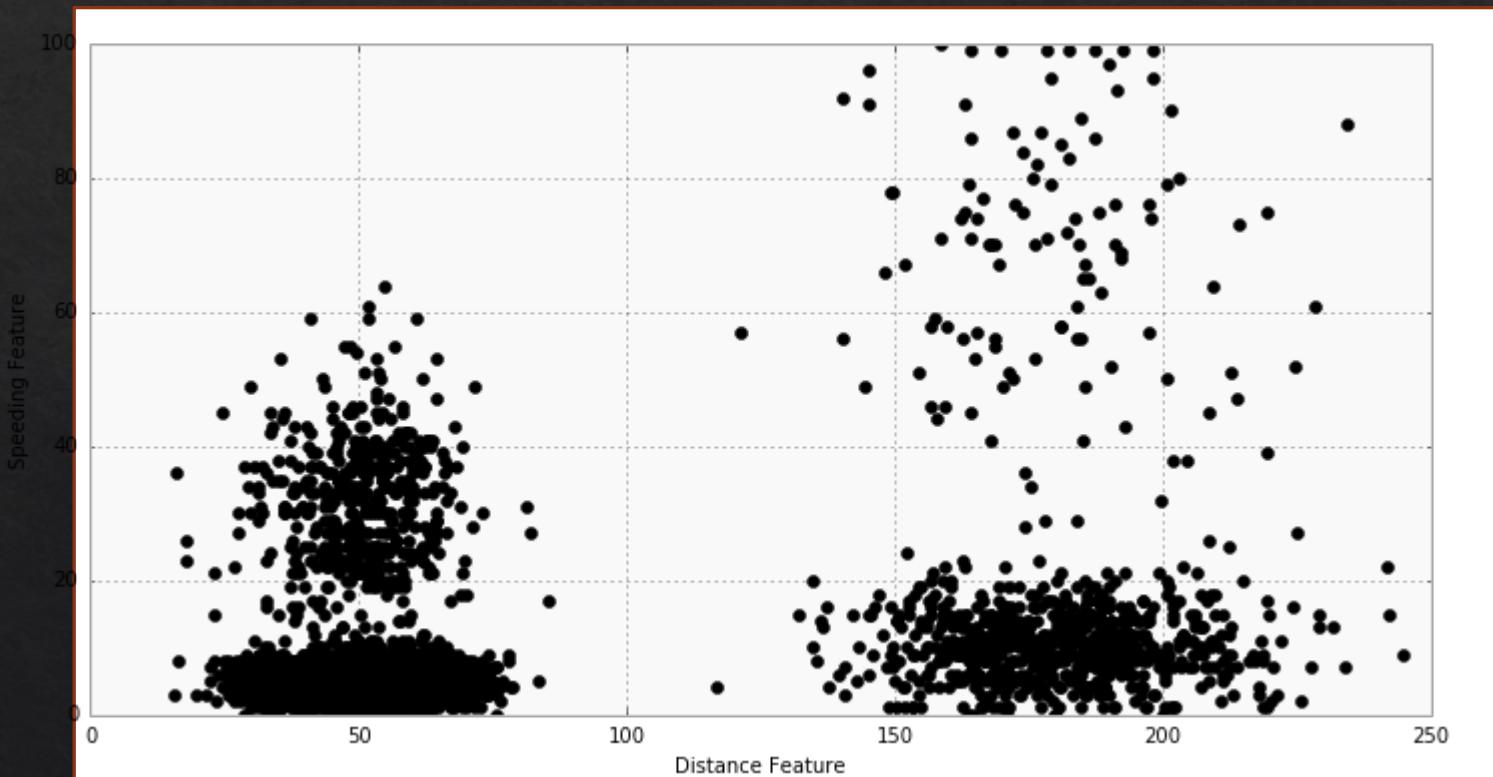


Usages

- Document Classification
 - Delivery Store Optimization
 - Identifying Crime Localities
 - Customer Segmentation
 - Automatic Clustering of IT Alerts
 - Insurance Fraud Detection
 - ...

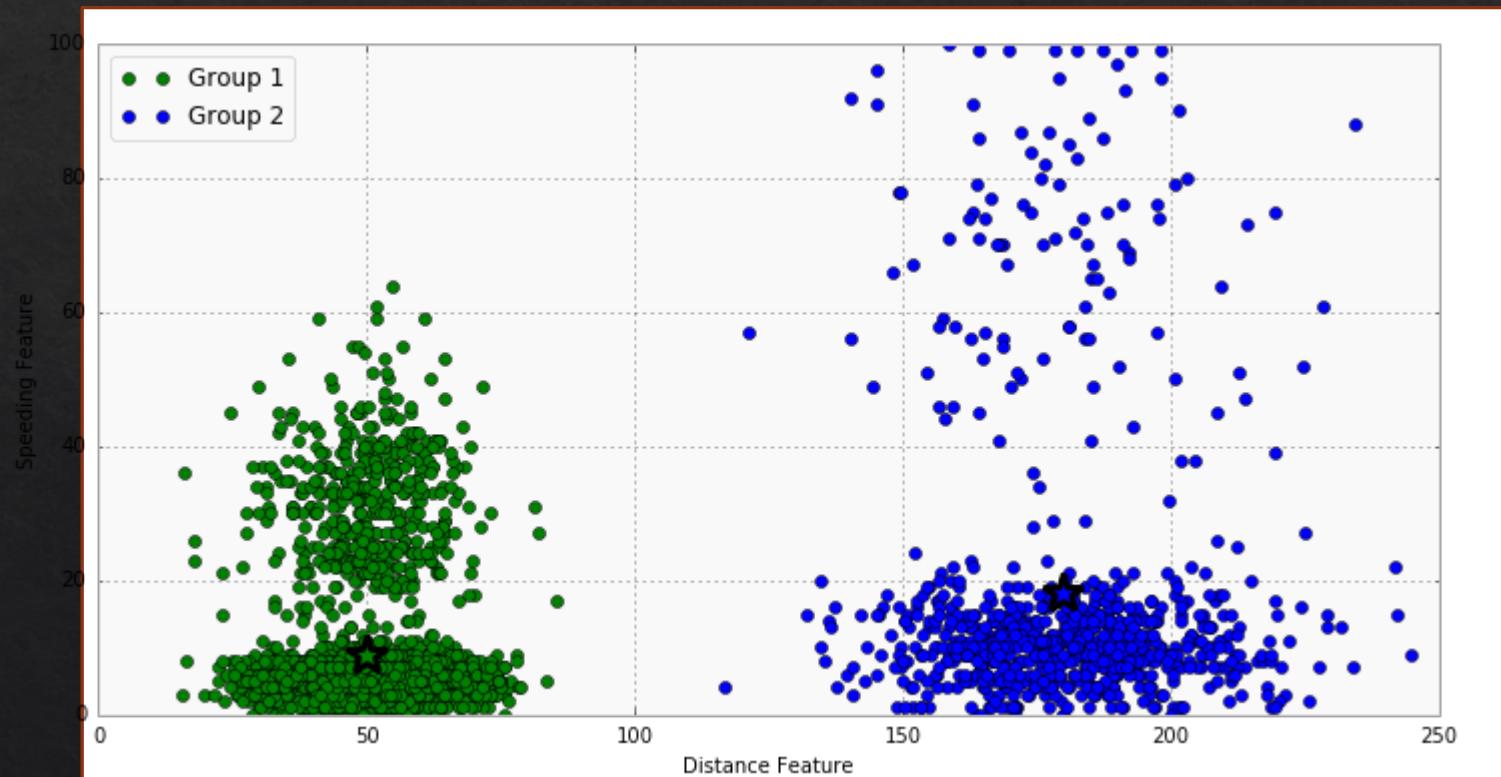
K-means Example

The chart below shows the dataset for 4,000 drivers, with the distance feature on the x-axis and speeding feature on the y-axis.



K-means when K=2

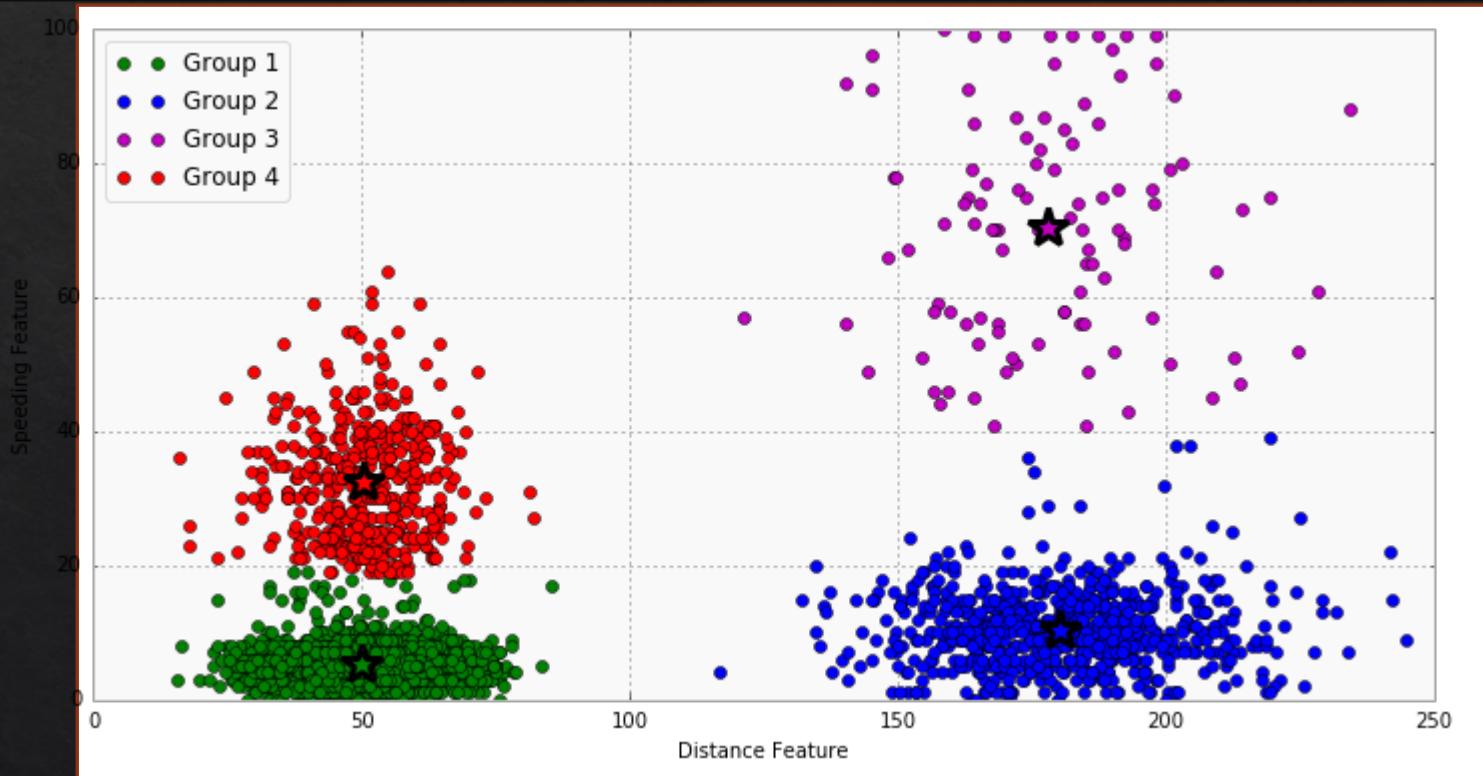
Using domain knowledge of the dataset, we can infer that Group 1 is urban drivers and Group 2 is rural drivers.



K-means when K=4

Speeding drivers have been separated from those who follow speed limits, in addition to the rural vs. urban divide.

The threshold for speeding is lower with the urban driver group than for the rural drivers, likely due to urban drivers spending more time in intersections and stop-and-go traffic.



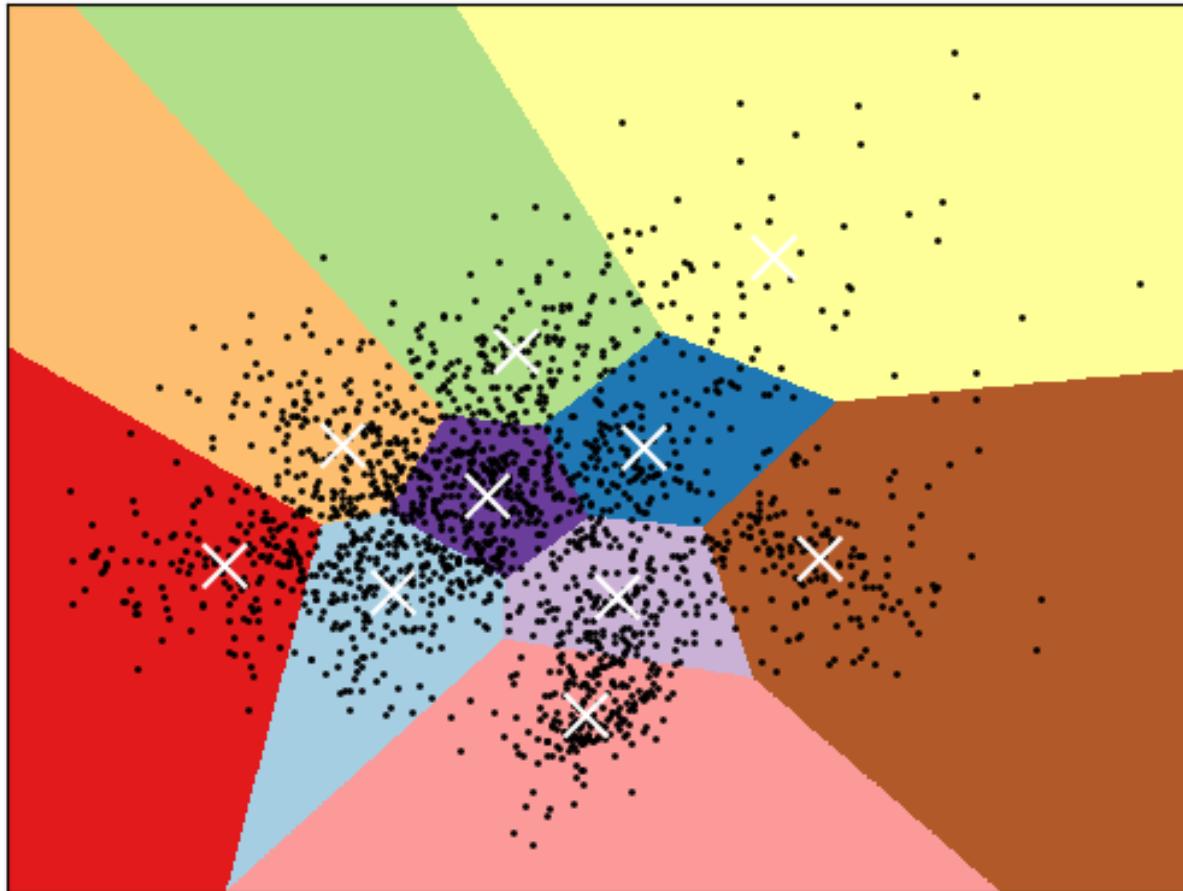
Feature Engineering

*Additional
Information:
not within the
scope of the
course*

- K-means clustering is a form of **unsupervised** machine learning
- Feature engineering is the process of using domain knowledge to choose **which data metrics to input as features** into a machine learning algorithm.
- Using **meaningful features that capture the variability of the data** is essential for the algorithm to find the naturally-occurring groups.
- Categorical data needs to be encoded.

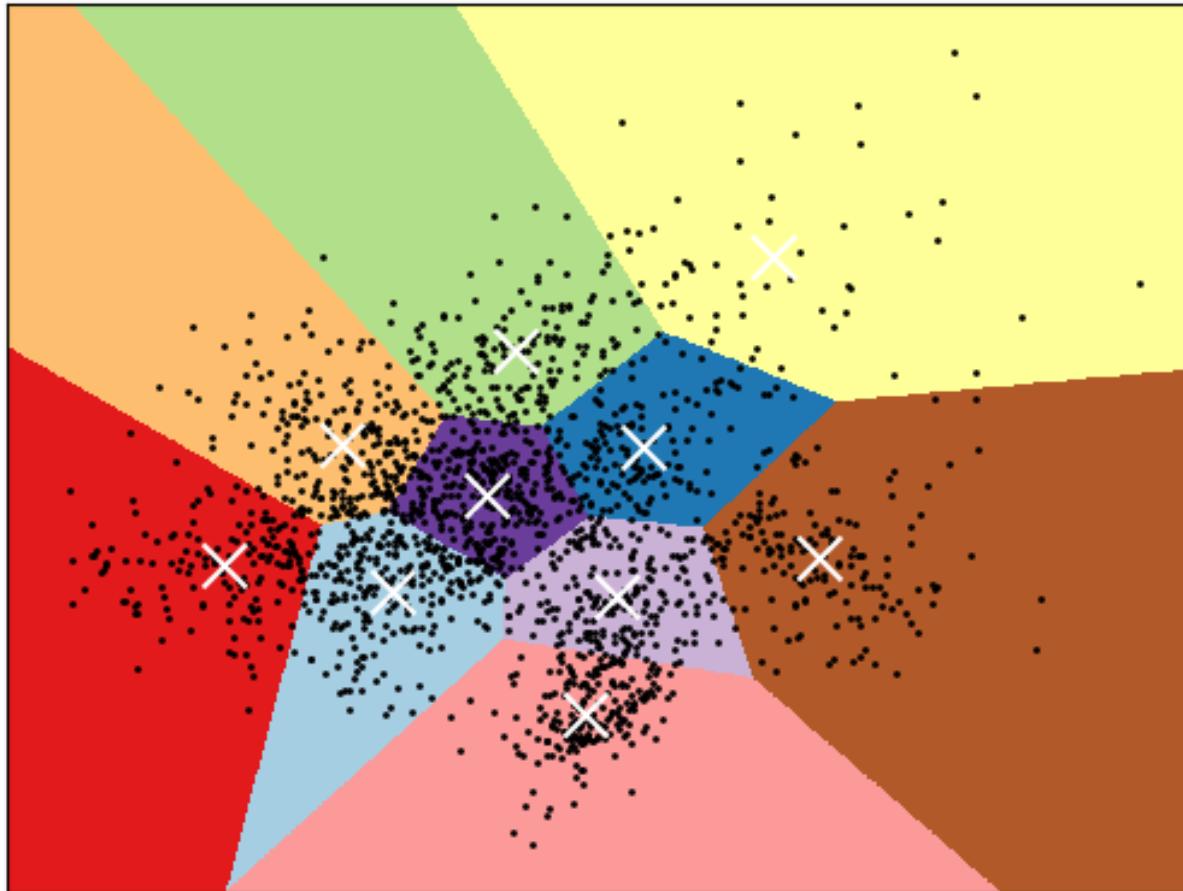
K-Means clustering on the handwritten digits data

K-means clustering on the digits dataset (PCA-reduced data)
Centroids are marked with white cross



The background is a **VORONOI DIAGRAM** and the cells are called **VORONOI CELLS**

K-means clustering on the digits dataset (PCA-reduced data)
Centroids are marked with white cross



Does K-means try to optimize any error?

- The K-means algorithm aims to minimize the within-cluster sum-of-squares:

$$\sum_{j=1}^k \sum_{x \in C_j} \|x - \mu_j\|^2$$

- Running time: $O(I * K * n * d)$
 - n : number of points
 - K : number of clusters
 - I : number of iterations
 - d : number of attributes

Does K-means try to optimize any error?

- The K-means algorithm aims to minimize the within-cluster sum-of-squares:

$$\sum_{j=1}^k \sum_{x \in C_j} \|x - \mu_j\|^2$$

- Running time: $O(IKnd)$, or linear if I,K,d are constant.

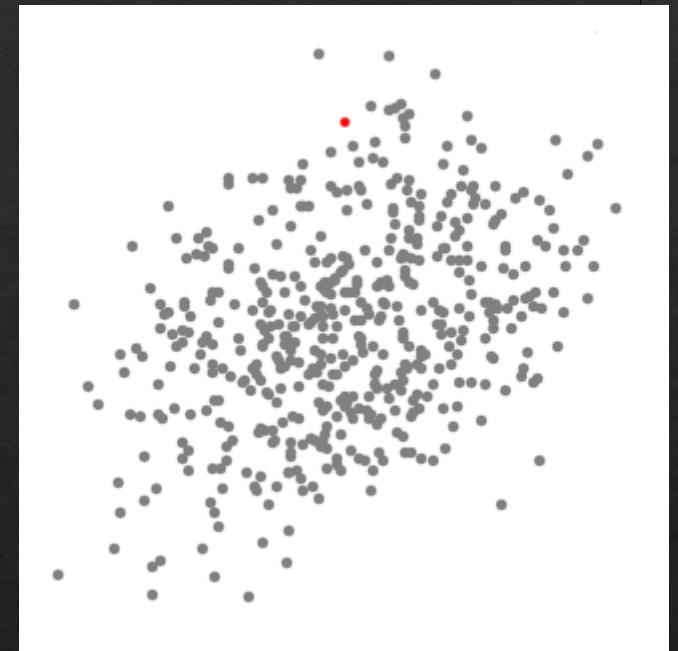
- n : number of points
- K : number of clusters
- I : number of iterations
- d : number of attributes

If d is large, apply
PCA before K-means

Mean-Shift Clustering

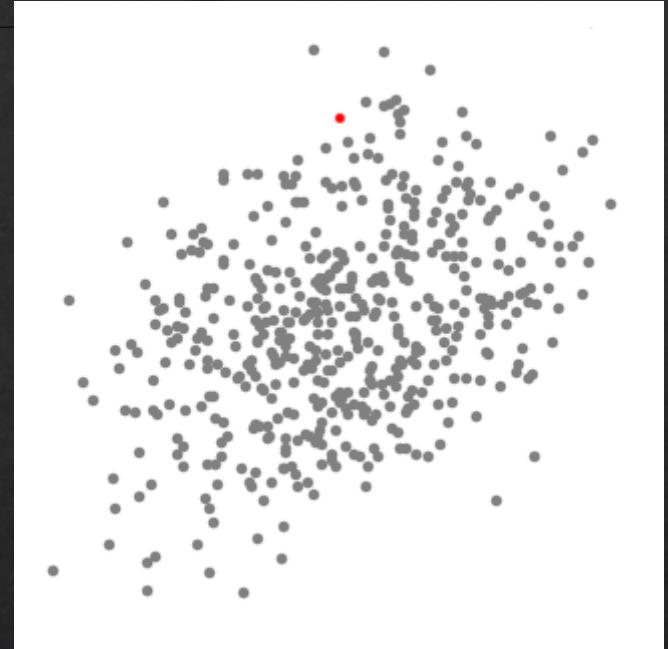
Mean-Shift Clustering

- First consider a set of points in two-dimensional space.
- Start with a circular sliding window centered at a point C (randomly selected) and having radius r determines the kernel.
- Mean-shift shifts this kernel iteratively to a higher density region on each step until convergence.



Mean-Shift Clustering

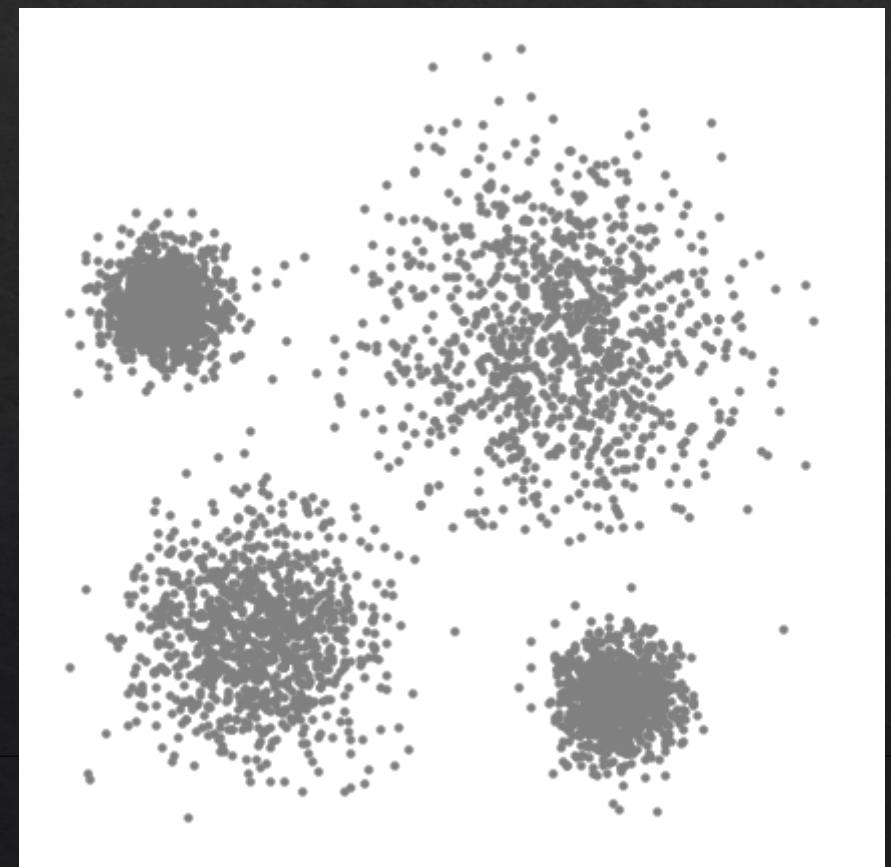
- How to Shift?
- shift the center point to the mean
of the points within the window/circle



By shifting to the mean of the points in the window it will gradually move towards areas of higher point density!

Mean-Shift Clustering

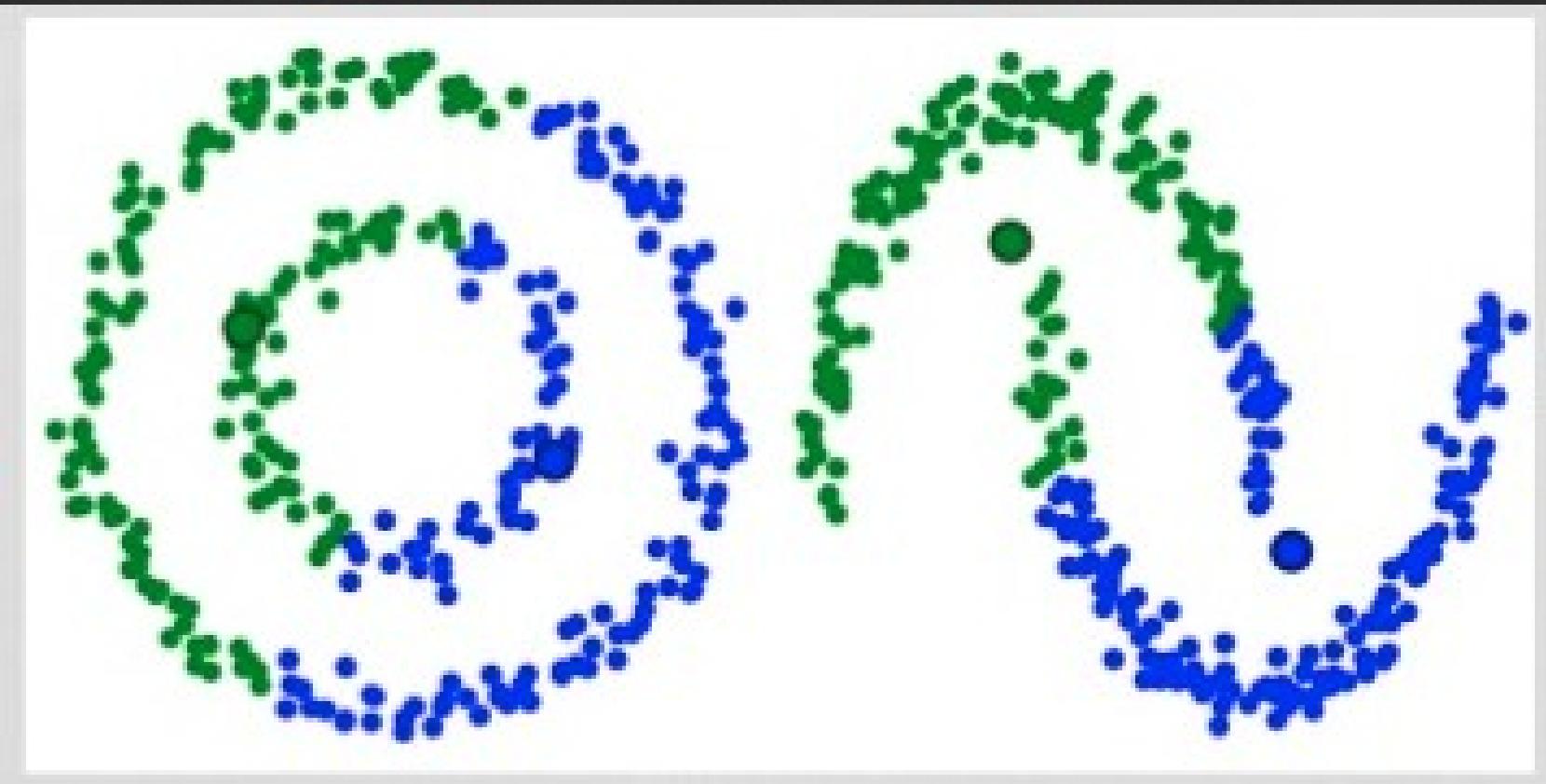
- This process is **done** with many sliding windows until all points lie within a window.
- When **multiple** sliding windows overlap the window containing the most points is preserved.
- The data points are then **clustered** according to the sliding window in which they reside.



Mean-Shift Clustering

- No need to guess the number of cluster centers!
- Slower than K-means --- Why?
- Selection of window size is difficult □

Failure Cases



Class Activity:

How to design a clustering algorithm to separate

the two

O-shapes or U-shapes in the picture?

Density-Based Spatial Clustering of Applications with Noise (DBSCAN)

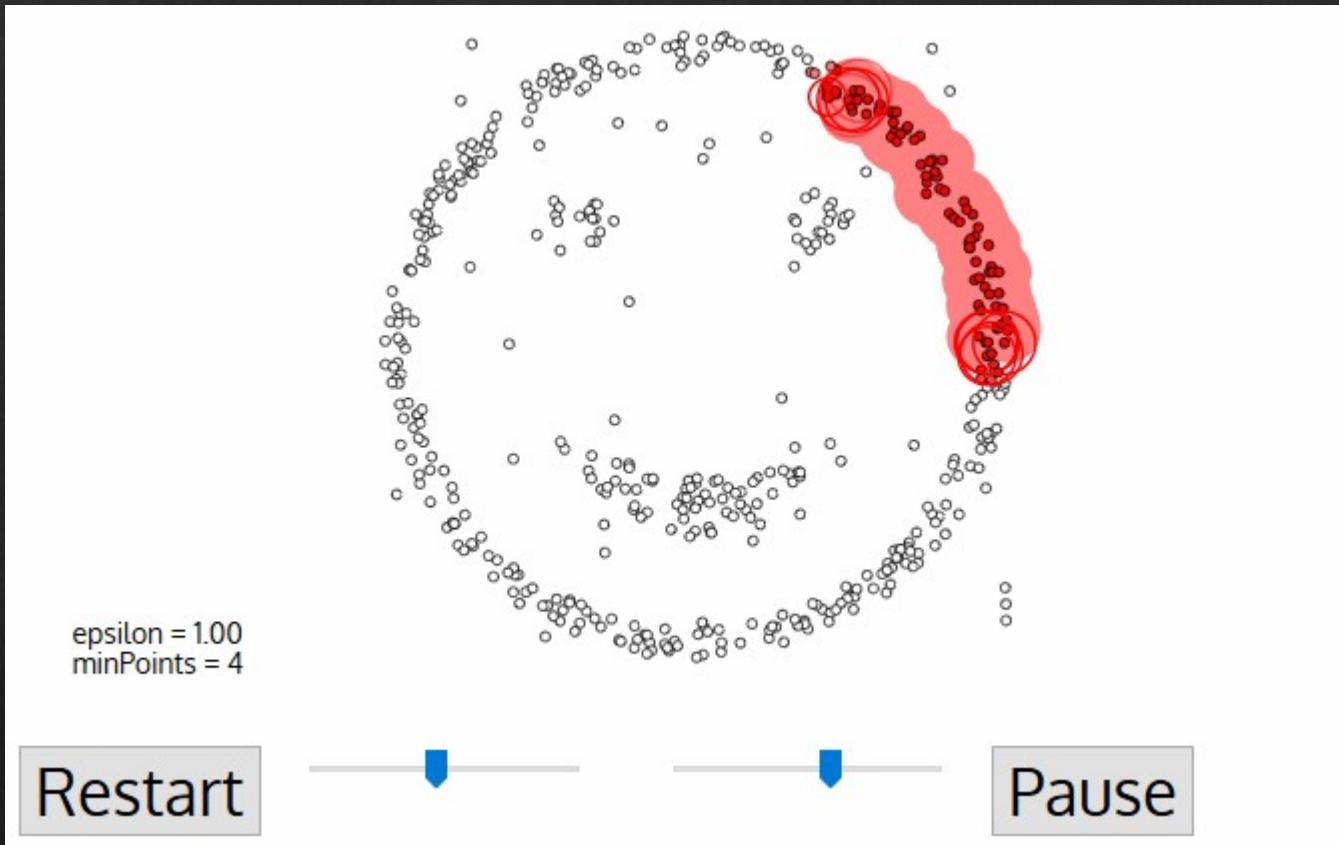
DBSCAN

- Start with an arbitrary starting data point X that has not been visited. The neighborhood of this point is extracted using a distance $\text{epsilon } \varepsilon$.
- If the number of points in the neighborhood $> \text{minPoints}$, then the clustering process starts and the current data point becomes the first point in the new cluster. Otherwise, this datapoint is just a noise.
- The points in the neighborhood of X is assigned the cluster of X and the process is repeated with these newly added points.

DBSCAN

- Once we're done with the current cluster, a new unvisited point is retrieved and processed, leading to the discovery of a further cluster or noise.
- This process repeats until all points are marked as visited.
- At the end of this, all points have been visited, each point well have been marked as either belonging to a cluster or being noise.

DBSCAN



<https://www.naftaliharris.com/blog/visualizing-dbscan-clustering/>

<https://www.naftaliharris.com/blog/visualizing-k-means-cluste>

DBSCAN

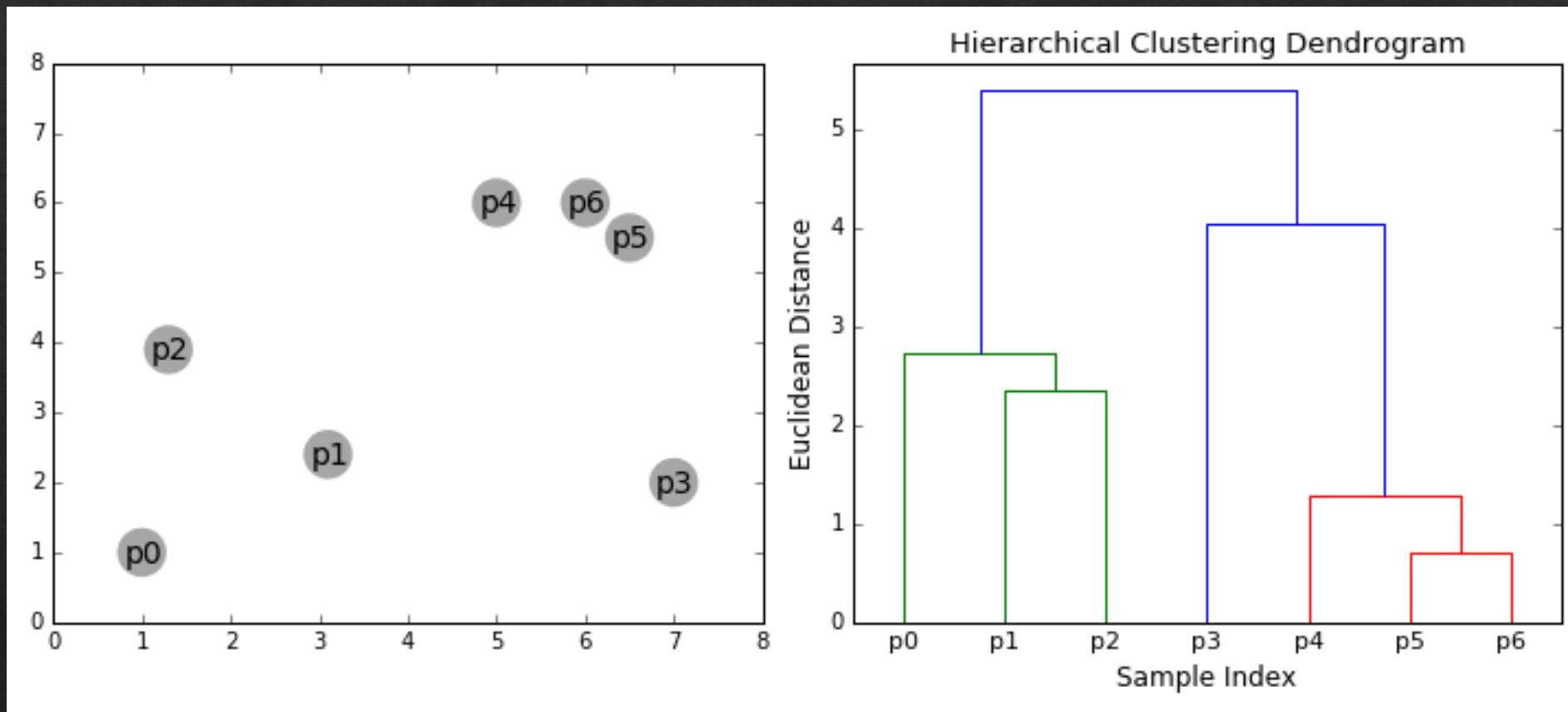
- Does not require a pre-set number of clusters
- Identifies outliers as noises unlike mean-shift which simply throws them into a cluster
- Able to find arbitrarily sized and arbitrarily shaped clusters
- Depending on parameters, maybe slower than K-means [requires query data structure for fast implementation]
- Doesn't perform as well as others when the clusters are of varying density. How to choose epsilon ϵ and minPoint?

Agglomerative Hierarchical Clustering

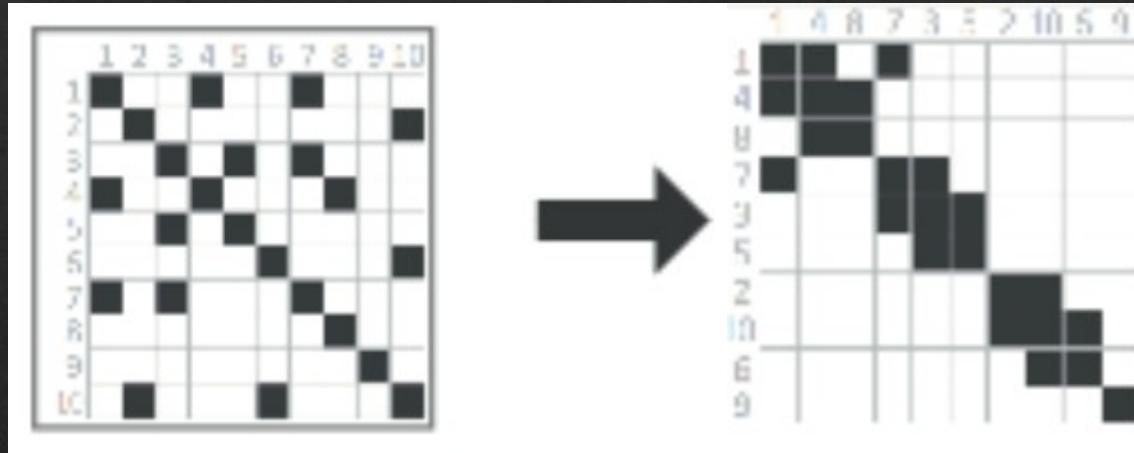
Hierarchical Clustering

- Start by treating each data point as a single cluster
- Choose a **distance metric** that measures the distance between two clusters
- On each iteration we **combine** two clusters into one. The two clusters to be combined are selected as those with the smallest distance
- **Repeat** until we have the expected number of clusters

Hierarchical Clustering Example

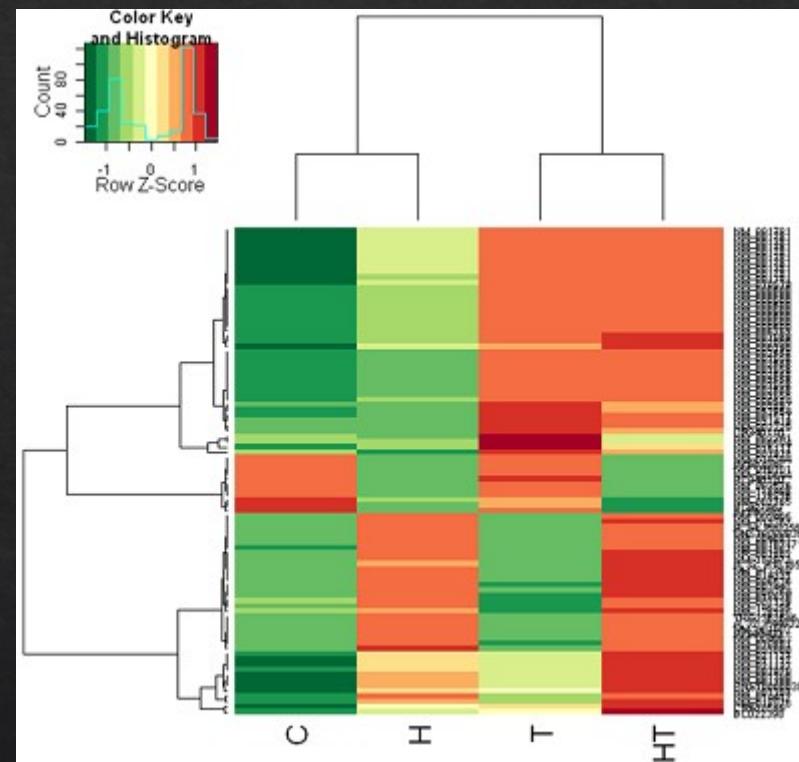
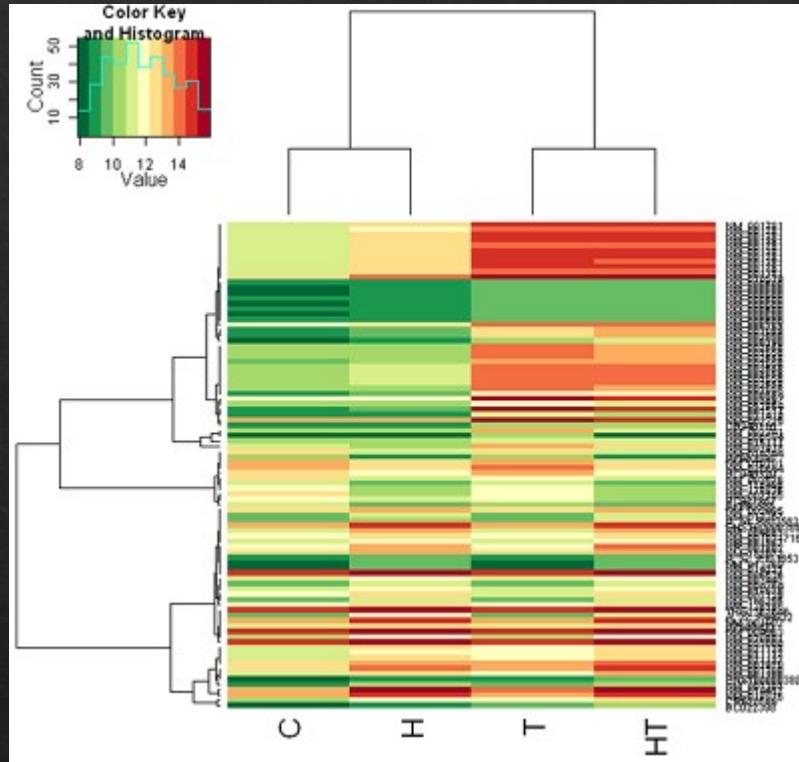


Good Clustering Makes a Difference



We can reorder the rows and columns to reveal cluster information

Good Clustering Makes a Difference

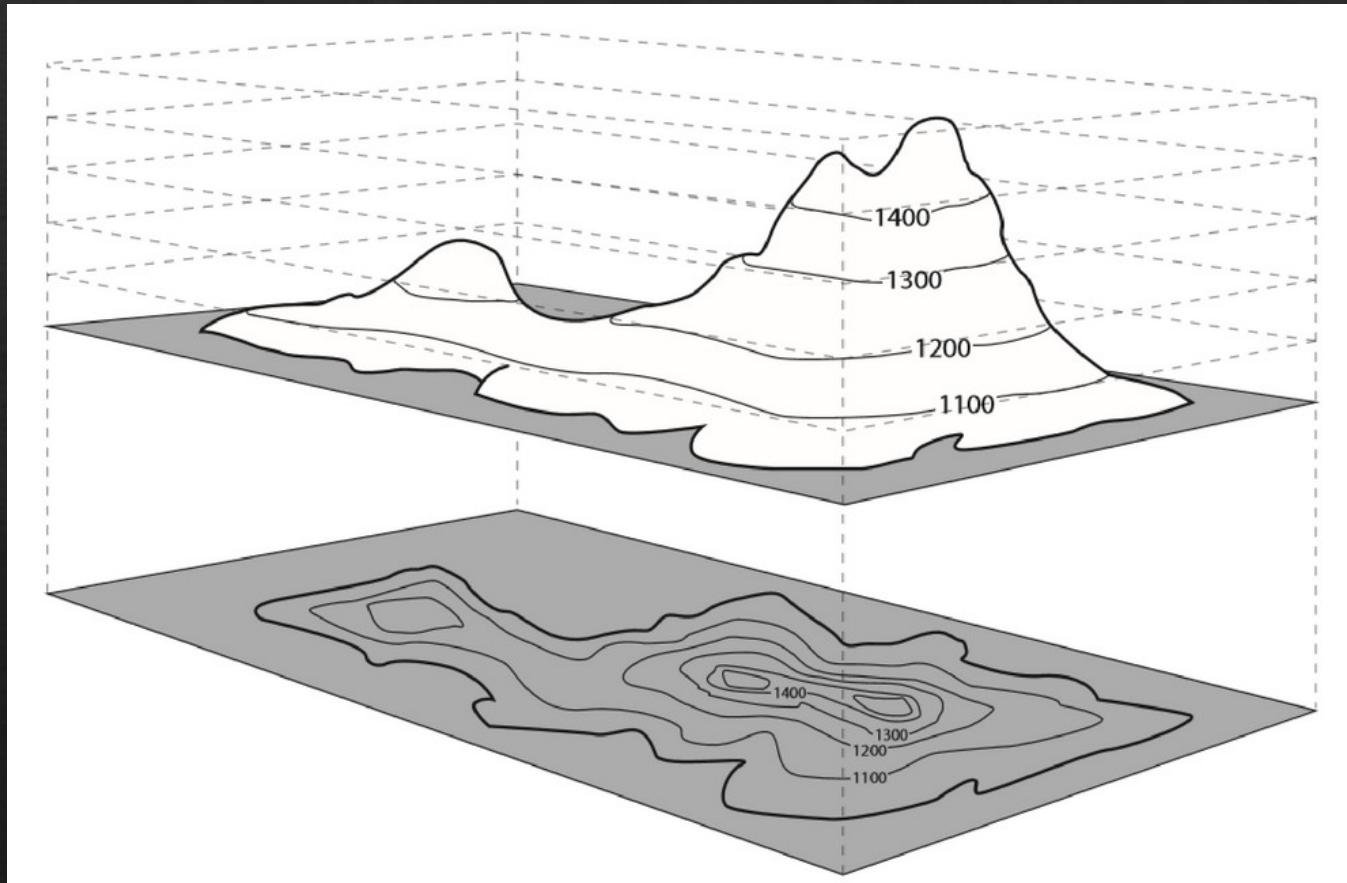


If the cell adjacencies are not important, then clustering on a heatmap may reveal interesting cluster information

Hierarchical Clustering

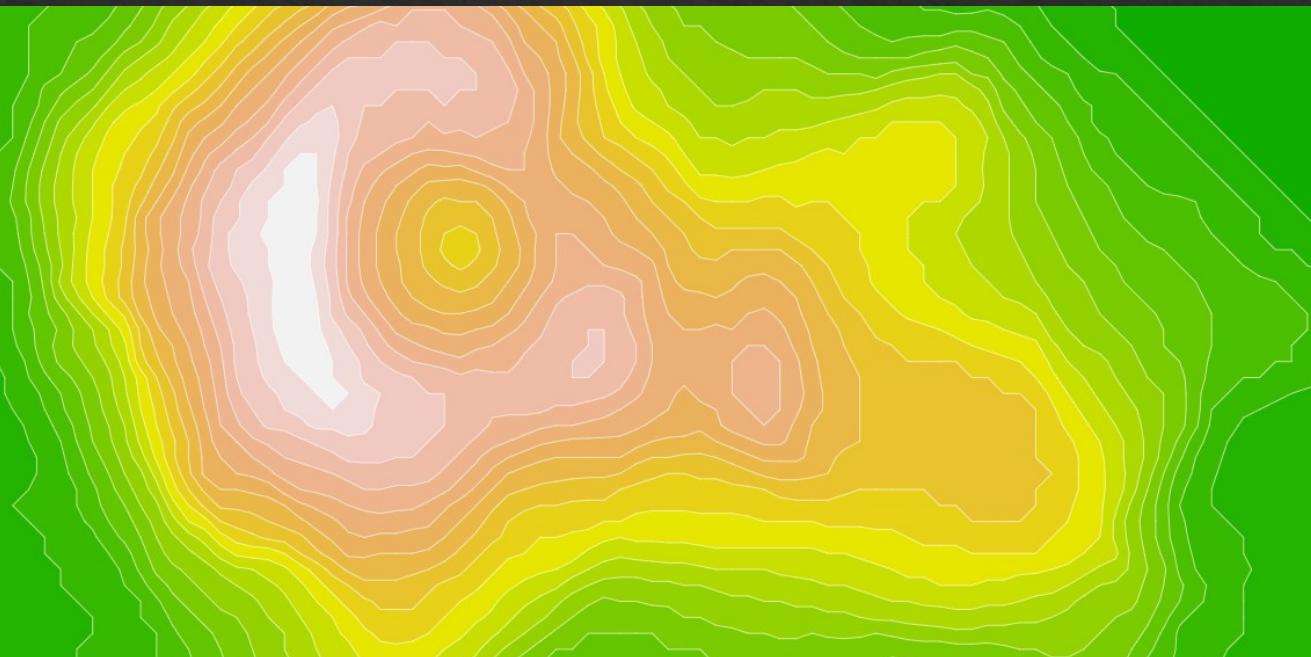
- Does not require a pre-set number of clusters
- Can recover hierarchical structure!
- Dendograms are great for visualization
- Slower – all pair distance computation [requires query data structure for fast implementation]

Contouring



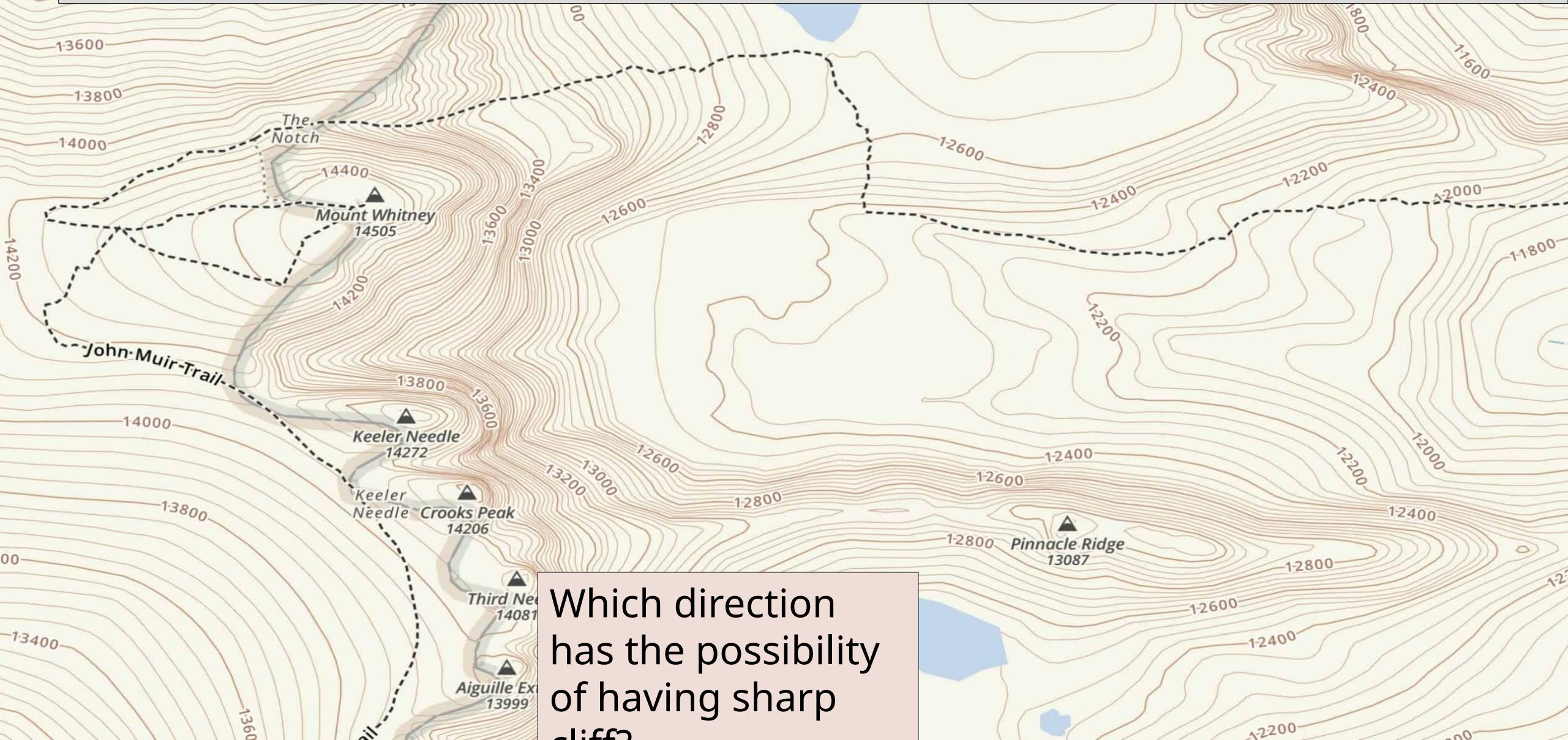
Contouring

Isoline – a line following a single data level, or isovalue



<https://bl.ocks.org/mbostock/4241134>

Locate the most challenging (steep slope) part in the dotted trail



Mount Whitney

nt in
U.S.

14494'

Keeler
Needle

Zone

Iceberg
Lake

MOUNTAINEER'S ROUTE

Pinnacle

Ridge

12250

12500

13500

11500

Scout

▼ Search

ex: Museums in New York, NY

[Get Directions](#) [History](#)

▼ Places

My Places

Sightseeing Tour

Make sure 3D Buildings layer is checked

GPS device

Created Wed Sep 2 09:25:58 2020

Temporary Places

CalTopo Scanned Topos Export (15/16)

GPS device

Created Tue Dec 1 13:21:27 2020

Tracks

Mt Whitney Hike

Mt Whitney Hike

Points

Path

Layers

Primary Database

Announcements

Borders and Labels

Places

Photos

Roads

3D Buildings

Weather

Gallery

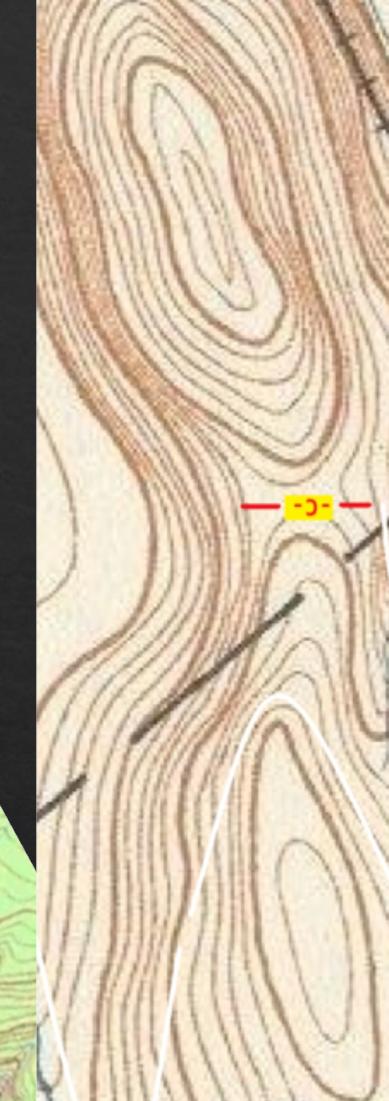
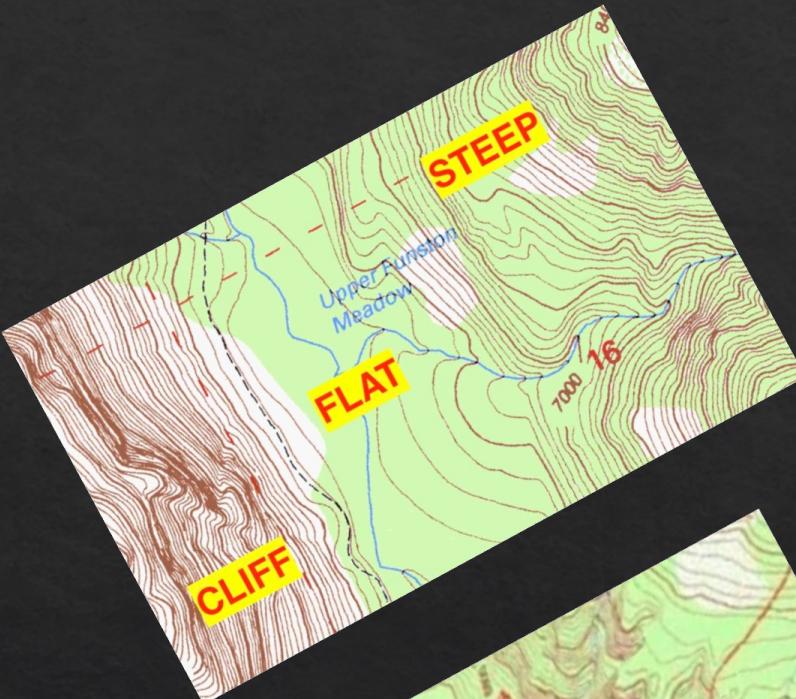
More

Terrain

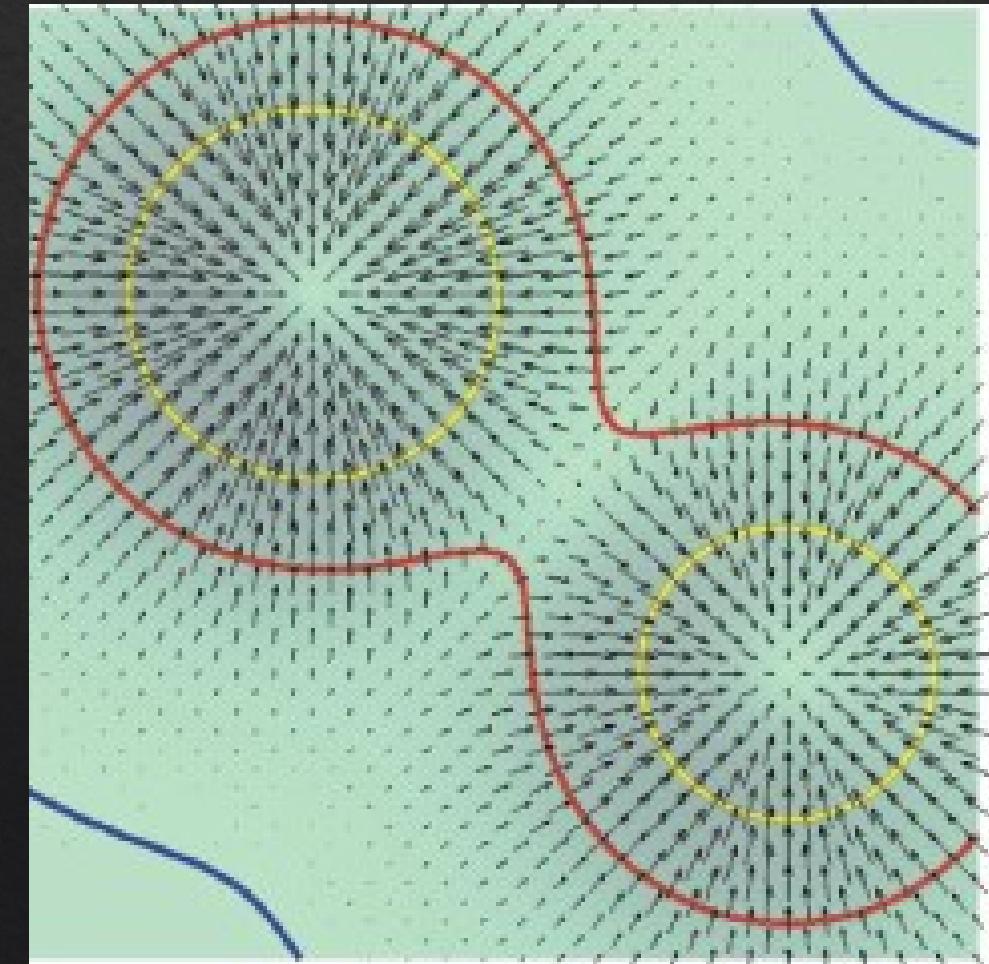
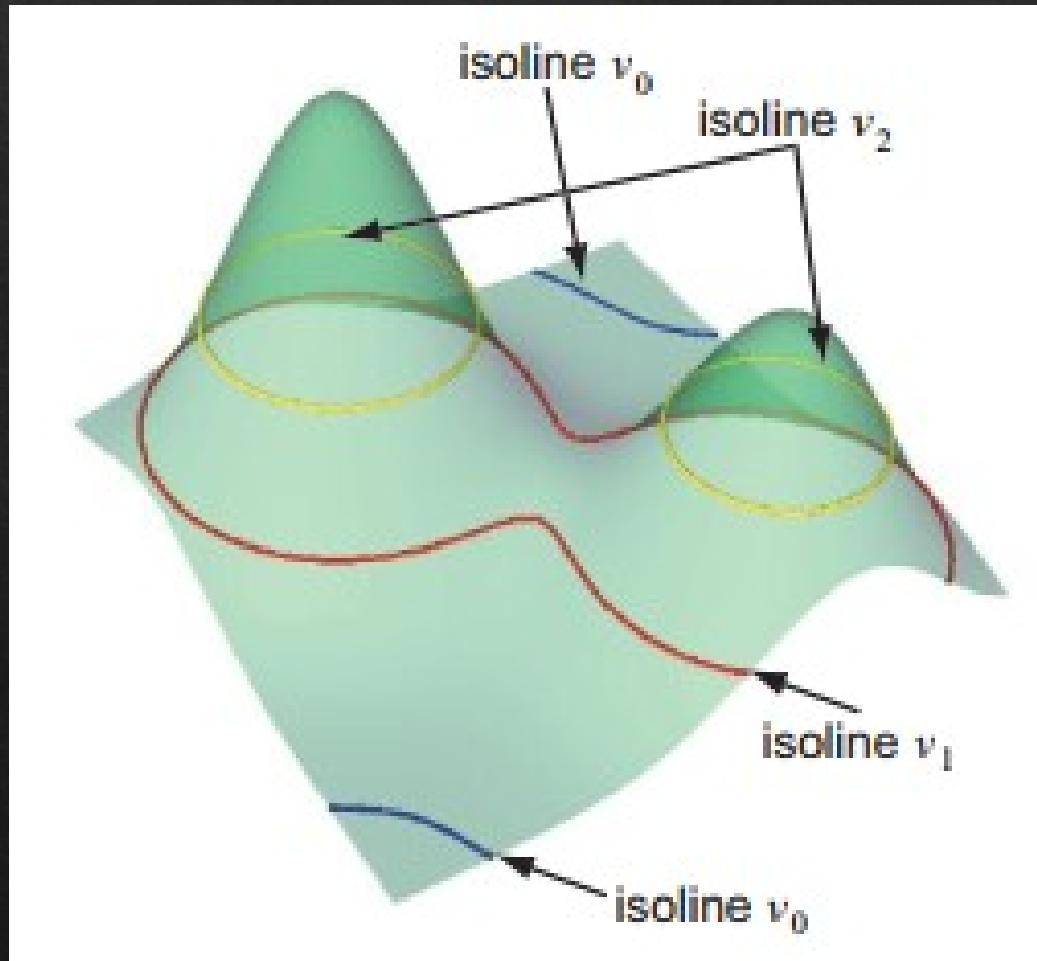




Class Activity:
Draw a hypothetical 2D contour plot for the given
setting



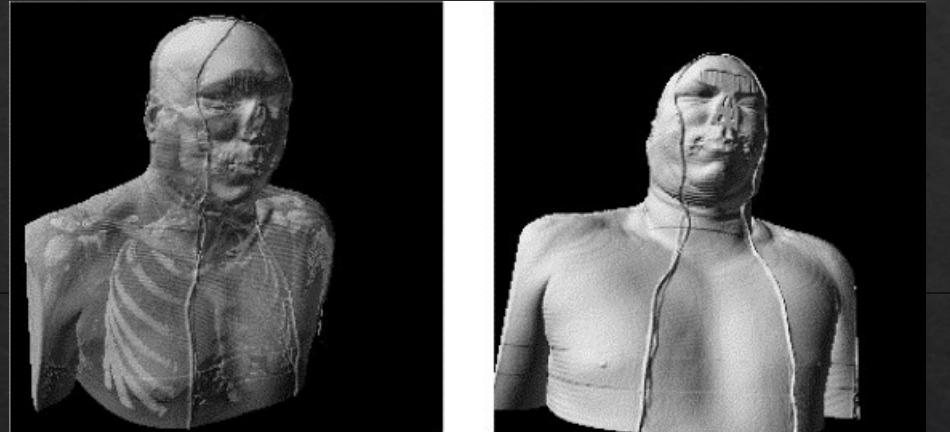
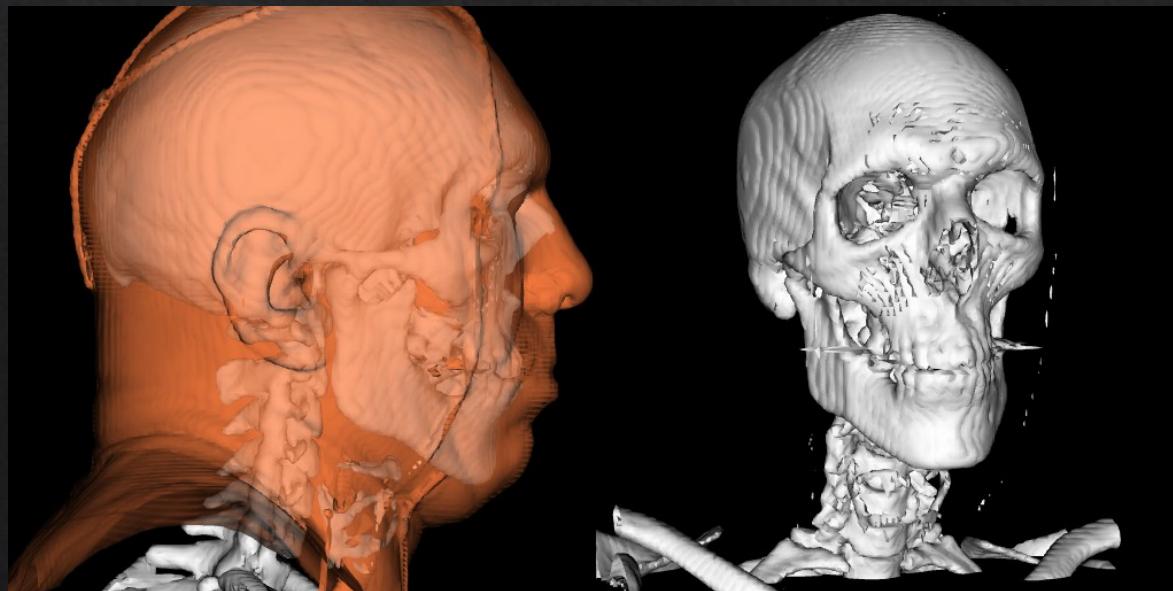
Contouring



Contouring

Isoline – a line following a single data level in 2D, or isovalue

Iso-surface – a surface following
a single data level in 3D



(a) two isosurfaces (skin+bone)

(b) one isosurface (skin)



(c) one isosurface (bone)

Marching Square

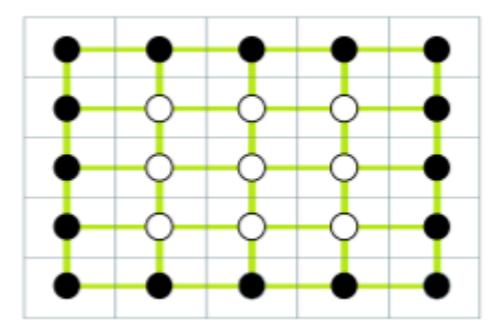


1	1	1	1	1
1	2	3	2	1
1	3	3	3	1
1	2	3	2	1
1	1	1	1	1

Threshold
with iso-value
→
Threshold 1.5

0	0	0	0	0
0	1	1	1	0
0	1	1	1	0
0	1	1	1	0
0	0	0	0	0

Binary image
to cells
→

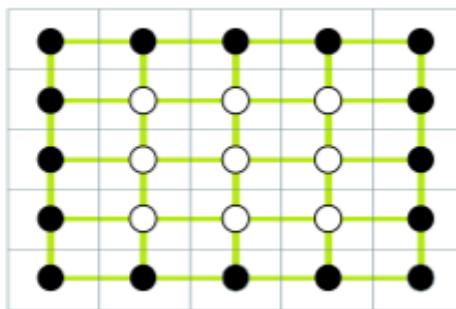


1	1	1	1	1
1	2	3	2	1
1	3	3	3	1
1	2	3	2	1
1	1	1	1	1

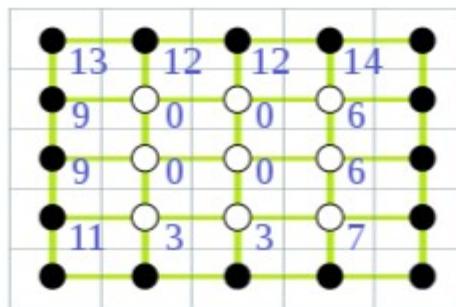
Threshold
with iso-value


0	0	0	0	0
0	1	1	1	0
0	1	1	1	0
0	1	1	1	0
0	0	0	0	0

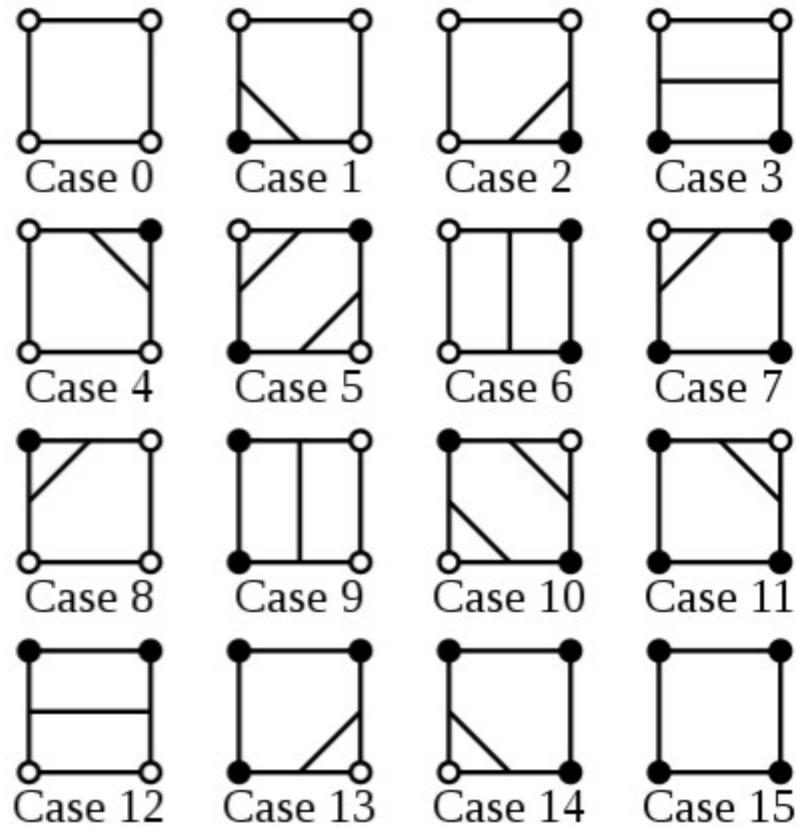
Binary image
to cells

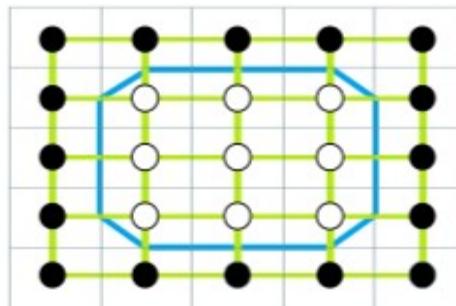
Give every cell a
number based on
which corners are
true/false

Look-up table contour lines



Look up the contour
lines in the database
and put them in
the cells

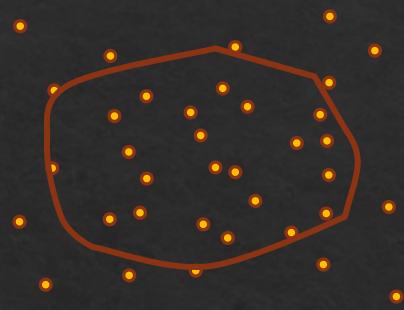
Look at the original
values and use linear
interpolation to
determine a
more accurate position
of all the line end-points


1	1	1	1	1
1	2	3	2	1
1	3	3	3	1
1	2	3	2	1
1	1	1	1	1

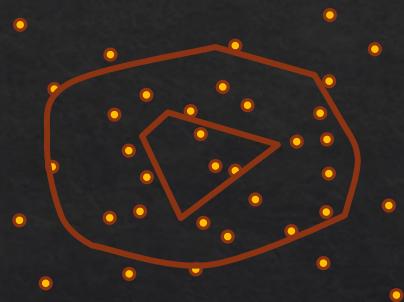
Marching Square



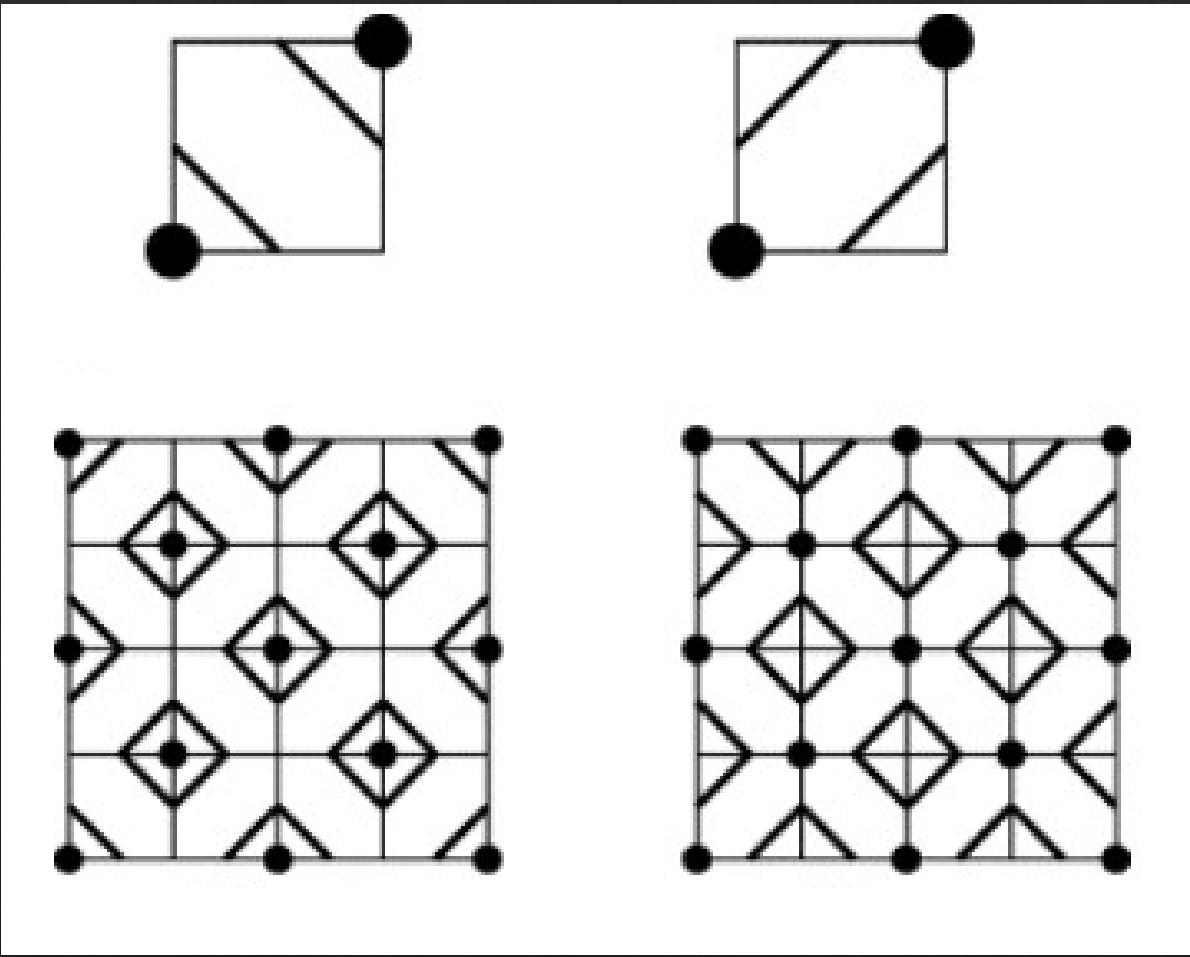
Threshold 1.5



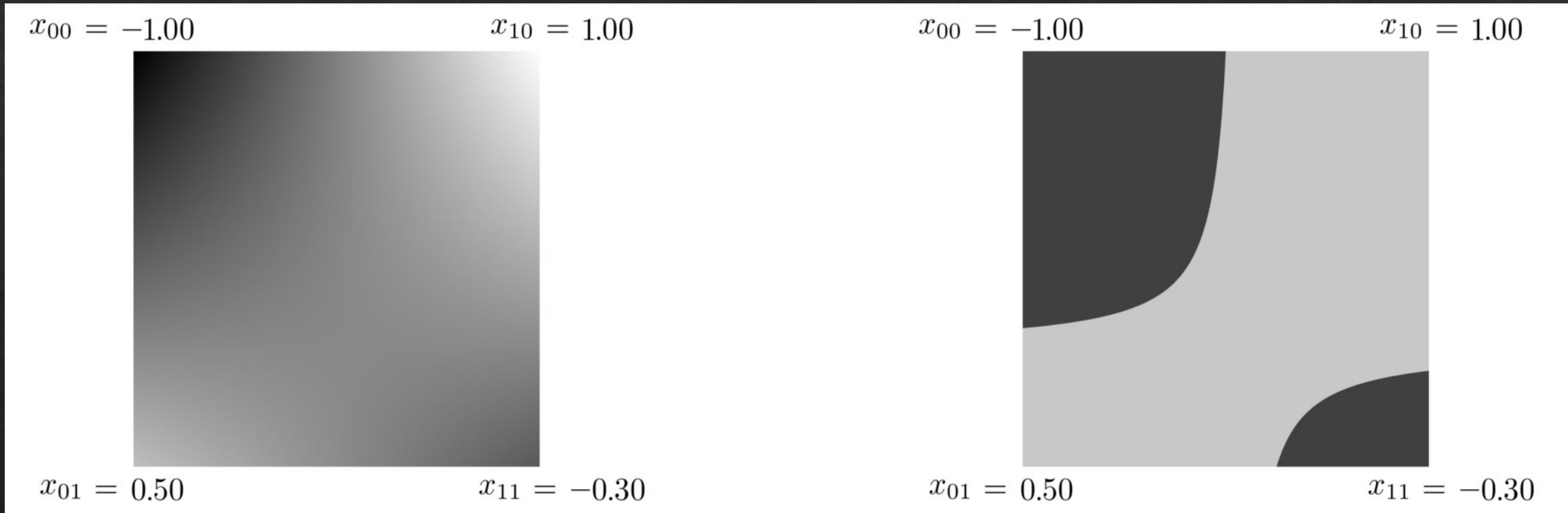
Threshold 2.5



Ambiguity

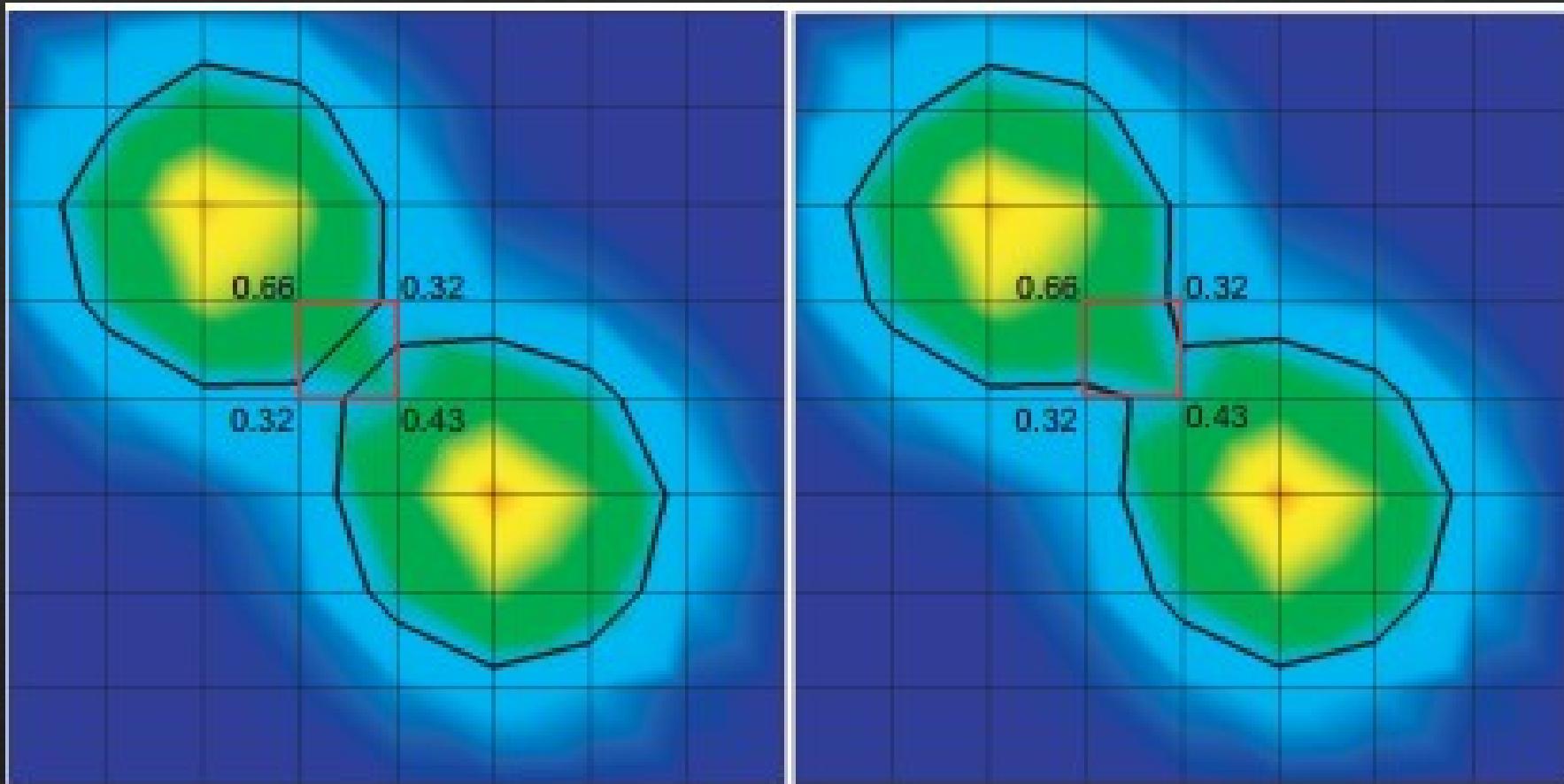


Resolving Ambiguity (interpolate)



Assume the isoline threshold is 0.
The bottom left and top right gets
connected if we make the points above 0
to white.

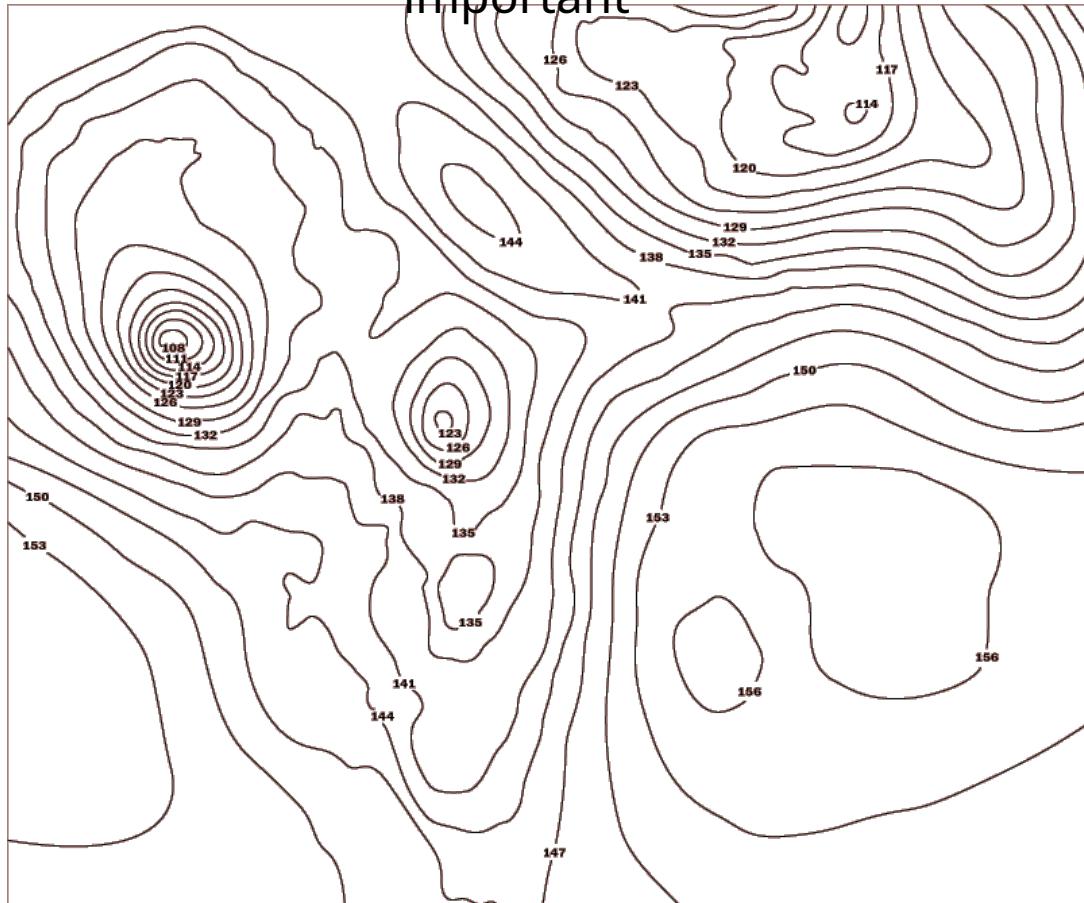
Ambiguity



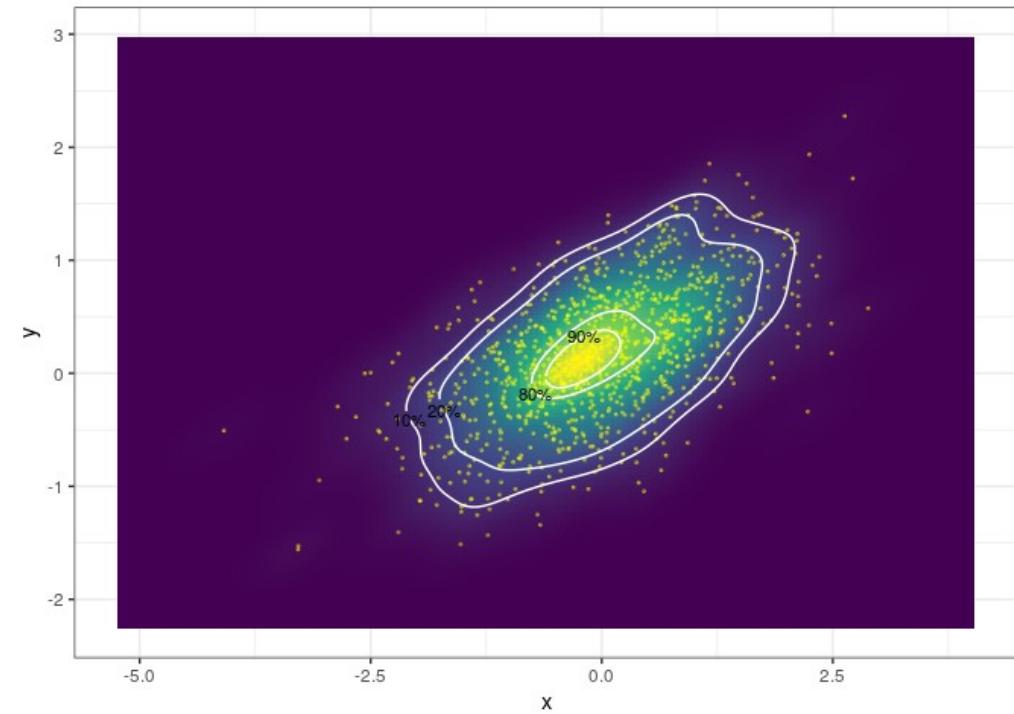
In practice, an implementation may choose either of the connection possibilities separately for each cell. If one has additional knowledge on the contour topology, e.g., that it must have a single connected component, this information can be used to discriminate between the two connection possibilities.

What threshold to use?

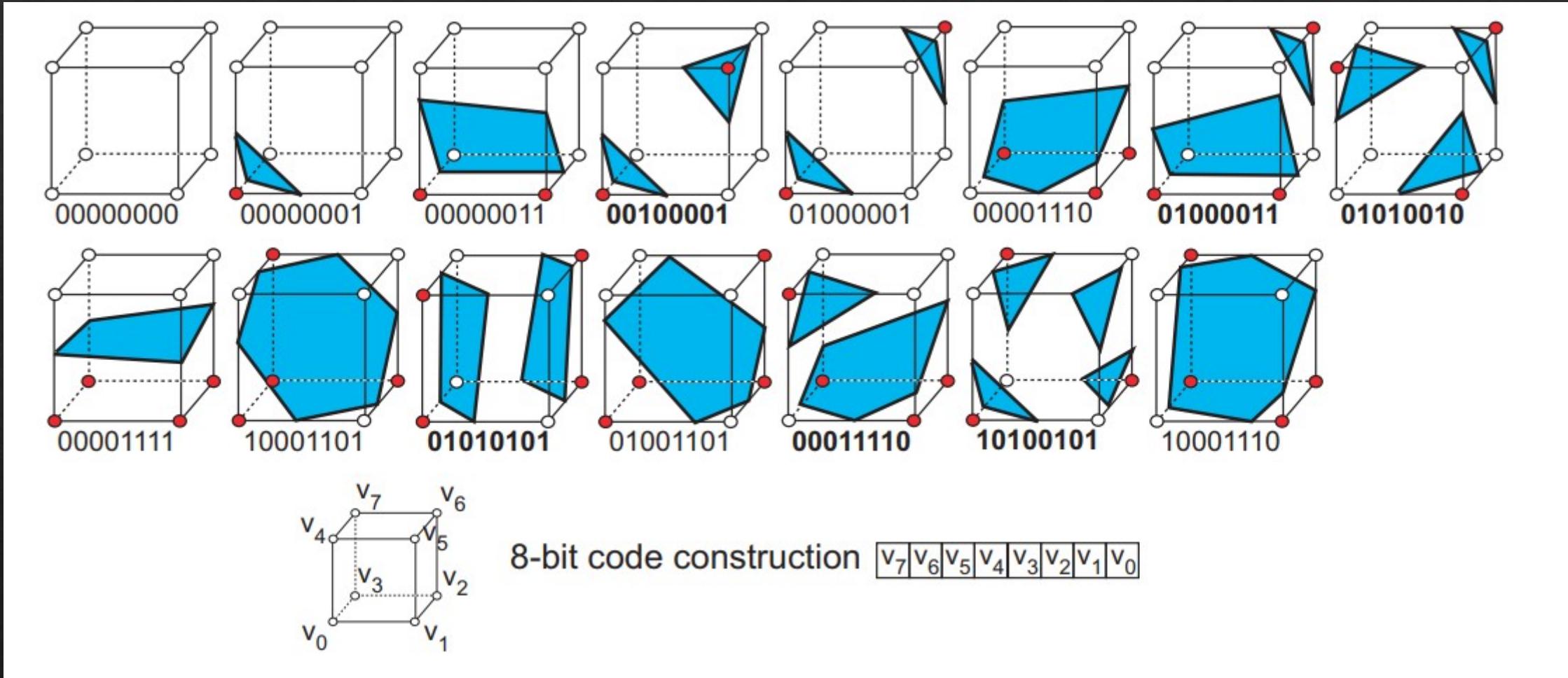
Use regular division when the gradient information is important



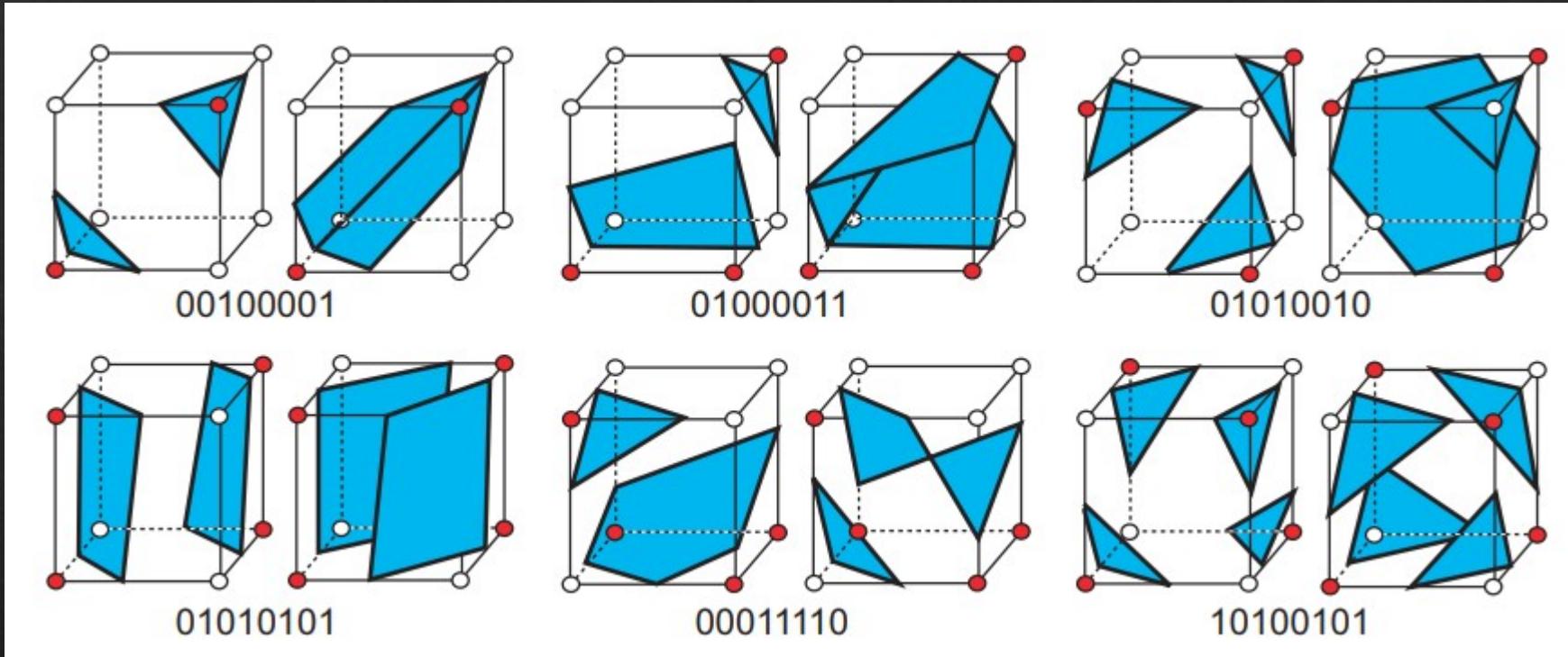
Use percentile when the frequency distribution over



Marching Cube



Ambiguity



Marching Cube in Practice



Two nested isosurfaces of
a tooth scan dataset