

# CPSC 340: Machine Learning and Data Mining

Density-Based Clustering

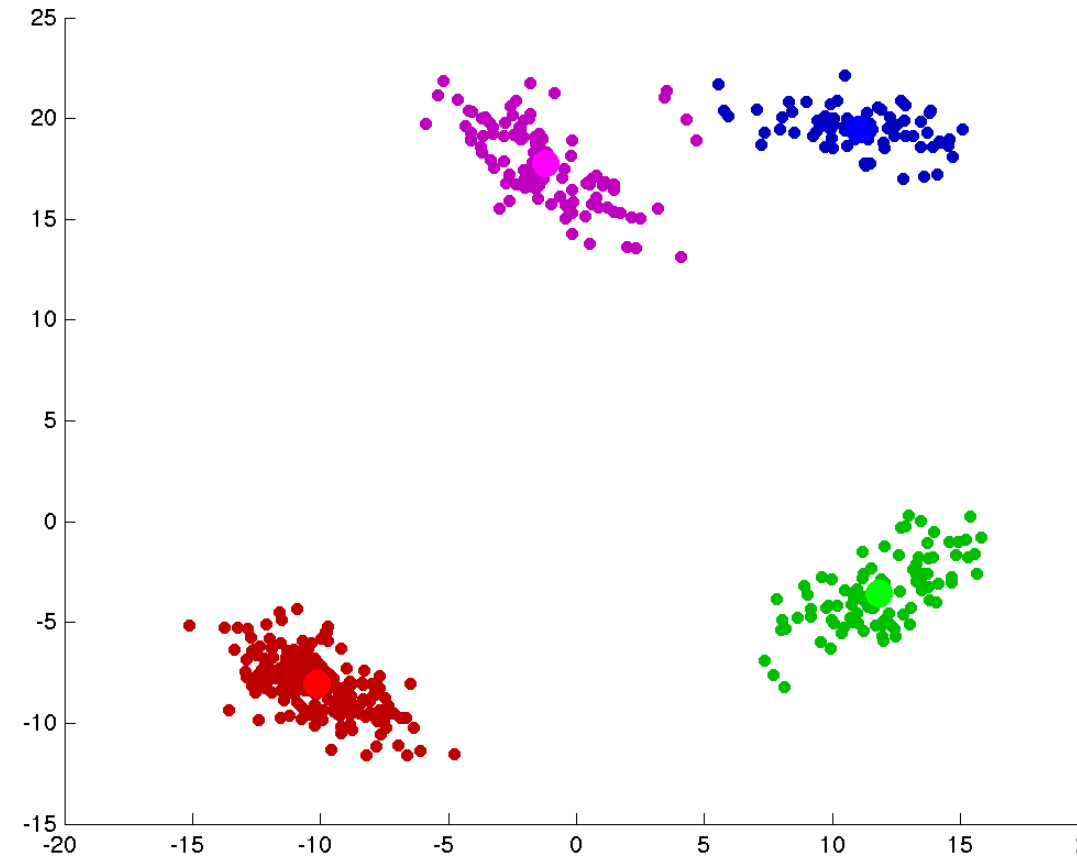
Fall 2017

# Admin

- **Assignment 1** is due Friday.
  - Needs Julie 0.6 (you can use JuliaBox if you can't get Julia/PyPlot working).
  - There was a bug in the decision tree predict function.
  - There was a minor error in the example\_knn.jl function.
  - Follow the assignment guidelines naming convention (a1.zip/a1.pdf).
- Assignment 0 grades posted on Connect?

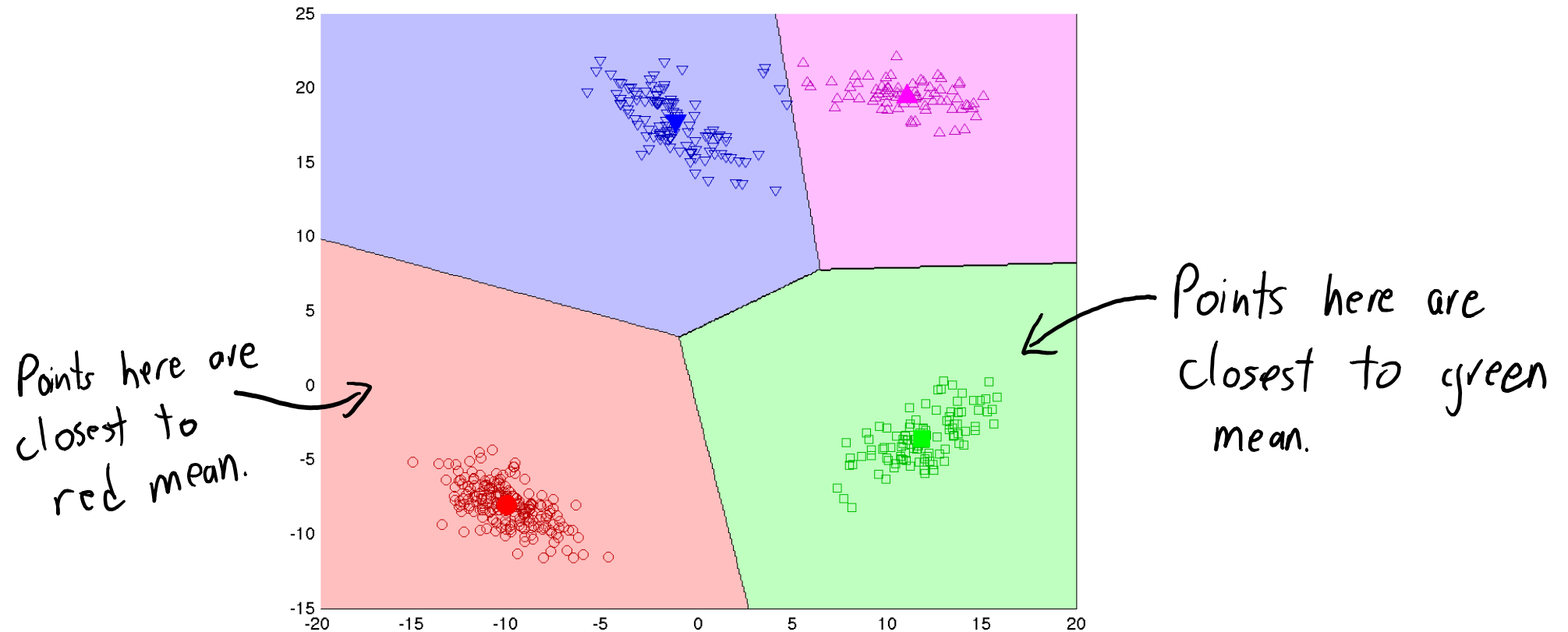
# Last Time: K-Means Clustering

- We want to **cluster** data:
  - Assign objects to groups.
- **K-means clustering**:
  - Define groups by “means”
  - Assigns objects to nearest mean.  
(And updates means during training.)
- Also used for **vector quantization**:
  - Use **means** as “prototypes” of groups.
- Issues with k-means:
  - Fast but sensitive to initialization.
  - Choosing ‘k’ is annoying.



# Shape of K-Means Clusters

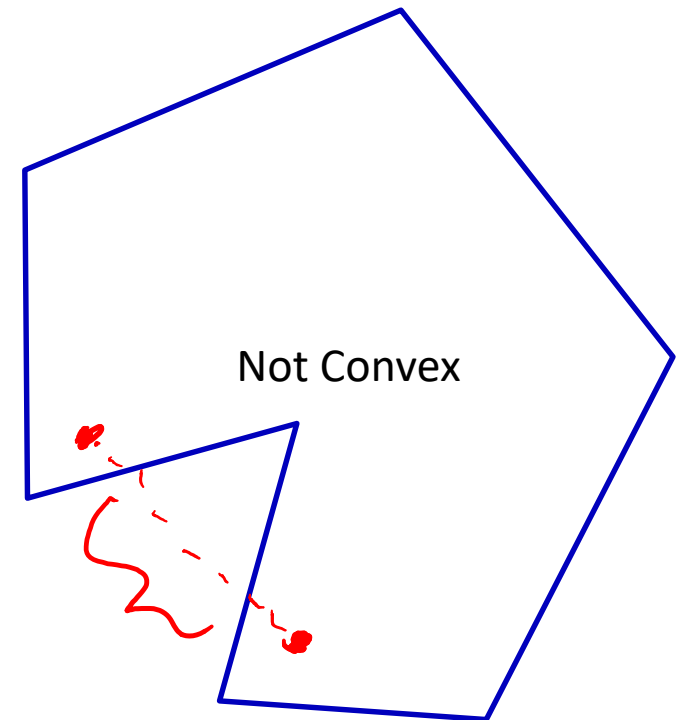
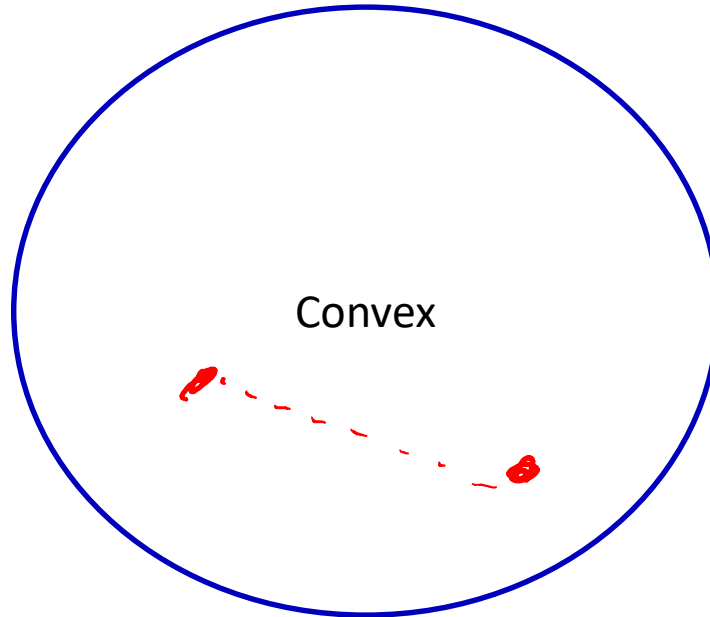
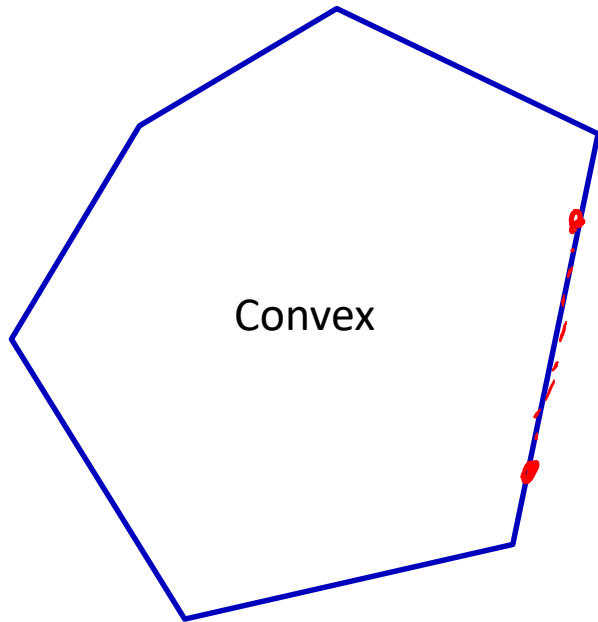
- K-means **partitions the space** based on the “closest mean”:



- Observe that the **clusters are convex** regions.

# Convex Sets

- A set is **convex** if **line between two points in the set stays in the set**.

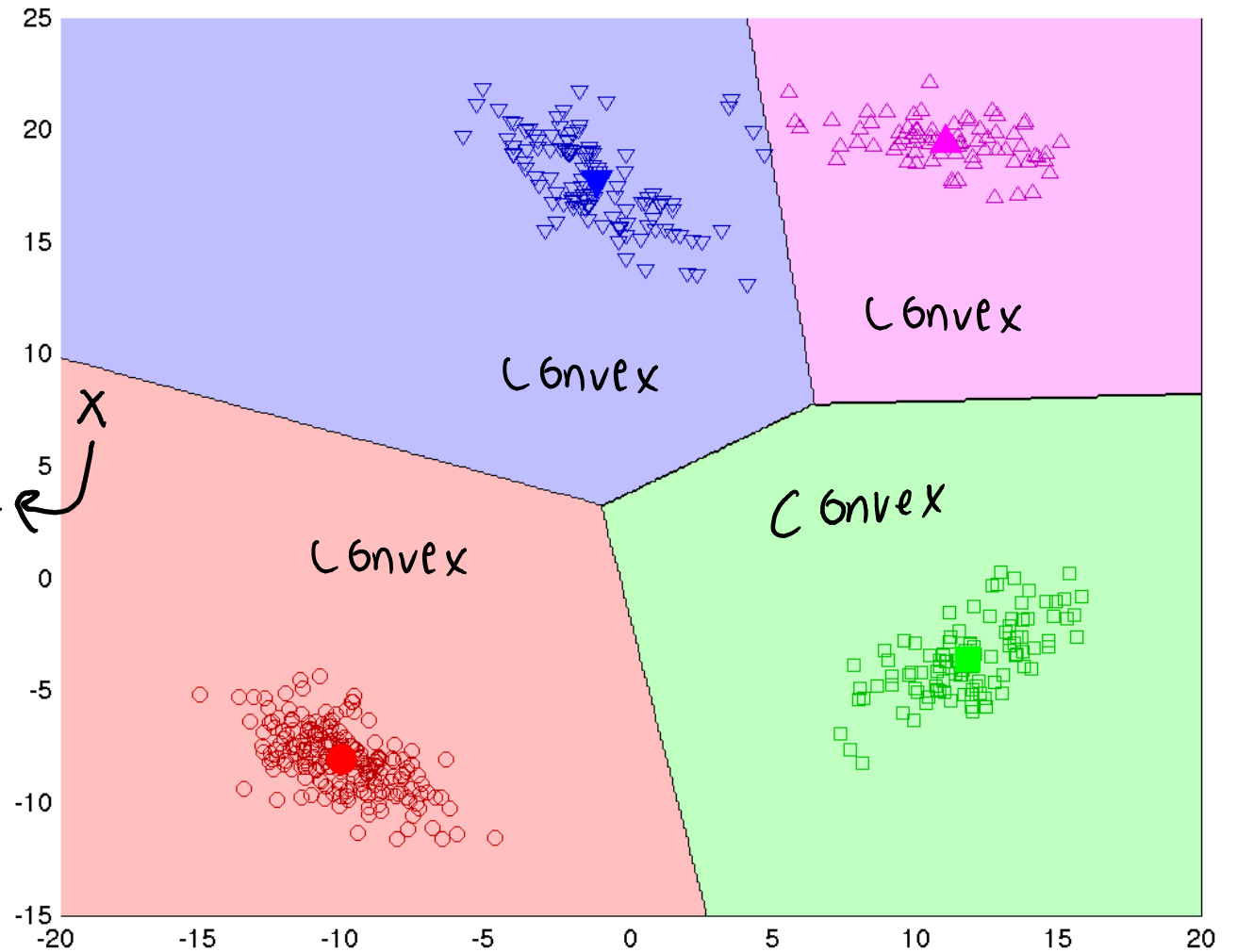


# Shape of K-Means Clusters

Issues with shape of k-means clusters:

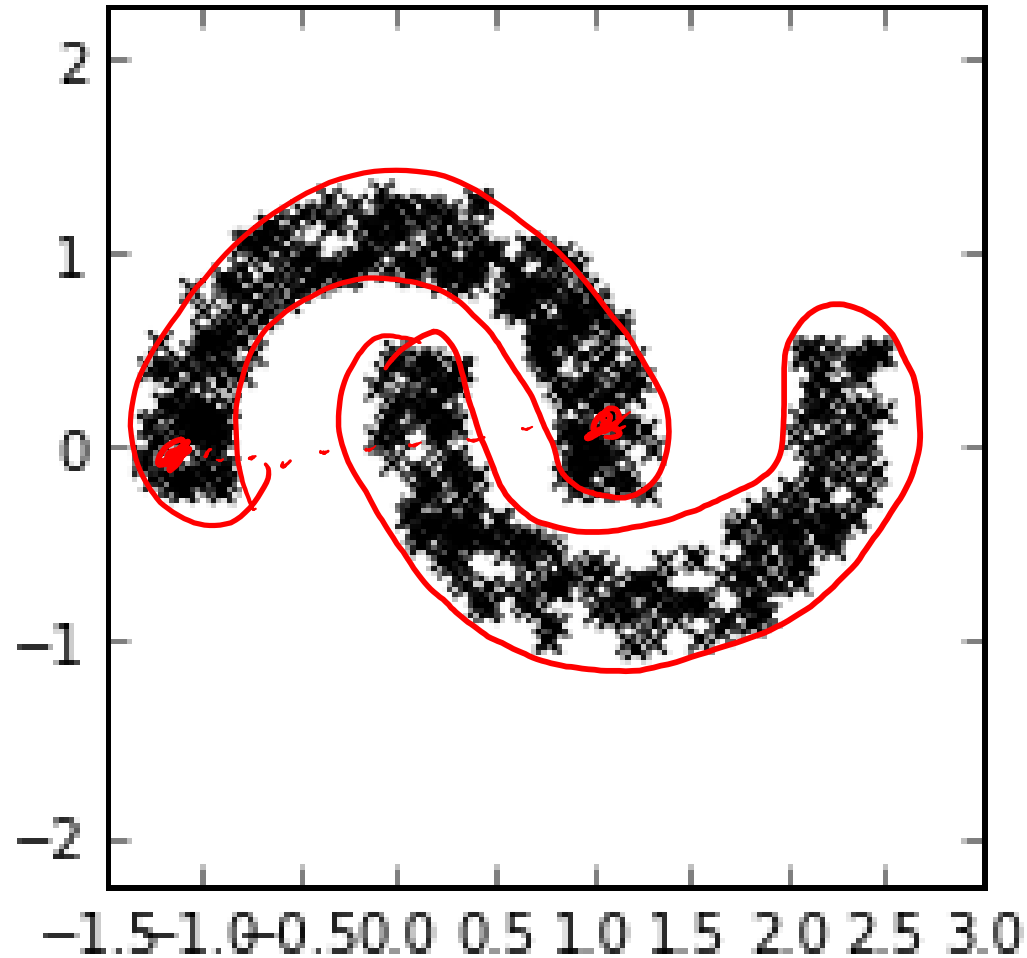
1. Clusters in the data might not be convex.

2. Does this point really belong in red cluster?

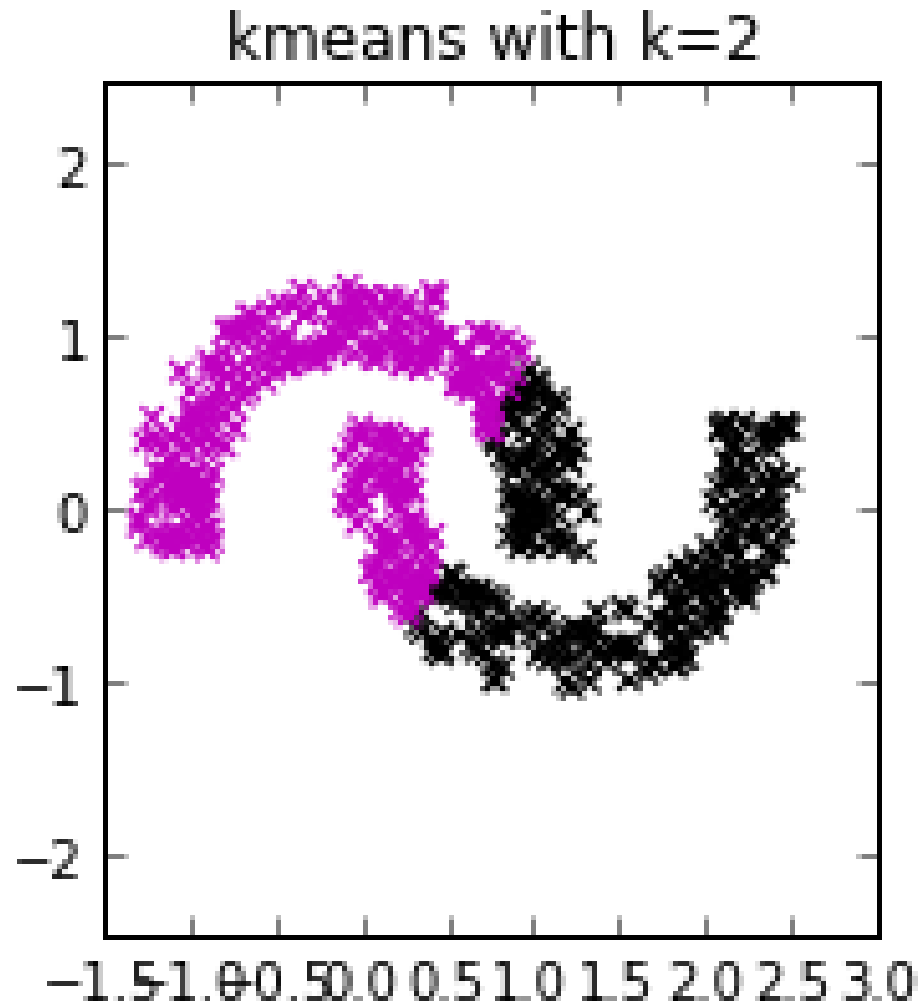


# K-Means with Non-Convex Clusters

Non-convex banana-shaped data points



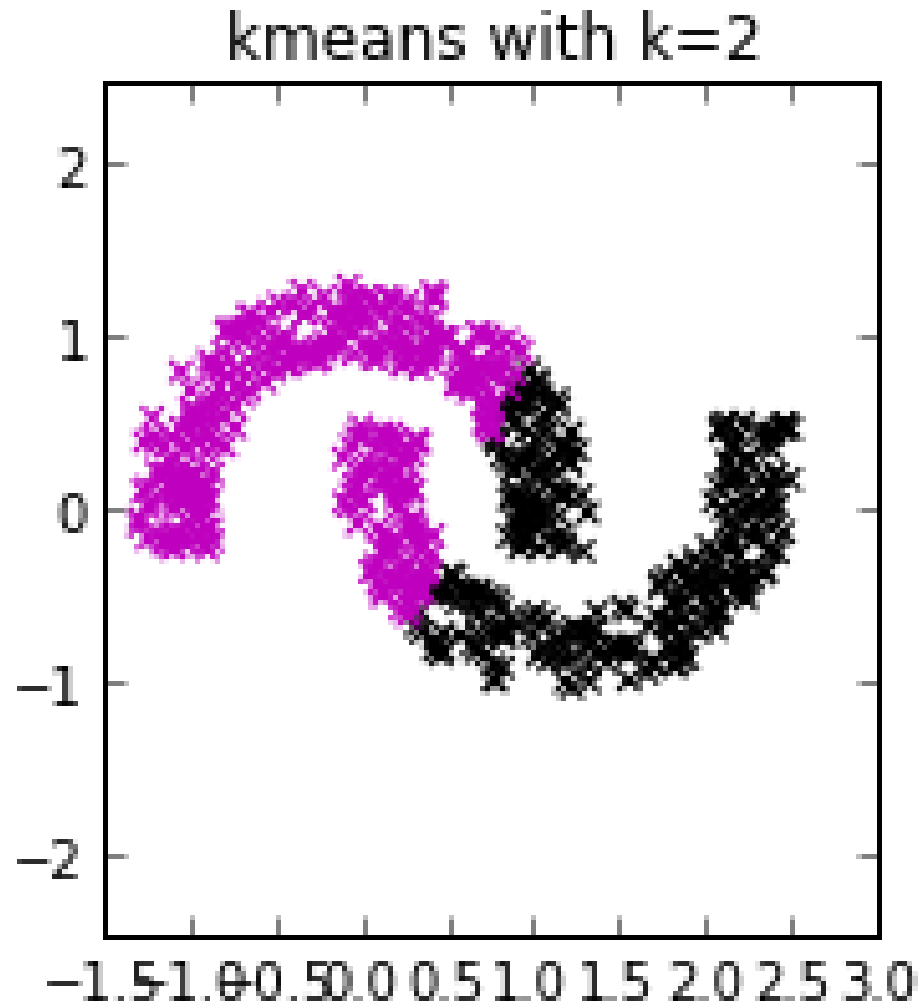
# K-Means with Non-Convex Clusters



K-means **cannot separate**  
non-convex clusters

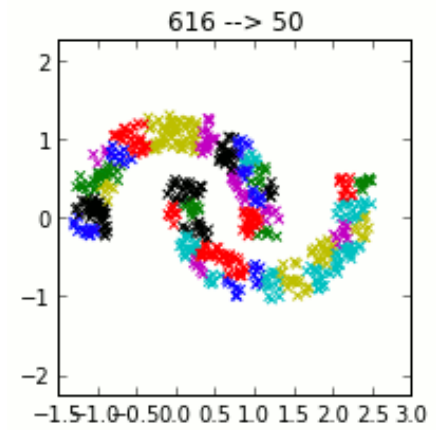


# K-Means with Non-Convex Clusters



K-means **cannot separate**  
non-convex clusters

Though over-clustering can help  
(next class)



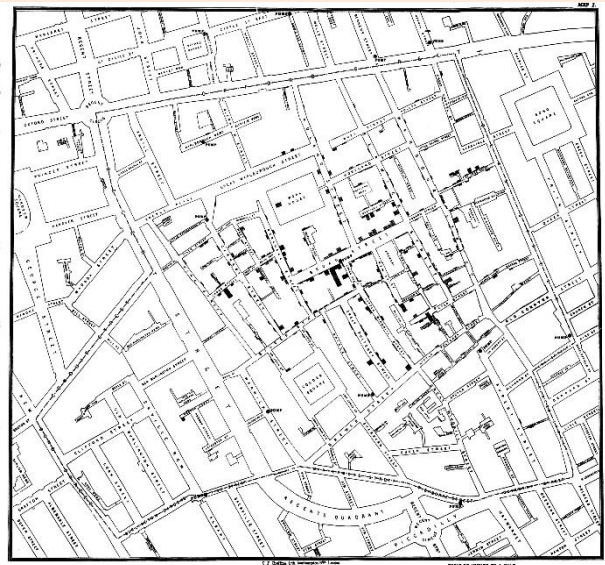
- wallpaper
- funny
- season 6
- winter is coming meme
- genderbent
- let it snow
- real life
- ghost
- wallpaper hd
- the watch meme
- his wolf
- game throne
- simpsons

Did you mean: [jon snow](#)



# John Snow and Cholera Epidemic

- John Snow's 1854 spatial histogram of deaths from cholera:



- Found cluster of cholera deaths around a particular water pump.
  - Went against airborne theory, but pump later found to be contaminated.
  - “Father” of epidemiology.



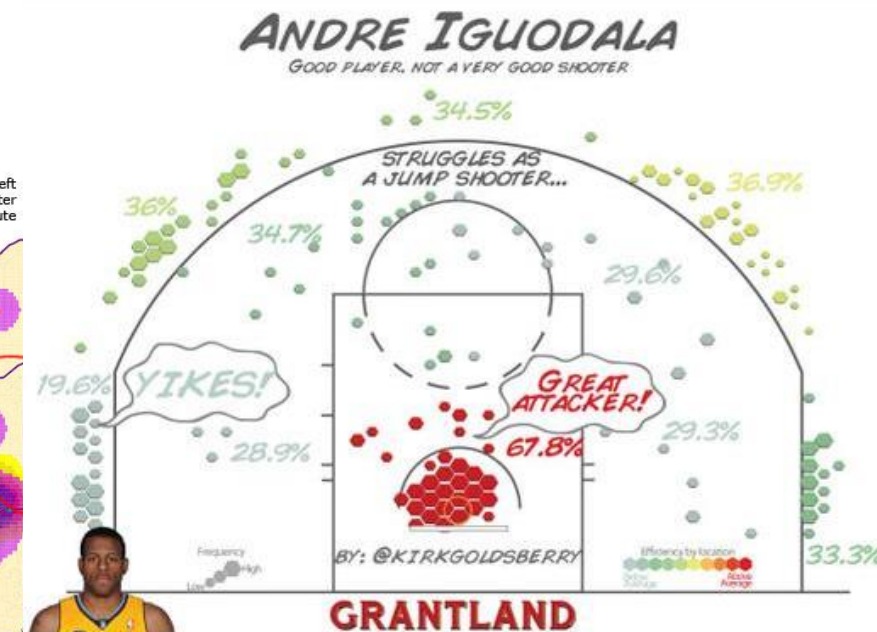
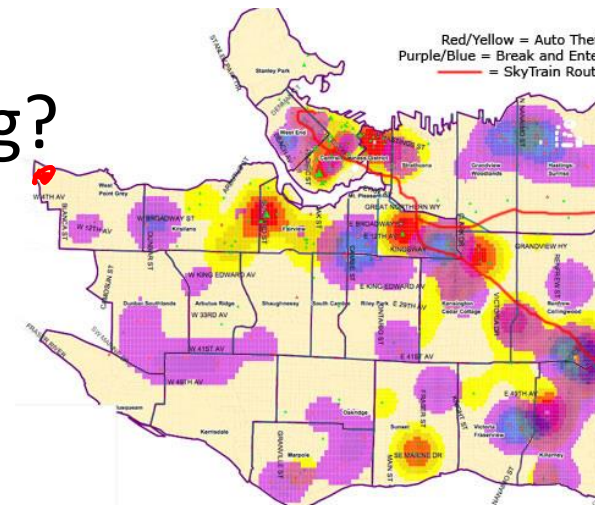
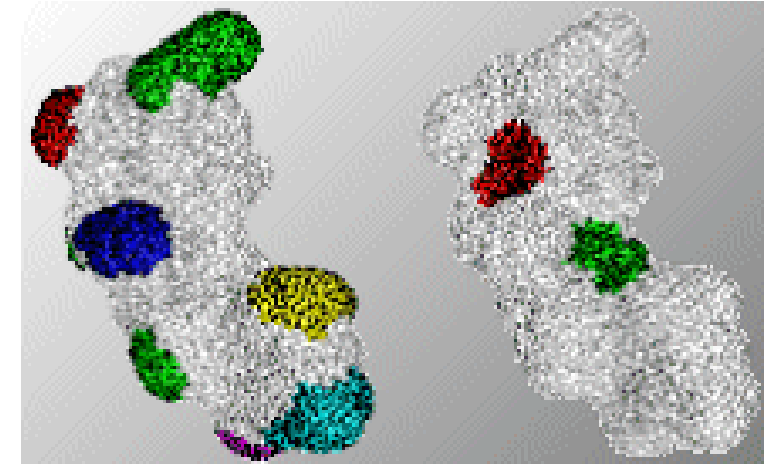
# Motivation for Density-Based Clustering

- **Density-based clustering:**
  - Clusters are **defined by “dense” regions.**
  - Objects in **non-dense regions don’t get clustered.**
    - Not trying to “partition” the space.
- Clusters can be **non-convex:**
  - Elephant clusters affected by vegetation, mountains, rivers, water access, etc.
- It’s a **non-parametric clustering** method:
  - No fixed number of clusters ‘k’.
  - Clusters can become more complicated with more data.



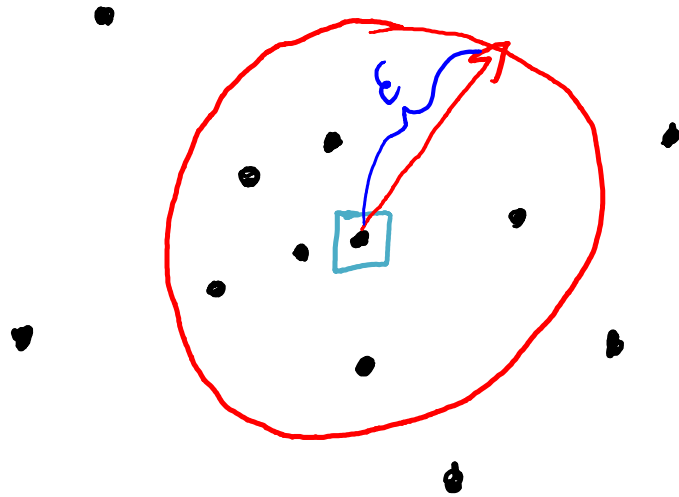
# Other Potential Applications

- Where are high crime regions of a city?
- Where should taxis patrol?
- Where does Iguodala make/miss shots?
- Which products are similar to this one?
- Which pictures are in the same place?
- Where can protein 'dock'?
- Where are people tweeting?



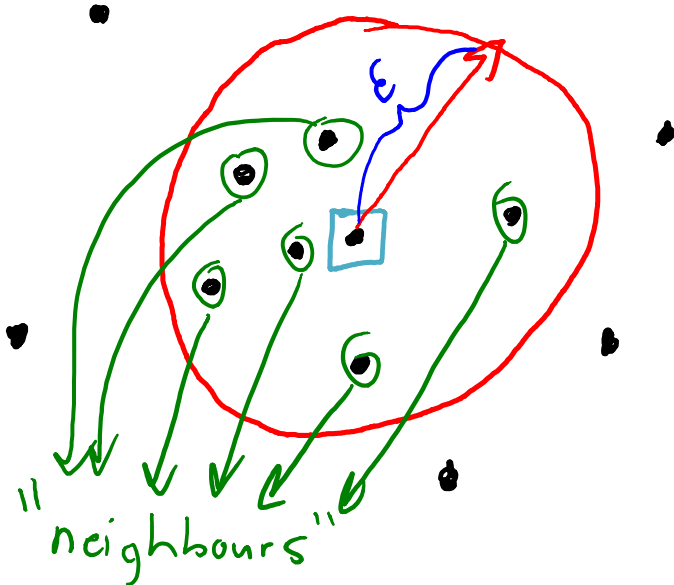
# Density-Based Clustering

- Density-based clustering algorithm (DBSCAN) has two hyperparameters:
  - Epsilon ( $\epsilon$ ): distance we use to decide if another point is a “neighbour”.



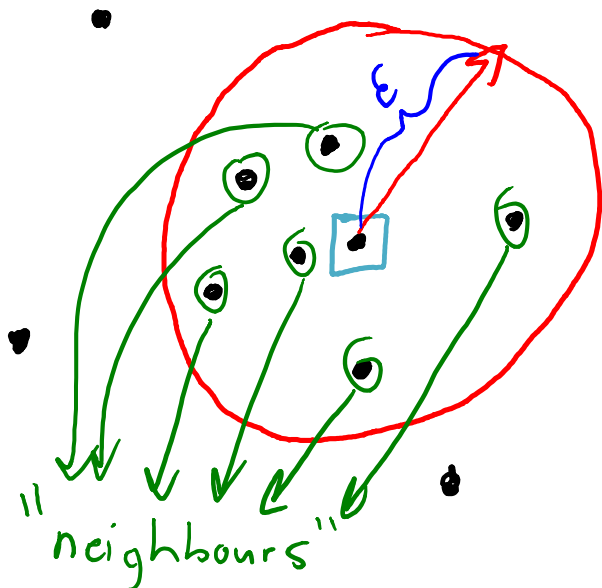
# Density-Based Clustering

- Density-based clustering algorithm (DBSCAN) has two hyperparameters:
  - Epsilon ( $\epsilon$ ): distance we use to decide if another point is a “neighbour”.



# Density-Based Clustering

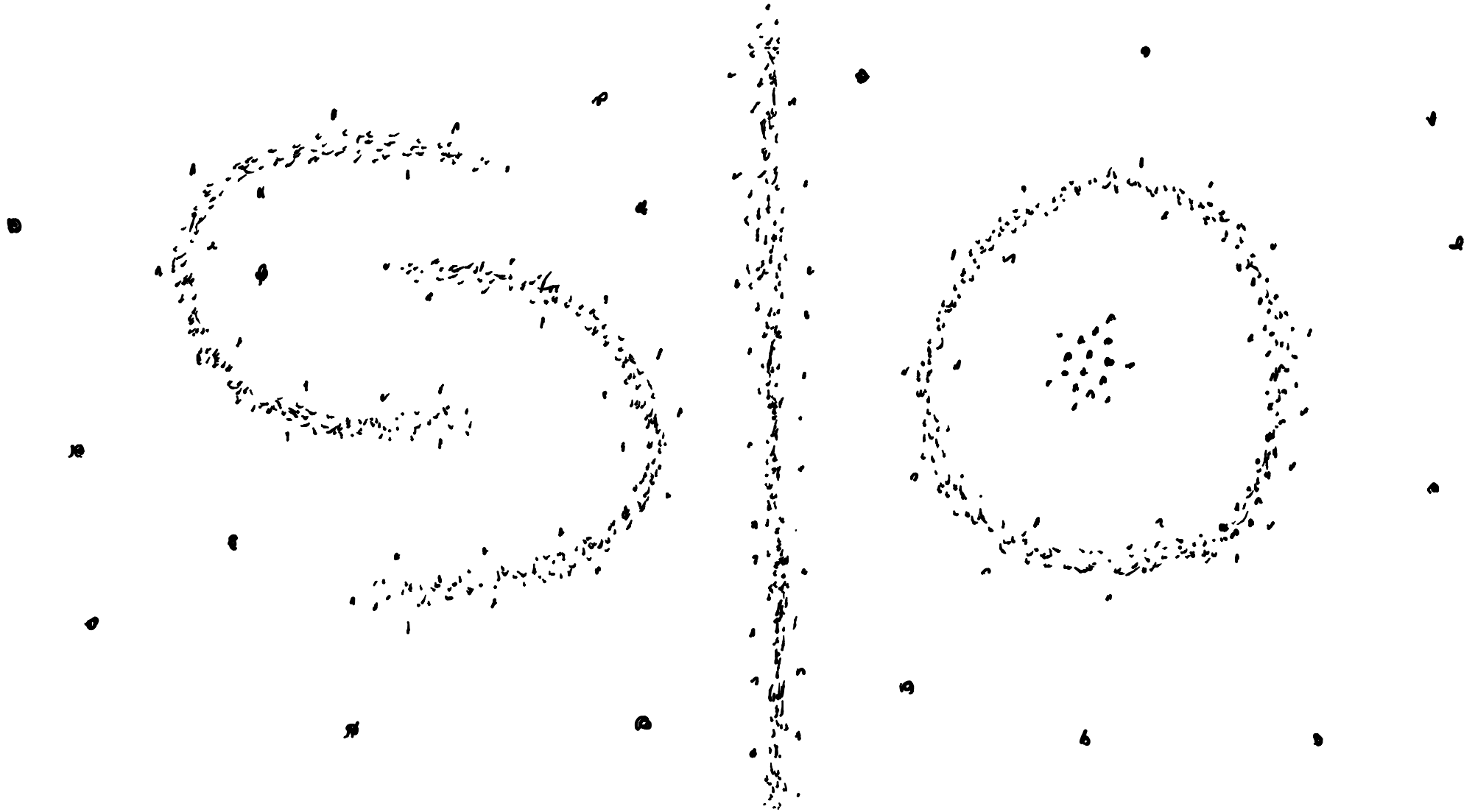
- **Density-based clustering** algorithm (DBSCAN) has two hyperparameters:
  - **Epsilon ( $\epsilon$ )**: distance we use to decide if another point is a “**neighbour**”.
  - **MinNeighbours**: **number of neighbours** needed to say a region is “dense”.
    - If you have at least minNeighbours “neighbours”, you are called a “**core**” point.



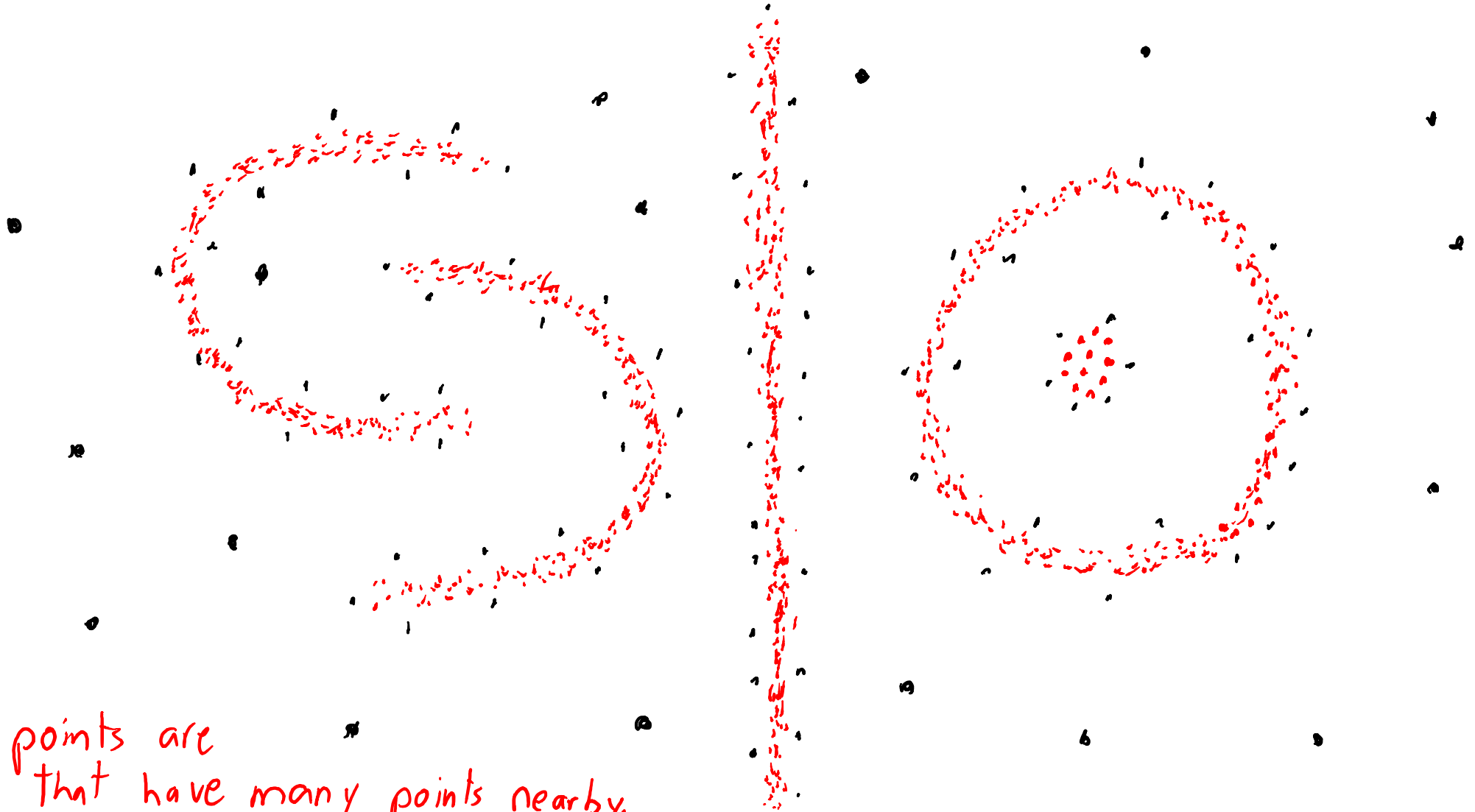
→ E.g., if  $\text{minNeighbours} = 3$   
then this is a “core”  
point since 6 points are  
“neighbours”



# Density-Based Clustering

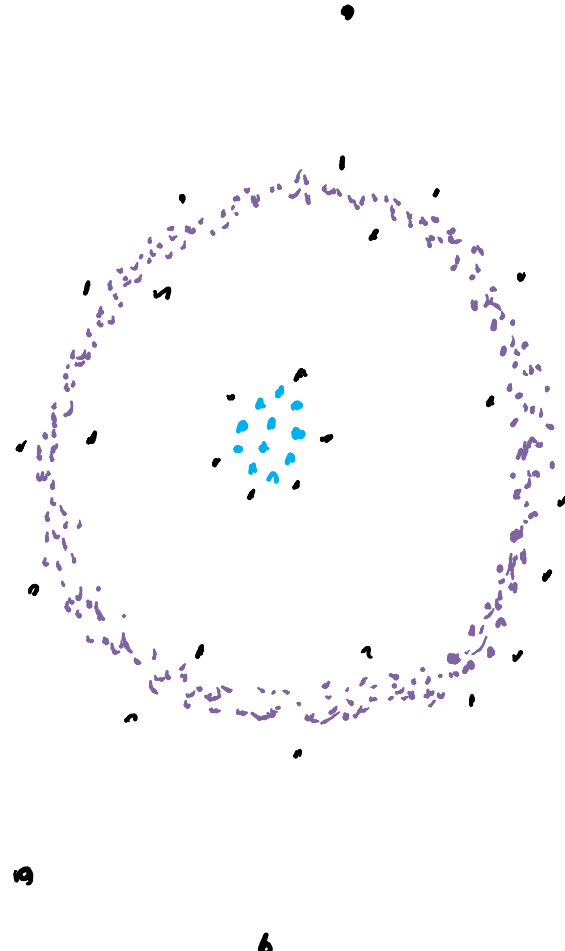
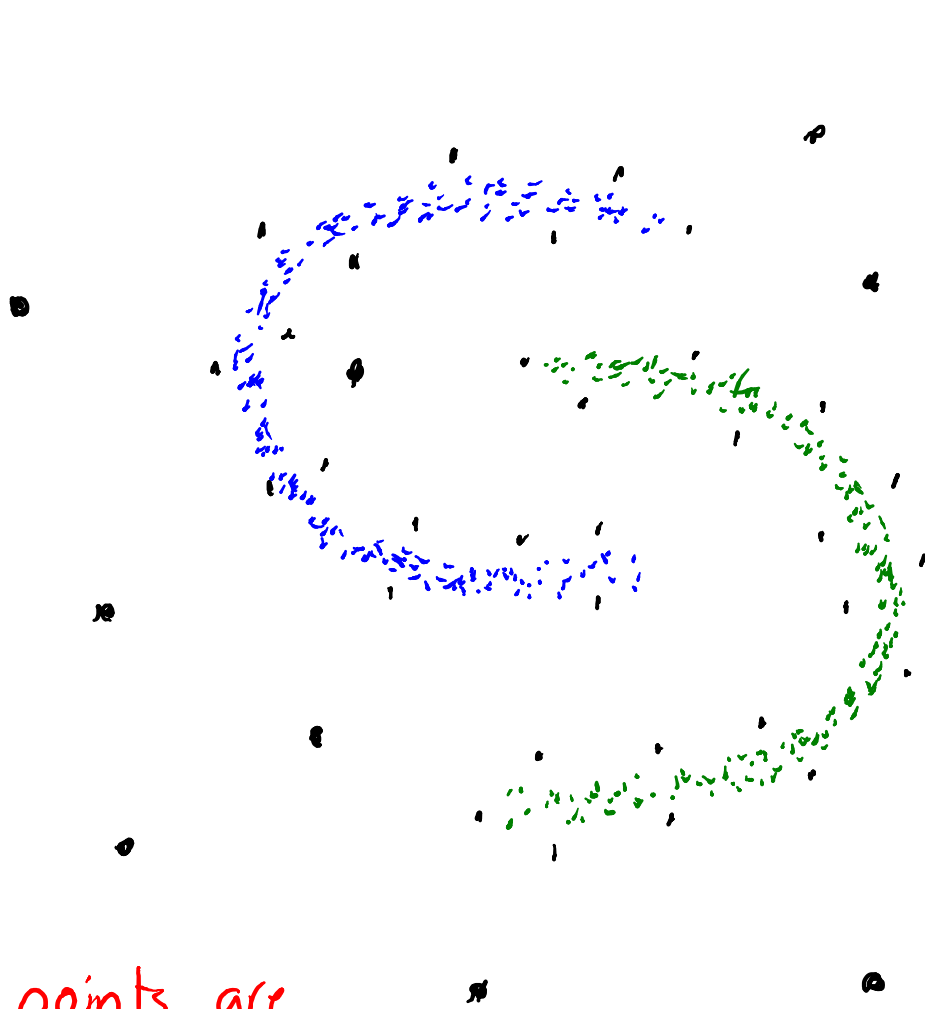


# Density-Based Clustering



"Core" points are points that have many points nearby.

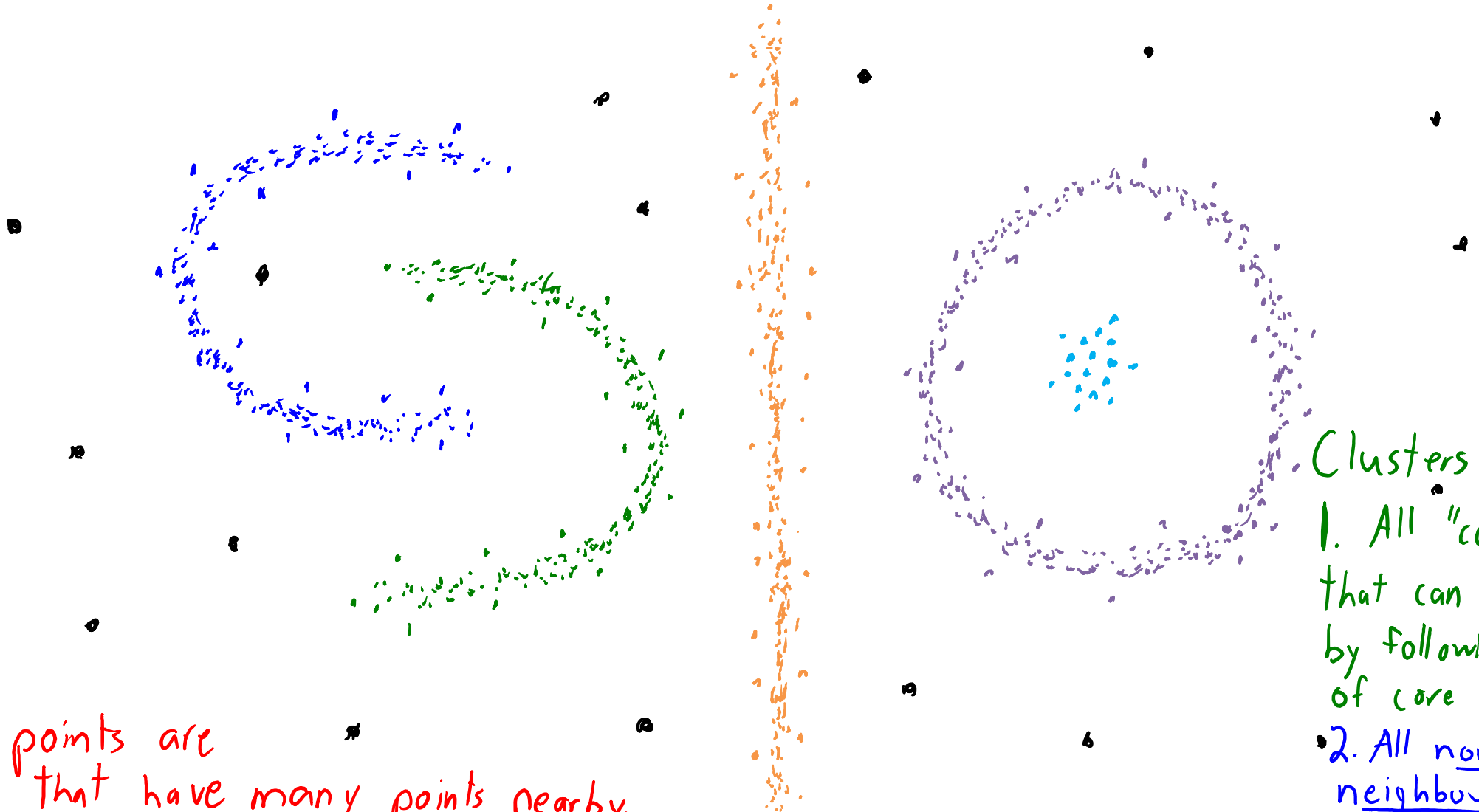
# Density-Based Clustering



Clusters contain:  
1. All "core" points that can be reached by following a sequence of core points.

"Core" points are points that have many points nearby.

# Density-Based Clustering

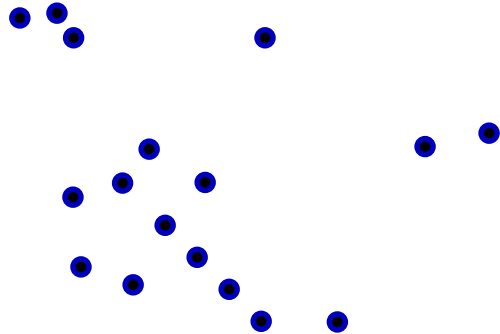


"Core" points are points that have many points nearby.

- Clusters contain:
1. All "core" points that can be reached by following a sequence of core points.
  2. All non-core neighbours of core points (boundary points)

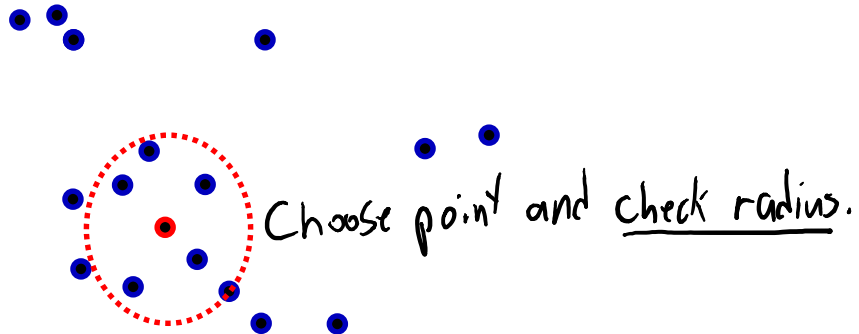
# Density-Based Clustering (Example)

- Each “core” point defines a cluster:
  - Consisting of “core” point and all its “neighbours”.
- Merge clusters if “core” points are “neighbours” of each other.



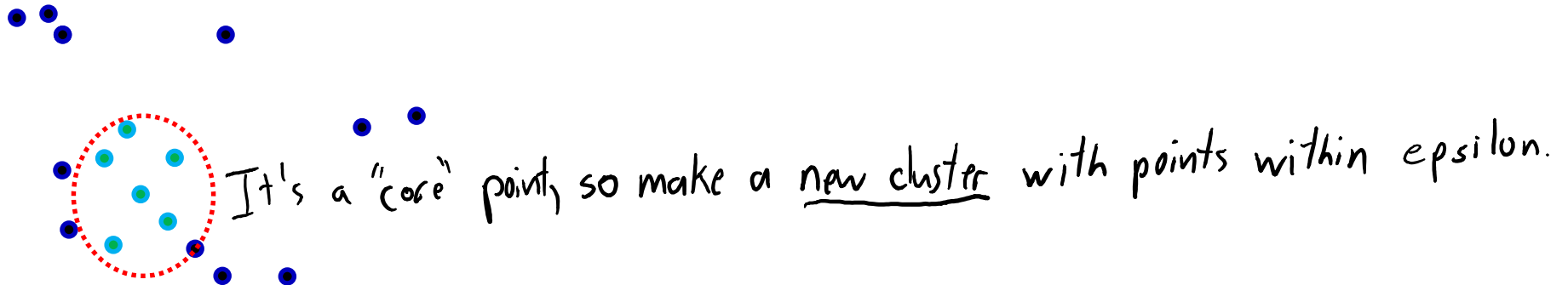
# Density-Based Clustering (Example)

- Each “core” point defines a cluster:
  - Consisting of “core” point and all its “neighbours”.
- Merge clusters if “core” points are “neighbours” of each other.



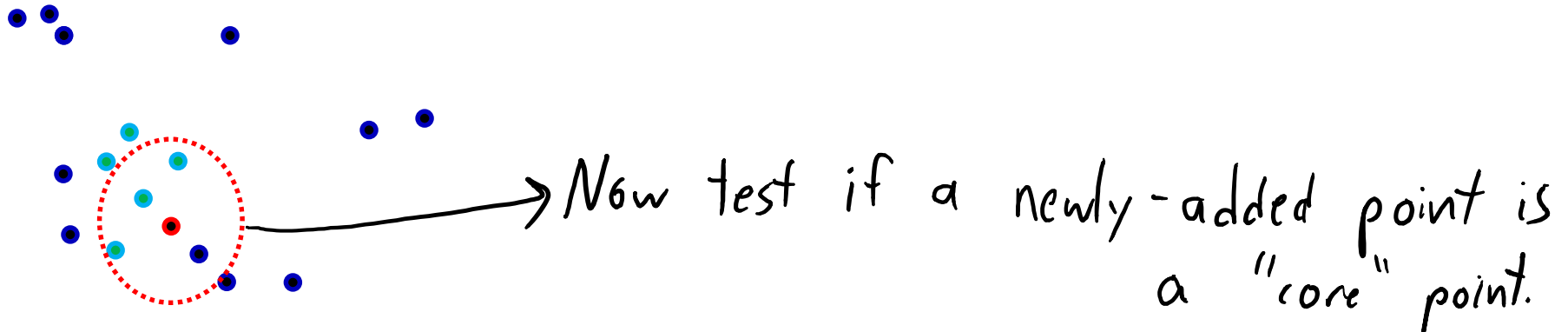
# Density-Based Clustering (Example)

- Each “core” point defines a cluster:
  - Consisting of “core” point and all its “neighbours”.
- Merge clusters if “core” points are “neighbours” of each other.



# Density-Based Clustering (Example)

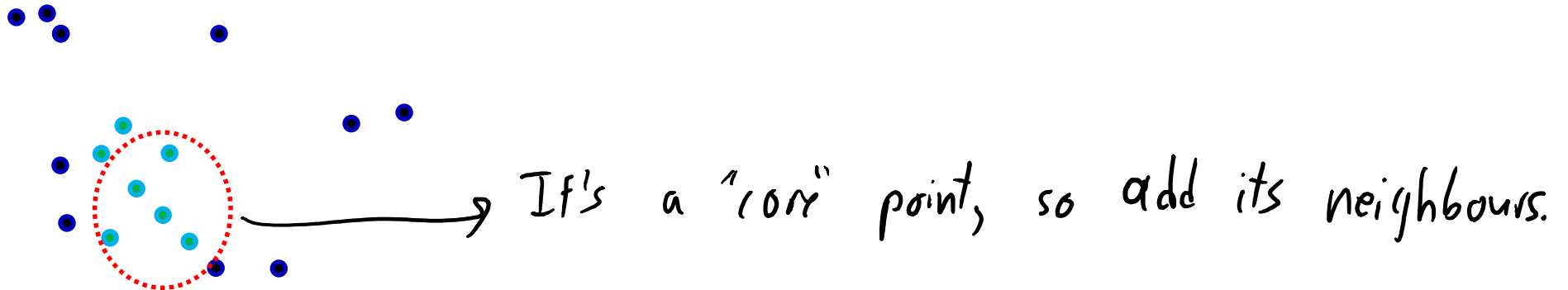
- Each “core” point defines a cluster:
  - Consisting of “core” point and all its “neighbours”.
- Merge clusters if “core” points are “neighbours” of each other.





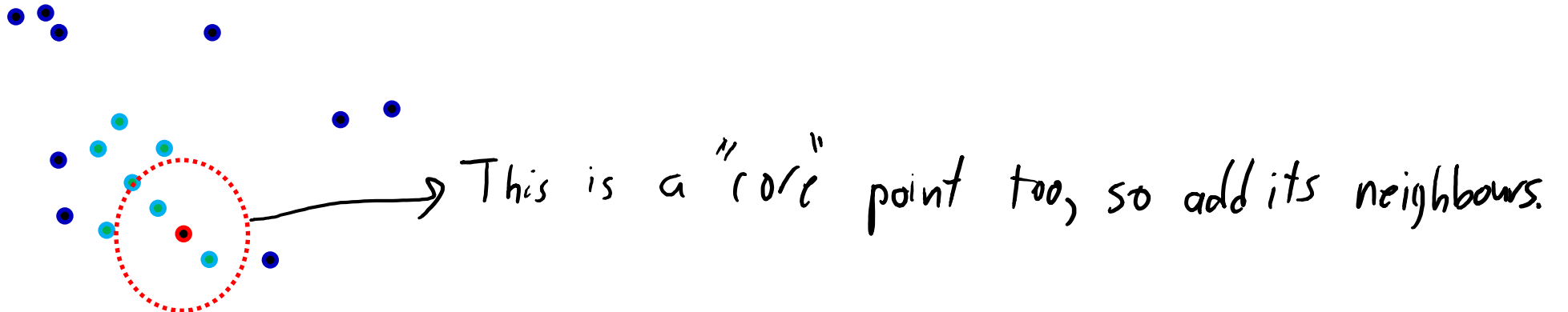
# Density-Based Clustering (Example)

- Each “core” point defines a cluster:
  - Consisting of “core” point and all its “neighbours”.
- Merge clusters if “core” points are “neighbours” of each other.



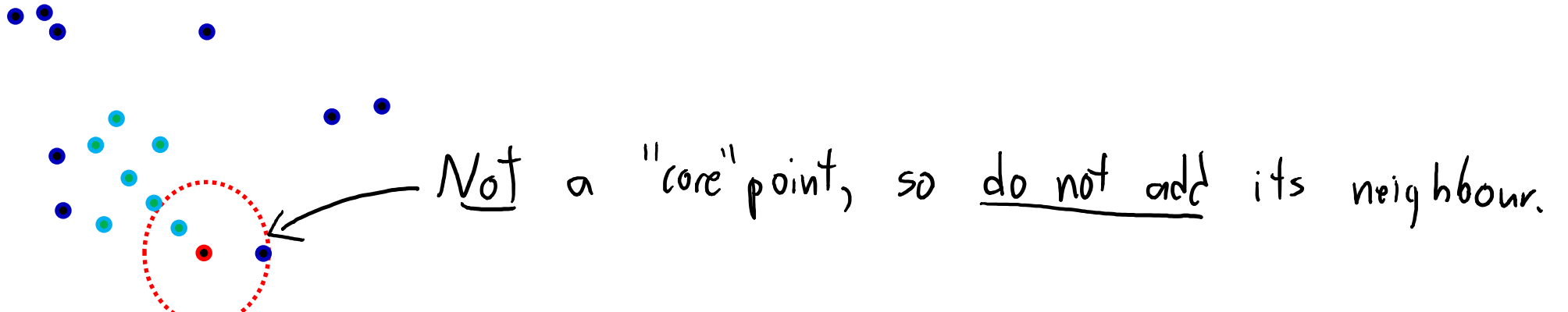
# Density-Based Clustering (Example)

- Each “core” point defines a cluster:
  - Consisting of “core” point and all its “neighbours”.
- Merge clusters if “core” points are “neighbours” of each other.



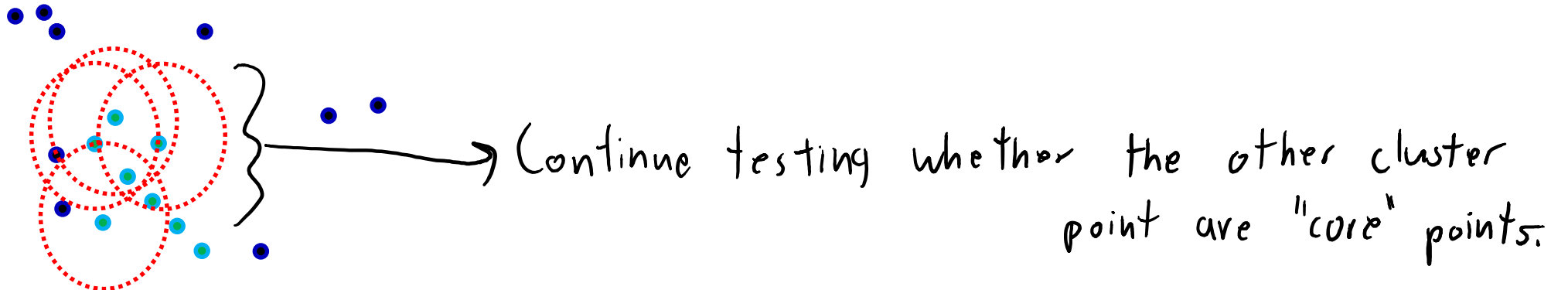
# Density-Based Clustering (Example)

- Each “core” point defines a cluster:
  - Consisting of “core” point and all its “neighbours”.
- Merge clusters if “core” points are “neighbours” of each other.



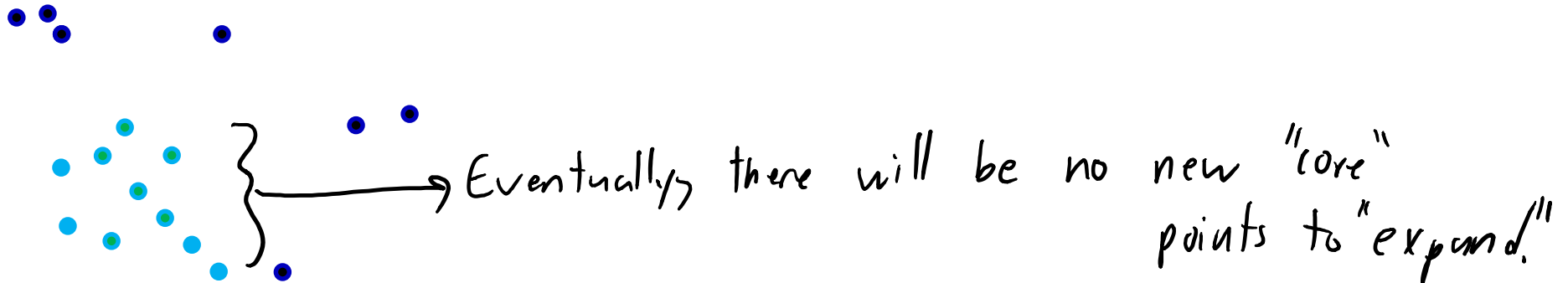
# Density-Based Clustering (Example)

- Each “core” point defines a cluster:
  - Consisting of “core” point and all its “neighbours”.
- Merge clusters if “core” points are “neighbours” of each other.



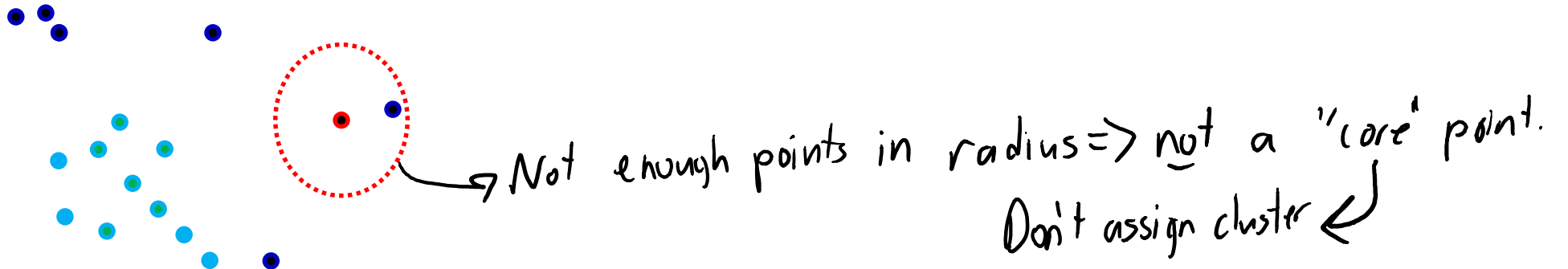
# Density-Based Clustering (Example)

- Each “core” point defines a cluster:
  - Consisting of “core” point and all its “neighbours”.
- Merge clusters if “core” points are “neighbours” of each other.



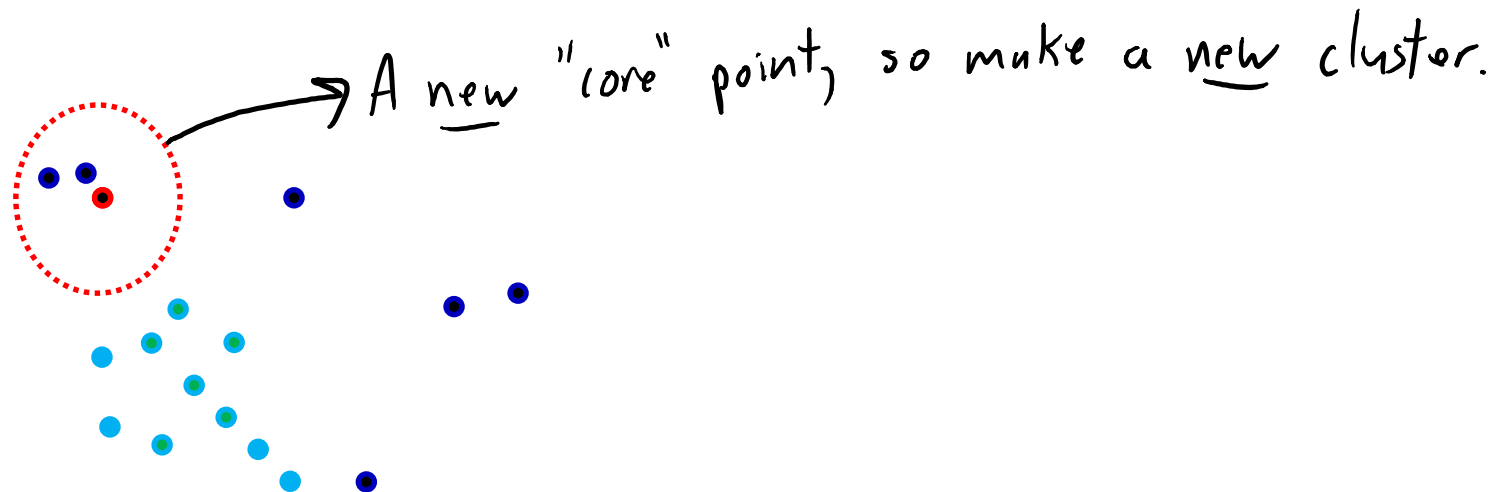
# Density-Based Clustering (Example)

- Each “core” point defines a cluster:
  - Consisting of “core” point and all its “neighbours”.
- Merge clusters if “core” points are “neighbours” of each other.



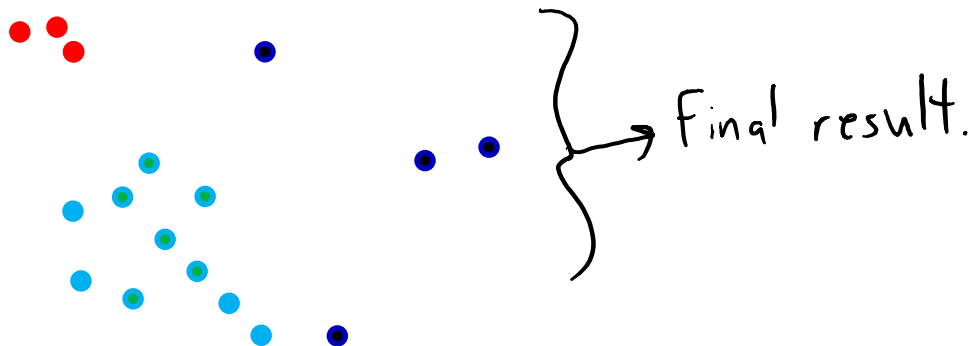
# Density-Based Clustering (Example)

- Each “core” point defines a cluster:
  - Consisting of “core” point and all its “neighbours”.
- Merge clusters if “core” points are “neighbours” of each other.



# Density-Based Clustering (Example)

- Each “core” point defines a cluster:
  - Consisting of “core” point and all its “neighbours”.
- Merge clusters if “core” points are “neighbours” of each other.



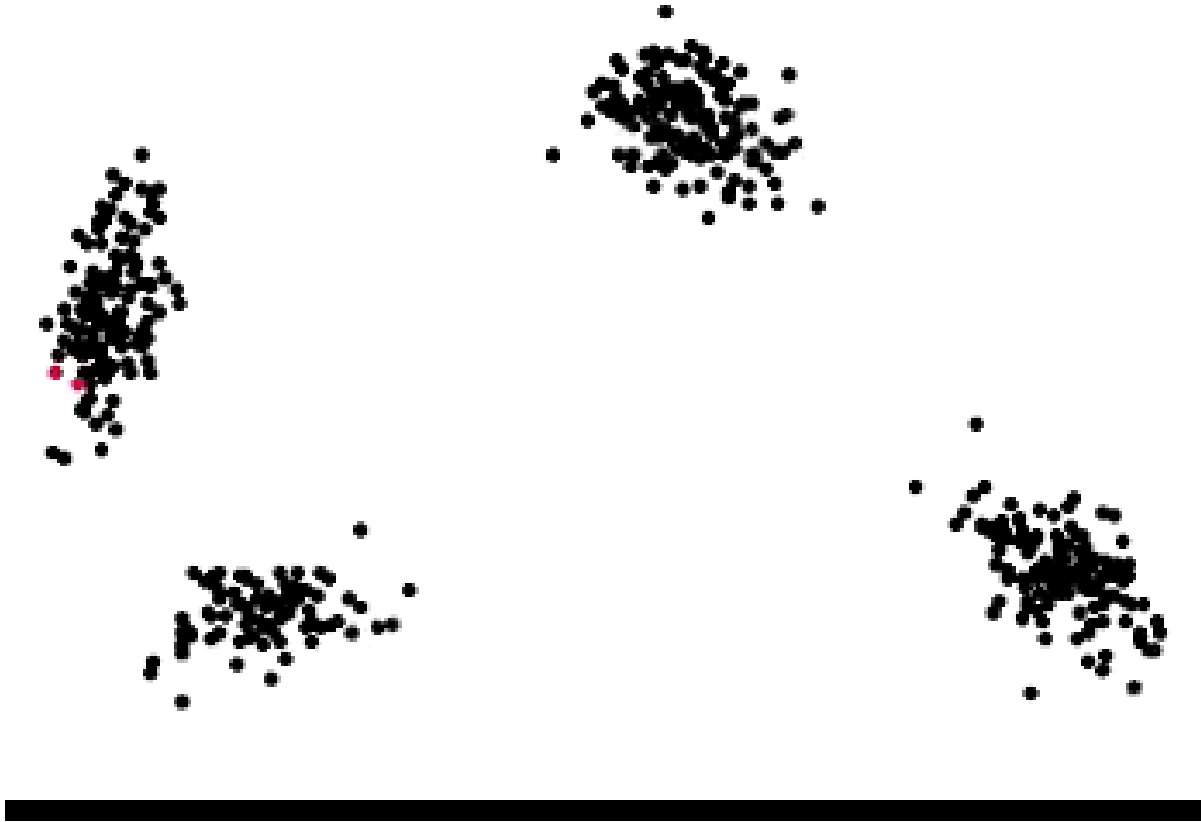


# Density-Based Clustering Pseudo-Code

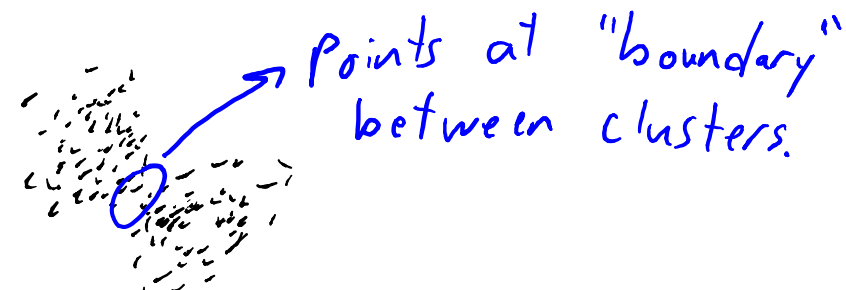
- For each example  $x_i$ :
  - If  $x_i$  is already assigned to a cluster, do nothing.
  - Test whether  $x_i$  is a ‘core’ point ( $\geq \text{minNeighbours}$  examples within ‘ $\epsilon$ ’).
    - If  $x_i$  is not core point, do nothing (this could be an outlier).
    - If  $x_i$  is a core point, “expand” cluster.
- “Expand” cluster function:
  - Assign all  $x_j$  within distance ‘ $\epsilon$ ’ of core point  $x_i$  to cluster.
  - For each newly-assigned neighbour  $x_j$  that is a core point, “expand” cluster.

# Density-Based Clustering in Action

---



# Density-Based Clustering Issues

- Some points are not assigned to a cluster.
  - Good or bad, depending on the application.
- Ambiguity of “non-core” (boundary) points:
  - Otherwise, not sensitive to initialization (except for boundary points).
- Sensitive to the choice of  $\epsilon$  and minNeighbours.
  - Need to compute distances to training points.
- If you get a new example, finding cluster is expensive.
  - Need to compute distances to training points.
- In high-dimensions, need a lot of points to ‘fill’ the space.

(pause)

# Ensemble Clustering

? question ☆

stop following 23 views

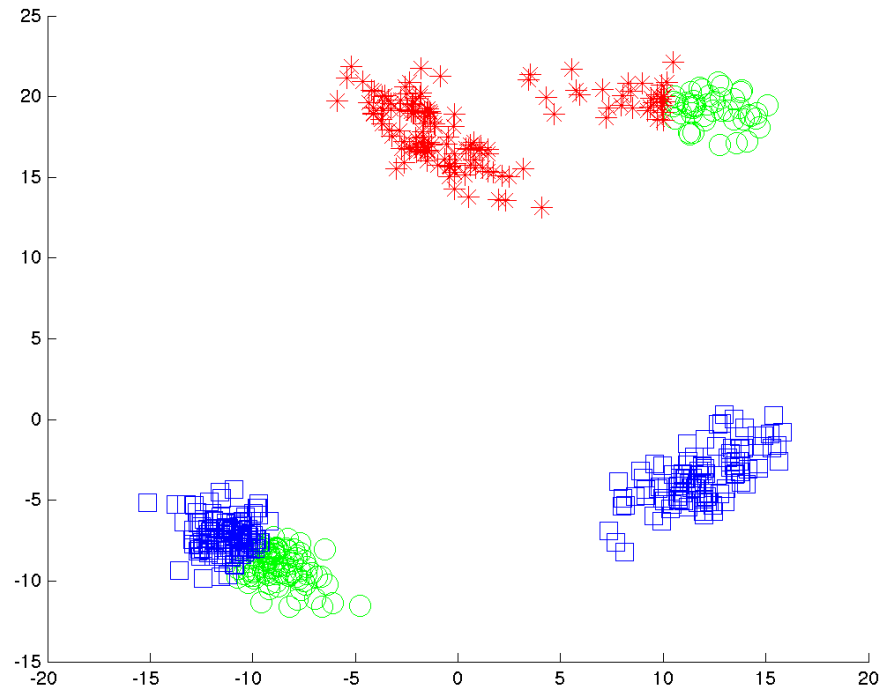
## Multiple random runs of K means

I was wondering how running K Means (original version, not K means ++ ) several times with random initializations can help us make an accurate model. K Means outputs the class labels of all the samples. We definitely can't use mode of all the labels it got in different runs because class labels from different runs don't make any sense when compared. We somehow have to see what points are coming in the same cluster in a lot of runs..I am not sure, how do we do it?

- We can consider **ensemble methods** for clustering.
  - “Consensus clustering”
- It's a good/important idea:
  - **Bootstrapping** is widely-used.
  - “Do clusters change if the data was slightly different?”
- But we **need to be careful** about how we combine models.

# Ensemble Clustering

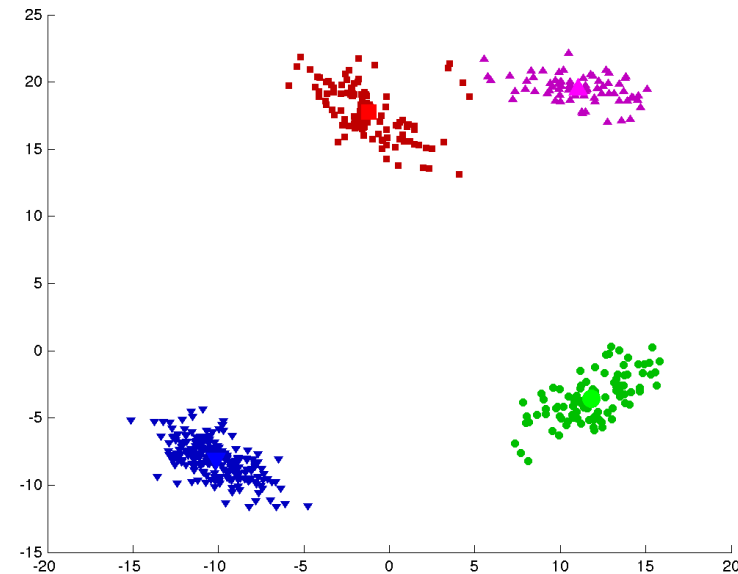
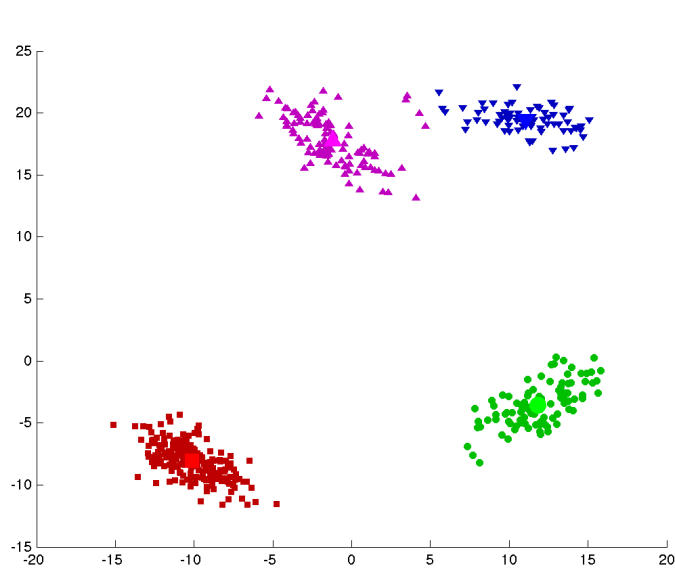
- E.g., run k-means 20 times and then cluster using the mode of each  $\hat{y}_i$ .
- Normally, averaging across models doing different things is good.



- But this is a bad ensemble method: **worse than k-means on its own.**

# Label Switching Problem

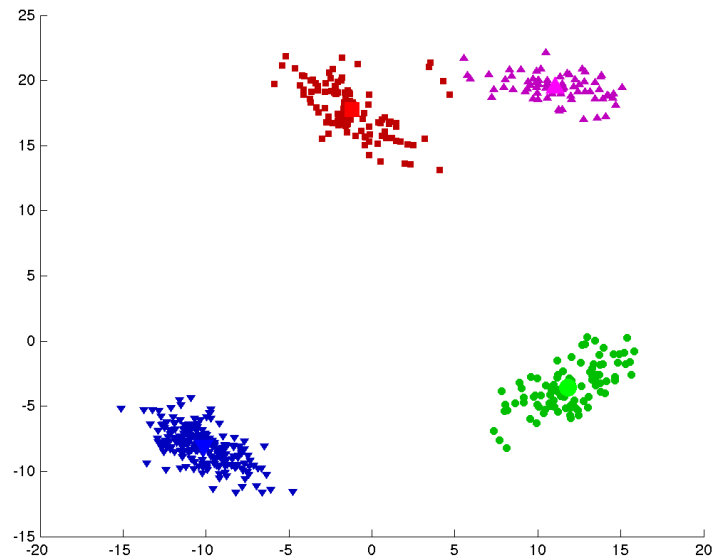
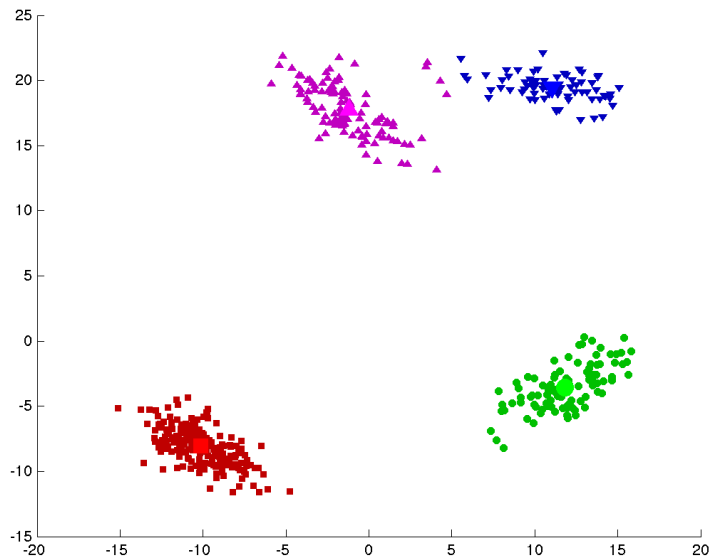
- This doesn't work because of “label switching” problem:
  - The cluster labels  $\hat{y}_i$  are meaningless.
  - We could get same clustering with permuted labels:



- All  $\hat{y}_i$  become equally likely as number of initializations increases.

# Addressing Label Switching Problem

- Ensembles can't depend on label "meaning":
  - Don't ask "is point  $x_i$  in red square cluster?", which is meaningless.
  - Ask "is point  $x_i$  in the same cluster as  $x_j$ ?", which is meaningful.



- Bonus slides give an example method ("UBClustering").



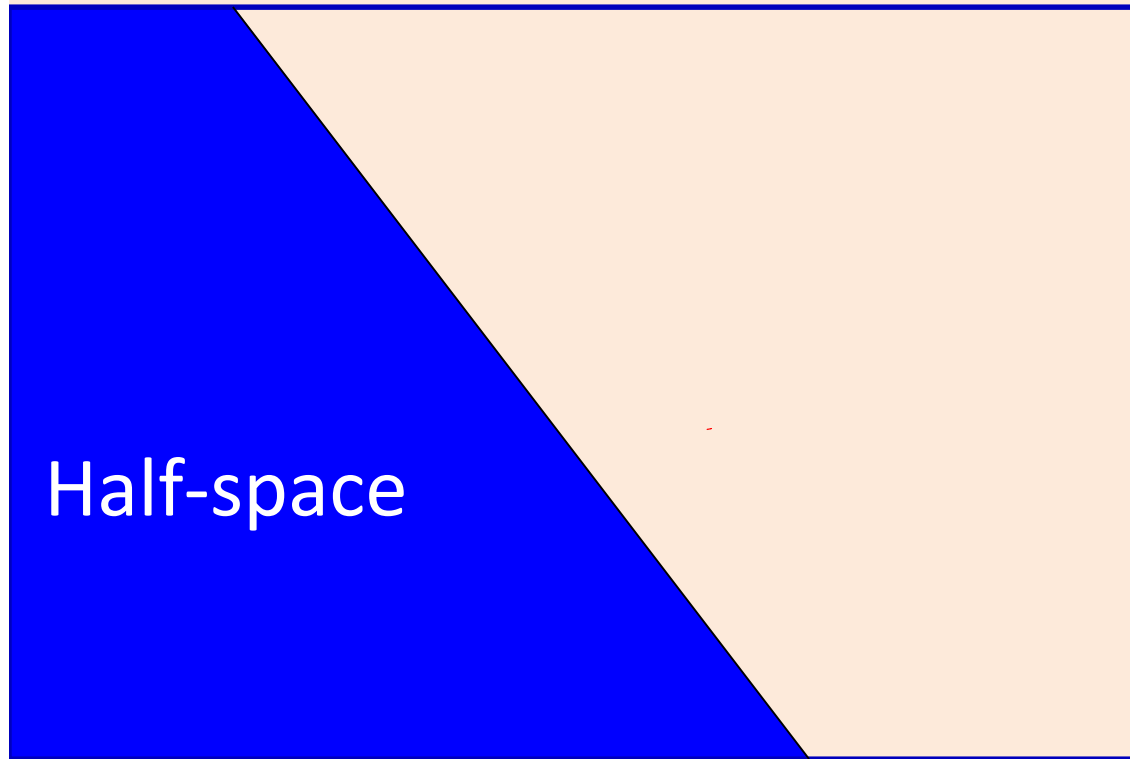
# Summary

- **Shape of K-means clusters:**
  - Partitions space into convex sets.
- **Density-based clustering:**
  - “Expand” and “merge” dense regions of points to find clusters.
  - Not sensitive to initialization or outliers.
  - Useful for finding non-convex connected clusters.
- **Ensemble clustering:** combines multiple clusterings.
  - Can work well but need to account for **label switching**.
- **Next time:**
  - Discovering the tree of life.

# Why are k-means clusters convex?

- K-means clusters are formed by the **intersection** of **half-spaces**.

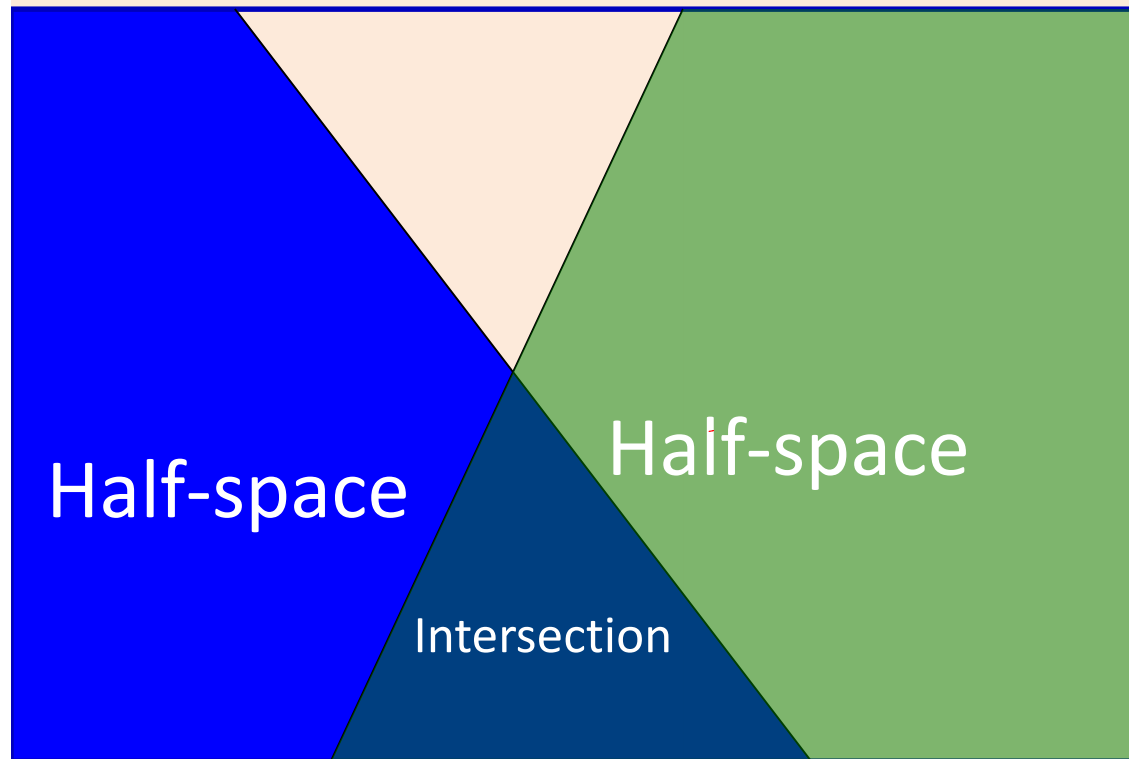
Half-space is Set of points satisfying a linear inequality, like  $\sum_{j=1}^d a_j x_j \leq b$



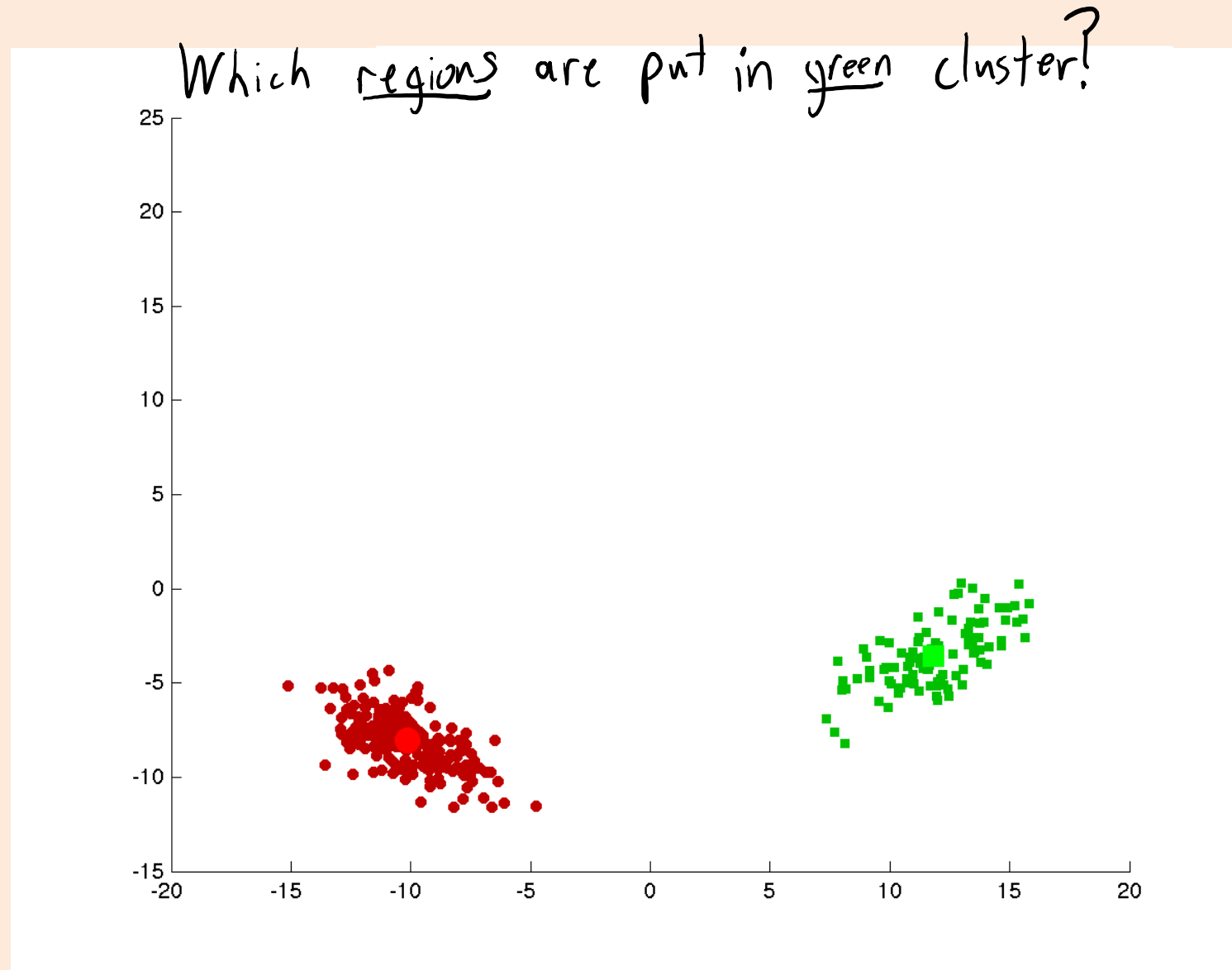
# Why are k-means clusters convex?

- K-means clusters are formed by the **intersection** of **half-spaces**.

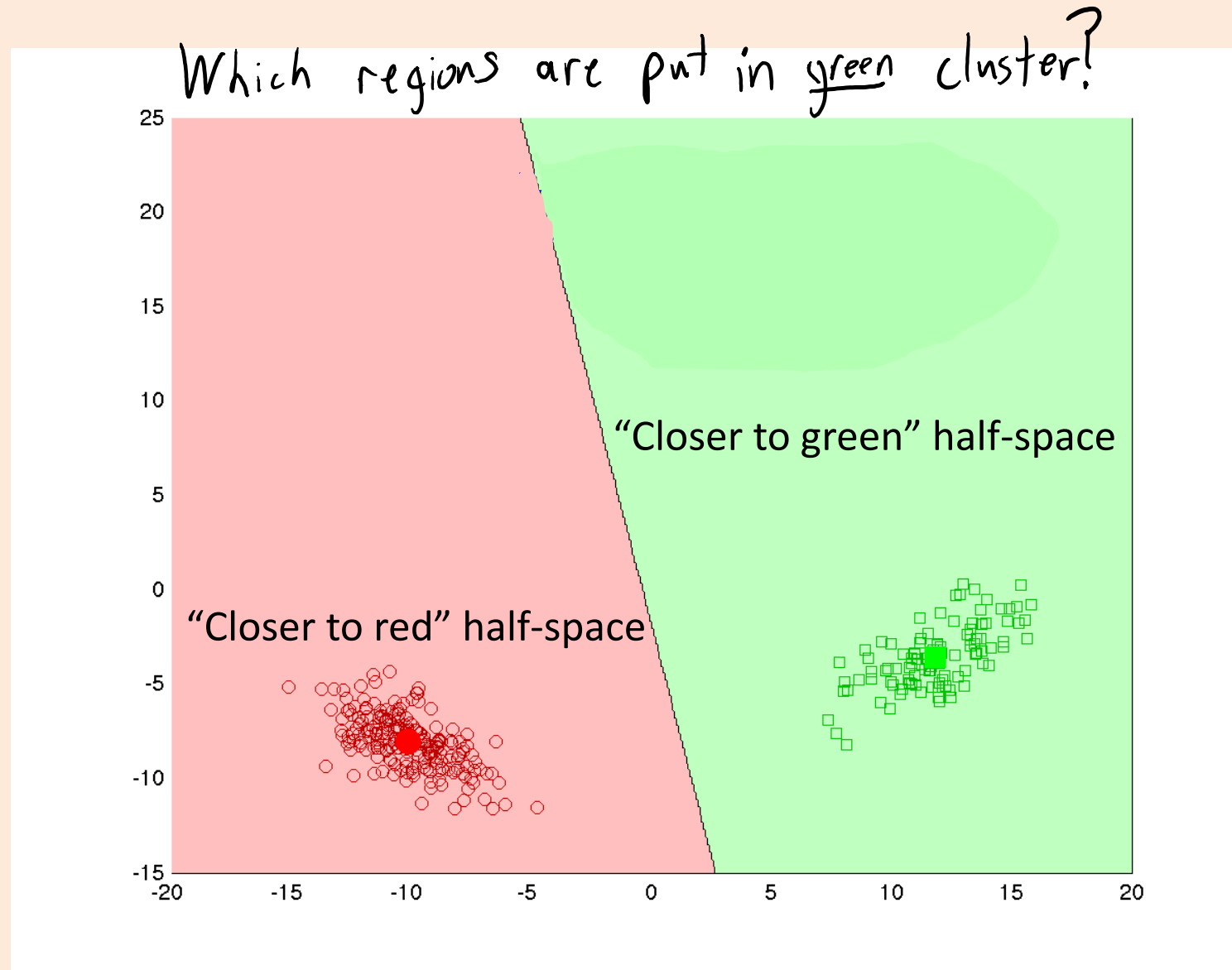
Half-space is Set of points satisfying a linear inequality, like  $\sum_{j=1}^d a_j x_j \leq b$



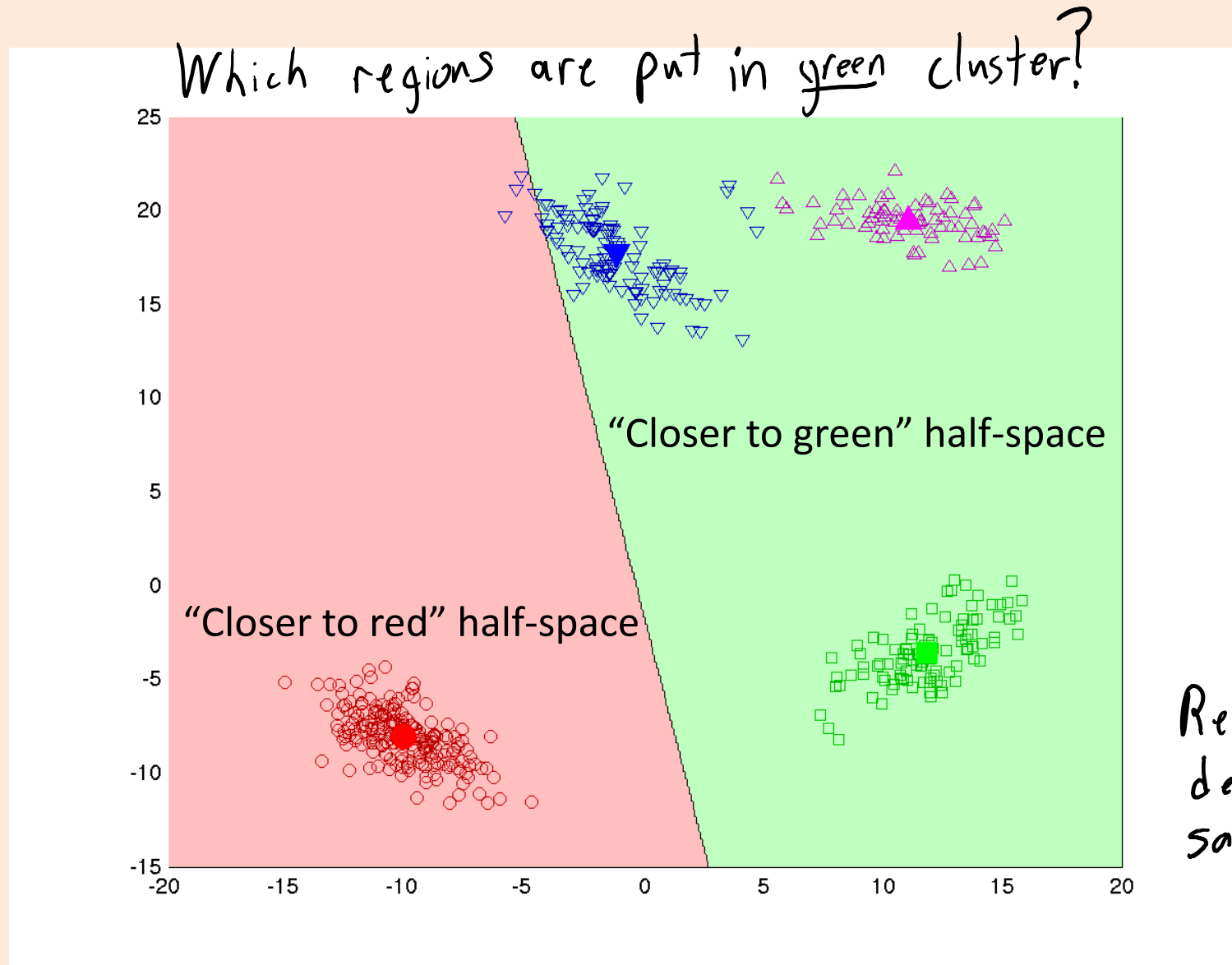
# Why are k-means clusters convex?



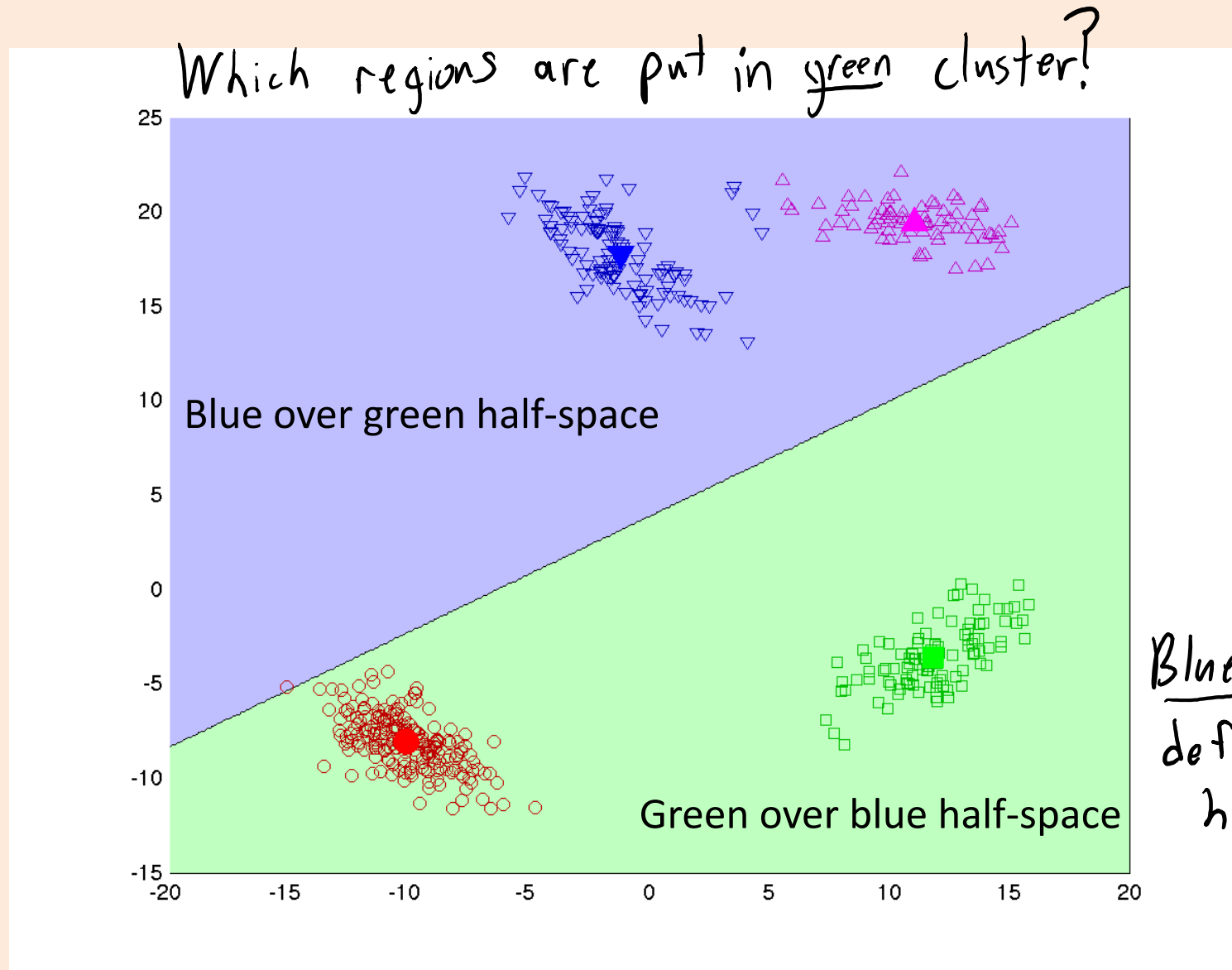
# Why are k-means clusters convex?



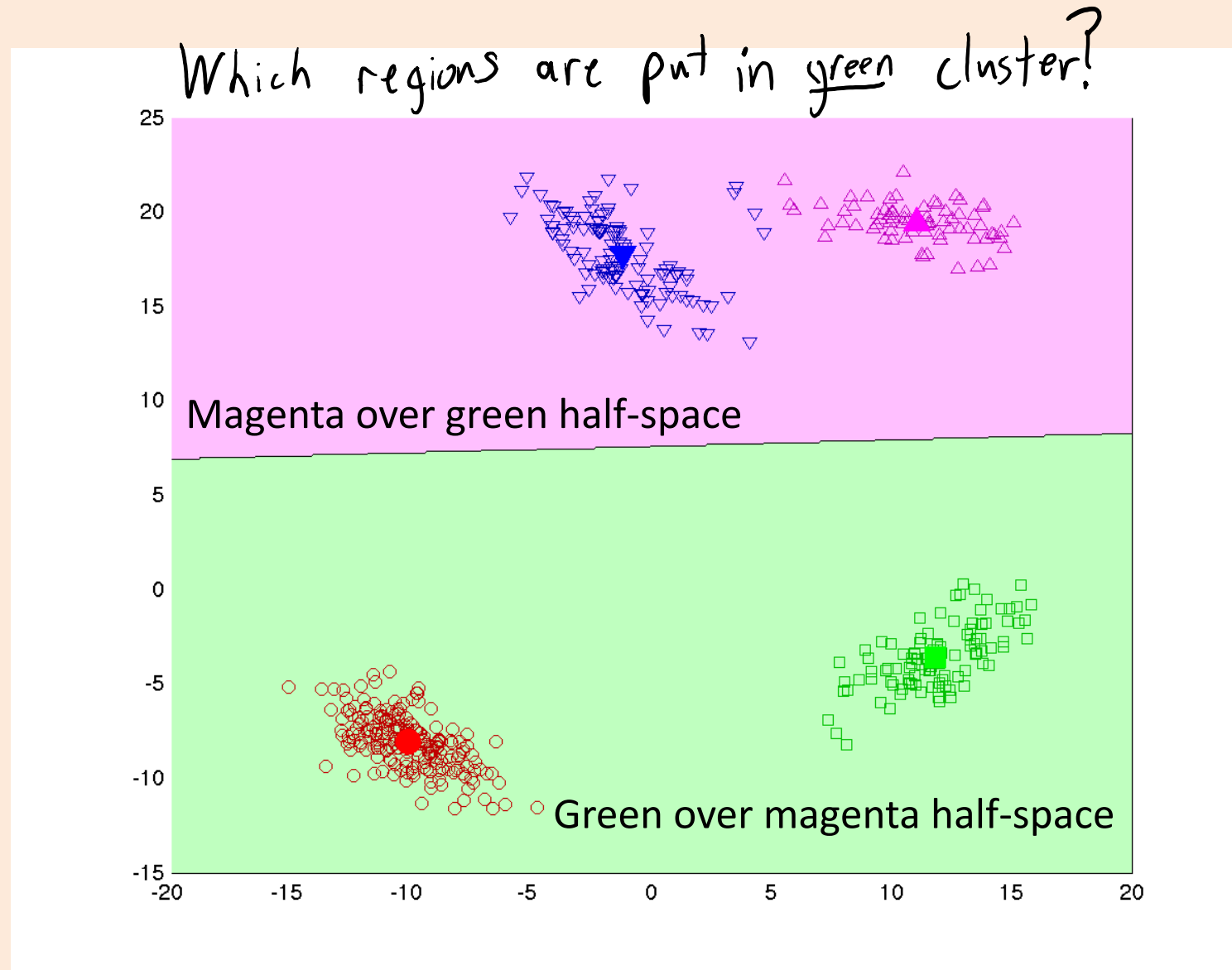
# Why are k-means clusters convex?



# Why are k-means clusters convex?

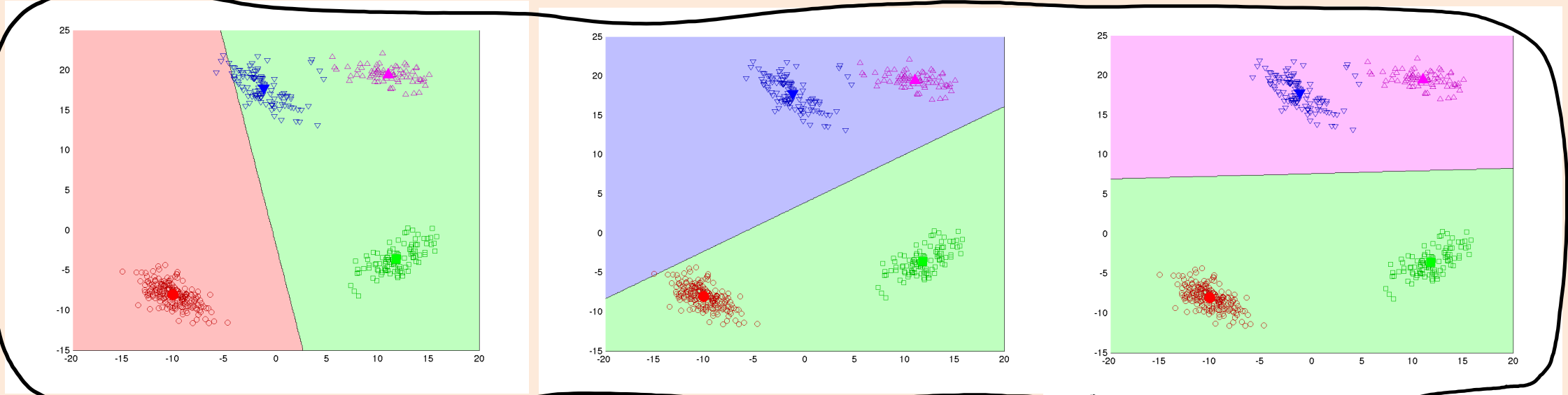


# Why are k-means clusters convex?



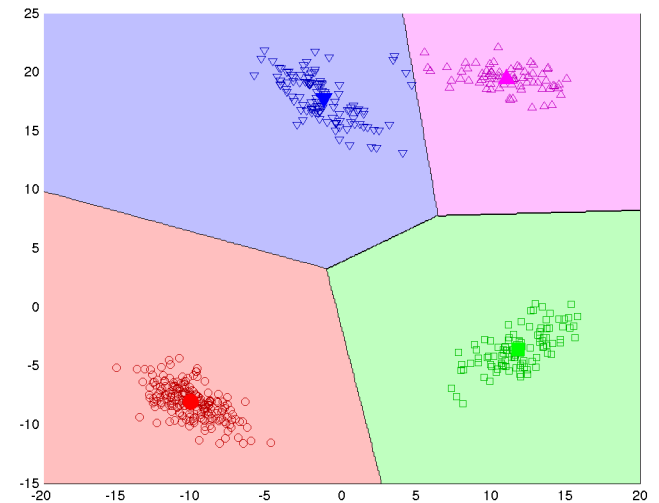


# Why are k-means clusters convex?



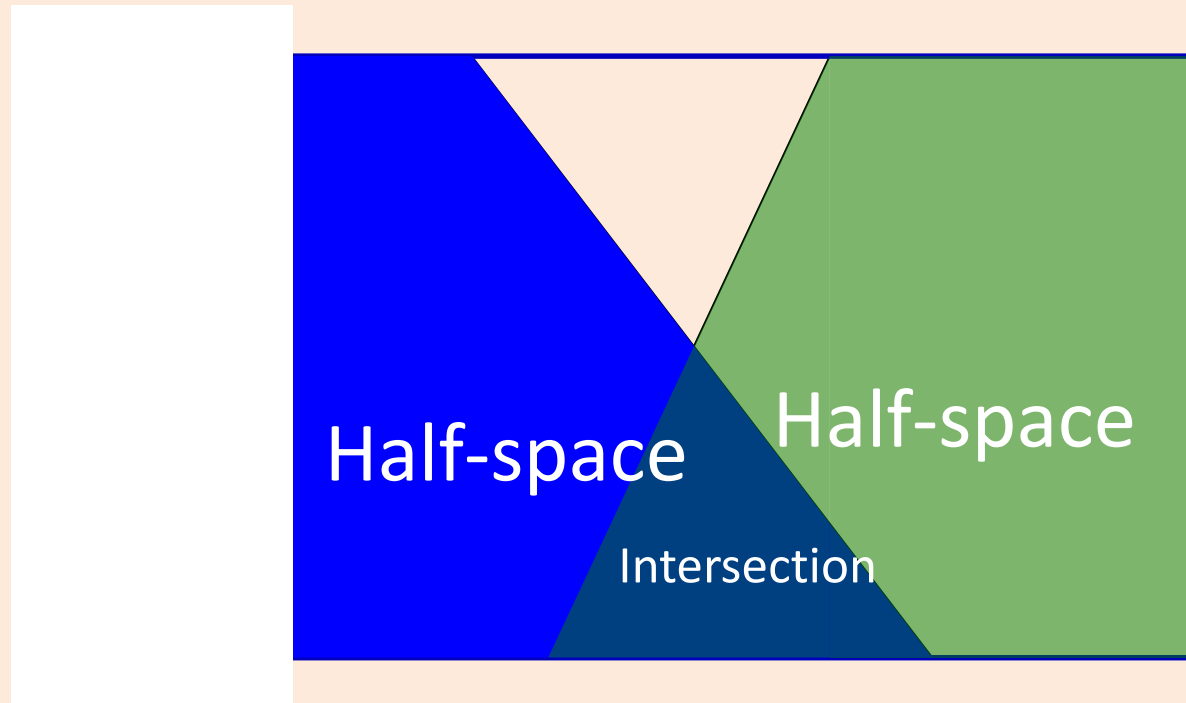
Green "cluster" is the intersection of these three half-spaces.

Here is what the four clusters look like:



# Why are k-means clusters convex?

- Half-spaces are convex sets.
- Intersection of convex sets is a convex set.
- So intersection of half-spaces is convex.



# Why are k-means clusters convex?

- Formal **proof** that "cluster 1" is convex (so all clusters are convex):

Let  $x_i$  and  $x_j$  be arbitrary points in cluster 1.

→ By def'n of cluster 1,  $\|x_i - w_1\| \leq \|x_i - w_c\|$  for all 'c' } equality  
 $\|x_j - w_1\| \leq \|x_j - w_c\|$  for all 'c' } for  $c=1$

→ Let  $x_m$  be an arbitrary point between  $x_i$  and  $x_j$ .

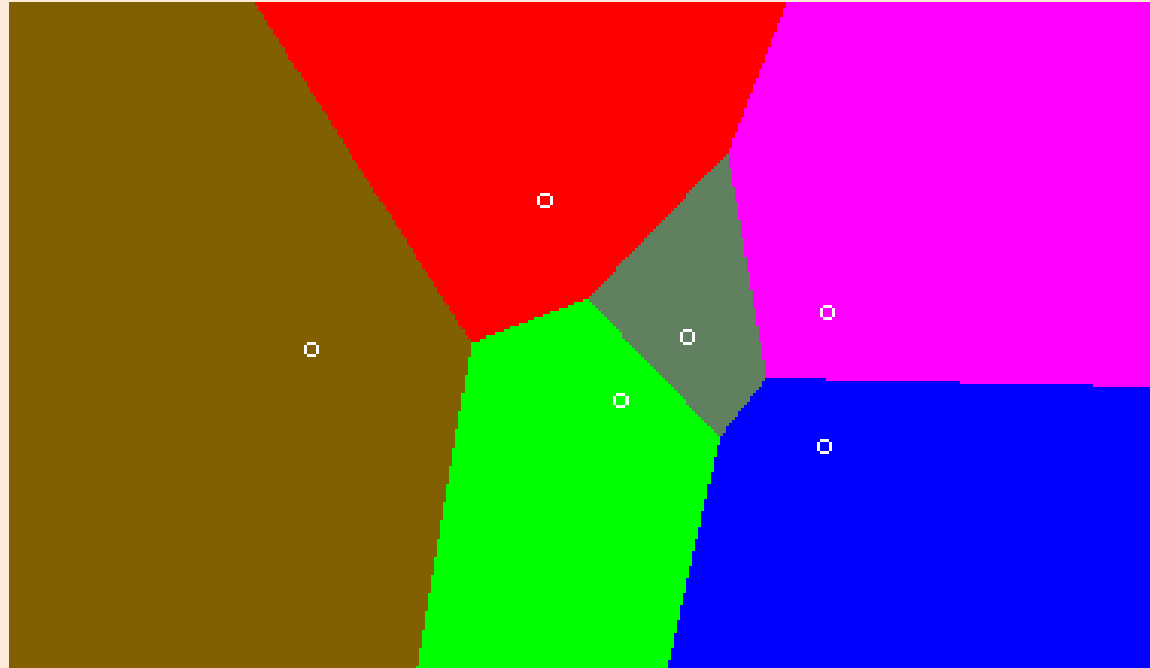
→ So we can write it as  $x_m = \theta x_i + (1-\theta)x_j$  for some  $\theta \in [0, 1]$

$$\begin{aligned} \text{Then } \|x_m - w_1\| &= \|\theta x_i + (1-\theta)x_j - (\theta w_1 + (1-\theta)w_1)\| && (w_1 = \theta w_1 + (1-\theta)w_1) \\ &\leq \|\theta x_i - \theta w_1\| + \|(1-\theta)x_j - (1-\theta)w_1\| && (\text{triangle inequality}) \end{aligned}$$

$$\begin{aligned} &= \theta \|x_i - w_1\| + (1-\theta) \|x_j - w_1\| && (\text{homogeneity of norms}) \\ \left. \begin{array}{l} x_i \text{ and } x_j \\ \text{are in cluster 1} \end{array} \right\} &\leq \theta \|x_i - w_c\| + (1-\theta) \|x_j - w_c\| = \|x_j - w_c\| && \text{so } x_m \text{ is in cluster 1!} \end{aligned}$$

# Voronoi Diagrams

- The k-means partition can be visualized as a [Voronoi diagram](#):

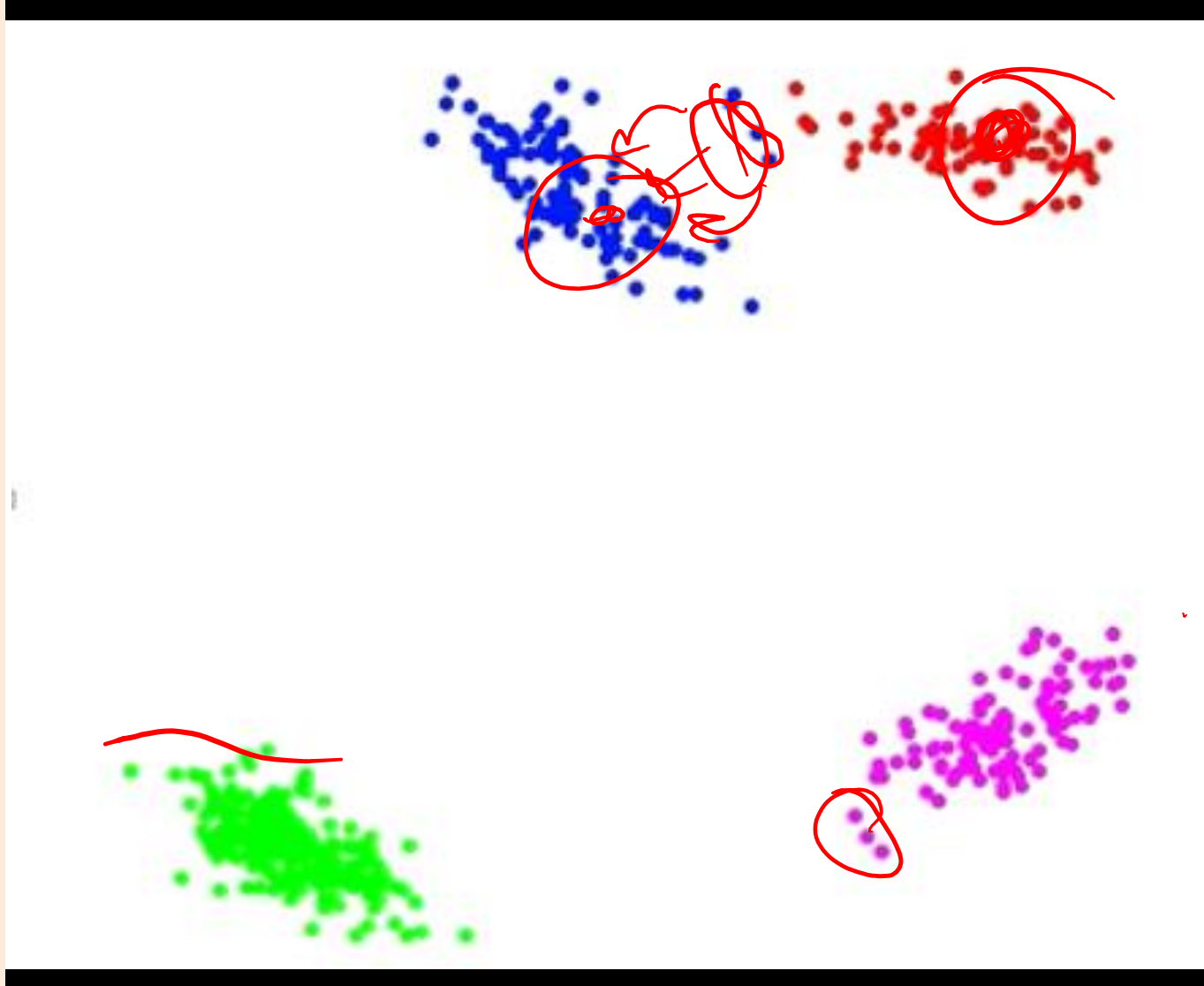


- Can be a useful visualization of “nearest available” problems.
  - E.g., [nearest tube station in London](#).

# UBClustering Algorithm

- Let's define a new ensemble clustering method: **UBClustering**.
  1. Run k-means with 'm' different random initializations.
  2. For each object i and j:
    - Count the number of times  $x_i$  and  $x_j$  are in the same cluster.
    - Define  $p(i,j) = \text{count}(x_i \text{ in same cluster as } x_j)/m$ .
  3. Put  $x_i$  and  $x_j$  in the same cluster if  $p(i,j) > 0.5$ .
- Like DBSCAN **merge clusters** in step 3 if i or j are already assigned.
  - You can implement this with a DBSCAN code (just changes "distance").
  - Each  $x_i$  has an  $x_j$  in its cluster with  $p(i,j) > 0.5$ .
  - Some points are not assigned to any cluster.

# UBClustering Algorithm



It looks like DBSCAN, but far-away points will be assigned to a cluster if they always appear in same cluster as other points.