

Unstructured Data Analysis

Lecture 8: Clustering (cont'd), topic modeling

George Chen

Last Time:

Automatically Choosing k , the Number of Clusters

Simple strategy:

(1) compute a score for each k you're willing to try

(2) use whichever k achieves the best score

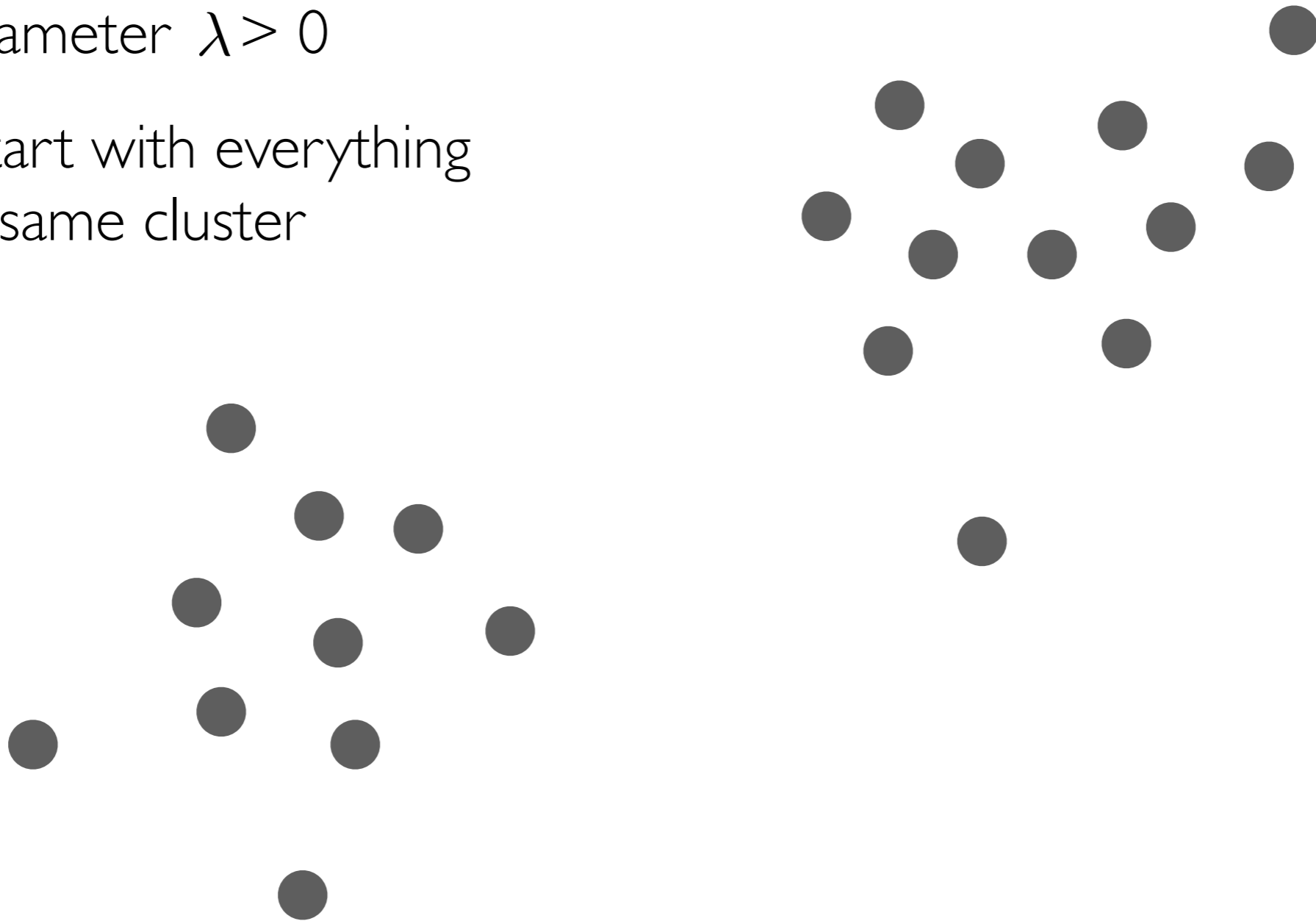
- There is no single best score function
- Fitting k-means/GMMs is in general *random*
(for example, in the CH index demo, if we don't set `random_state`, then the CH indices computed will be different every time we run the code, and the best k could change!)

There are other clustering methods that do not require specifying the number of clusters (e.g., DP-means, DP-GMM, many variants of hierarchical clustering, density-based clustering)

DP-means

Step 0. Pick concentration parameter $\lambda > 0$

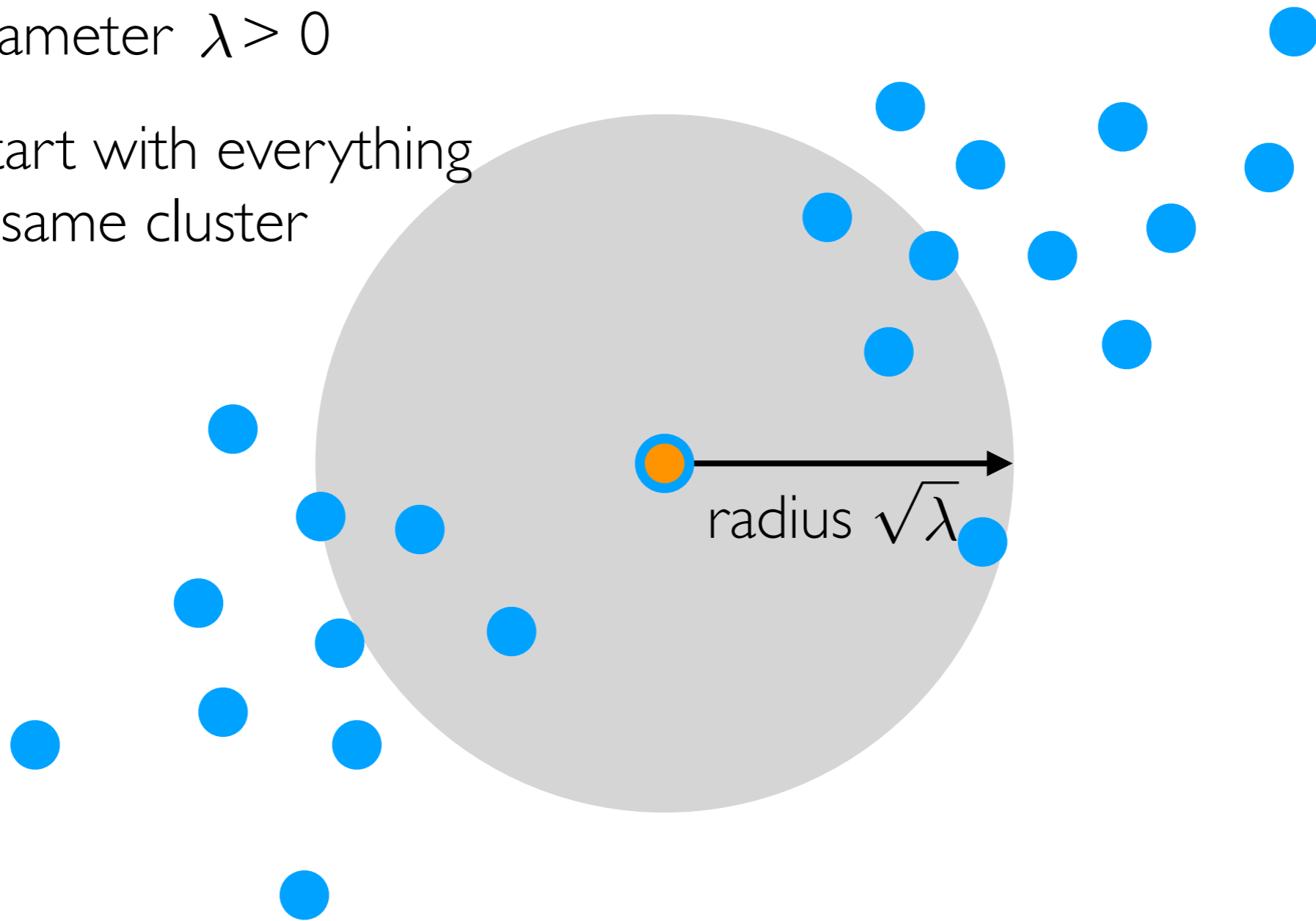
Step 1. Start with everything in same cluster



DP-means

Step 0. Pick concentration parameter $\lambda > 0$

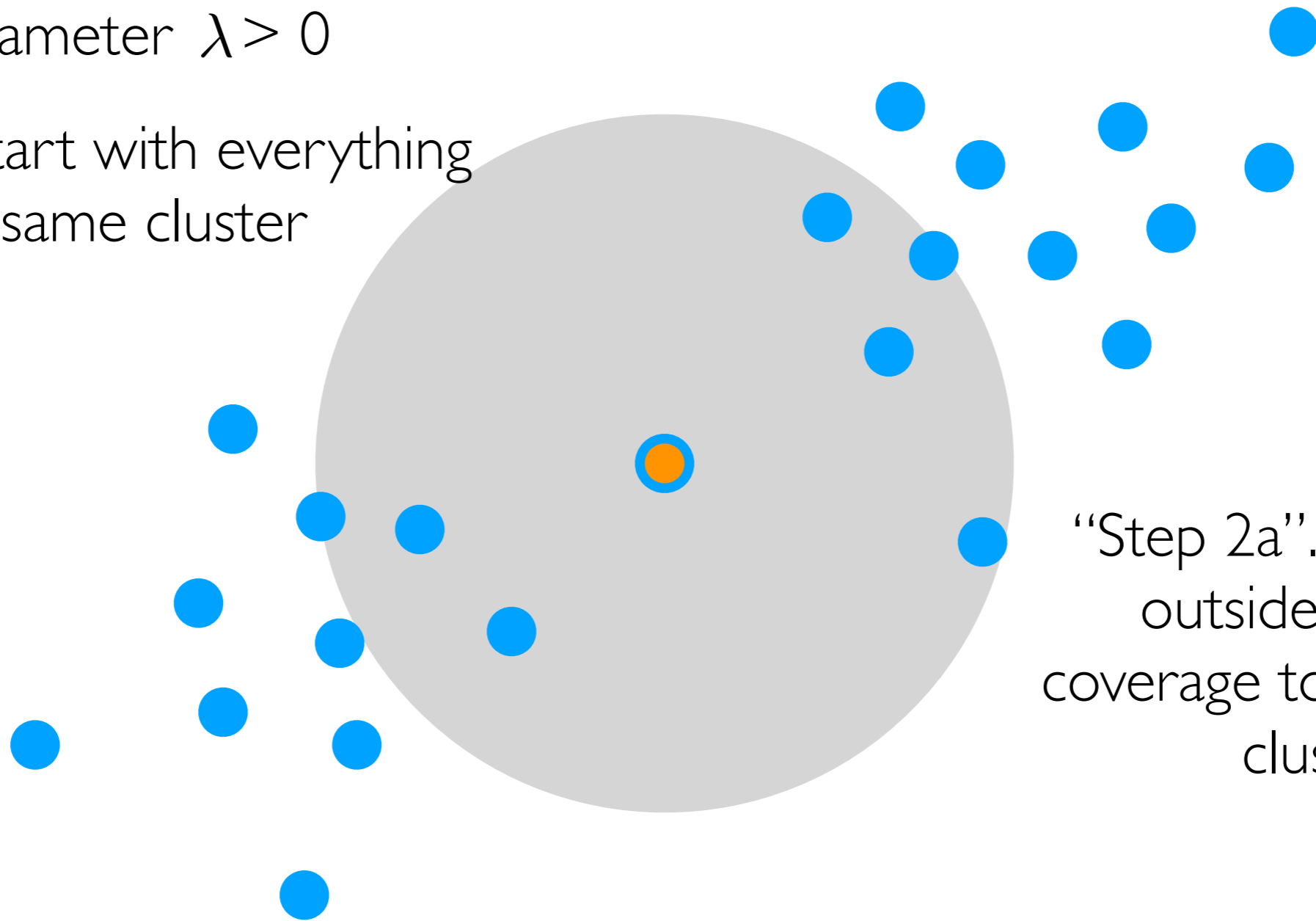
Step 1. Start with everything in same cluster



DP-means

Step 0. Pick concentration parameter $\lambda > 0$

Step 1. Start with everything in same cluster

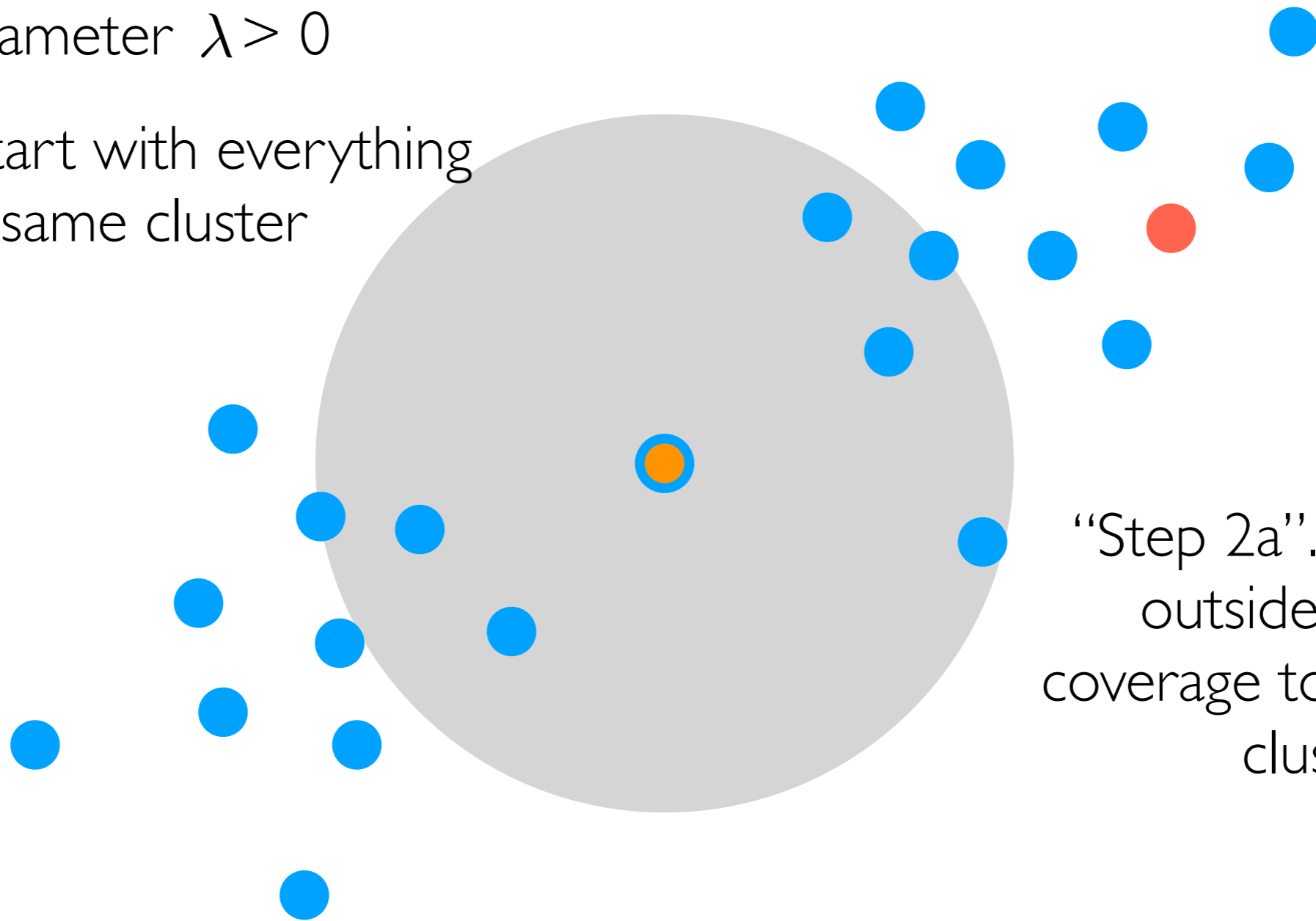


“Step 2a”. Pick point outside of gray coverage to make new cluster

DP-means

Step 0. Pick concentration parameter $\lambda > 0$

Step 1. Start with everything in same cluster

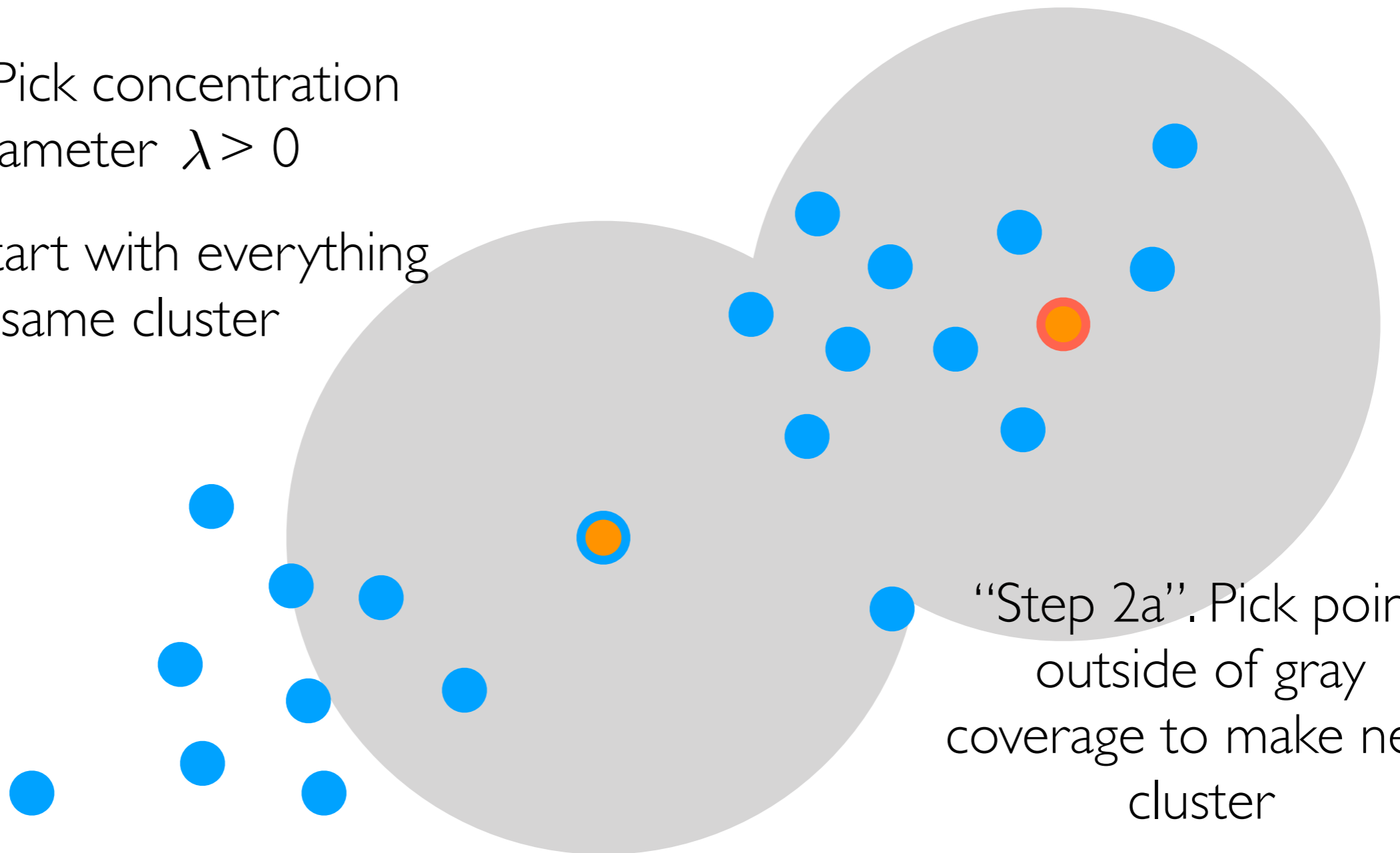


“Step 2a”. Pick point outside of gray coverage to make new cluster

DP-means

Step 0. Pick concentration parameter $\lambda > 0$

Step 1. Start with everything in same cluster



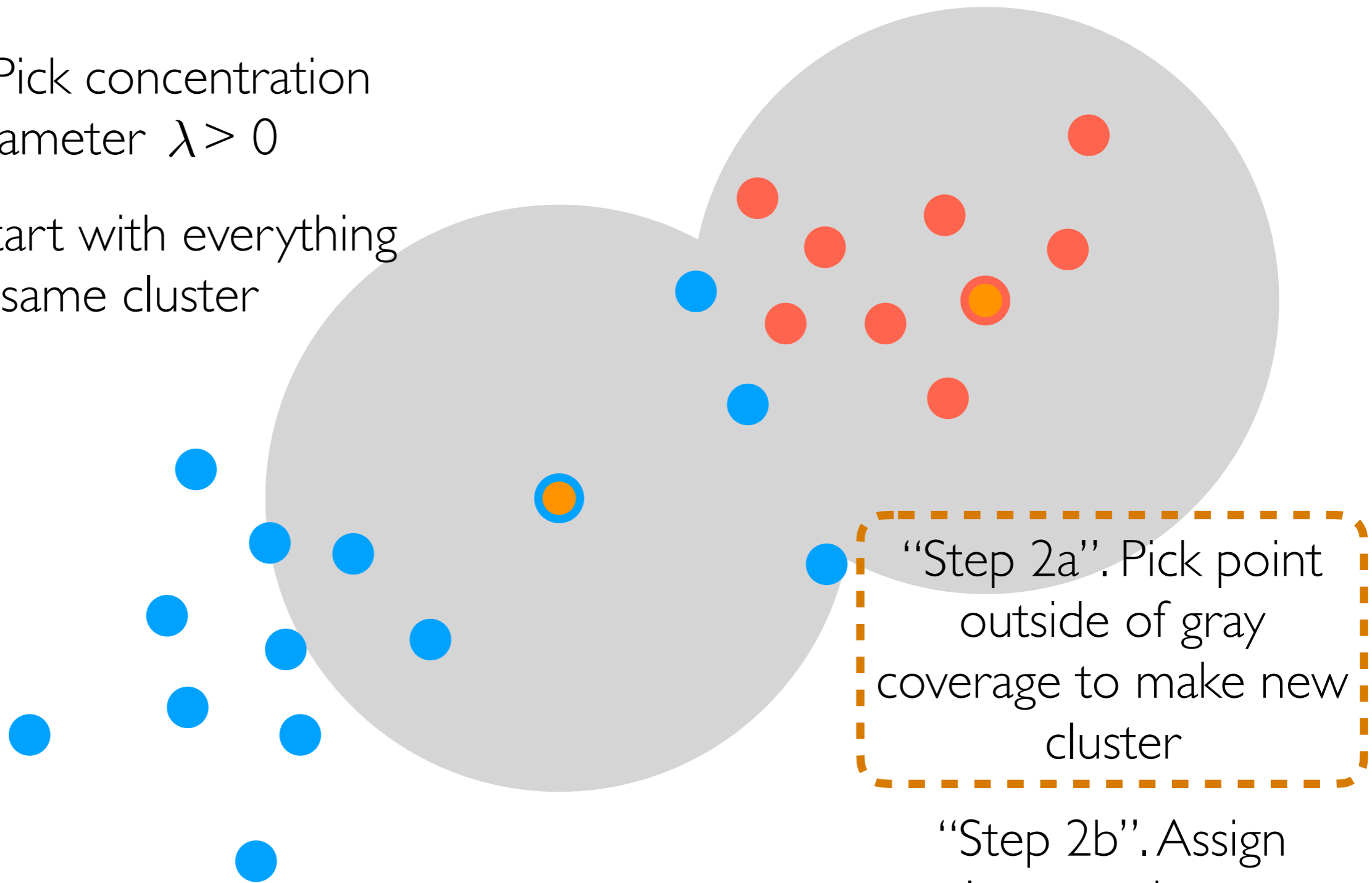
“Step 2a”. Pick point outside of gray coverage to make new cluster

“Step 2b”. Assign closest points to current clusters

DP-means

Step 0. Pick concentration parameter $\lambda > 0$

Step 1. Start with everything in same cluster



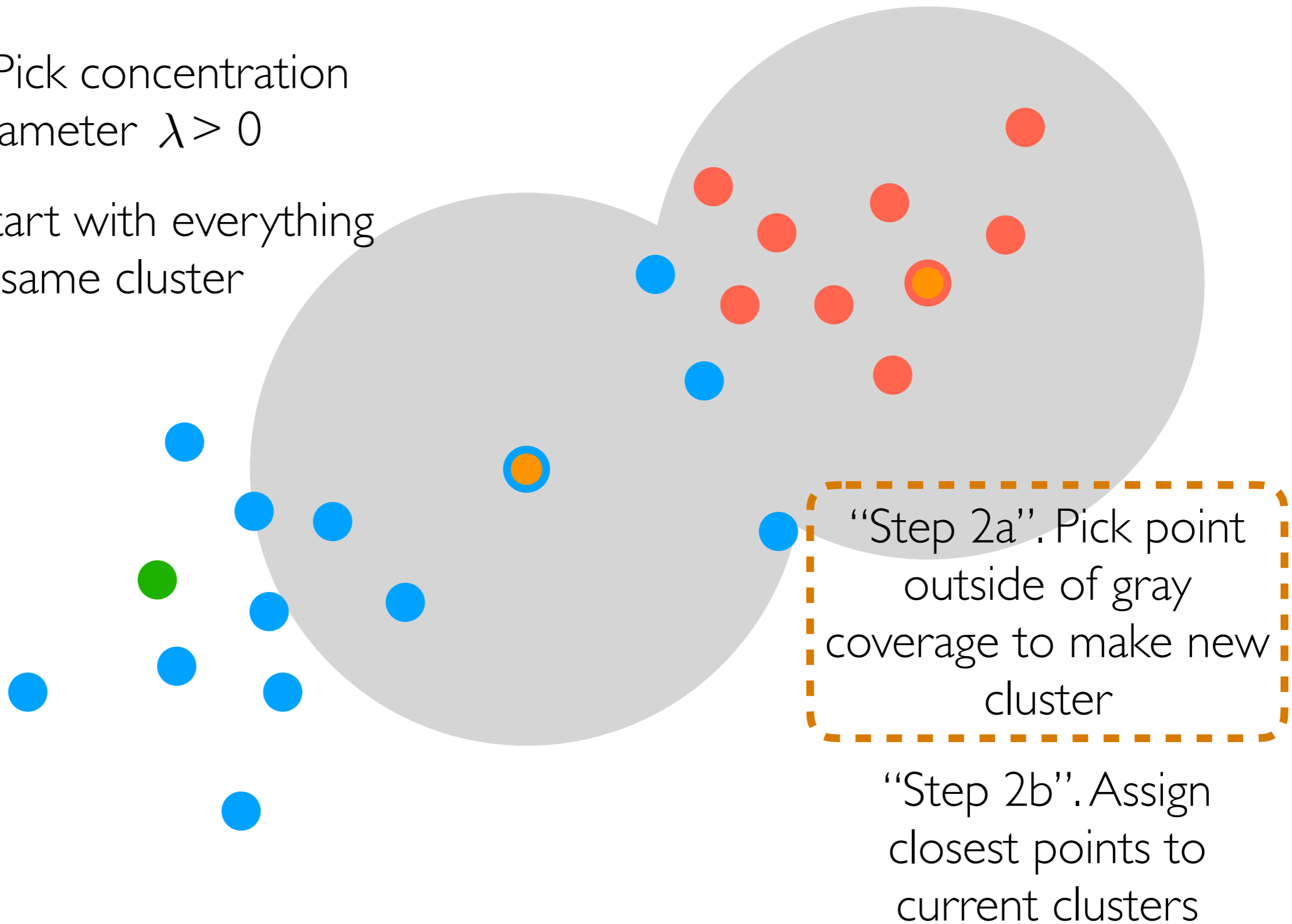
"Step 2a". Pick point outside of gray coverage to make new cluster

"Step 2b". Assign closest points to current clusters

DP-means

Step 0. Pick concentration parameter $\lambda > 0$

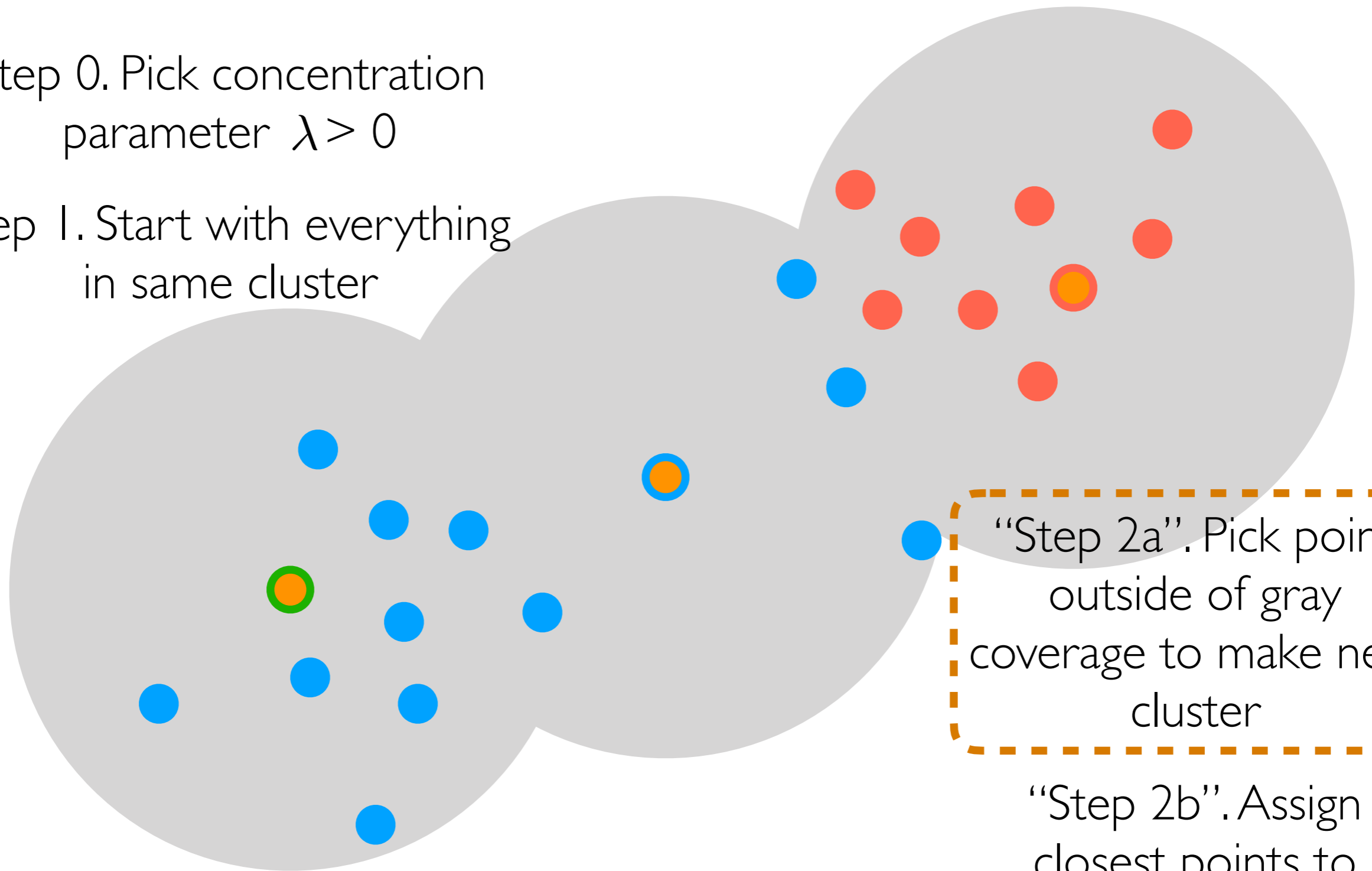
Step 1. Start with everything in same cluster



DP-means

Step 0. Pick concentration parameter $\lambda > 0$

Step 1. Start with everything in same cluster



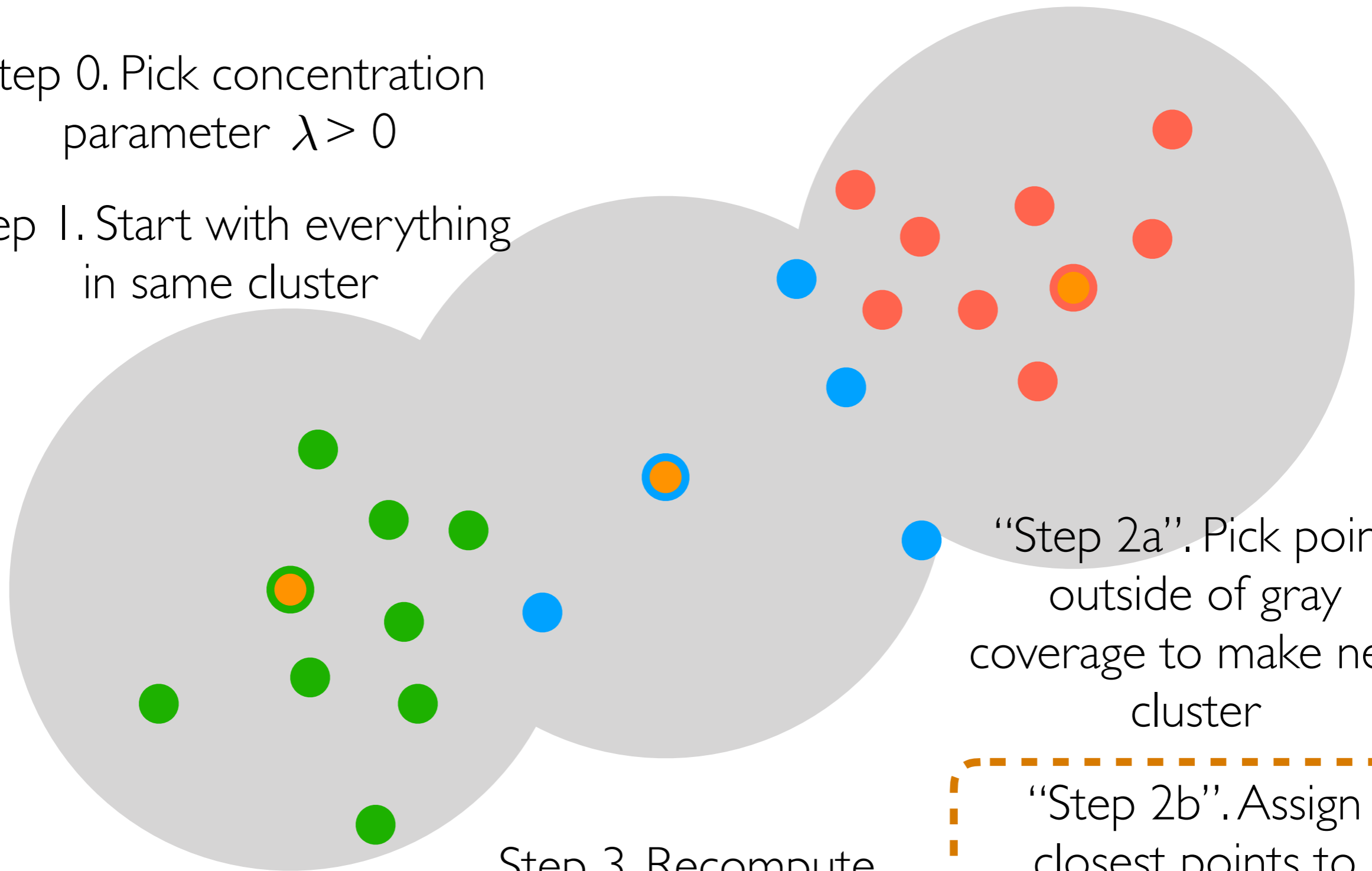
“Step 2a”. Pick point outside of gray coverage to make new cluster

“Step 2b”. Assign closest points to current clusters

DP-means

Step 0. Pick concentration parameter $\lambda > 0$

Step 1. Start with everything in same cluster



“Step 2a”. Pick point outside of gray coverage to make new cluster

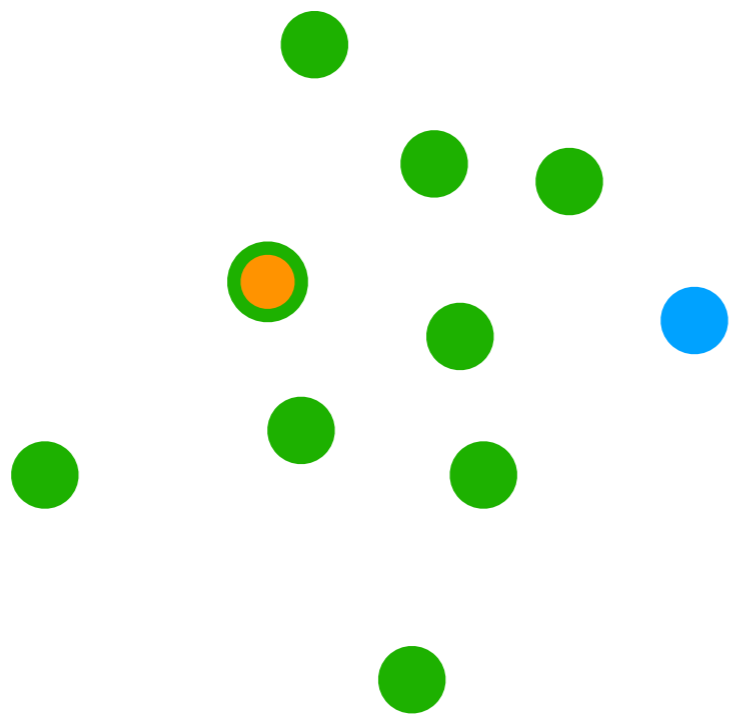
“Step 2b”. Assign closest points to current clusters

Step 3. Recompute cluster centers

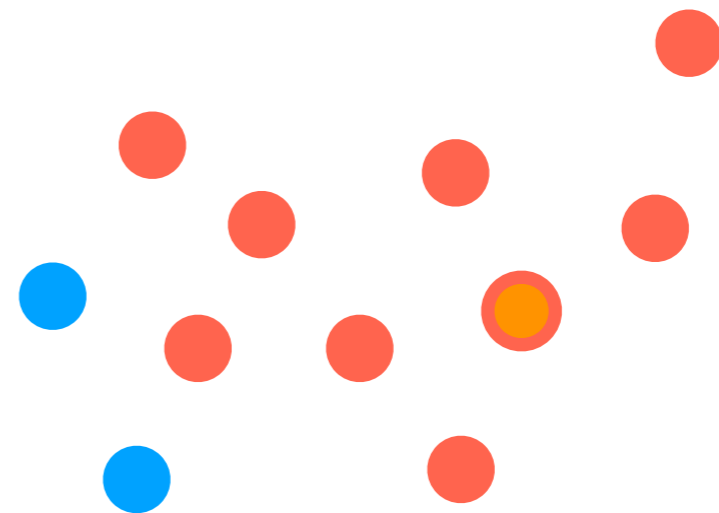
DP-means

Step 0. Pick concentration parameter $\lambda > 0$

Step 1. Start with everything in same cluster



Step 3. Recompute cluster centers



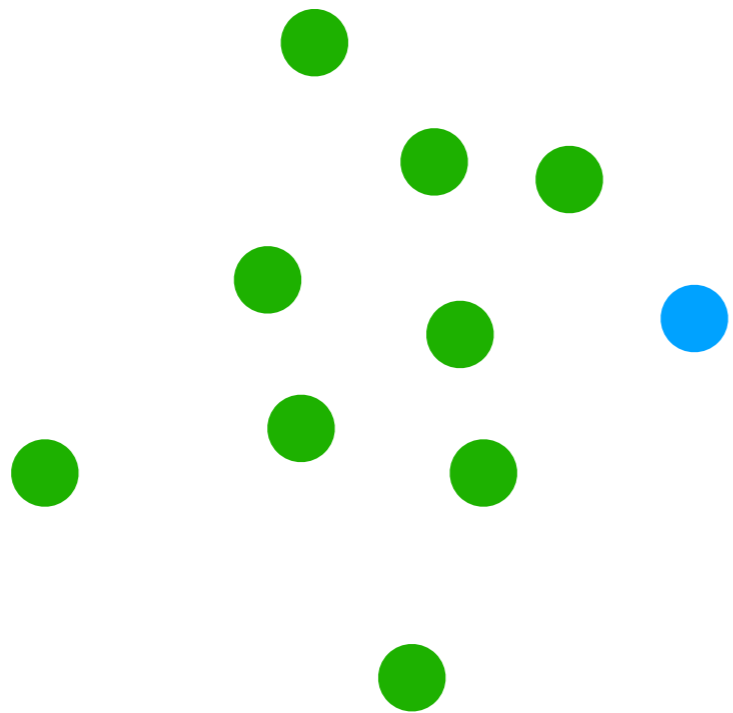
“Step 2a”. Pick point outside of gray coverage to make new cluster

“Step 2b”. Assign closest points to current clusters

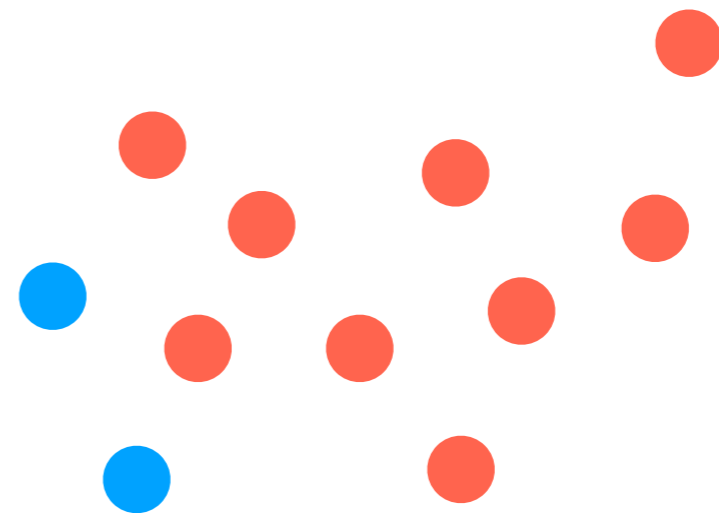
DP-means

Step 0. Pick concentration parameter $\lambda > 0$

Step 1. Start with everything in same cluster



Step 3. Recompute cluster centers



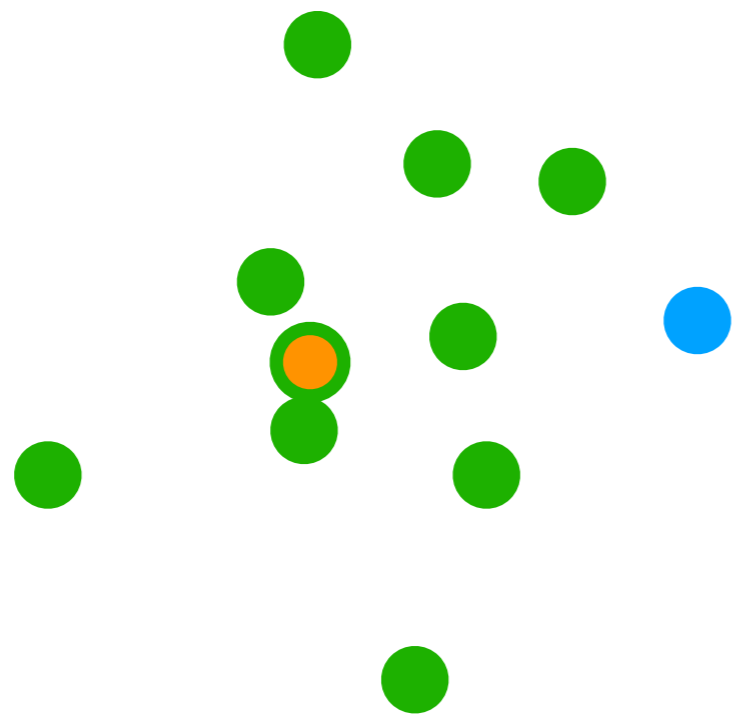
“Step 2a”. Pick point outside of gray coverage to make new cluster

“Step 2b”. Assign closest points to current clusters

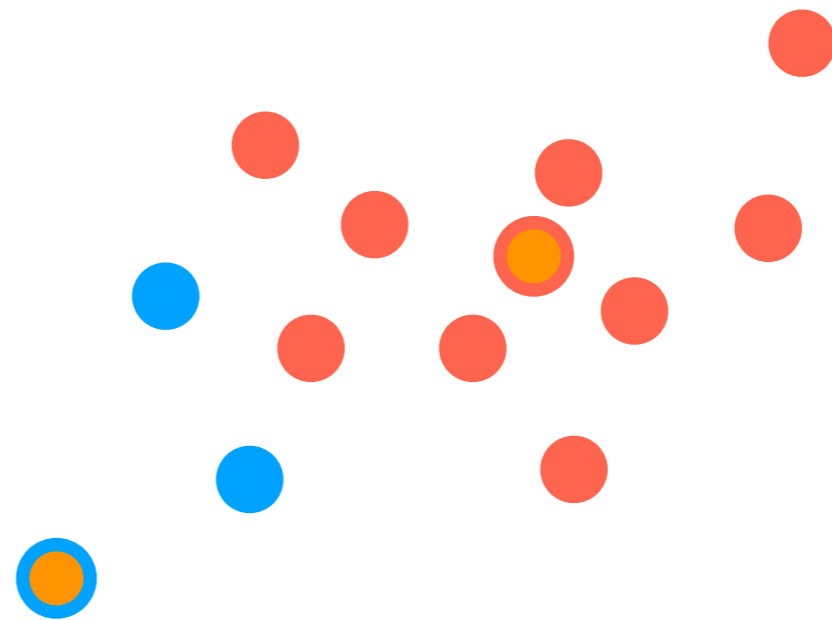
DP-means

Step 0. Pick concentration parameter $\lambda > 0$

Step 1. Start with everything in same cluster



Step 3. Recompute cluster centers



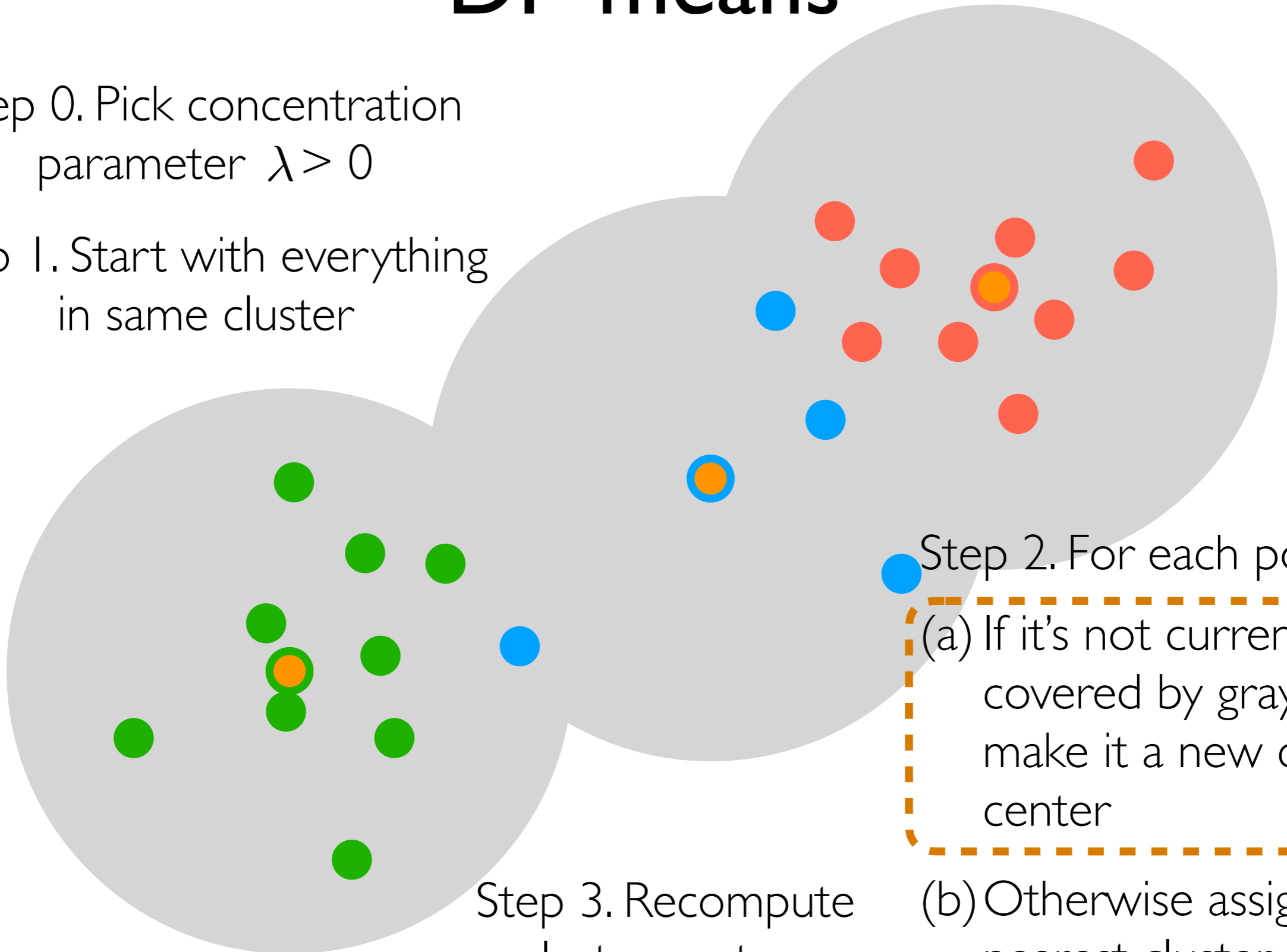
“Step 2a”. Pick point outside of gray coverage to make new cluster

“Step 2b”. Assign closest points to current clusters

DP-means

Step 0. Pick concentration parameter $\lambda > 0$

Step 1. Start with everything in same cluster



Step 2. For each point:

(a) If it's not currently covered by gray balls, make it a new cluster center

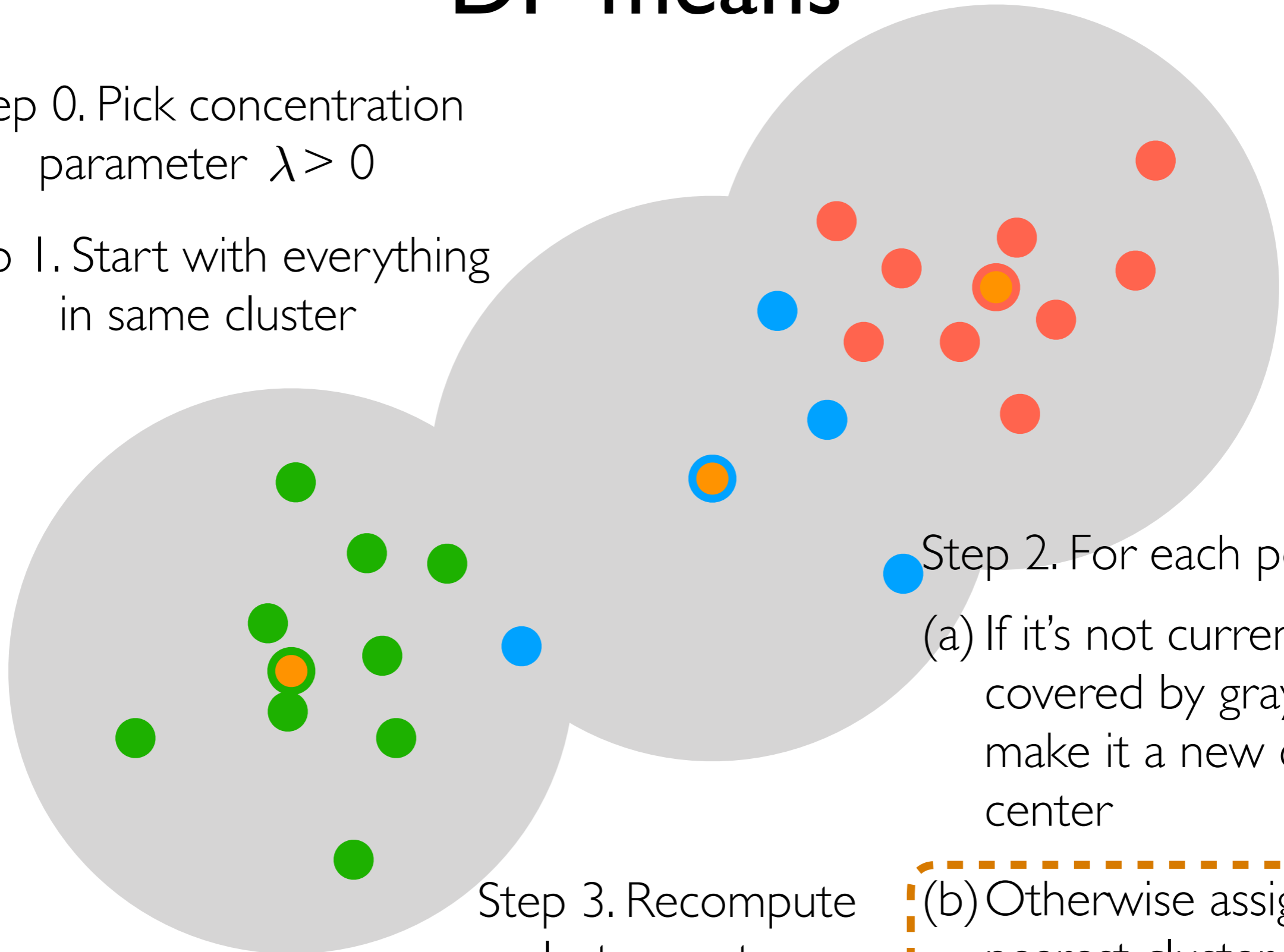
(b) Otherwise assign it to nearest cluster

Step 3. Recompute cluster centers

DP-means

Step 0. Pick concentration parameter $\lambda > 0$

Step 1. Start with everything in same cluster



Step 2. For each point:
(a) If it's not currently covered by gray balls, make it a new cluster center

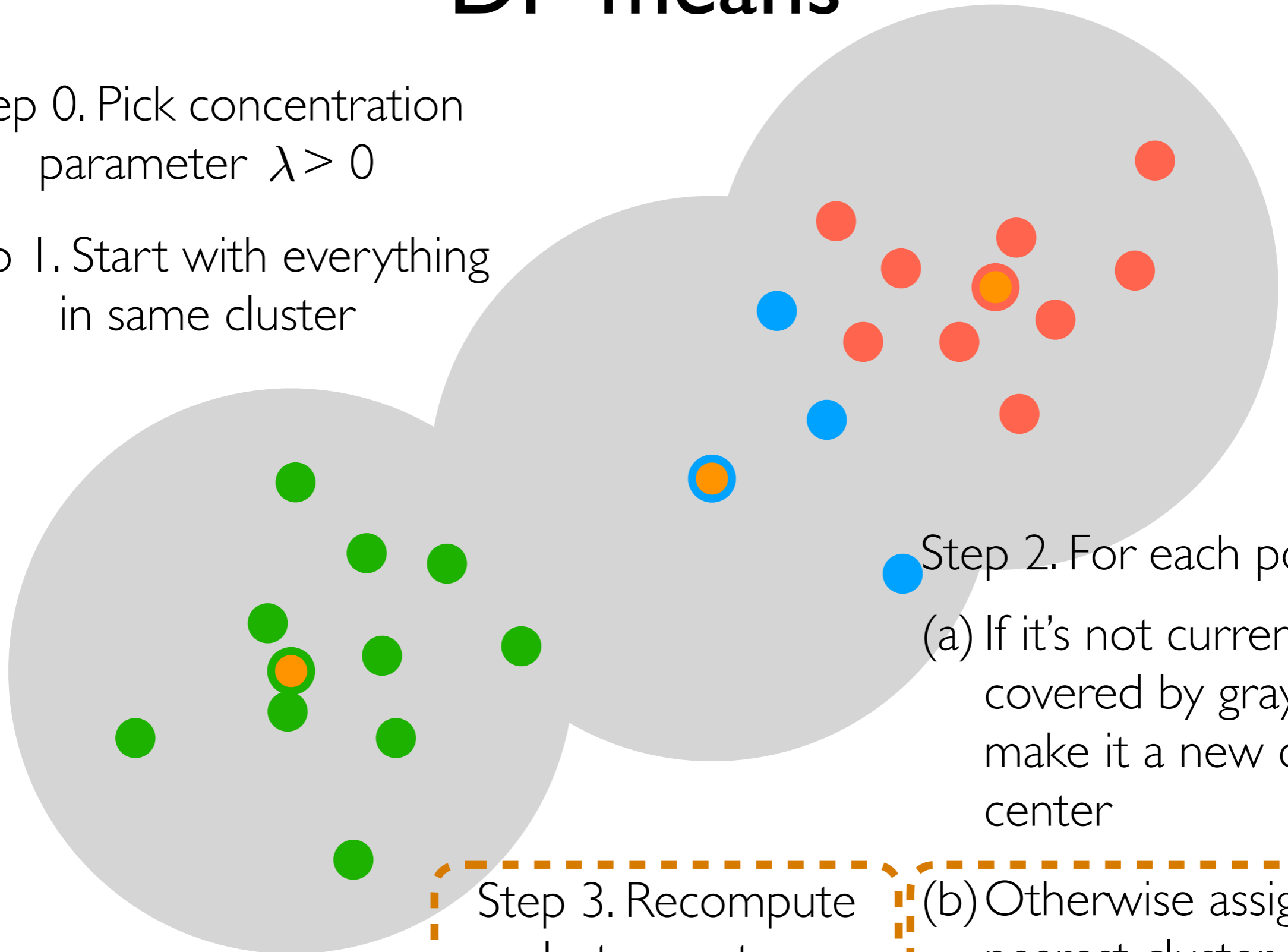
(b) Otherwise assign it to nearest cluster

Step 3. Recompute cluster centers

DP-means

Step 0. Pick concentration parameter $\lambda > 0$

Step 1. Start with everything in same cluster



Step 2. For each point:
(a) If it's not currently covered by gray balls, make it a new cluster center

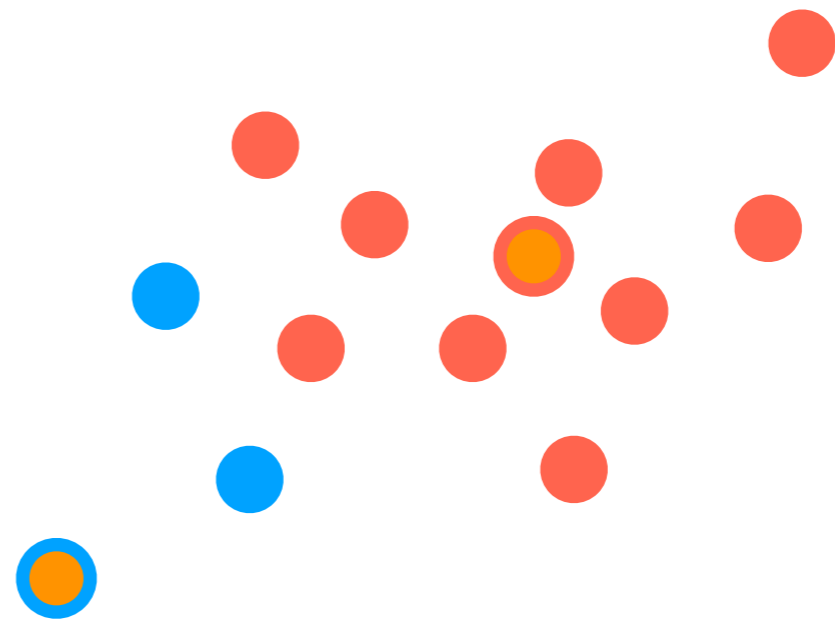
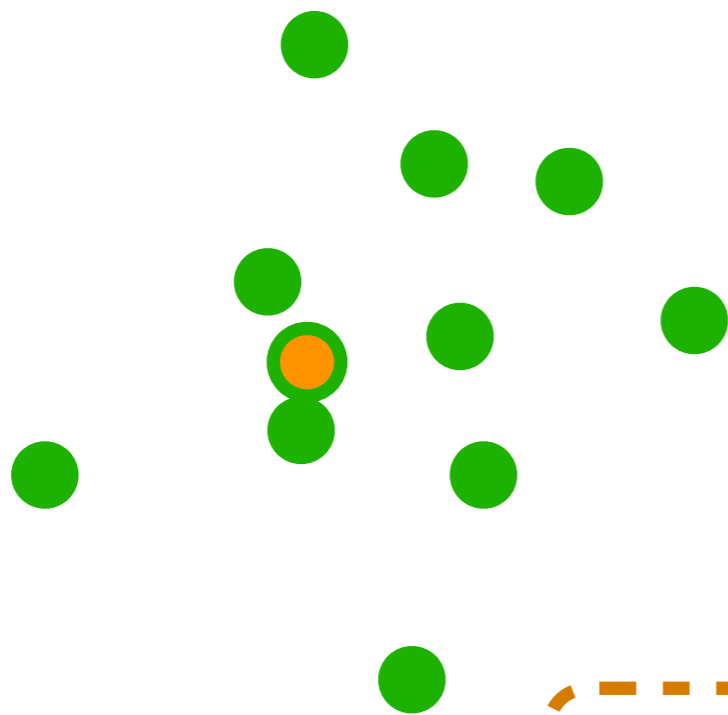
Step 3. Recompute cluster centers

(b) Otherwise assign it to nearest cluster

DP-means

Step 0. Pick concentration parameter $\lambda > 0$

Step 1. Start with everything in same cluster



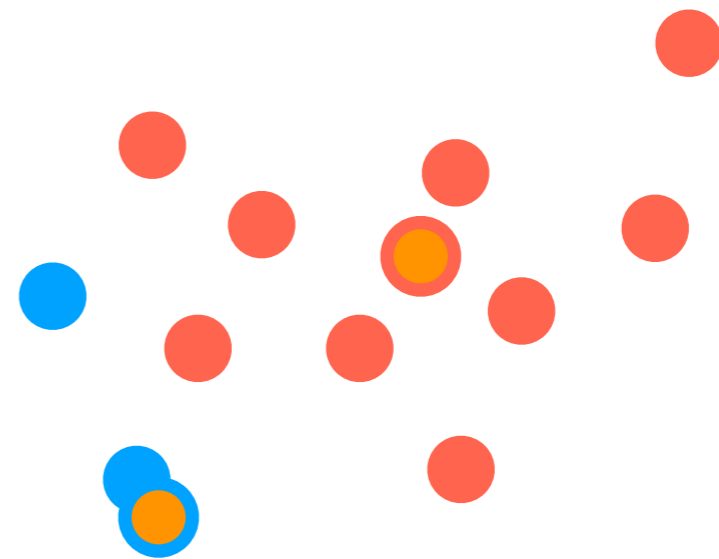
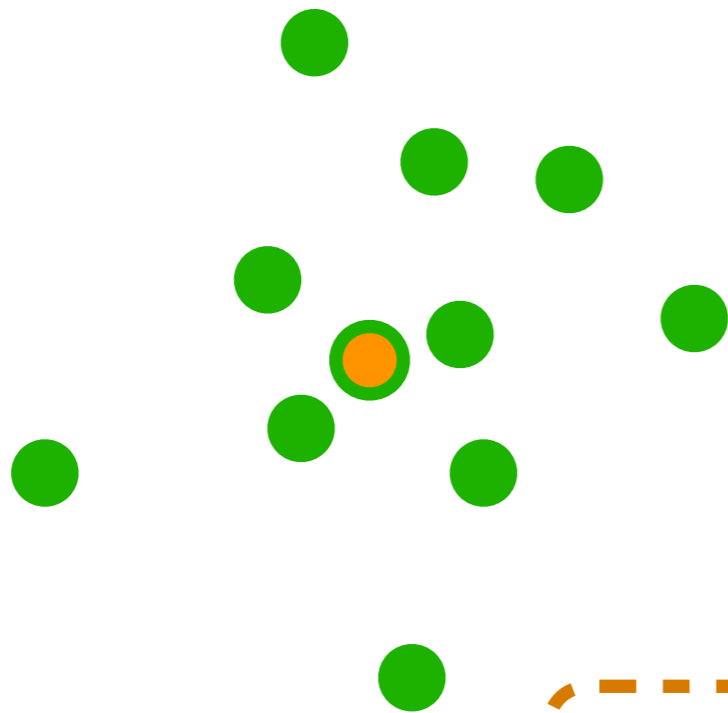
- Step 2. For each point:
- (a) If it's not currently covered by gray balls, make it a new cluster center
 - (b) Otherwise assign it to nearest cluster

Step 3. Recompute cluster centers

DP-means

Step 0. Pick concentration parameter $\lambda > 0$

Step 1. Start with everything in same cluster



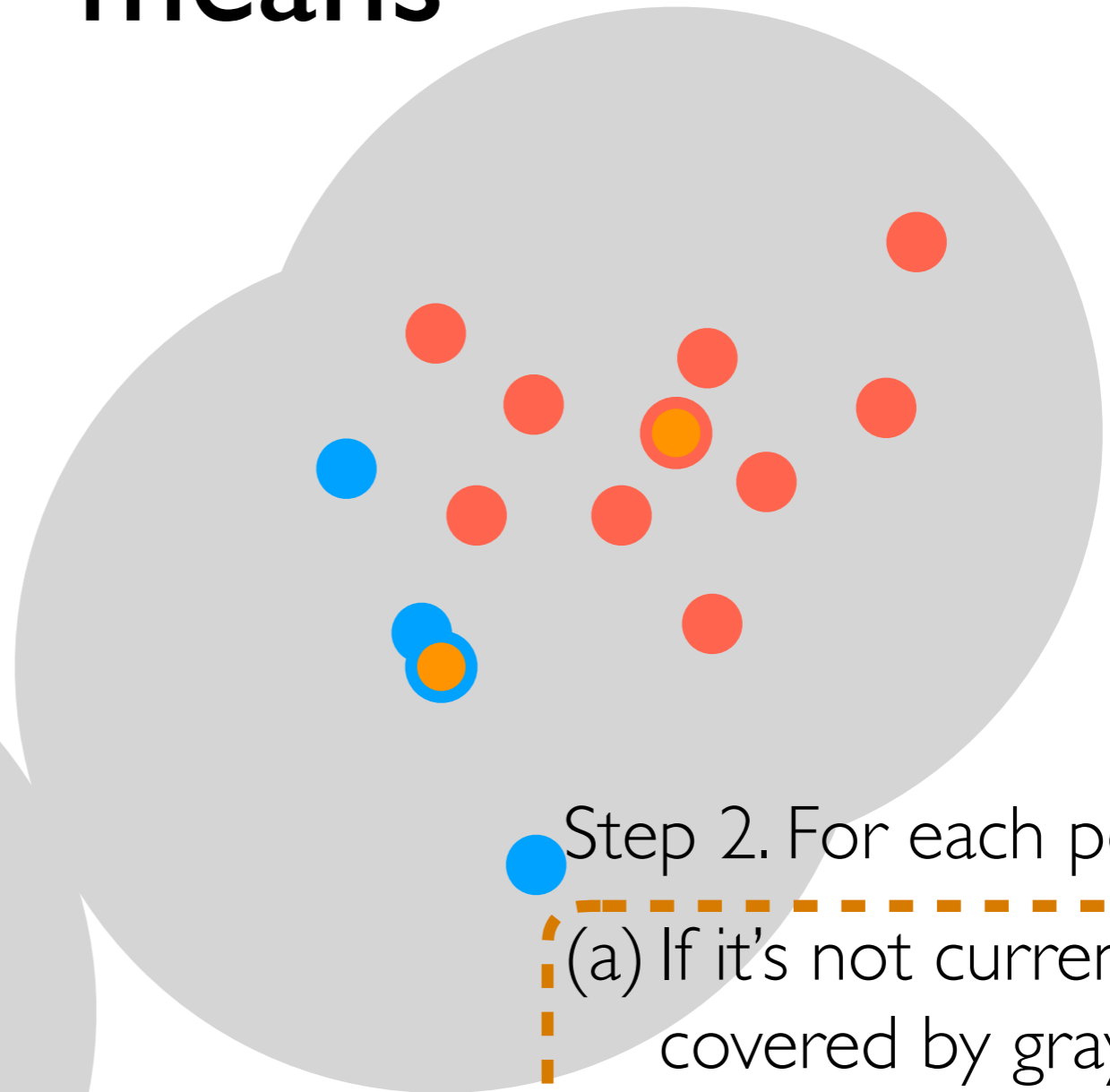
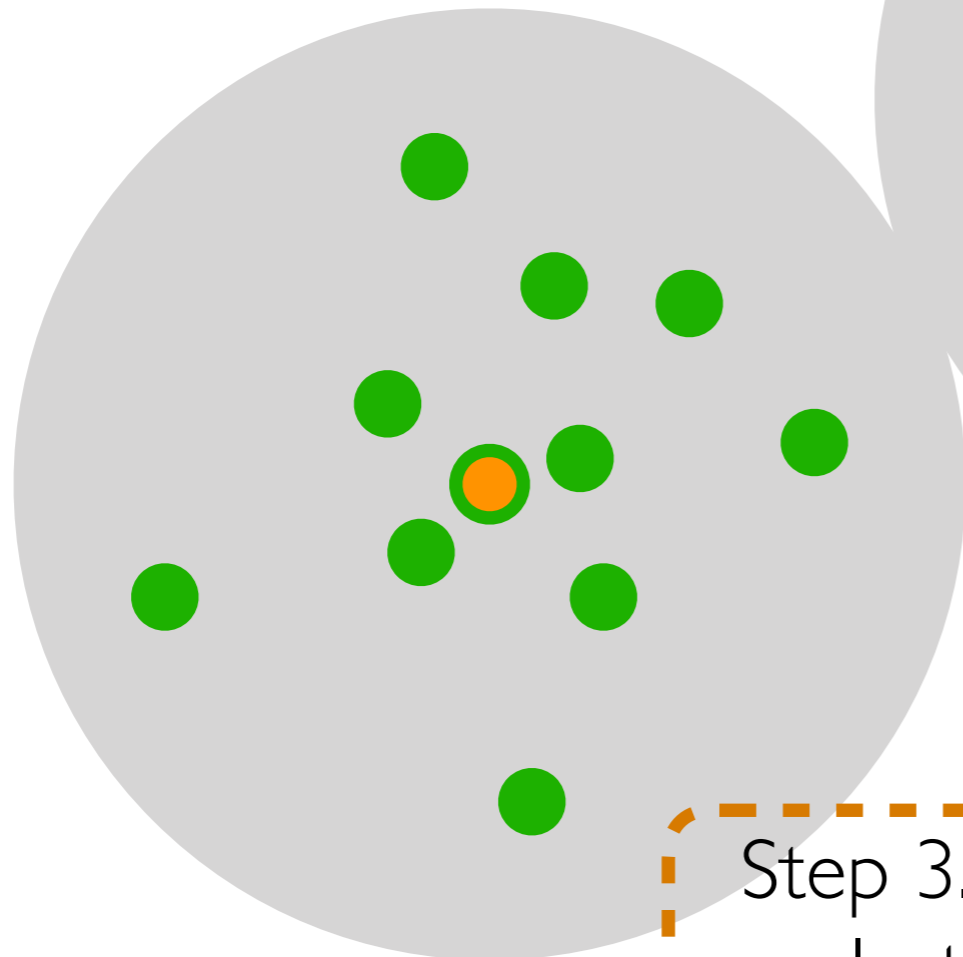
- Step 2. For each point:
- (a) If it's not currently covered by gray balls, make it a new cluster center
 - (b) Otherwise assign it to nearest cluster

Step 3. Recompute cluster centers

DP-means

Step 0: Pick concentration parameter $\lambda > 0$

Step 1: Start with everything in same cluster



Step 2. For each point:

(a) If it's not currently covered by gray balls, make it a new cluster center

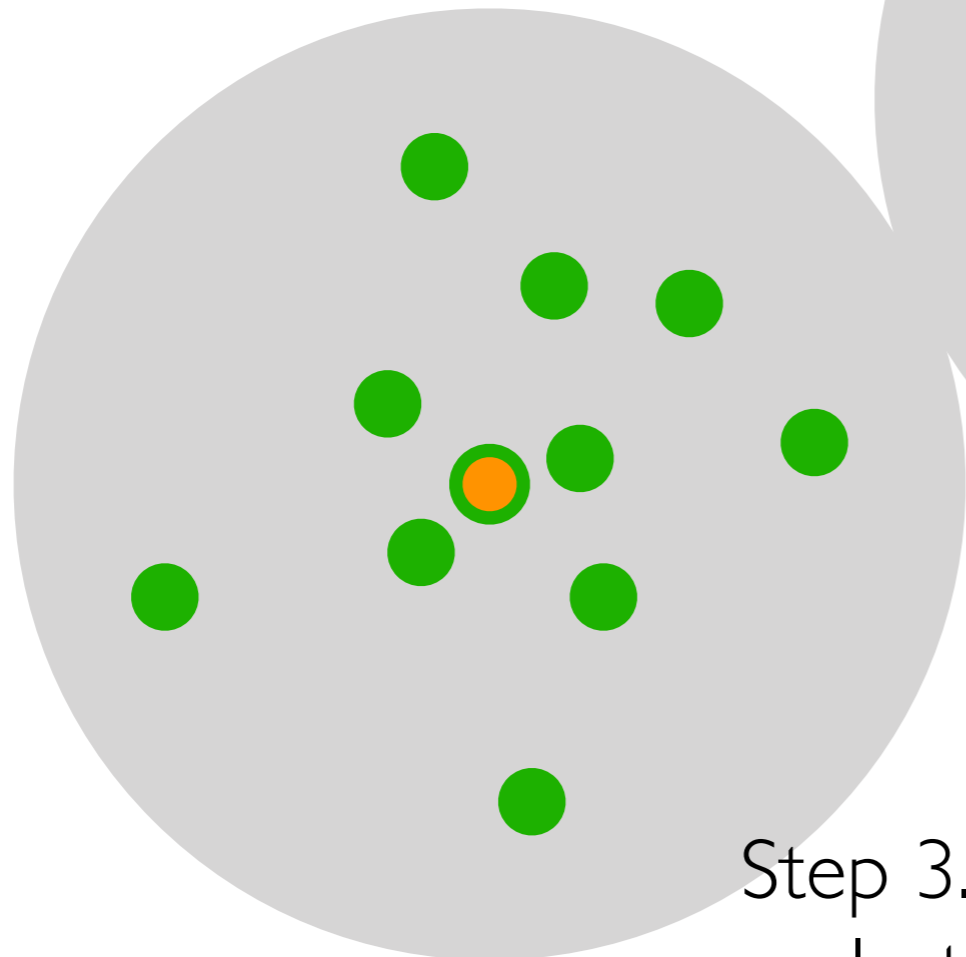
(b) Otherwise assign it to nearest cluster

Step 3. Recompute cluster centers

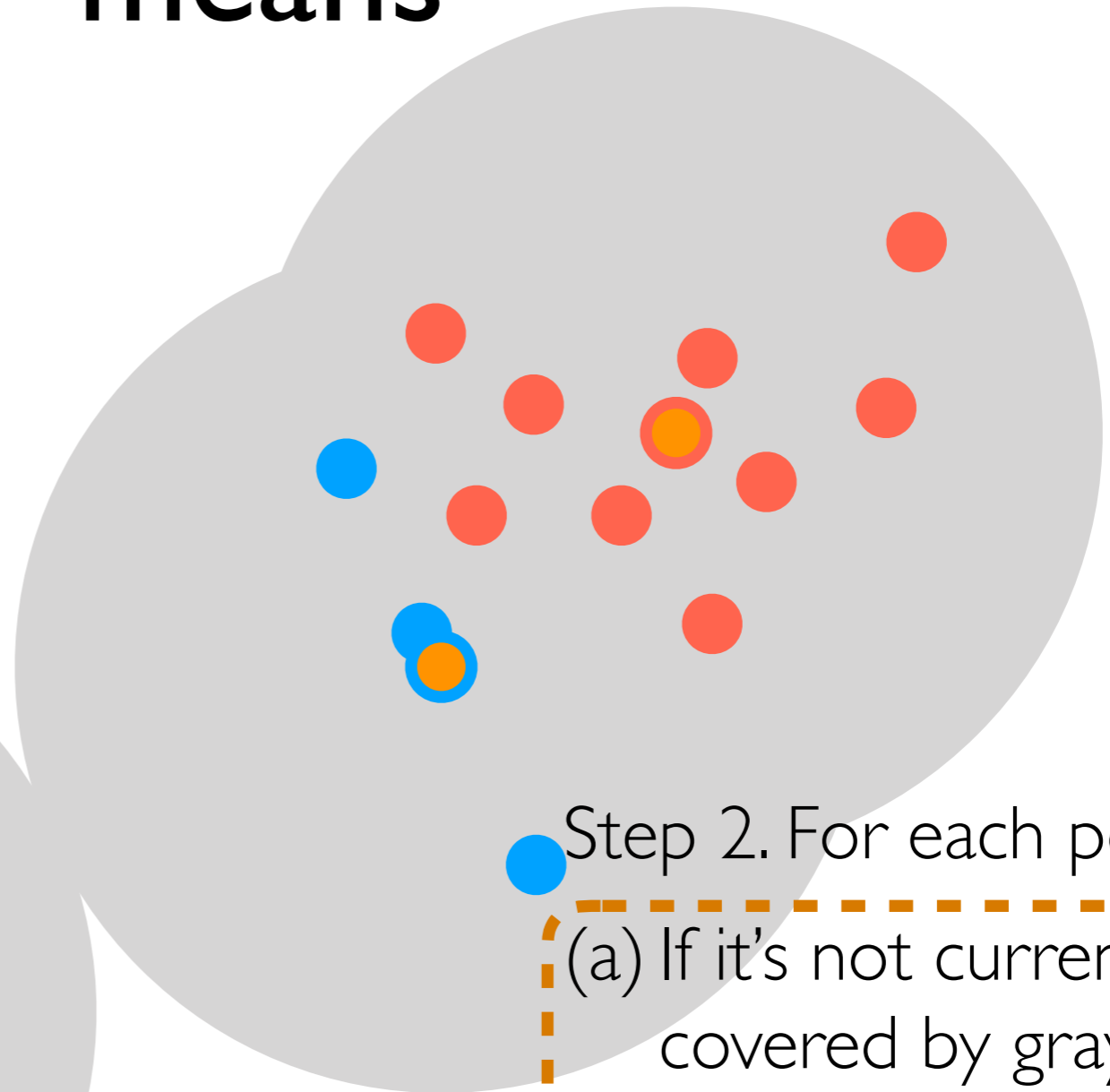
DP-means

Step 0: Pick concentration parameter $\lambda > 0$

Step 1: Start with everything in same cluster



Step 3. Recompute cluster centers



Step 2. For each point:

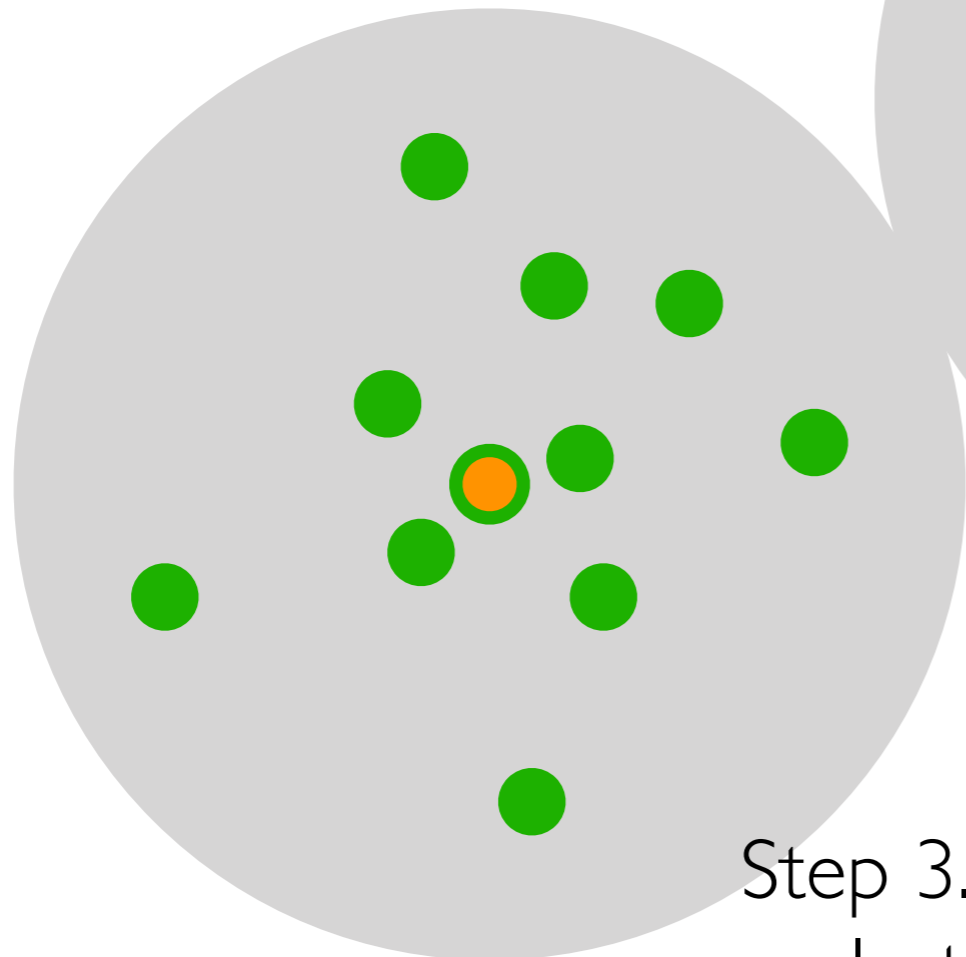
(a) If it's not currently covered by gray balls, make it a new cluster center

(b) Otherwise assign it to nearest cluster

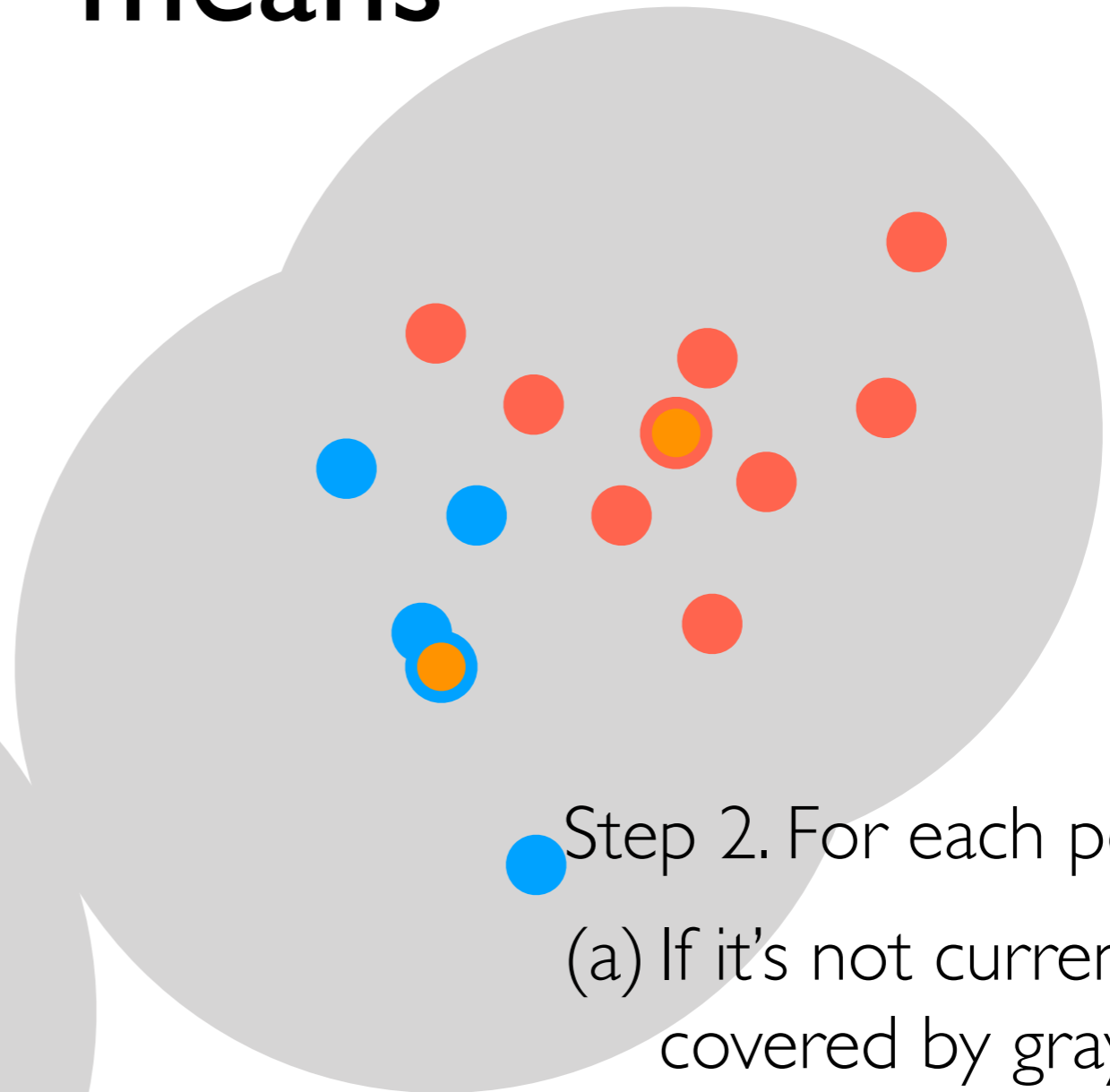
DP-means

Step 0: Pick concentration parameter $\lambda > 0$

Step 1: Start with everything in same cluster



Step 3. Recompute cluster centers



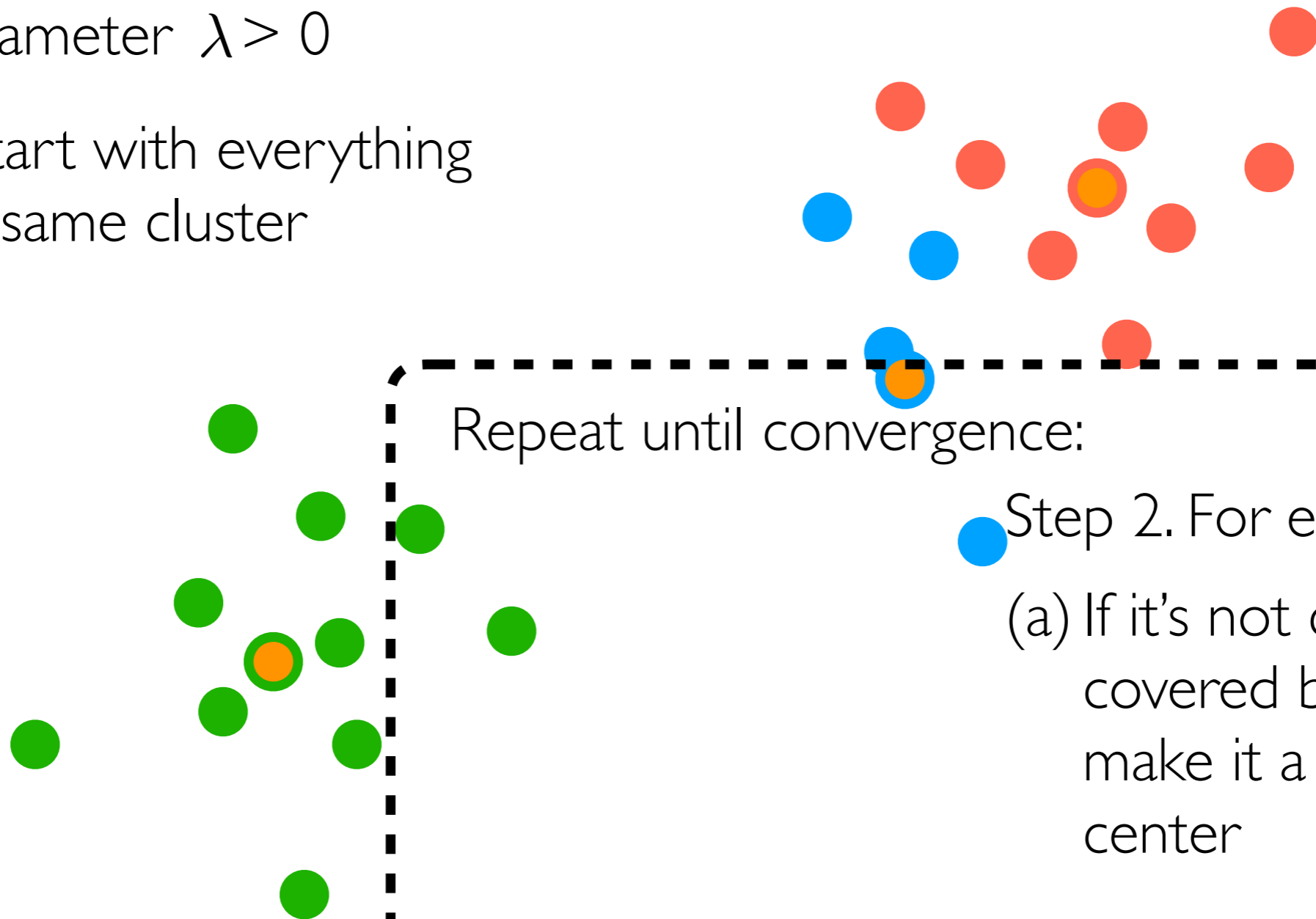
Step 2. For each point:
(a) If it's not currently covered by gray balls, make it a new cluster center

(b) Otherwise assign it to nearest cluster

DP-means

Step 0: Pick concentration parameter $\lambda > 0$

Step 1: Start with everything in same cluster



Repeat until convergence:

Step 2. For each point:

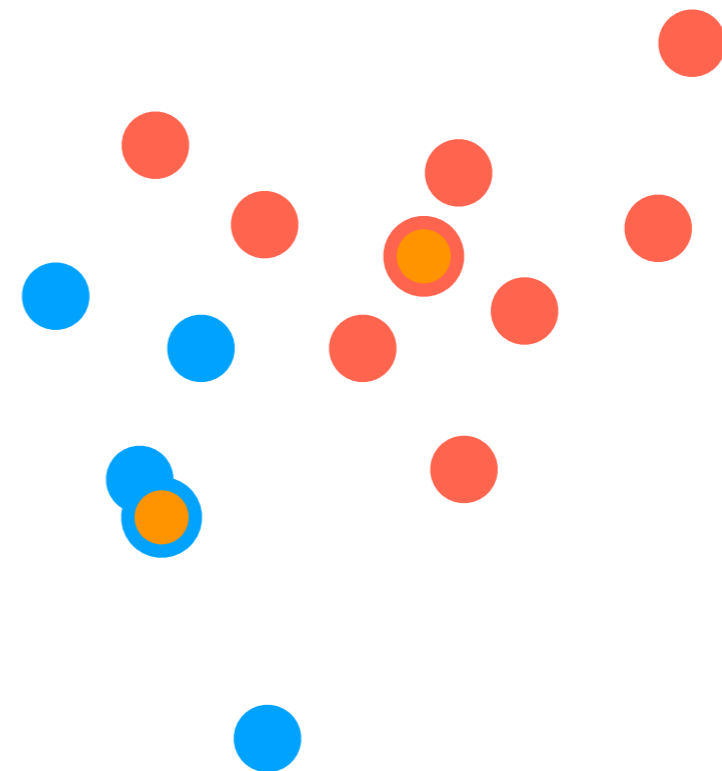
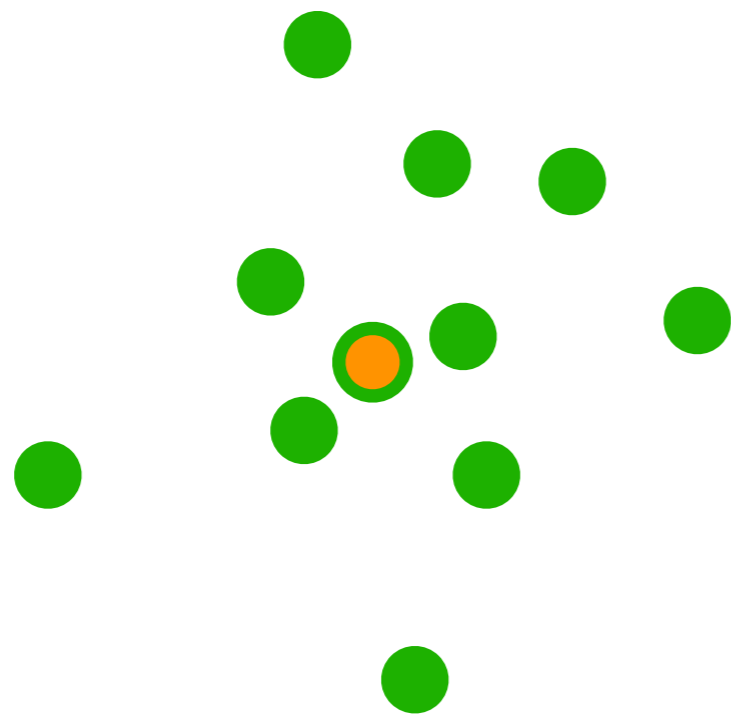
(a) If it's not currently covered by gray balls, make it a new cluster center

(b) Otherwise assign it to nearest cluster

Step 3. Recompute cluster centers

DP-means

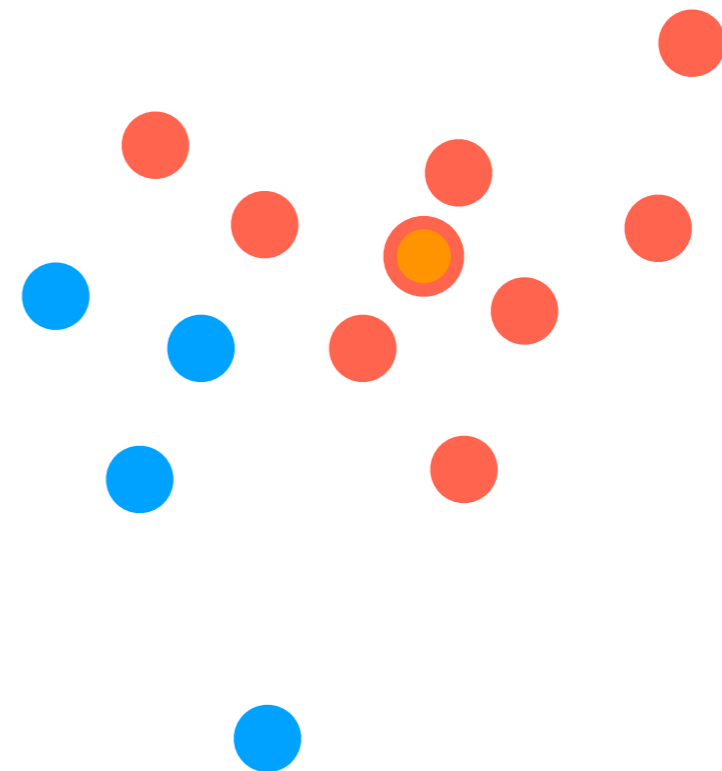
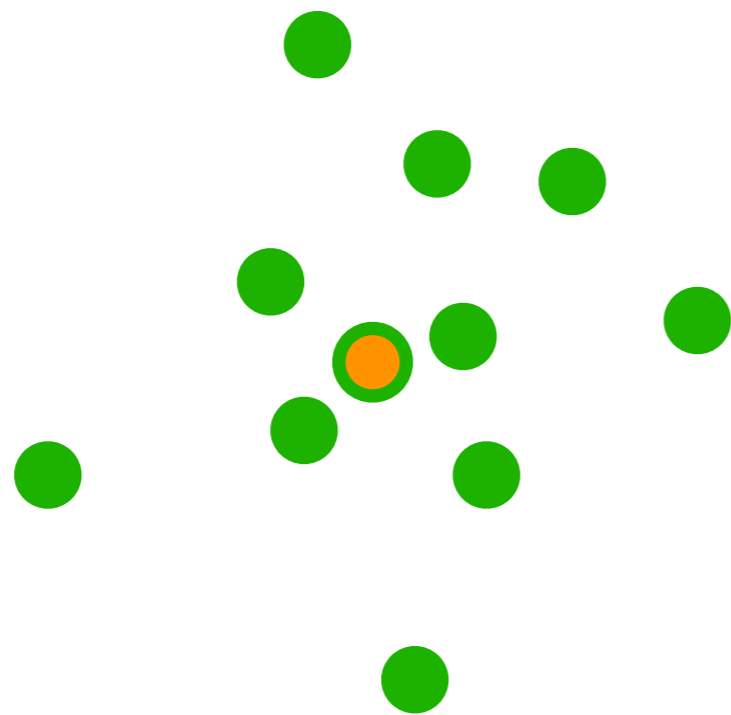
DP-means can produce a few extra small clusters



In practice: can reassign points in small clusters to bigger clusters

DP-means

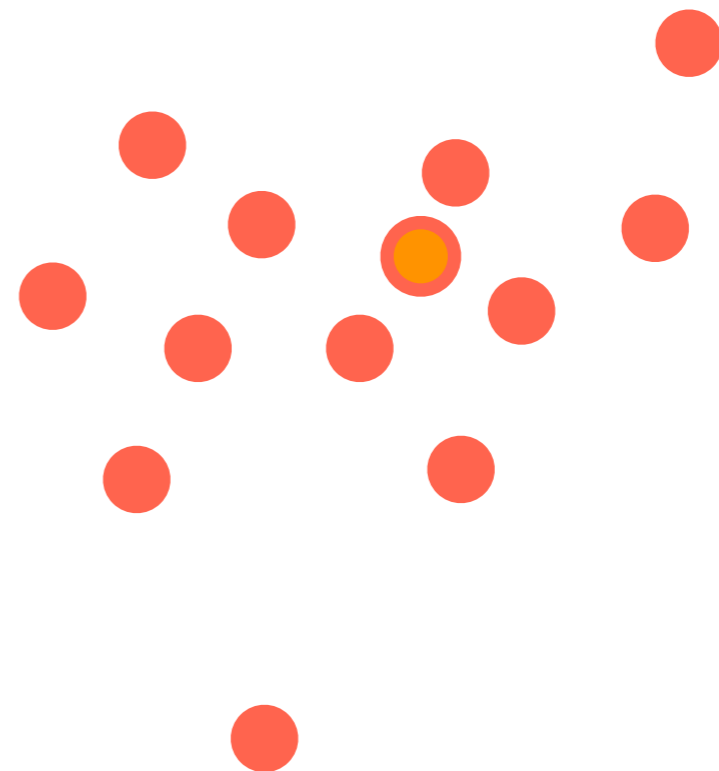
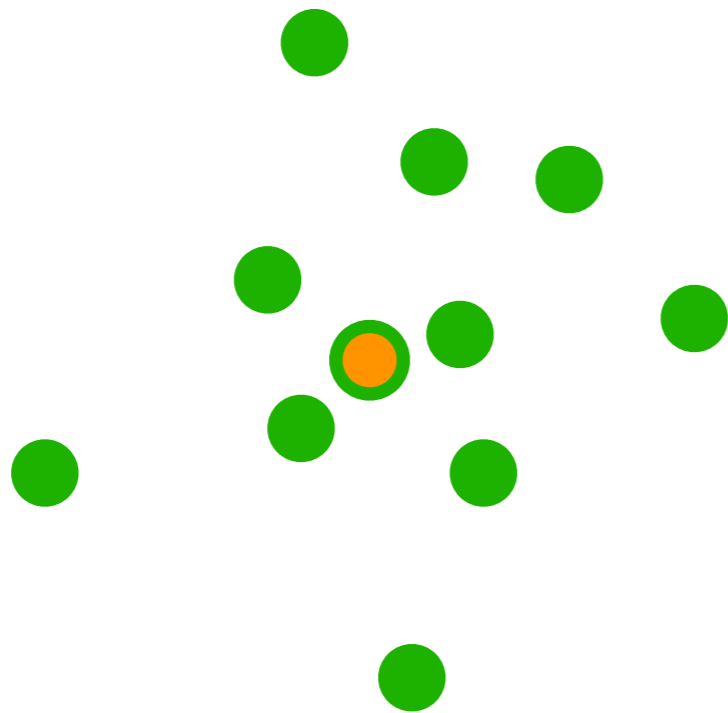
DP-means can produce a few extra small clusters



In practice: can reassign points in small clusters to bigger clusters

DP-means

DP-means can produce a few extra small clusters

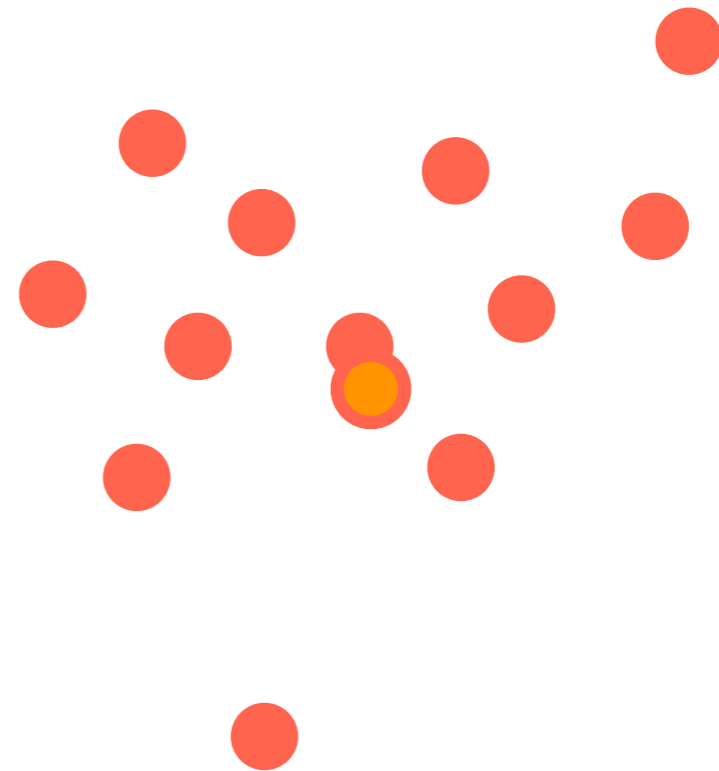
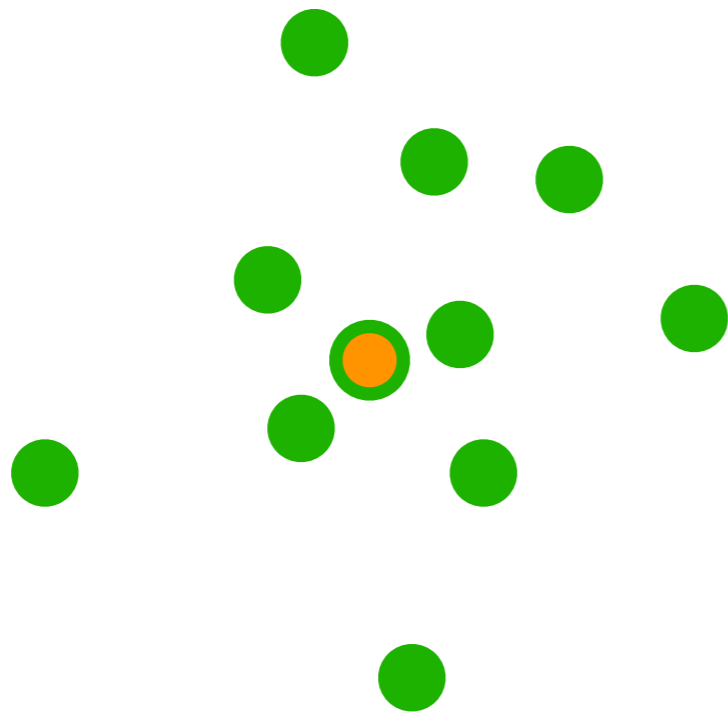


In practice: can reassign points in small clusters to bigger clusters

Can recompute cluster centers

DP-means

DP-means can produce a few extra small clusters



In practice: can reassign points in small clusters to bigger clusters

Can recompute cluster centers

DP-means

- Big picture: DP-means has a parameter controlling the size (radius) of clusters rather than number of clusters
- If your problem can more naturally be phrased as having cluster sizes that should not be too large, can use DP-means instead of k-means

Real example. *Satellite image analysis of rural India to find villages*

Each cluster is a village: don't know how many villages there are total but rough upper bound on radius of village can be specified

→ DP-means can provide a decent solution!

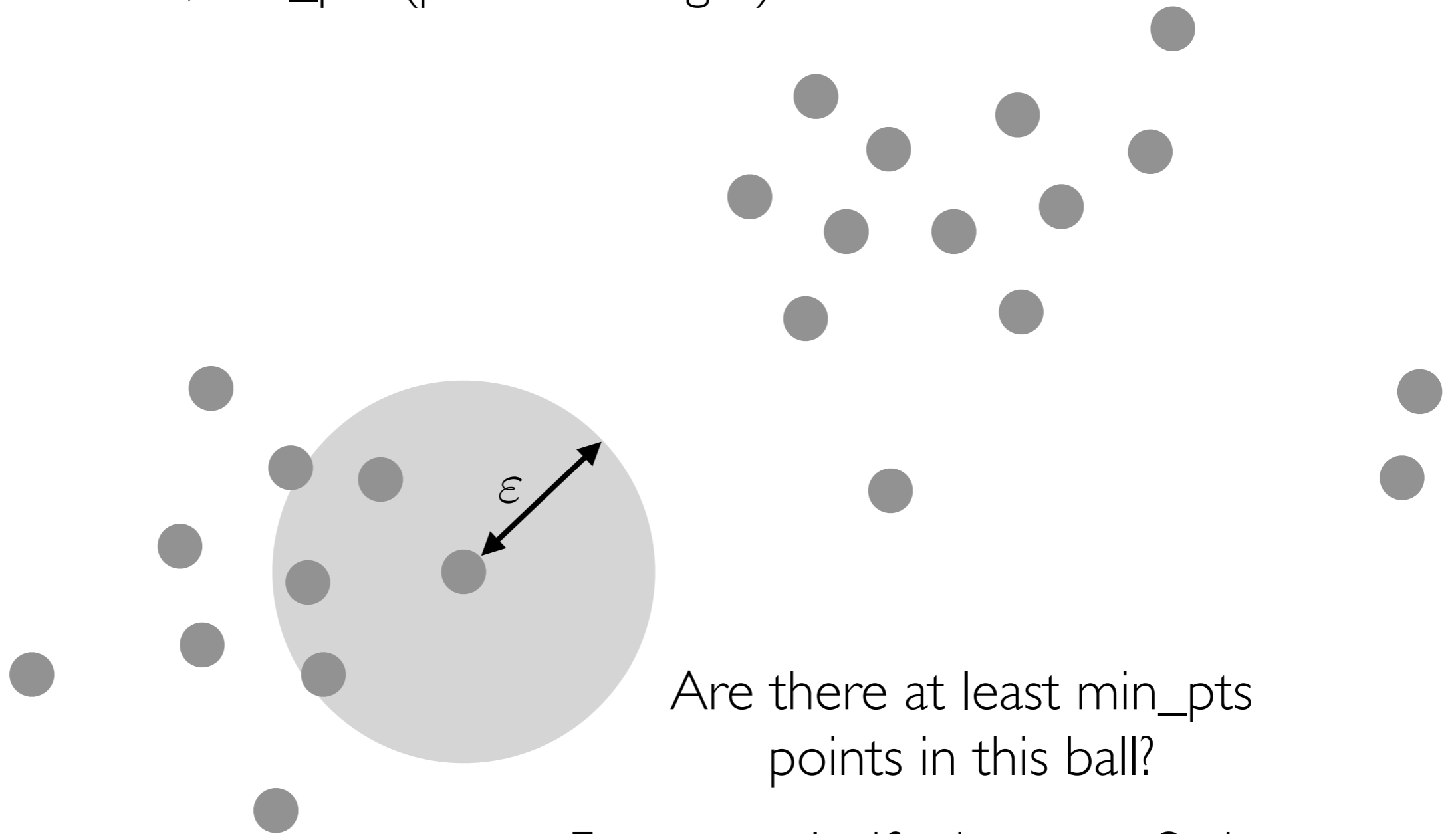
Technical remark: k-means approximates learning a GMM, DP-means approximates learning what's called a *Dirichlet-Process GMM*

Density-based Clustering

DBSCAN

Let's choose $\text{min_pts} = 3$

Pick radius $\varepsilon > 0$, min_pts (positive integer)



Are there at least min_pts points in this ball?

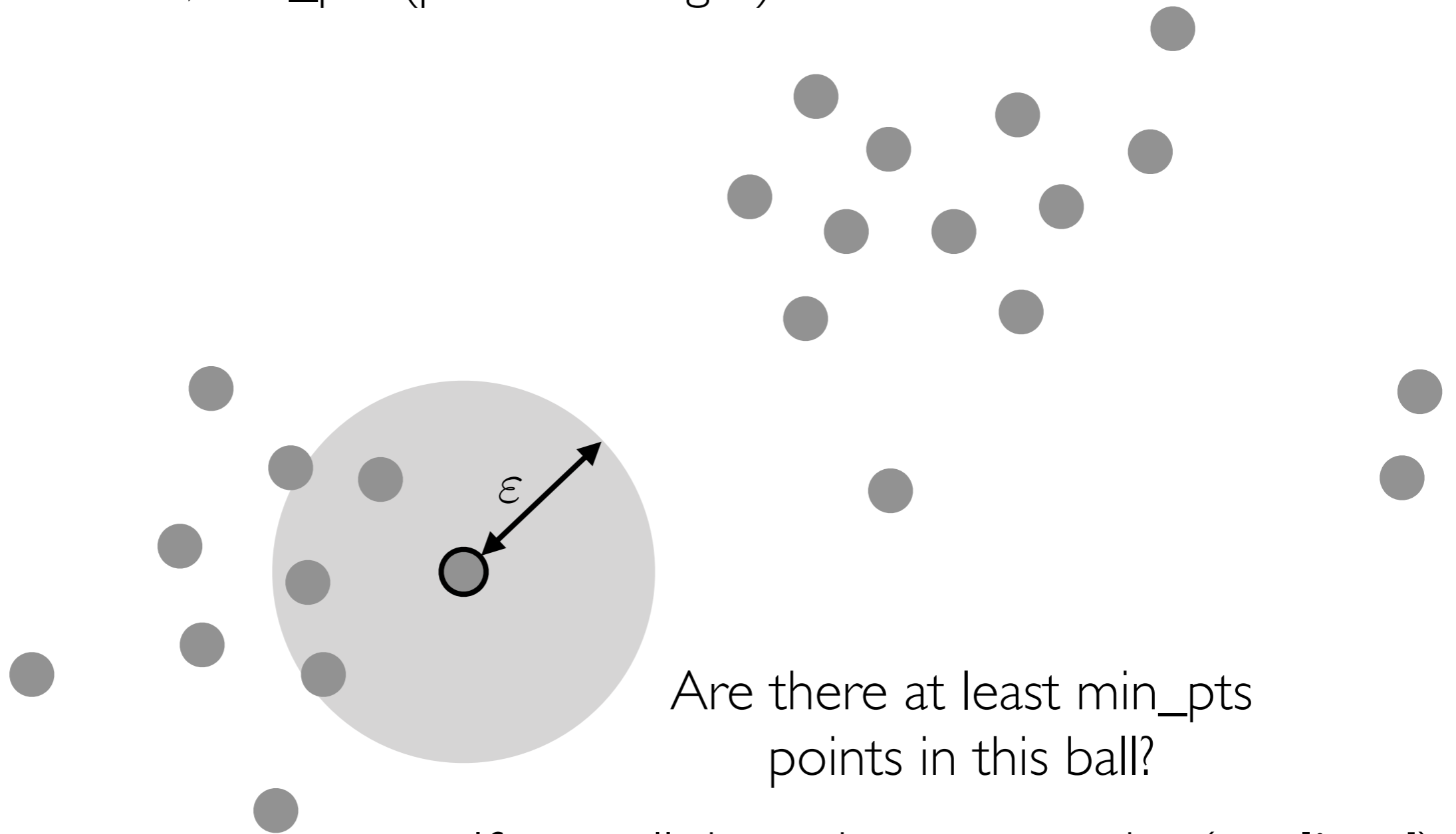
For example, if $\text{min_pts} = 3$, then yes

For example, if $\text{min_pts} = 10$, then no

DBSCAN

Let's choose $\text{min_pts} = 3$

Pick radius $\varepsilon > 0$, min_pts (positive integer)



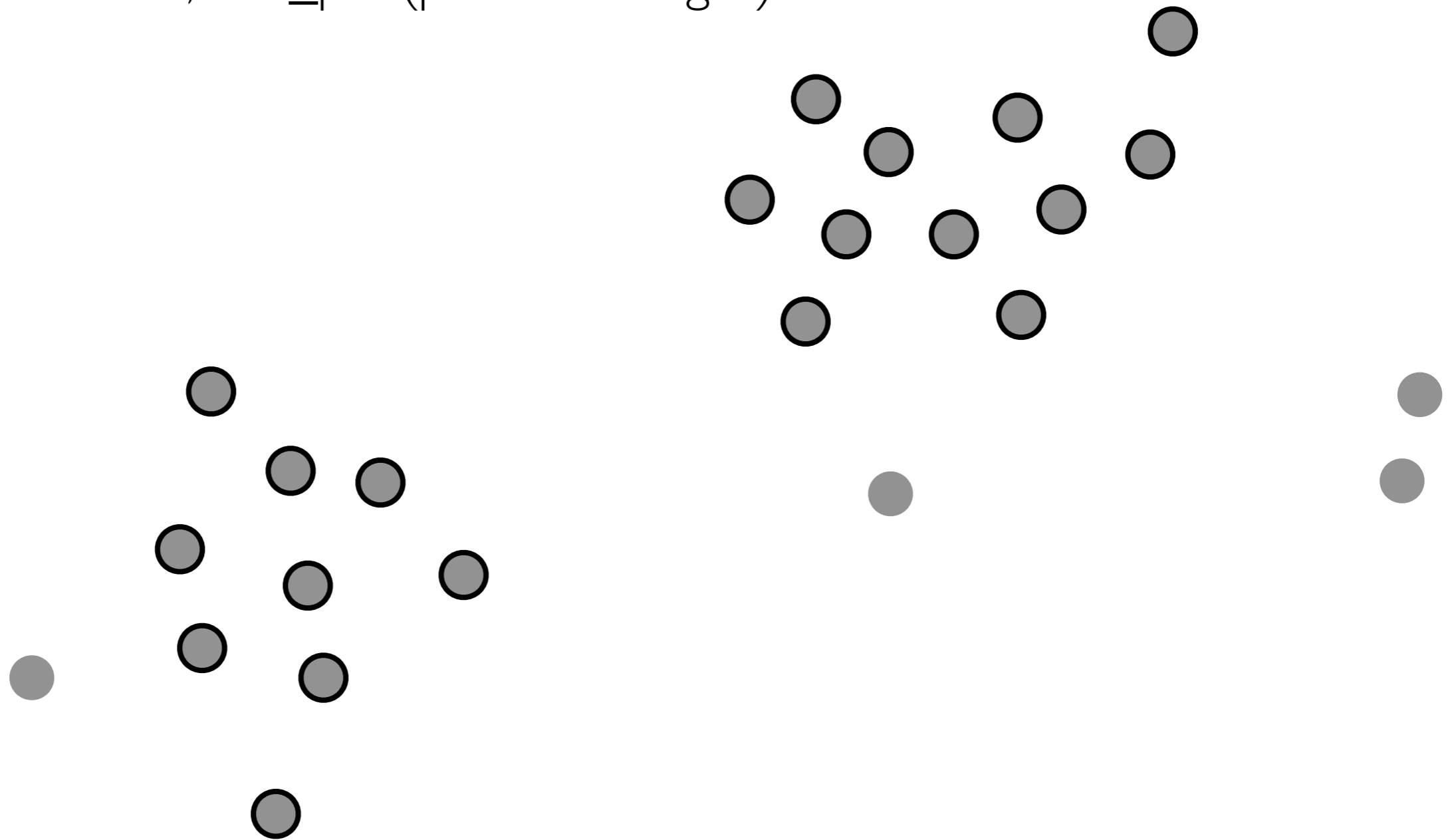
Are there at least min_pts points in this ball?

If yes: call the point a core point (**outlined**)

DBSCAN

Let's choose $\text{min_pts} = 3$

Pick radius $\varepsilon > 0$, min_pts (positive integer)

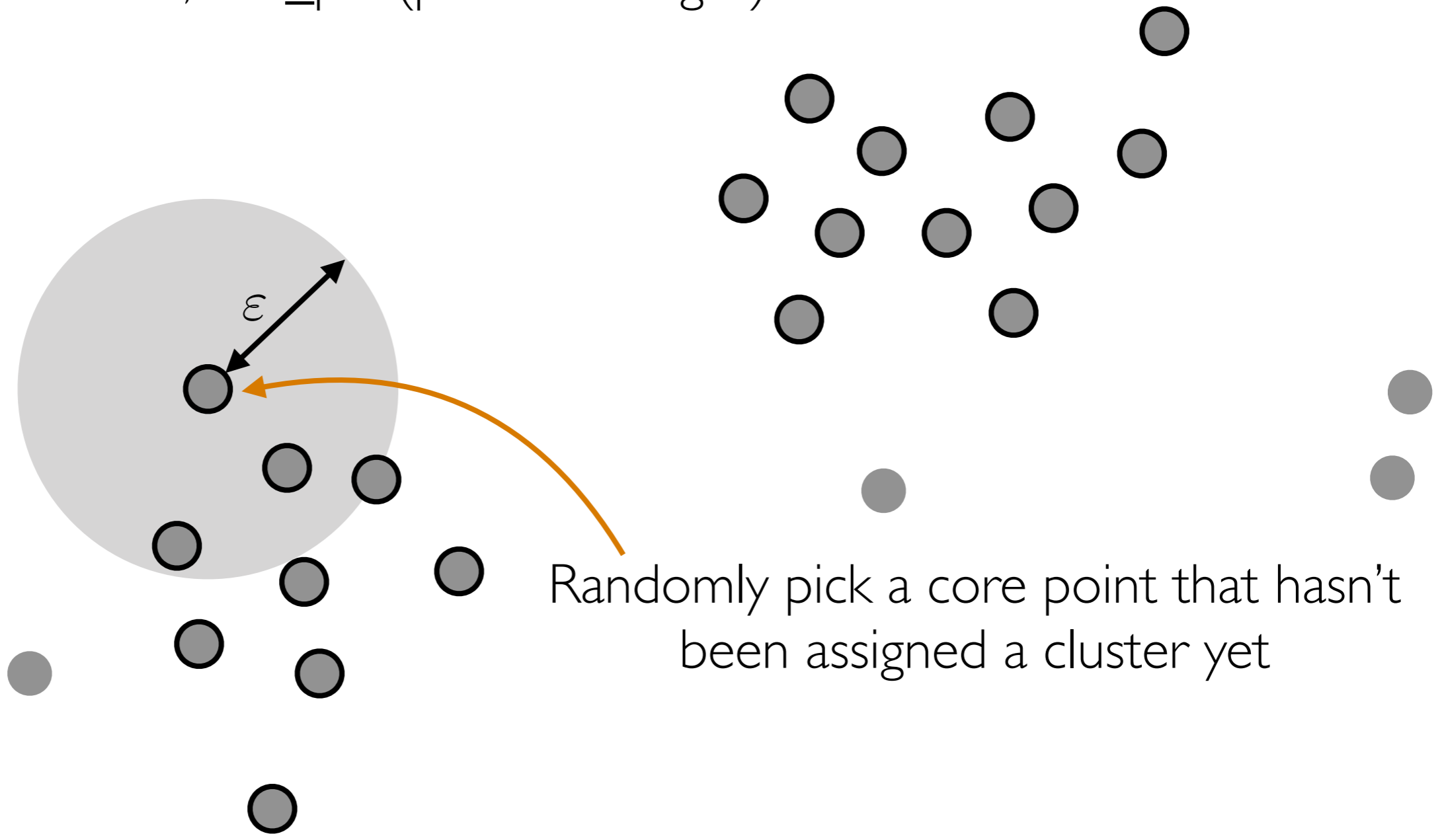


Core points (**outlined**)

DBSCAN

Let's choose $\text{min_pts} = 3$

Pick radius $\epsilon > 0$, min_pts (positive integer)



Randomly pick a core point that hasn't been assigned a cluster yet

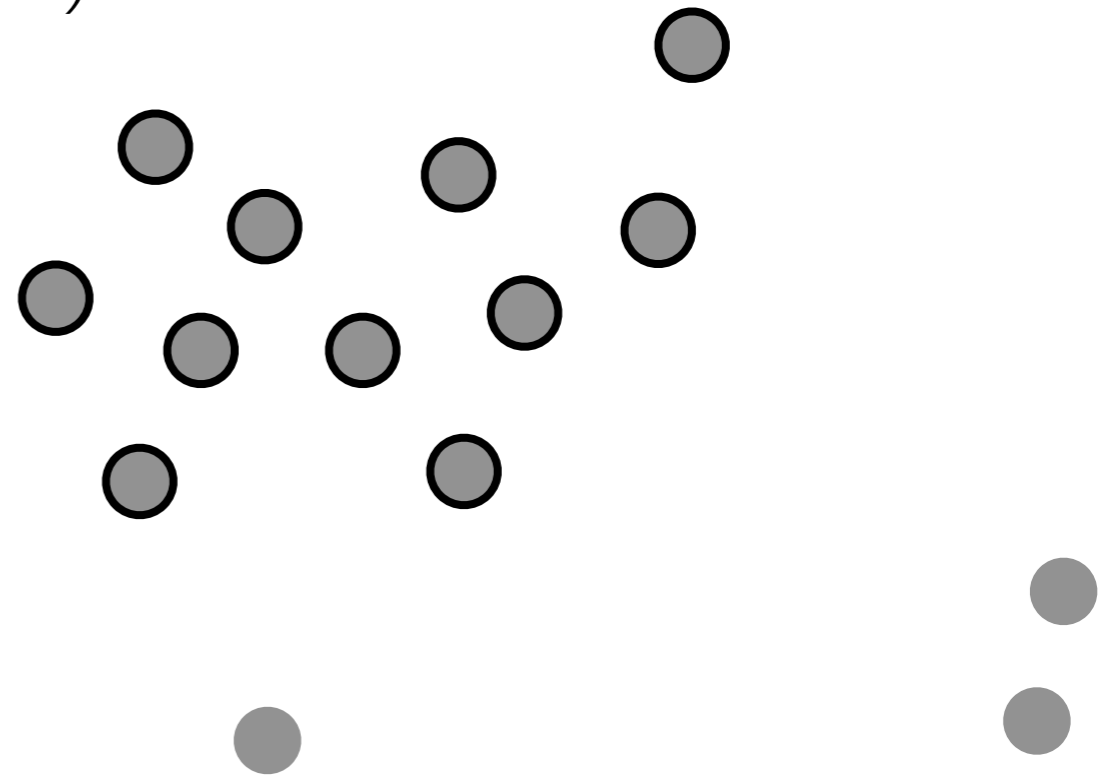
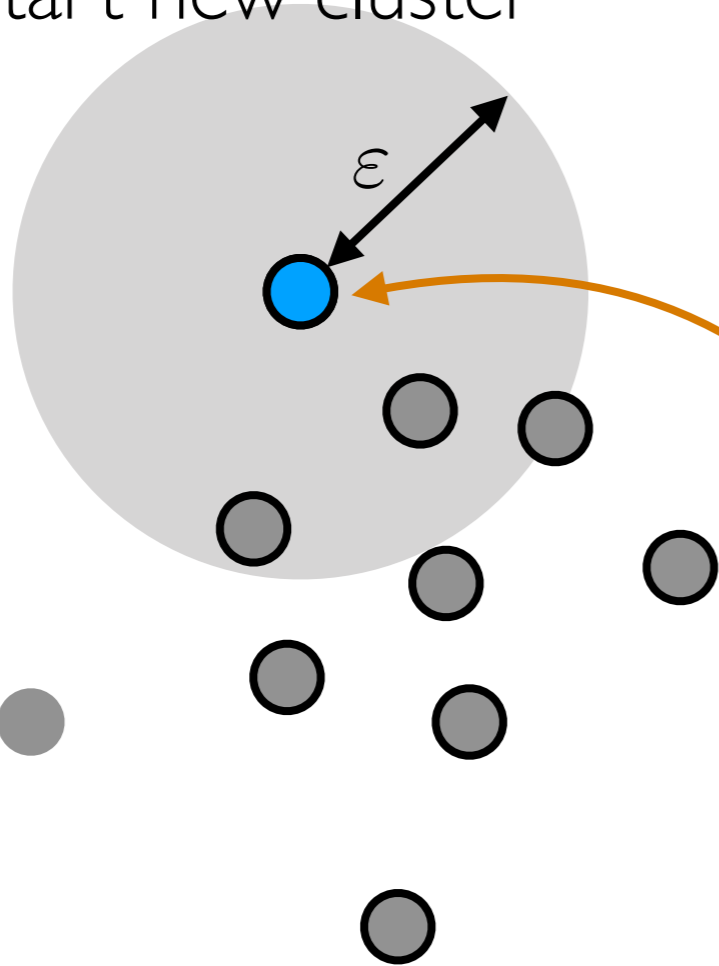
Core points (**outlined**)

DBSCAN

Let's choose $\text{min_pts} = 3$

Pick radius $\epsilon > 0$, min_pts (positive integer)

Start new cluster



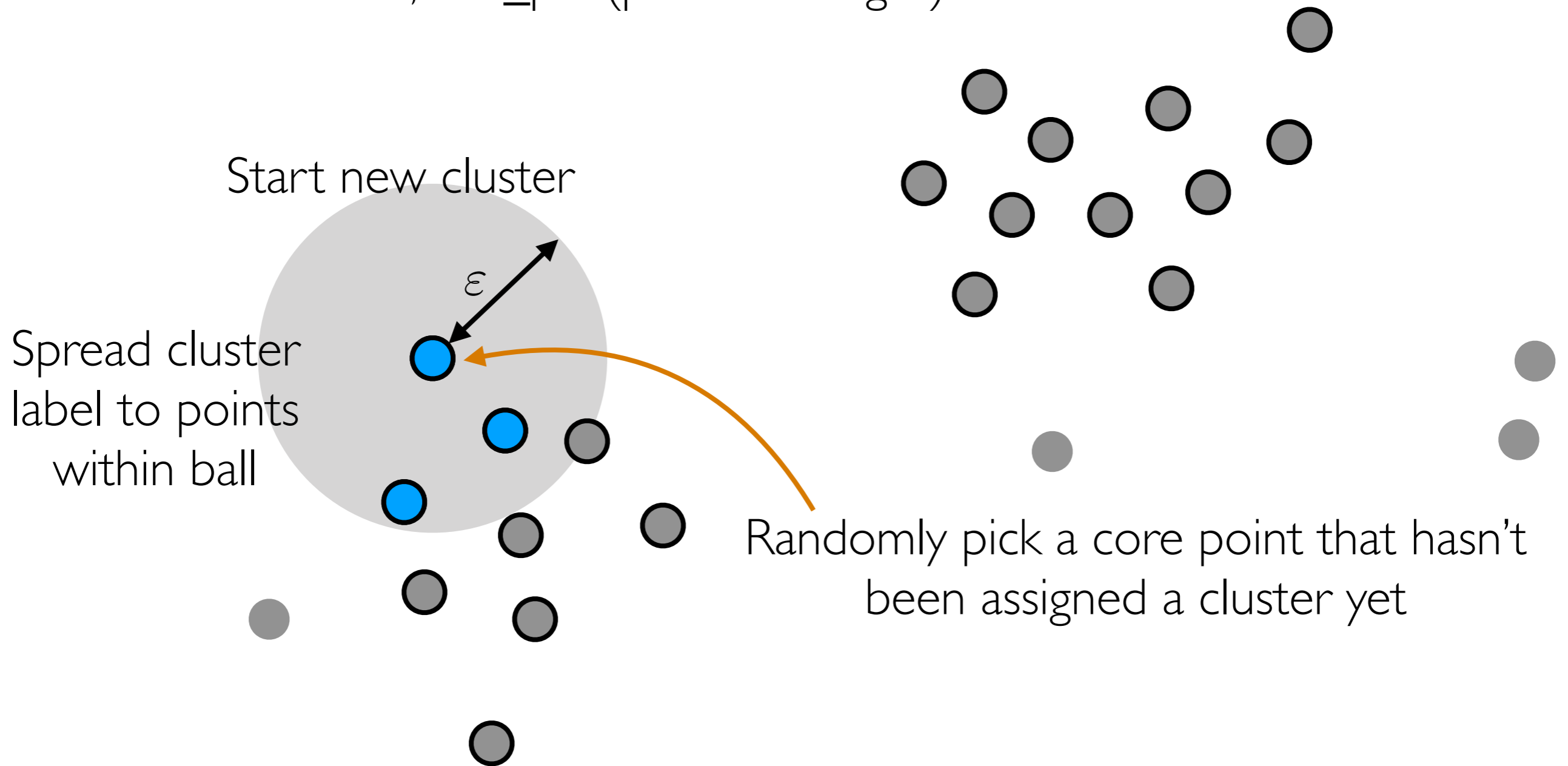
Randomly pick a core point that hasn't been assigned a cluster yet

Core points (**outlined**)

DBSCAN

Let's choose $\text{min_pts} = 3$

Pick radius $\epsilon > 0$, min_pts (positive integer)

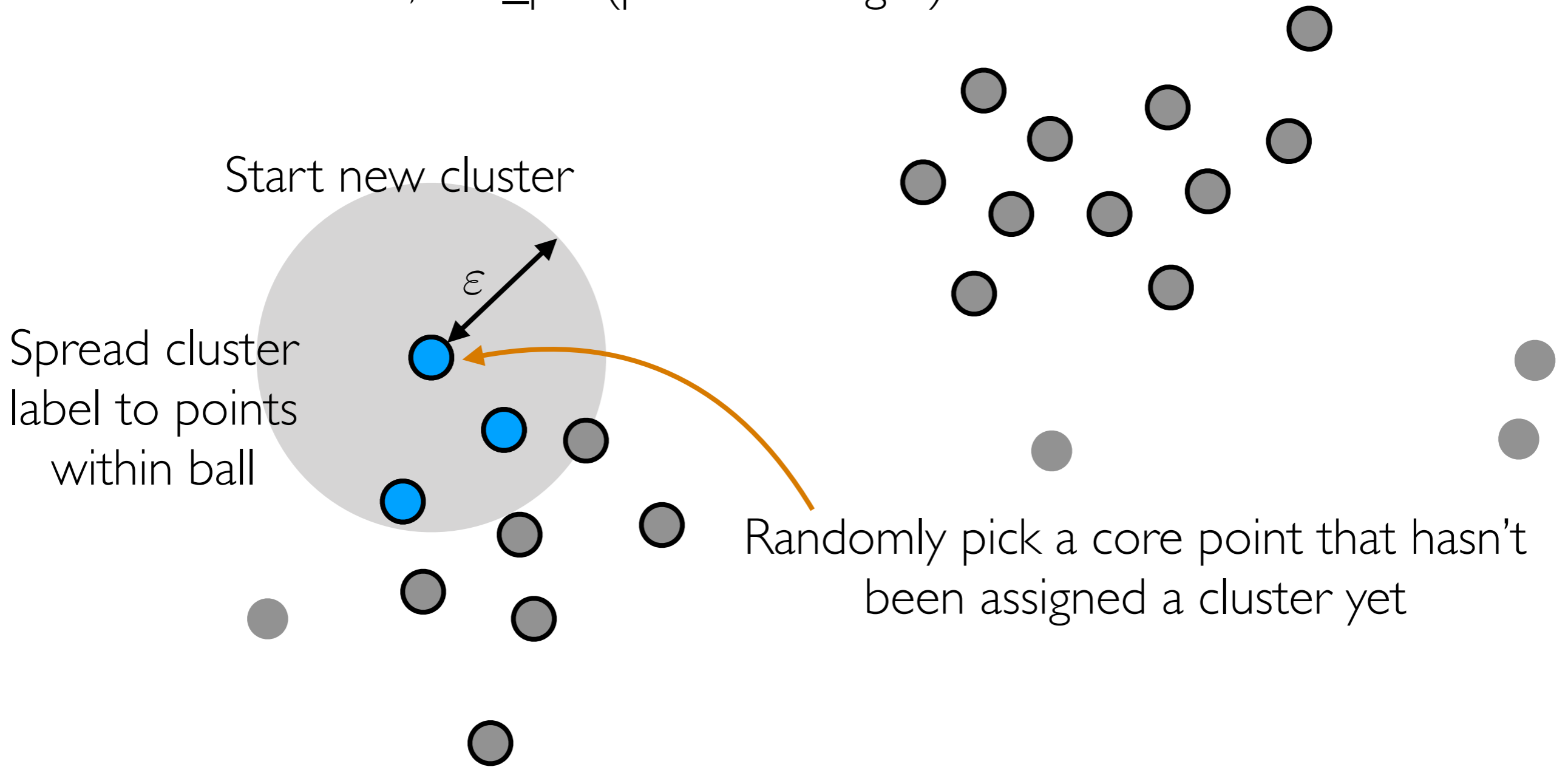


Core points (**outlined**)

DBSCAN

Let's choose $\text{min_pts} = 3$

Pick radius $\epsilon > 0$, min_pts (positive integer)



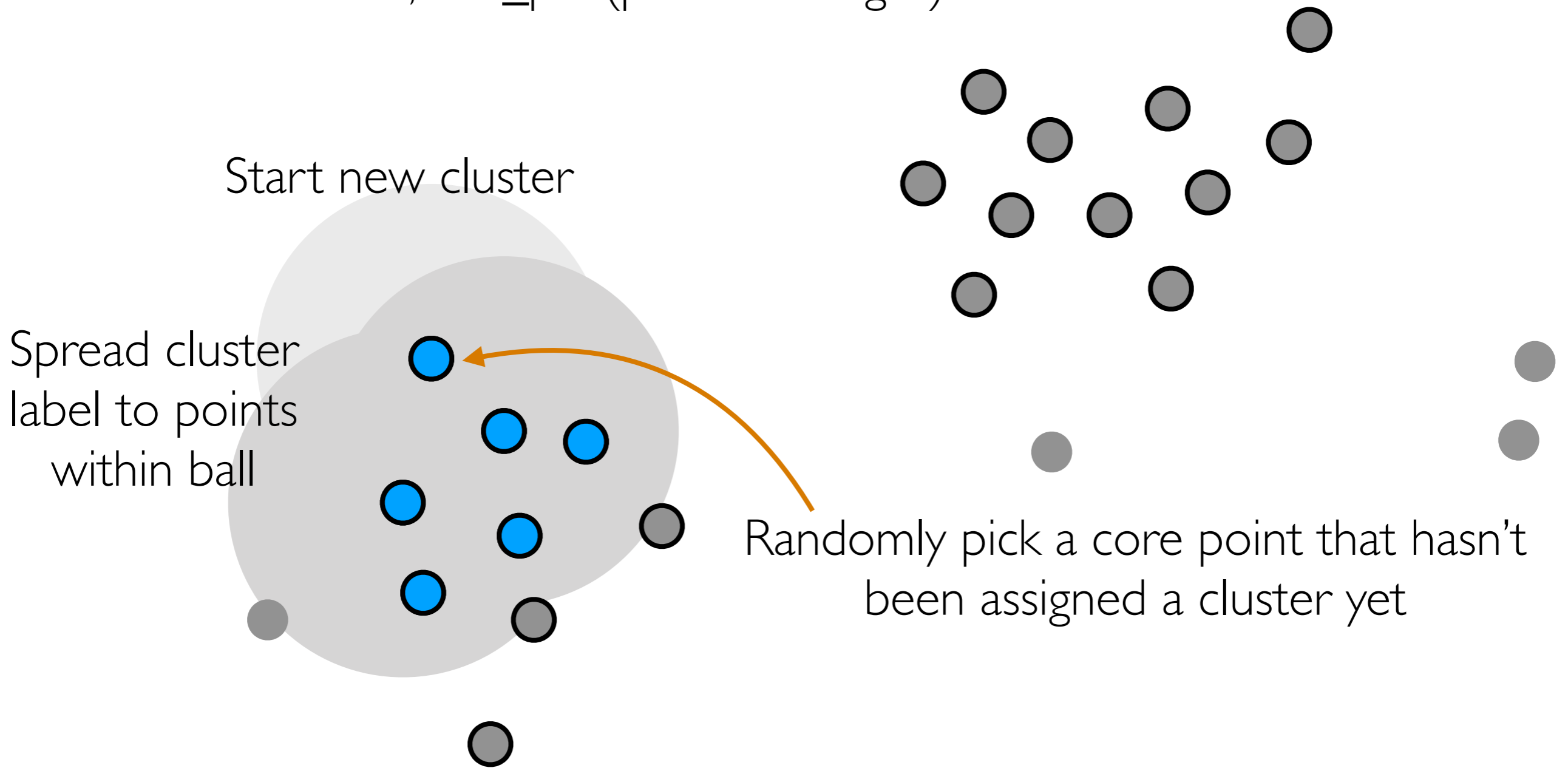
For core points within the ball:
spread the cluster label

Core points (**outlined**)

DBSCAN

Let's choose $\text{min_pts} = 3$

Pick radius $\epsilon > 0$, min_pts (positive integer)



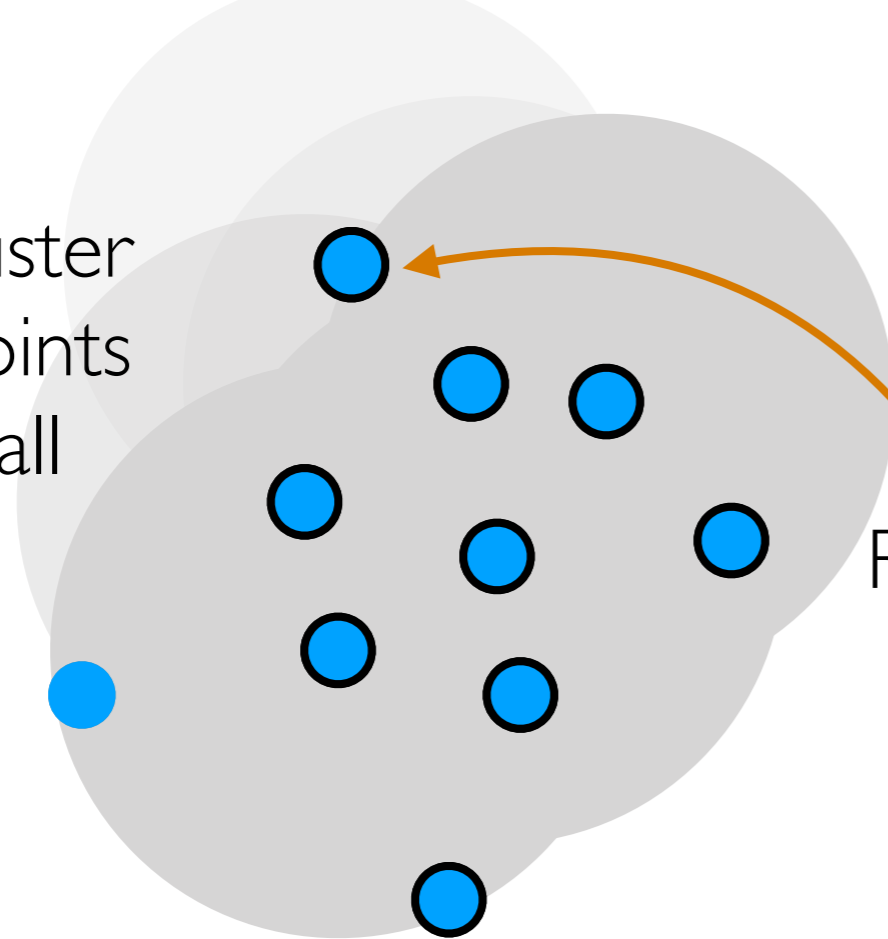
DBSCAN

Let's choose $\text{min_pts} = 3$

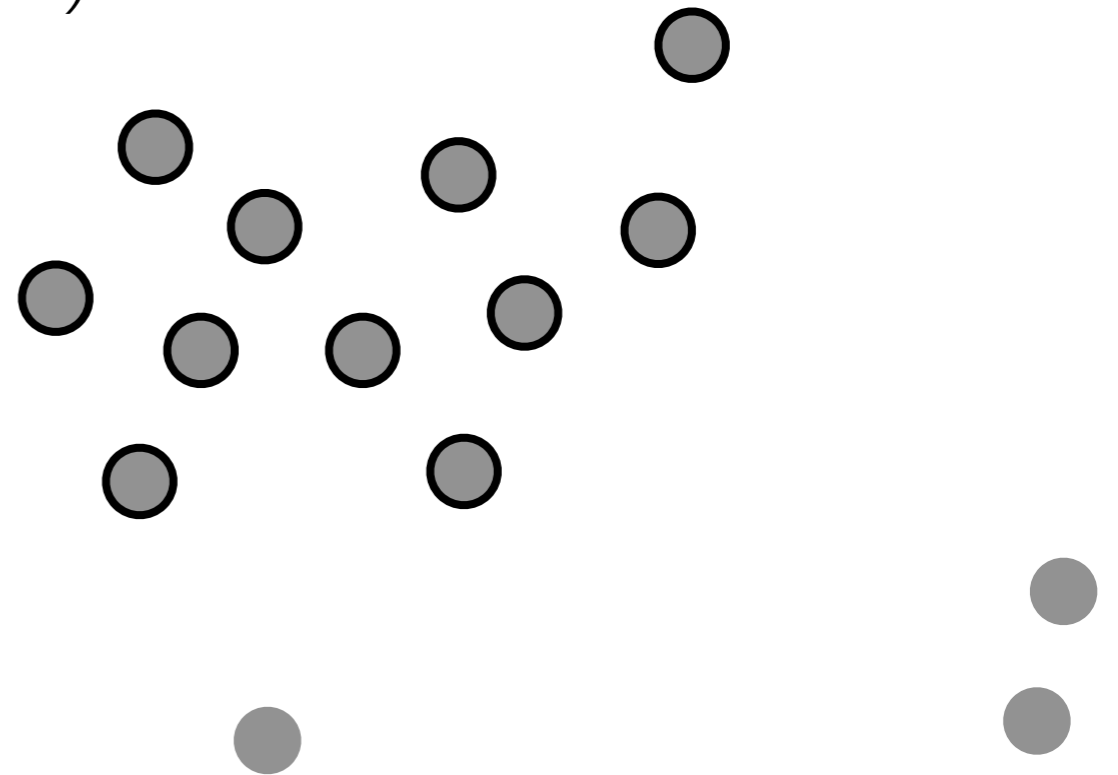
Pick radius $\epsilon > 0$, min_pts (positive integer)

Start new cluster

Spread cluster label to points within ball



Randomly pick a core point that hasn't been assigned a cluster yet



For core points within the ball:
spread the cluster label

Core points (**outlined**)

DBSCAN

Let's choose $\text{min_pts} = 3$

Pick radius $\epsilon > 0$, min_pts (positive integer)

Start new cluster

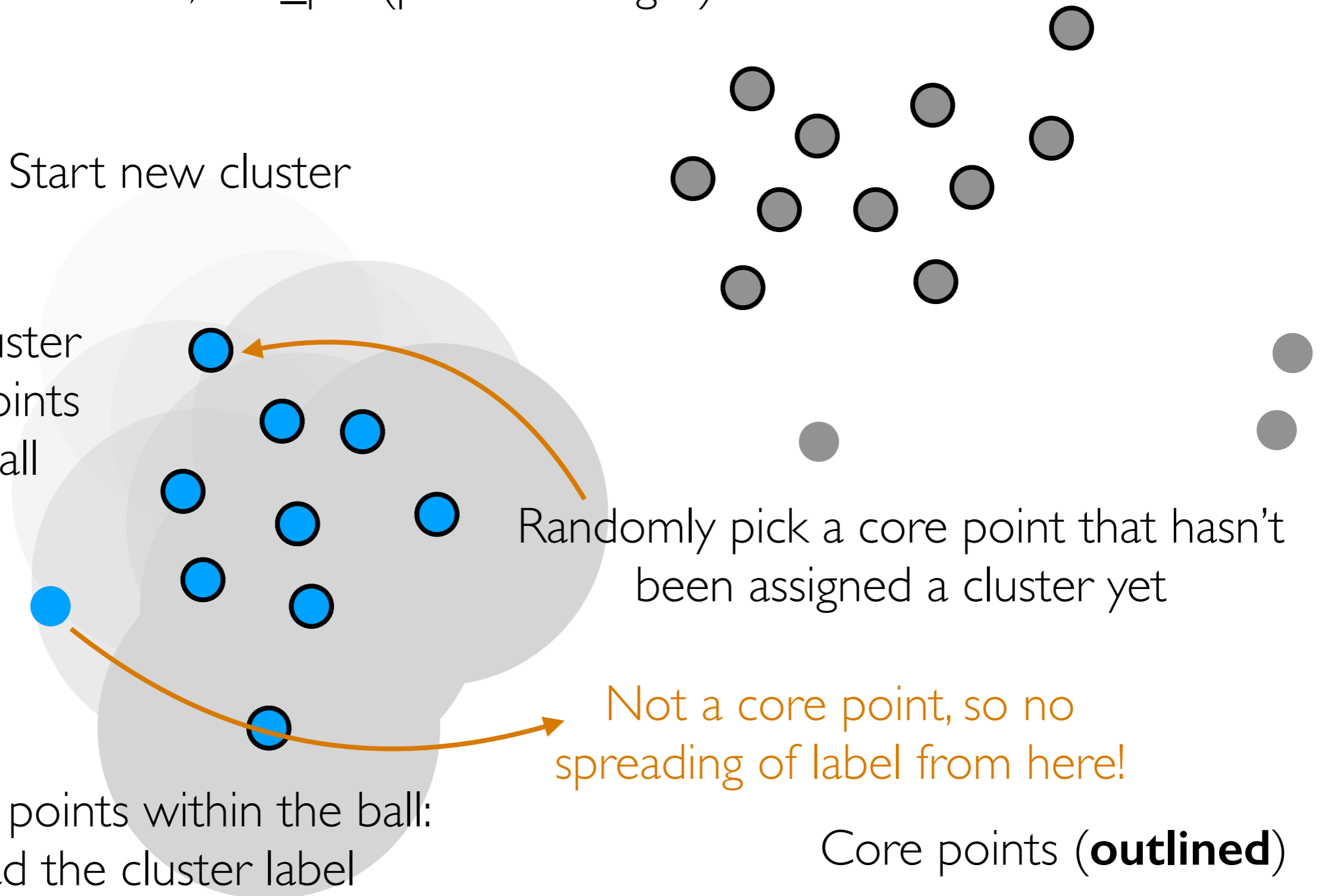
Spread cluster label to points within ball

Randomly pick a core point that hasn't been assigned a cluster yet

Not a core point, so no spreading of label from here!

For core points within the ball: spread the cluster label

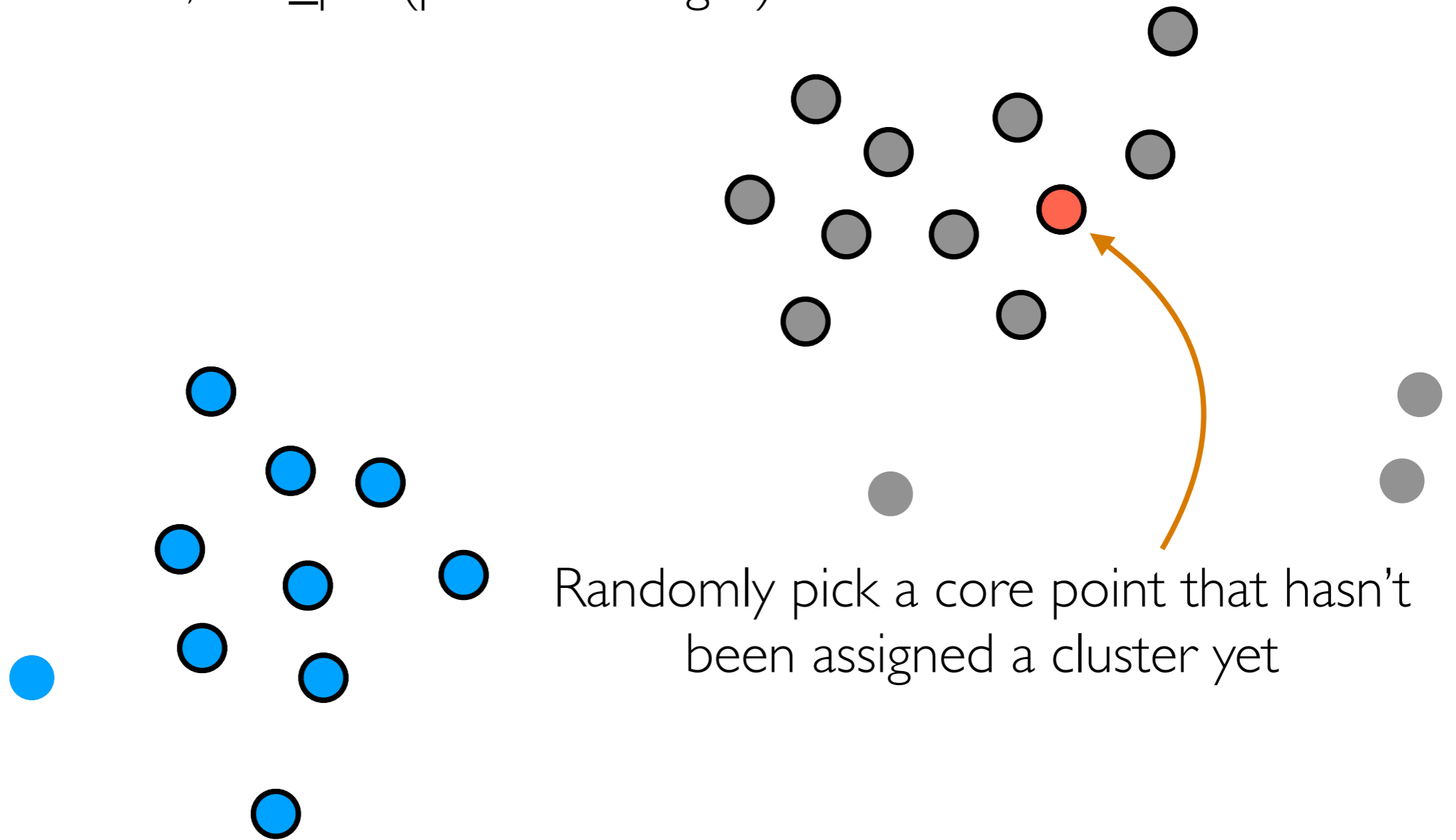
Core points (**outlined**)



DBSCAN

Let's choose $\text{min_pts} = 3$

Pick radius $\epsilon > 0$, min_pts (positive integer)

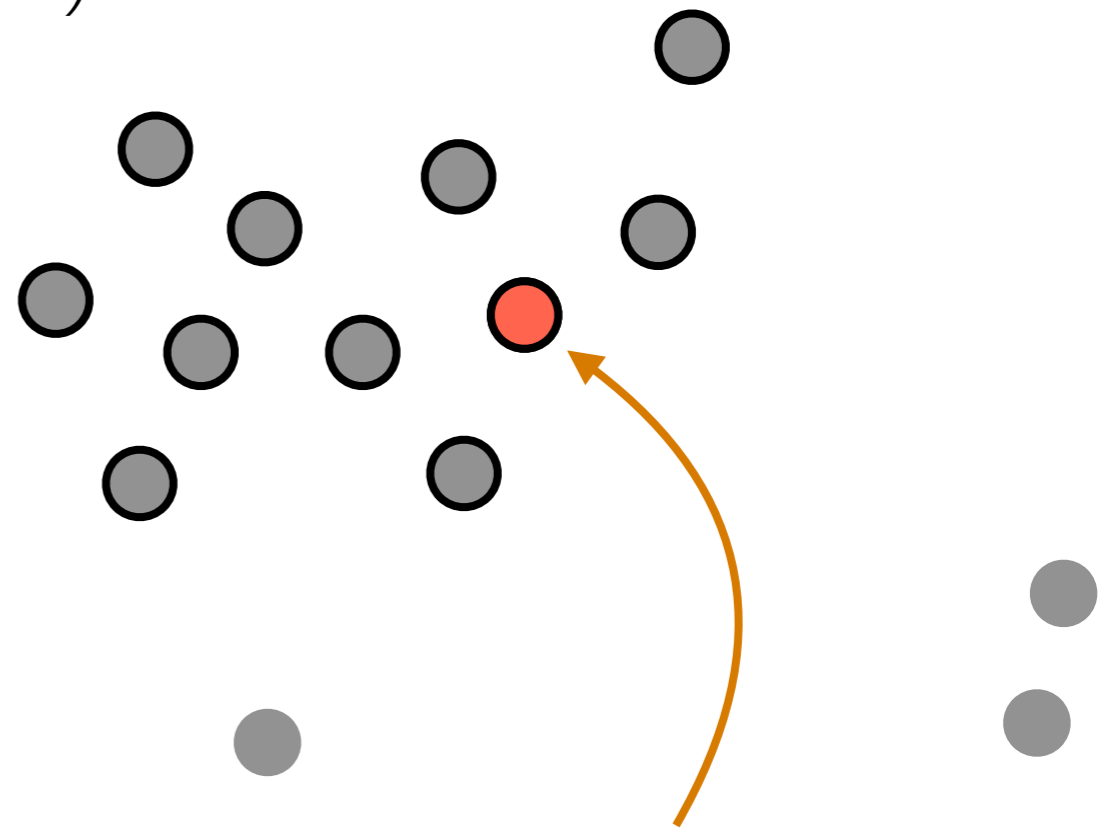
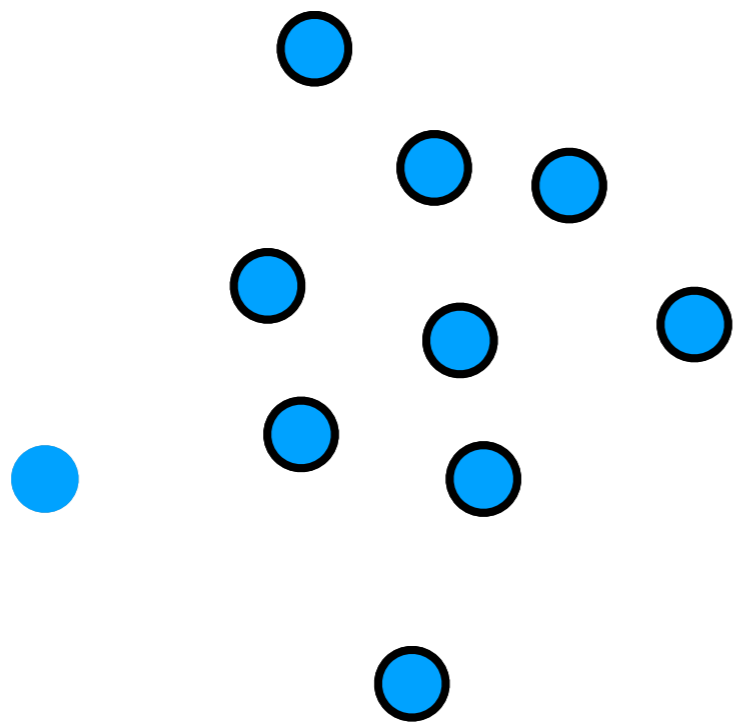


DBSCAN

Let's choose $\text{min_pts} = 3$

Pick radius $\epsilon > 0$, min_pts (positive integer)

Repeat "virus-spreading" like cluster label spreading; again, no spreading starting from non-core points



Randomly pick a core point that hasn't been assigned a cluster yet

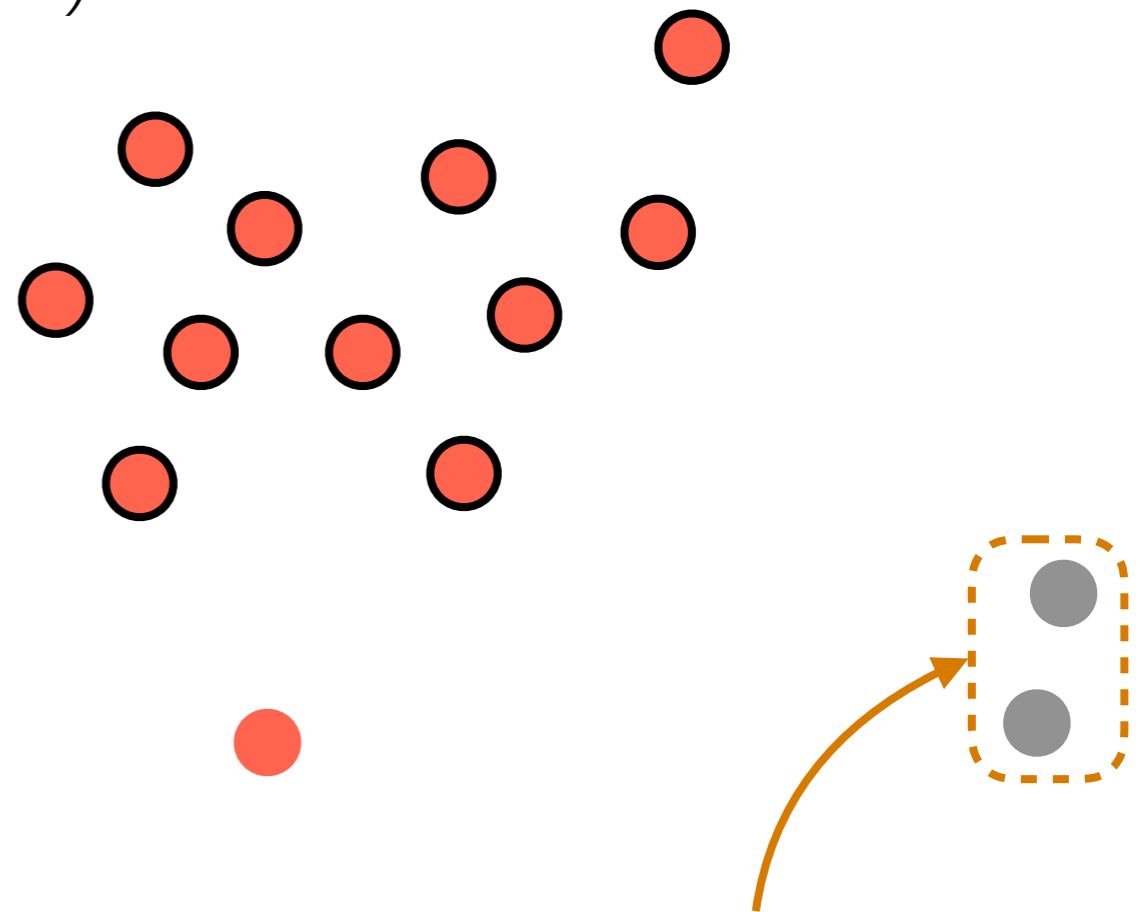
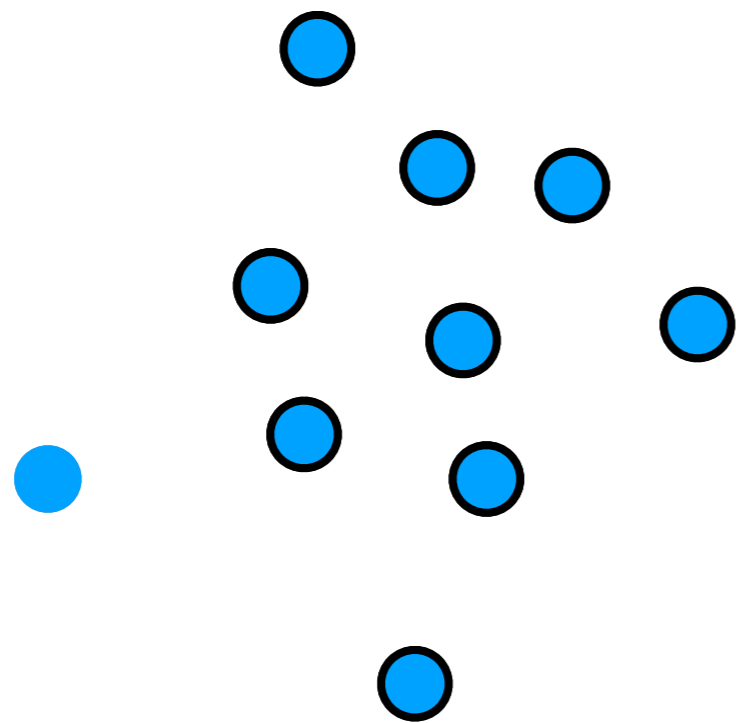
Core points (**outlined**)

DBSCAN

Let's choose $\text{min_pts} = 3$

Pick radius $\epsilon > 0$, min_pts (positive integer)

Repeat “virus-spreading” like cluster label spreading; again, no spreading starting from non-core points



Some points might actually *not* get clustered by DBSCAN and are declared as **outliers!**

Core points (**outlined**)

Some Last Remarks on Clustering

Demo for DP-GMM & DBSCAN are at the end of prev. lecture's demo

What about clustering unstructured data?

- Covered this Friday April 8 in CMU Adelaide's 95-865 recitation
- I'll get the recording for you to watch
(since CMU Pittsburgh has Spring Carnival, CMU Adelaide doesn't)

Important takeaway: ultimately, *you* have to decide on which clustering method and number of clusters make sense for your data

- After you run a clustering algorithm, make visualizations to interpret the clusters *in the context of your application!*
- Do not just blindly rely on numerical metrics (e.g., CH index)

If you can set up a prediction task, then you can use the prediction task to guide the clustering

Is clustering structure enough?

(Flashback) GMM with k Clusters

Cluster 1

Probability of generating a point from cluster 1 = π_1

Gaussian mean = μ_1

Gaussian covariance = Σ_1

...

Cluster k

Probability of generating a point from cluster k = π_k

Gaussian mean = μ_k

Gaussian covariance = Σ_k

How to generate points from this GMM:

1. Flip biased k -sided coin (the sides have probabilities π_1, \dots, π_k)

2. Let \mathbf{Z} be the side that we got (it is some value $1, \dots, k$)

3. Sample 1 point from the Gaussian for cluster \mathbf{Z}

Each data point has a single true cluster assignment \mathbf{Z}
& is generated from the Gaussian for cluster \mathbf{Z}

In reality, a data point could have “mixed” membership and belong to multiple clusters

We show a way to model this for text data

A text document can be about multiple **topics**

Each topic consists of a bunch of regularly co-occurring words
(example topics: “sports”, “medicine”, “movies”, “finance”)

Latent Dirichlet Allocation (LDA)

- A generative model
- Input: “document-word” matrix, and pre-specified # topics k

		Word			
		1	2	...	d
Document	1	Each row is a feature vector representing a raw counts histogram!			
	2				
	⋮				
	n				

i -th row, j -th column: # times word j appears in doc i

- Output: what the k topics are (details on this shortly)

LDA Generative Model Example

		Topic	
		weather	food
Document	Alice's text	0.1	0.9
	Bob's text	0.5	0.5

		Word			
		cold	hot	apple	pie
Topic	weather	0.3	0.7	0.0	0.0
	food	0.1	0.3	0.5	0.1

Each word in Alice's text is generated by:

1. Flip 2-sided coin for Alice
2. If weather: flip 4-sided coin for weather
If food: flip 4-sided coin for food

LDA Generative Model Example

		Topic	
		weather	food
Document	Alice's text	0.1	0.9
	Bob's text	0.5	0.5

		Word			
		cold	hot	apple	pie
Topic	weather	0.3	0.7	0.0	0.0
	food	0.1	0.3	0.5	0.1

Each word in Bob's text is generated by:

1. Flip 2-sided coin for Bob
2. If weather: flip 4-sided coin for weather
If food: flip 4-sided coin for food

LDA Generative Model Example

		Topic	
		weather	food
Document	Alice's text	0.1	0.9
	Bob's text	0.5	0.5

		Word			
		cold	hot	apple	pie
Topic	weather	0.3	0.7	0.0	0.0
	food	0.1	0.3	0.5	0.1

Each word in doc i is generated by:

1. Flip 2-sided coin for doc i
2. If weather: flip 4-sided coin for weather
If food: flip 4-sided coin for food

“Learning the topics” means figuring out these 4-sided coin probabilities

LDA Generative Model



LDA models each word in document i to be generated as:

- Randomly choose a topic \mathbf{Z} (use topic distribution for doc i)
- Randomly choose a word (use word distribution for topic \mathbf{Z})

LDA

- A generative model
- Input: “document-word” matrix, and pre-specified # topics k

		Word			
		1	2	...	d
Document	1	Each row is a feature vector representing a raw counts histogram!			
	2				
	\vdots				
	n				

i -th row, j -th column: # times word j appears in doc i

- Output: the k topics' distribution of words

LDA

Demo

An Alternative Feature Vector Representation for Text: TF-IDF

Intuition: words that appear in more documents are likely less useful (same intuition as stop words!) — let's *downweight* these words!



i -th row, j -th column: # times word j appears in doc i

multiply TF by $\log \frac{1}{\mathbb{P}(\text{document has word } j)} = \log \frac{1}{\left(\frac{\# \text{ documents that have word } j}{n}\right)}$

$= \log \frac{n}{\# \text{ documents that have word } j}$


Hack: add 1 to numerator & 1 to denominator

An Alternative Feature Vector Representation for Text: TF-IDF

Intuition: words that appear in more documents are likely less useful (same intuition as stop words!) — let's *downweight* these words!

There are many TF-IDF variants! (Lots of hacks!)

Document

	Word			
	1	2	...	d
1				
2				
⋮				
n				

Term frequency (TF)

i -th row, j -th column: # times word j appears in doc i

$$\times \left[1 + \log \frac{1 + n}{1 + \# \text{ documents that have word } j} \right]$$

Default TF-IDF weighting in sklearn

sklearn's default behavior further normalizes each row to have Euclidean norm 1

TF-IDF is in your HW2
(usage is similar to CountVectorizer from the demo)