

Clustering

The k -Means Algorithm

Cary Goltermann

Galvanize

2017

Supervised vs. Unsupervised Learning

Clustering

- Intuition
- Definition

k-Means Algorithm

- Pseudocode
- Centroid Initialization
- Stopping Criteria
- Step-through
- Evaluation
- Problems
- Choosing k

Supervised

- Have a target / label that we model.
- Models look like functions that take in data and create prediction.
- Have an error metric that we can use to compare models.

Supervised

- Have a target / label that we model.
- Models look like functions that take in data and create prediction.
- Have an error metric that we can use to compare models.

Unsupervised

- No labels → no target!
- No stark error metric to compare models with.
- It's easy to be wrong, but it's hard to prove you're right.
- Trying to uncover/
discover hidden structure in our data.

Supervised vs. Unsupervised Learning

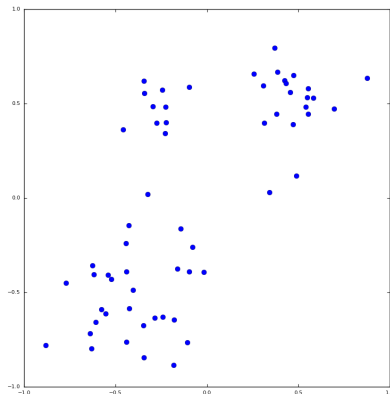
Clustering

- Intuition
- Definition

k-Means Algorithm

- Pseudocode
- Centroid Initialization
- Stopping Criteria
- Step-through
- Evaluation
- Problems
- Choosing *k*

What Is a Cluster?



- How many clusters do you see?
- What makes something a cluster?
- What makes something not a cluster?

Supervised vs. Unsupervised Learning

Clustering

- Intuition
- Definition

k-Means Algorithm

- Pseudocode
- Centroid Initialization
- Stopping Criteria
- Step-through
- Evaluation
- Problems
- Choosing *k*

Defining “Cluster”

- A partition of the dataset - not necessarily crisp.
- A strong internal similarity - small intra/within cluster distance.
- A strong external dissimilarity - large extra cluster distance.

Supervised vs. Unsupervised Learning

Clustering

- Intuition
- Definition

k-Means Algorithm

- Pseudocode
- Centroid Initialization
- Stopping Criteria
- Step-through
- Evaluation
- Problems
- Choosing k

The algorithm in all it's glory:

- ① Initialize centroids.
- ② While stopping condition not met:
 - ① Find closest centroid to each point.
 - ② Move centroids to the average of all the points closest to them.

The algorithm in all it's glory:

- ① Initialize centroids.
- ② While stopping condition not met:
 - ① Find closest centroid to each point.
 - ② Move centroids to the average of all the points closest to them.

This training algorithm may look pretty simple...

The algorithm in all it's glory:

- ① Initialize centroids.
- ② While stopping condition not met:
 - ① Find closest centroid to each point.
 - ② Move centroids to the average of all the points closest to them.

This training algorithm may look pretty simple... and that's because it is.

Supervised vs. Unsupervised Learning

Clustering

- Intuition
- Definition

k-Means Algorithm

- Pseudocode
- Centroid Initialization
- Stopping Criteria
- Step-through
- Evaluation
- Problems
- Choosing k

Centroid Initialization

- The simplest way to do this is to randomly choose k points from your data and make their locations your initial centroid locations.

Centroid Initialization

- The simplest way to do this is to randomly choose k points from your data and make their locations your initial centroid locations.
- Another straightforward method is to randomly assign each data point a number $1-k$, and start the initialize the k^{th} centroid to the average of the points with the k^{th} label (in each dimension).

A more advanced centroid initialization method, known as *k*-Means++, chooses well spread initial centroids.

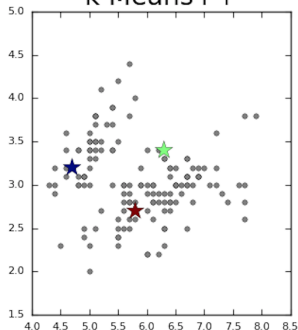
→ `sklearn: init='k-means++'`, set as default.

k-Means++ follows the procedure:

- 1 Choose the first centroid to be the location of a data point chosen at random.
- 2 For each remaining centroid, choose the location of a data point with probability proportional to its squared distance from the point's closest existing centroid (points further from existing centroids have higher probability of being chosen as the next centroid).

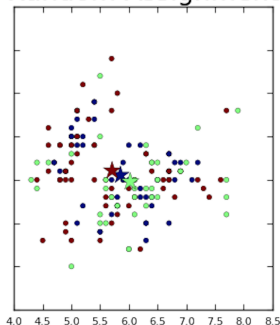
Initialization - Visual Comparison

k-Means++



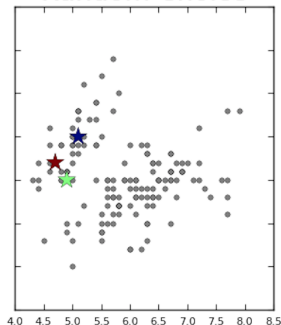
More even spread to start with.

Random Assignment



All start close to the center.

Random Choice



Who the eff knows... could be anything!

Supervised vs. Unsupervised Learning

Clustering

- Intuition
- Definition

k-Means Algorithm

- Pseudocode
- Centroid Initialization
- **Stopping Criteria**
- Step-through
- Evaluation
- Problems
- Choosing *k*

Stopping Criteria

We can update...

- for a pre-specified number of iterations.
→ `sklearn: max_iter=1000`.

Stopping Criteria

We can update...

- for a pre-specified number of iterations.
→ `sklearn: max_iter=1000`.
- until the centroids don't change at all - may take a ton of iterations.

Stopping Criteria

We can update...

- for a pre-specified number of iterations.
→ `sklearn: max_iter=1000`.
- until the centroids don't change at all - may take a ton of iterations.
- until the centroids don't move very much - takes fewer iterations.
→ `sklearn: tol=0.0001`, for tolerance of “how much”.

Supervised vs. Unsupervised Learning

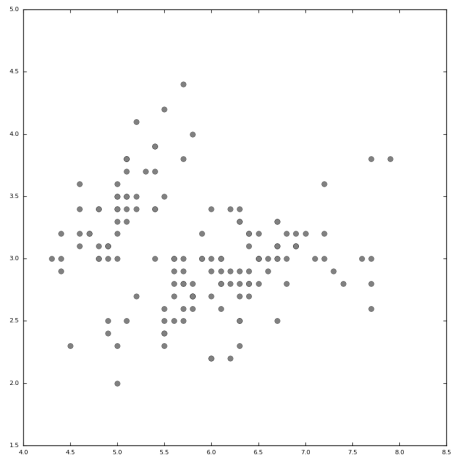
Clustering

- Intuition
- Definition

k-Means Algorithm

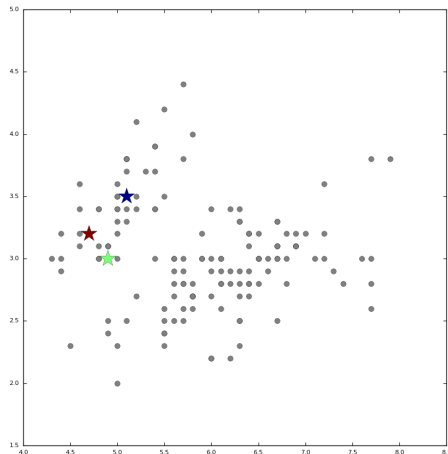
- Pseudocode
- Centroid Initialization
- Stopping Criteria
- **Step-through**
- Evaluation
- Problems
- Choosing *k*

Step-by-step Execution: DATA!!



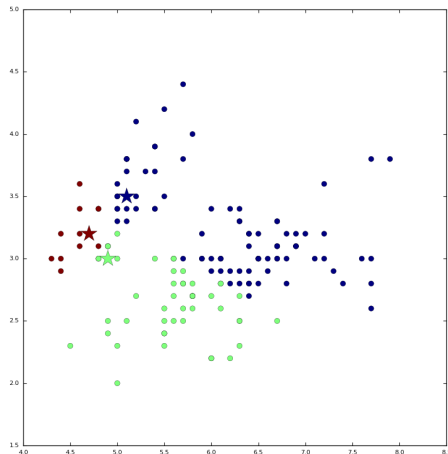
Step-by-step Execution: Initialize

- ① Initialize centroids.
- ② While not stopping condition:
 - ① Assign points to centroid
 - ② Move centroids to new average location



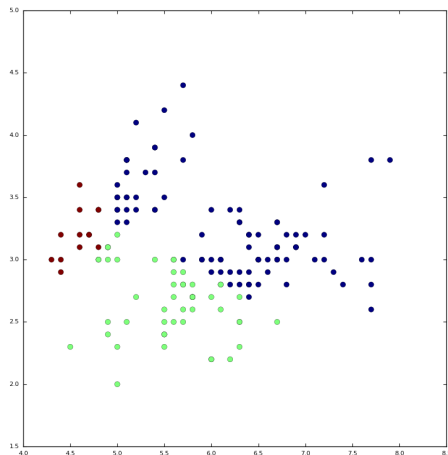
Step-by-step Execution: Iteration 1 - Step 1

- 1 Initialize centroids.
- 2 While not stopping condition:
 - 1 Assign points to centroid
 - 2 Move centroids to new average location



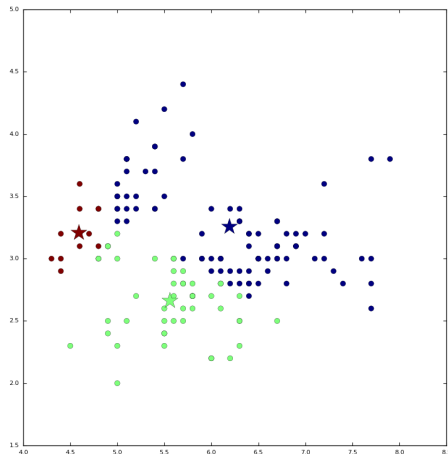
Step-by-step Execution: Iteration 1 - Prep Step 2

- 1 Initialize centroids.
- 2 While not stopping condition:
 - 1 Assign points to centroid
 - 2 Move centroids to new average location



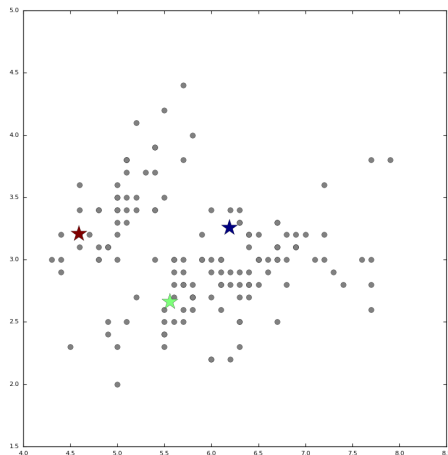
Step-by-step Execution: Iteration 1 - Step 2

- ① Initialize centroids.
- ② While not stopping condition:
 - ① Assign points to centroid
 - ② Move centroids to new average location



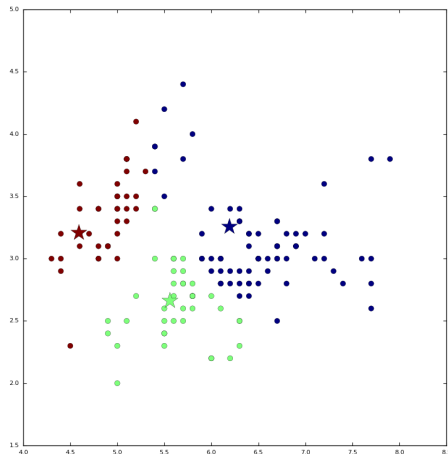
Step-by-step Execution: Iteration 2 - Prep Step 1

- 1 Initialize centroids.
- 2 While not stopping condition:
 - 1 Assign points to centroid
 - 2 Move centroids to new average location



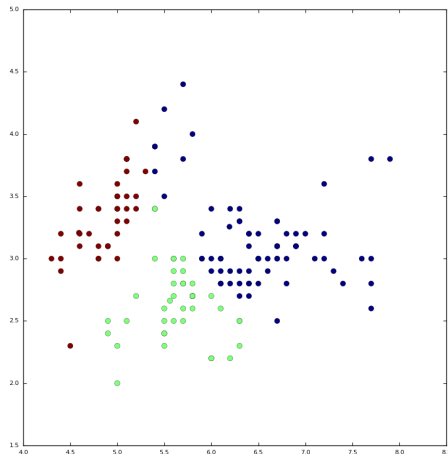
Step-by-step Execution: Iteration 2 - Step 1

- 1 Initialize centroids.
- 2 While not stopping condition:
 - 1 Assign points to centroid
 - 2 Move centroids to new average location



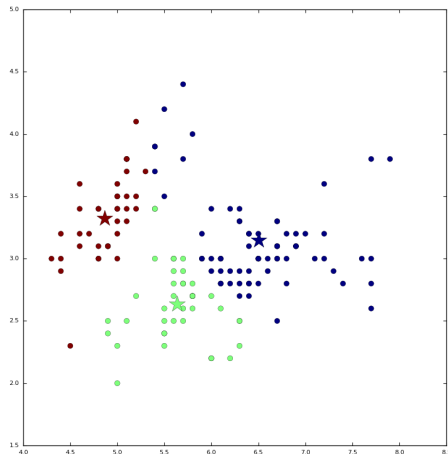
Step-by-step Execution: Iteration 2 - Prep Step 2

- 1 Initialize centroids.
- 2 While not stopping condition:
 - 1 Assign points to centroid
 - 2 Move centroids to new average location



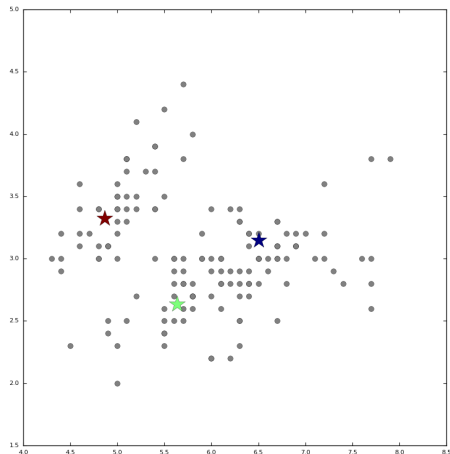
Step-by-step Execution: Iteration 2 - Step 2

- 1 Initialize centroids.
- 2 While not stopping condition:
 - 1 Assign points to centroid
 - 2 Move centroids to new average location



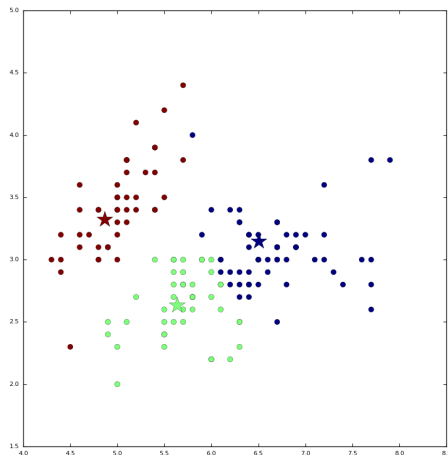
Step-by-step Execution: Iteration 3 - Prep Step 1

- ① Initialize centroids.
- ② While not stopping condition:
 - ① Assign points to centroid
 - ② Move centroids to new average location



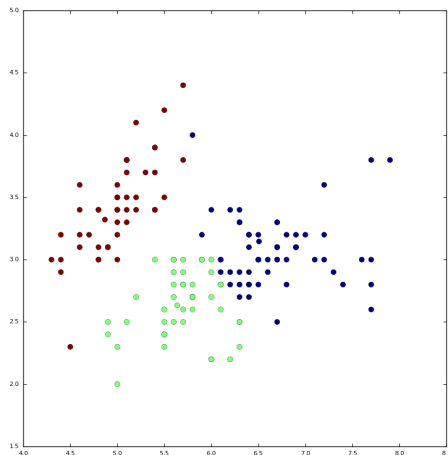
Step-by-step Execution: Iteration 3 - Step 1

- ① Initialize centroids.
- ② While not stopping condition:
 - ① Assign points to centroid
 - ② Move centroids to new average location



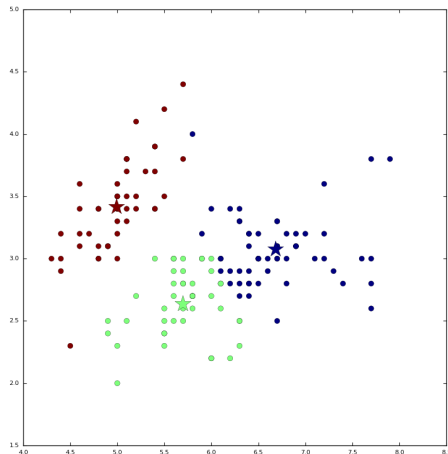
Step-by-step Execution: Iteration 3 - Prep Step 2

- 1 Initialize centroids.
- 2 While not stopping condition:
 - 1 Assign points to centroid
 - 2 Move centroids to new average location



Step-by-step Execution: Iteration 3 - Step 2

- 1 Initialize centroids.
- 2 While not stopping condition:
 - 1 Assign points to centroid
 - 2 Move centroids to new average location



Supervised vs. Unsupervised Learning

Clustering

- Intuition
- Definition

k-Means Algorithm

- Pseudocode
- Centroid Initialization
- Stopping Criteria
- Step-through
- **Evaluation**
- Problems
- Choosing *k*

Evaluating k -Means

- How can we quantify how “good” our clustering is?

Evaluating k -Means

- How can we quantify how “good” our clustering is?
- A good measure should quantify how similar things are in a cluster.

Evaluating k -Means

- How can we quantify how “good” our clustering is?
- A good measure should quantify how similar things are in a cluster.
- The metric that we will use is called intra-cluster or within cluster variance:

$$WCV = \sum_{k=1}^K \frac{1}{|C_k|} \sum_{i,i' \in C_k} \sum_{j=1}^p (x_{ij} - x_{i'j})^2$$

Supervised vs. Unsupervised Learning

Clustering

- Intuition
- Definition

k-Means Algorithm

- Pseudocode
- Centroid Initialization
- Stopping Criteria
- Step-through
- Evaluation
- **Problems**
- Choosing *k*

- Centroids that are “discovered” will likely be different depending on initialization.

- Centroids that are “discovered” will likely be different depending on initialization.
 - Run algorithm more than once and choose the run that yields the smallest intra-cluster variance.

- Centroids that are “discovered” will likely be different depending on initialization.
 - Run algorithm more than once and choose the run that yields the smallest intra-cluster variance.
- k -Means is highly dependent on distance as a metric.

- Centroids that are “discovered” will likely be different depending on initialization.
 - Run algorithm more than once and choose the run that yields the smallest intra-cluster variance.
- k -Means is highly dependent on distance as a metric.
 - Normalize features before clustering.

- Centroids that are “discovered” will likely be different depending on initialization.
 - Run algorithm more than once and choose the run that yields the smallest intra-cluster variance.
- k -Means is highly dependent on distance as a metric.
 - Normalize features before clustering.
 - Have to think about the curse of dimensionality.

Supervised vs. Unsupervised Learning

Clustering

- Intuition
- Definition

k -Means Algorithm

- Pseudocode
- Centroid Initialization
- Stopping Criteria
- Step-through
- Evaluation
- Problems
- Choosing k

Choosing k

Unsupervised

Choosing k is HARD!!! It usually takes some work and you're never quite sure if you're "right".

Unsupervised

Choosing k is HARD!!! It usually takes some work and you're never quite sure if you're "right".

There are a number of ways you can go about choosing k :

- Domain knowledge
- Elbow method
- Silhouette score
- GAP Statistic

- Looks at the total amount of within-cluster sum of squares (WCSS) across all the clusters for different values of k .

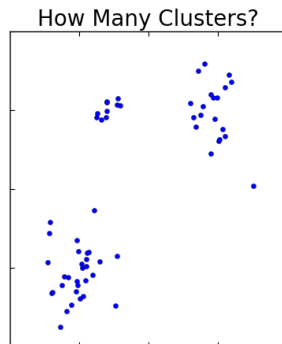
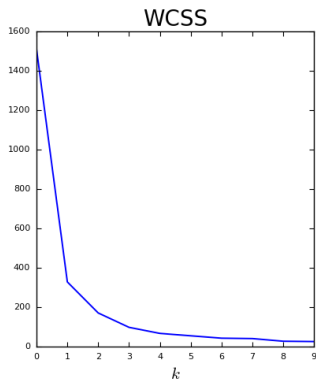
$$WCSS = \sum_{k=1}^K \sum_{i, i' \in C_k} \sum_{j=1}^p (x_{ij} - x_{i'j})^2$$

- Looks at the total amount of within-cluster sum of squares (WCSS) across all the clusters for different values of k .

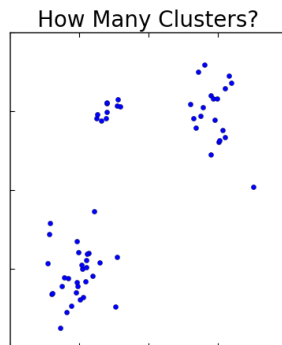
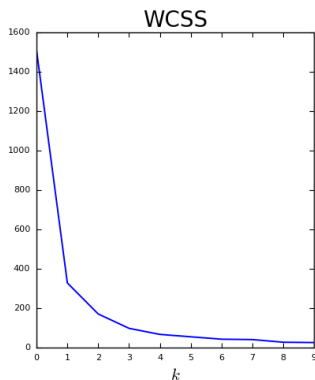
$$WCSS = \sum_{k=1}^K \sum_{i,i' \in C_k} \sum_{j=1}^p (x_{ij} - x_{i'j})^2$$

- Chooses the k such that adding one more cluster doesn't decrease the WCSS by much more. Leads us to look for an elbow in the k vs. WCSS plot.

Elbow Method

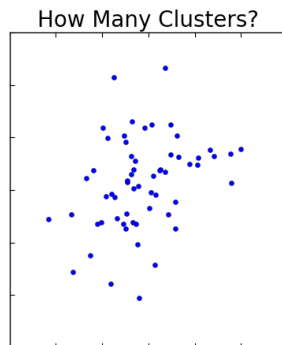
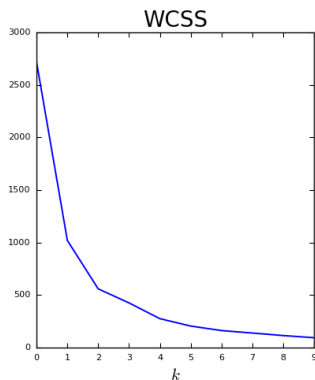


Elbow Method

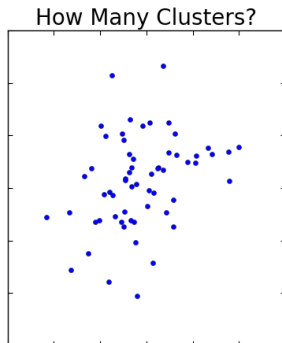
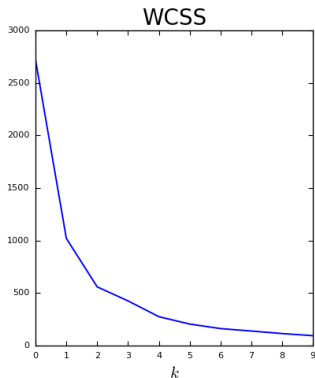


Question: Do you think the elbow will always be so obvious?

Elbow Method - Not Always So Clear



Elbow Method - Not Always So Clear



Question: How is this related to the curse of dimensionality?