

# **INTRO to DATA SCIENCE**

## **LECTURE 11: K-MEANS CLUSTERING**

# I. CLUSTER ANALYSIS

	<i>continuous</i>	<i>categorical</i>
<i>supervised</i>	<i>regression</i>	<i>classification</i>
<i>unsupervised</i>	<i>dimension reduction</i>	<i>clustering</i>

---

<i><b>supervised</b></i>	<i><b>making predictions</b></i>
<i><b>unsupervised</b></i>	<i><b>discovering patterns</b></i>

<i><b>supervised</b></i>	<i><b>labeled examples</b></i>
<i><b>unsupervised</b></i>	<i><b>no labeled examples</b></i>

*Q: What is a cluster?*

*Q: What is a cluster?*

*A: A group of **similar data points**.*

*Q: What is a cluster?*

*A: A group of **similar** data points.*

*The concept of similarity is central to the definition of a cluster, and therefore to cluster analysis.*



*Q: What is the purpose of cluster analysis?*

*Q: What is the purpose of cluster analysis?*

*A: To enhance our understanding of a dataset by dividing the data into groups.*

*Q: What is the purpose of cluster analysis?*

*A: To enhance our understanding of a dataset by dividing the data into groups.*

*Clustering provides a layer of abstraction from individual data points.*

*The goal is to extract and enhance the natural structure of the data (not to impose arbitrary structure!)*

## People You May Know



**Kamal Kumar**

1 mutual friend  
[Add to My Friends](#)



**MrsI F**

1 mutual friend  
[Add to My Friends](#)



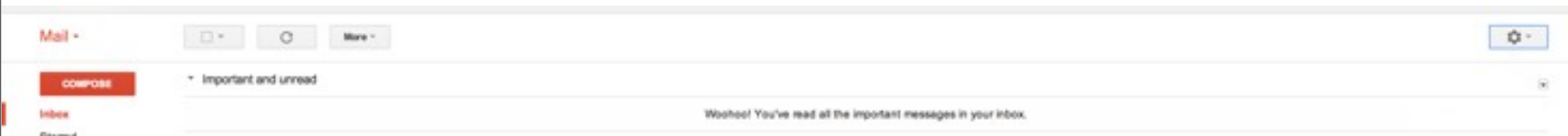
**Imran  
Memmedov**

1 mutual friend  
[Add to My Friends](#)



**Rick Cruz**

1 mutual friend  
[Add to My Friends](#)



## Priority Inbox: Unsupervised Learning

Group mails into groups and decide which group represents important mails



Roll over image to zoom in



See 1 customer image

Share your own customer images

## Star Trek [Blu-ray] (2009)

Chris Pine (Actor), Zachary Quinto (Actor), J.J. Abrams (Director) | Rated: PG-13 | Format: Blu-ray  
[View details](#) [\(2,040 customer reviews\)](#)

List Price: ~~\$22.66~~

Price: **\$9.99 & FREE Shipping** on orders over \$25. [Details](#)

You Save: \$12.99 (57%)

In Stock.

Ships from and sold by Amazon.com. Gift-wrap available.

**Want it Wednesday, June 5?** Order within 20 hrs 11 mins and choose **One-Day Shipping** at checkout. [Details](#)

**23 new** from \$9.98 **18 used** from \$9.78 **1 collectible** from \$49.99

Watch Instantly with <b>amazon instant video</b>		<b>Rent</b>	<b>Buy</b>
Star Trek (2009)		\$2.99	\$9.99
Other Formats & Versions			
Blu-ray	1-Disc Version	Amazon Price	New from Used from
DVD	Single-Disc Edition	\$8.49	\$5.65 \$3.29



This week only, save up to 62% on [Fingerpuppet: The Complete Series](#) in our Deal of the Week. Offer ends June 8, 2013. [Learn more](#)

### Frequently Bought Together



Price for all three: **\$66.47**

[Add all three to Cart](#)

[Add all three to Wish List](#)

Some of these items ship sooner than the others. [Show details](#)

- This item:** Star Trek [Blu-ray] ~ Chris Pine Blu-ray **\$9.99**
- Star Trek Into Darkness (Blu-ray 3D + Blu-ray + DVD + Digital Copy) ~ Chris Pine Blu-ray **\$24.99**
- Iron Man 3 (Two-Disc Blu-ray / DVD + Digital Copy) ~ Robert Downey Jr. Blu-ray **\$31.49**

### What Other Items Do Customers Buy After Viewing This Item?

Star Trek Into Darkness (Blu-ray 3D + Blu-ray + DVD + Digital Copy) ~ Chris Pine Blu-ray  
[View details](#) (199)  
**\$24.99**

Star Trek: Original Motion Picture Collection (Star Trek I, II, III, IV, V, VI + The Captain's Summit Bonus Disc) [Blu-ray] ~ William Shatner Blu-ray  
[View details](#) (571)  
**\$53.56**

Sin City (Two-Disc Theatrical & Recut, Extended, and Unrated Versions) [Blu-ray] ~ Jessica Alba Blu-ray  
[View details](#) (933)  
**\$4.99**

Quantity:

Yes, I want **FREE Two-Day Shipping** with **Amazon Prime**

[Add to Cart](#)

or

[Sign in](#) to turn on 1-Click ordering.

[Add to Wish List](#)

### Sell Us Your Item

For up to a **\$2.60 Gift Card**

[Trade in](#)

[Learn more](#)

### More Buying Choices

TechShowMe [Add to Cart](#)  
**\$19.99 & FREE Shipping** on orders over \$25. [Details](#)

**43 used & new** from \$9.78

Have one to sell? [Sell on Amazon](#)

Share [Email](#) [Facebook](#) [Twitter](#) [Pinterest](#)

*Q: How do you solve a clustering problem?*

*Q: How do you solve a clustering problem?*

*A: Think of a cluster as a “potential class”; then the solution to a clustering problem is to programatically determine these classes.*



*Q: How do you solve a clustering problem?*

*A: Think of a cluster as a “potential class”; then the solution to a clustering problem is to programatically determine these classes.*

*The real purpose of clustering is data exploration, so a solution is anything that contributes to your understanding.*

# II. K-MEANS CLUSTERING

*Q: What is  $k$ -means clustering?*

*Q: What is  $k$ -means clustering?*

*A: A **greedy** learner that **partitions** a data set into  $k$  clusters.*

*Q: What is  $k$ -means clustering?*

*A: A **greedy** learner that **partitions** a data set into  $k$  clusters.*

*greedy – captures local structure (depends on initial conditions)*

*Q: What is  $k$ -means clustering?*

*A: A **greedy** learner that **partitions** a data set into  $k$  clusters.*

*greedy – captures local structure*

*partition – each point belongs to exactly one cluster*

*Q: How are these partitions determined?*

*Q: How are these partitions determined?*

*A: Each point is assigned to the cluster with the nearest **centroid**.*



*Q: How are these partitions determined?*

*A: Each point is assigned to the cluster with the nearest **centroid**.*

*centroid – the mean of the data points in a cluster*

*→ requires continuous (vector-like) features*

*→ highlights iterative nature of algorithm*

*One important point to keep in mind is that partitions are not scale-invariant!*

*One important point to keep in mind is that partitions are not scale-invariant!*

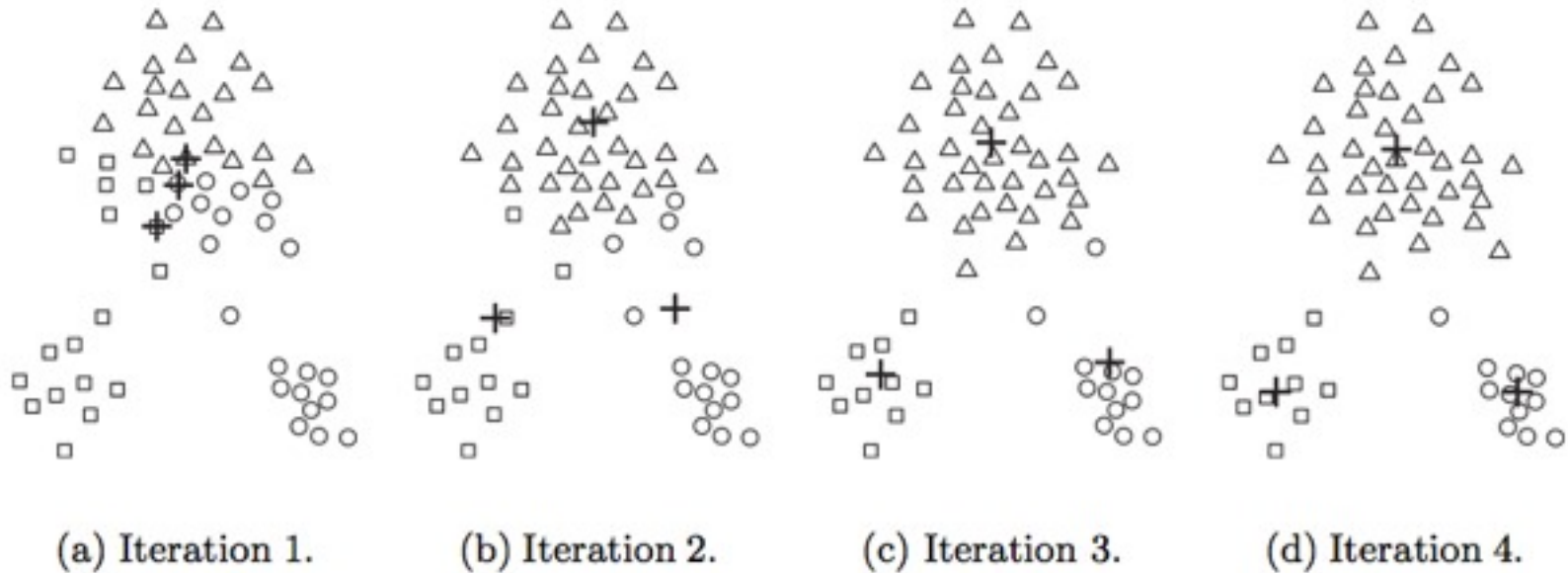
*This means that the same data can yield very different clustering results depending on the scale and the units used.*

*One important point to keep in mind is that partitions are not scale-invariant!*

*This means that the same data can yield very different clustering results depending on the scale and the units used.*

*Therefore it's important to think about your data representation before applying a clustering algorithm.*

- 1) *choose  $k$  initial centroids (note that  $k$  is an input)*
- 2) *for each point:*
  - *find distance to each centroid*
  - *assign point to nearest centroid*
- 3) *recalculate centroid positions*
- 4) *repeat steps 2-3 until stopping criteria met*



**Figure 8.3.** Using the K-means algorithm to find three clusters in sample data.

source: <http://www-users.cs.umn.edu/~kumar/dmbook/ch8.pdf>

*K-means is algorithmically pretty efficient (time & space complexity is linear in number of records).*

*Q: How do you choose the initial centroid positions?*



*Q: How do you choose the initial centroid positions?*

*A: There are several options:*

*Q: How do you choose the initial centroid positions?*

*A: There are several options:*

- randomly (but may yield divergent behavior)*

*Q: How do you choose the initial centroid positions?*

*A: There are several options:*

- randomly (but may yield divergent behavior)*
- perform alternative clustering task, use resulting centroids as initial k-means centroids*

*Q: How do you determine which centroid is the nearest?*

*Q: How do you determine which centroid is the nearest?*

*The “nearness” criterion is determined by the similarity/distance measure we discussed earlier.*

*Q: How do you determine which centroid is the nearest?*

*The “nearness” criterion is determined by the similarity/distance measure we discussed earlier.*

*This measure makes quantitative inference possible.*

*There are a number of different similarity measures to choose from, and in general the right choice depends on the problem.*

*There are a number of different similarity measures to choose from, and in general the right choice depends on the problem.*

*For data that takes values in  $\mathbb{R}^n$ , the typical choice is the **Euclidean distance**:*

$$d(x, y) = \sqrt{\sum (x_i - y_i)^2}$$



*There are a number of different similarity measures to choose from, and in general the right choice depends on the problem.*

*For data that takes values in  $\mathbb{R}^n$ , the typical choice is the **Euclidean distance**:*

$$d(x, y) = \sqrt{\sum (x_i - y_i)^2}$$

*We can express different semantics about our data through the choice of metric.*

*The matrix whose entries  $D_{ij}$  contain the values  $d(x, y)$  for all  $x$  and  $y$  is called the **distance matrix**.*

*The matrix whose entries  $D_{ij}$  contain the values  $d(x, y)$  for all  $x$  and  $y$  is called the **distance matrix**.*

*The distance matrix contains all of the information we know about the dataset.*

*The matrix whose entries  $D_{ij}$  contain the values  $d(x, y)$  for all  $x$  and  $y$  is called the **distance matrix**.*

*The distance matrix contains all of the information we know about the dataset.*

*For this reason, it's really the choice of metric that determines the definition of a cluster.*

*Q: How do we recompute the positions of the centroids at each iteration of the algorithm?*

*Q: How do we recompute the positions of the centroids at each iteration of the algorithm?*

*A: By optimizing an **objective function** that tells us how “good” the clustering is.*

*Q: How do we recompute the positions of the centroids at each iteration of the algorithm?*

*A: By optimizing an **objective function** that tells us how “good” the clustering is.*

*The iterative part of the algorithm (recomputing centroids and reassigning points to clusters) explicitly tries to minimize this objective function.*

*Ex: Using the Euclidean distance measure, one typical objective function is the **sum of squared errors** from each point  $x$  to its centroid  $c_i$ :*

$$SSE = \sum_{i=1}^K \sum_{x \in C_i} d(x, c_i)^2$$



*Ex: Using the Euclidean distance measure, one typical objective function is the **sum of squared errors** from each point  $x$  to its centroid  $c_i$ :*

$$SSE = \sum_{i=1}^K \sum_{x \in C_i} d(x, c_i)^2$$

*Given two clusterings, we will prefer the one with the lower SSE since this means the centroids have converged to better locations (a better local optimum).*

*We iterate until some stopping criteria are met; in general, suitable convergence is achieved in a small number of steps.*

# **III. CLUSTER VALIDATION**

---

<i><b>supervised</b></i> <i><b>unsupervised</b></i>	<i><b>test out your predictions</b></i> <i><b>...</b></i>
--	--

*In general,  $k$ -means will converge to a solution and return a partition of  $k$  clusters, even if no natural clusters exist in the data.*

*In general,  $k$ -means will converge to a solution and return a partition of  $k$  clusters, even if no natural clusters exist in the data.*

*We will look at two validation metrics useful for partitional clustering, **cohesion and separation**.*

**Cohesion** *measures clustering effectiveness within a cluster.*

$$\hat{C}(C_i) = \sum_{x \in C_i} d(x, c_i)$$

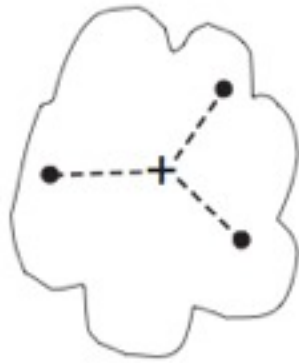
**Cohesion** *measures clustering effectiveness within a cluster.*

$$\hat{C}(C_i) = \sum_{x \in C_i} d(x, c_i)$$

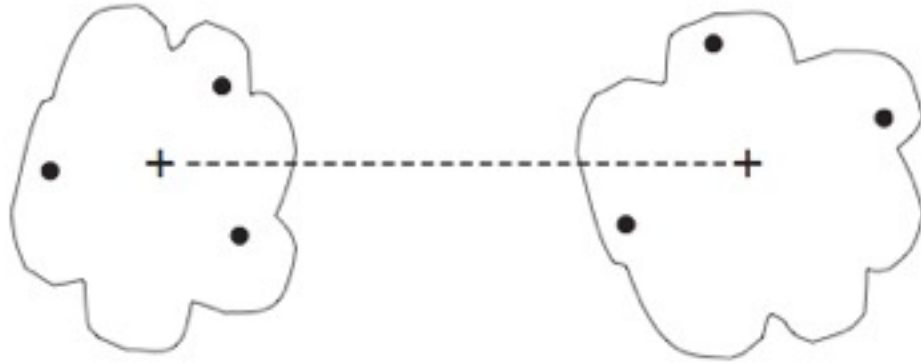
**Separation** *measures clustering effectiveness between clusters.*

$$\hat{S}(C_i, C_j) = d(c_i, c_j)$$





(a) Cohesion.



(b) Separation.

**Figure 8.28.** Prototype-based view of cluster cohesion and separation.

source: <http://www-users.cs.umn.edu/~kumar/dmbook/ch8.pdf>

*We can turn these values into overall measures of clustering validity by taking a weighted sum over clusters:*

$$\hat{V}_{total} = \sum_1^K w_i \hat{V}(C_i)$$

*Here  $V$  can be cohesion, separation, or some function of both.*

*We can turn these values into overall measures of clustering validity by taking a weighted sum over clusters:*

$$\hat{V}_{total} = \sum_1^K w_i \hat{V}(C_i)$$

*Here  $V$  can be cohesion, separation, or some function of both.*

*The weights can all be set to 1 (best for  $k$ -means), or proportional to the cluster masses (the number of points they contain).*

*Cluster validation measures can be used to identify clusters that should be split or merged, or to identify individual points with disproportionate effect on the overall clustering.*

*One useful application of cluster validation is to determine the best number of clusters for your dataset.*

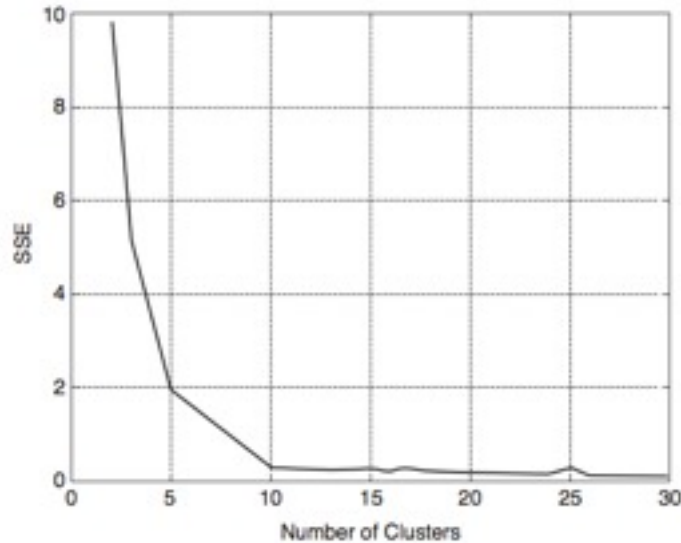
*One useful application of cluster validation is to determine the best number of clusters for your dataset.*

*Q: How would you do this?*

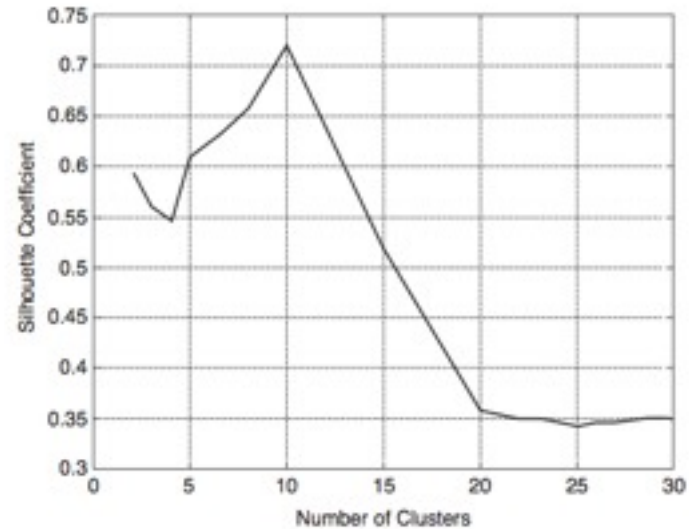
*One useful application of cluster validation is to determine the best number of clusters for your dataset.*

*Q: How would you do this?*

*A: By computing the overall SSE for different values of  $k$ .*



**Figure 8.32.** SSE versus number of clusters for the data of Figure 8.29.



**Figure 8.33.** Average silhouette coefficient versus number of clusters for the data of Figure 8.29.

source: <http://www-users.cs.umn.edu/~kumar/dmbook/ch8.pdf>



*Ultimately, cluster validation and clustering in general are suggestive techniques that rely on human interpretation to be meaningful.*

---

**INTRO TO DATA SCIENCE**

---

**EX: K-MEANS CLUSTERING**

Friday, April 4, 14