

APACHE SPARK ADVANCED

DS8003 – MGT OF BIG DATA AND TOOLS

Ryerson University

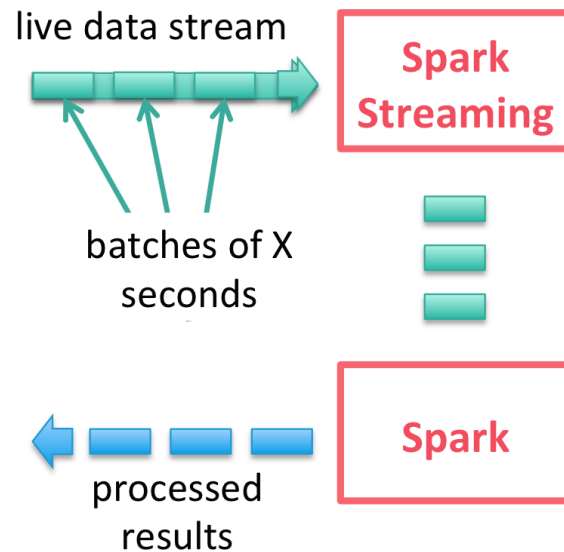
Instructor: Kanchana Padmanabhan

Streaming

2

Run a streaming computation as a series of very small deterministic batch jobs

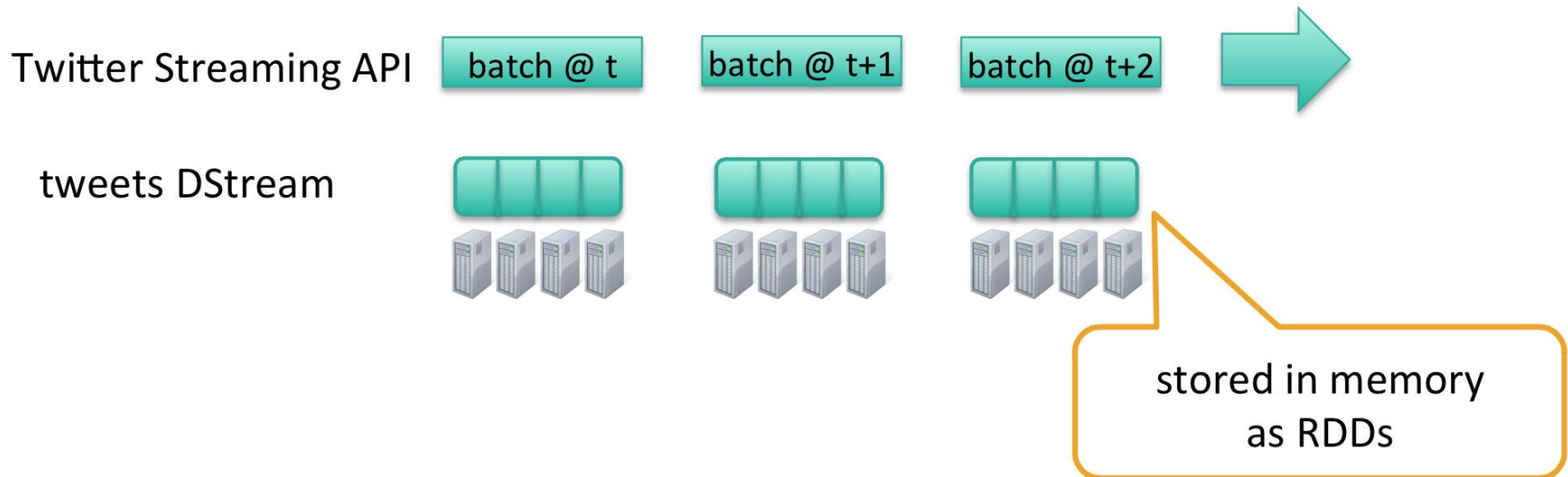
- Chop up the live stream into batches of X seconds (1/2 second or 1 second)
- Spark treats each batch of data as RDD
- Run same code in batch mode and streaming



Streaming - DStream

3

- Discretized Stream (Dstream)
 - ▣ Represents a stream of data
 - ▣ Implemented as sequence of RDDs
 - ▣ Created from streaming input sources
 - ▣ Created by applying transformations



Streaming - Code

4

1. `from pyspark import SparkContext`
2. `from pyspark.streaming import StreamingContext`
3. `conf = SparkConf().setAppName("Testing Spark Streaming")`
4. `sc = SparkContext(conf = conf)` [Required only if code is in python file.
Not required in pyspark console]
5. `stream = StreamingContext(sc, 1)`
6. `lines = stream.socketTextStream("localhost", 9999)`
7. `counts = lines.flatMap(lambda line: line.split(" ")).map(lambda word: (word, 1)).reduceByKey(lambda a, b: a+b)`
8. `counts.pprint()`
9. `stream.start()`
10. `stream.awaitTermination()`
11. [Code is same as **streamingExample.py**]

Testing your streaming code

5

- Example details
here: <http://spark.apache.org/docs/latest/streaming-programming-guide.html>
- Open two connections to VirtualBox (you will open two consoles)
- One first console type “spark-submit streamingExample.py”
- On second console type command “nc -lk 9999”. Start typing some text on the console screen
- The text will get processed by the spark program running on the first console

DataFrames - SQLContext

6

- ❑ Spark SQL is a Spark module for structured data processing
- ❑ The DataFrame provides easier access to data, because it looks conceptually like a table
- ❑ A DataFrame is a *Dataset* organized into named columns. It is conceptually equivalent to a table in a relational database or a data frame in R/Python
- ❑ An extension SQLContext is HiveContext
- ❑ SqlContext can also load JSON's into DataFrames
 - ▣ Uses “null” for missing values

Creating a DataFrame

7

- ❑ `from pyspark.sql import SQLContext`
- ❑ `sqlContext = SQLContext(sc)`
- ❑ `v = sqlContext.createDataFrame([
 ("a", "Alice", 34),
 ("b", "Bob", 36),
 ("c", "Charlie", 30),
], ["id", "name", "age"])`
- ❑ `v.select("id").show()`
- ❑ `v.groupBy("age").count().show()`

Machine Learning - Mllib

8

- Mllib is Spark's machine learning (ML) library. Its goal is to make practical machine learning scalable and easy.
- It consists of common learning algorithms and utilities, including classification, regression, clustering, collaborative filtering, dimensionality reduction
- Spark excels at iterative computation, enabling Mllib to run fast. At the same time, we care about algorithmic performance: Mllib contains high-quality algorithms that leverage iteration, and can yield better results than the one-pass approximations sometimes used on MapReduce.
- Supposedly, running times or up to 100x faster than Hadoop MapReduce

<http://spark.apache.org/docs/latest/mllib-guide.html>

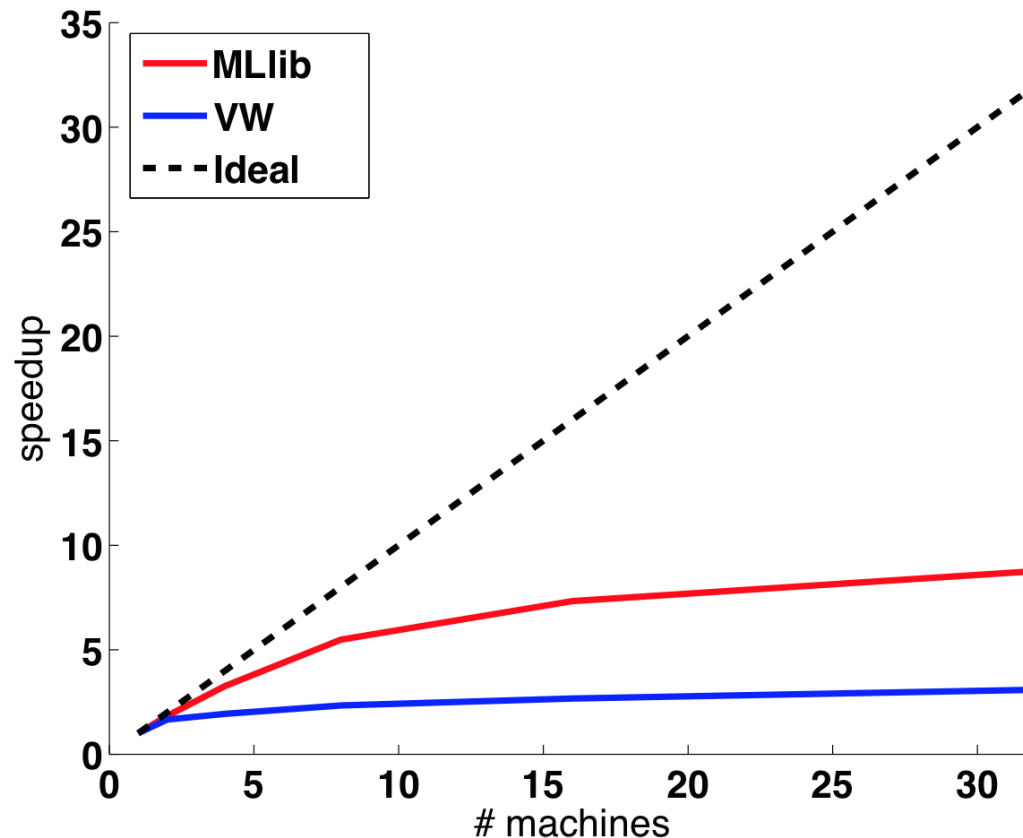
[https://databricks-](https://databricks-training.s3.amazonaws.com/slides/Spark_Summit_MLlib_070214_v2.pdf)

[training.s3.amazonaws.com/slides/Spark_Summit_MLlib_070214_v2.pdf](https://databricks-training.s3.amazonaws.com/slides/Spark_Summit_MLlib_070214_v2.pdf)

Mllib - scaling

9

Logistic Regression



Classification

10

- Classification takes a set of data already divided into predefined groups and searches for patterns in the data that differentiate those groups.
- The discovered patterns then can be used to classify other data where the right group designation is unknown (though other attributes may be known).
- Logistic regression – as a classifier. It is widely used to predict a binary response.

<http://www.britannica.com/technology/data-mining>

Logistic Regression – As a Classifier

11

- Models relationship between set of variables
 - ▣ – dichotomous such as rain (yes/no)
 - ▣ – categorical (department, country)
 - ▣ – continuous (age, salary, weight, height...)
- Binary outcome (Y) variable (sentiment, 1 = positive, and 0 = negative)

<https://www.nemoursresearch.org/open/StatClass/January2011/Class8.pdf>

http://www.ncss.com/wp-content/themes/ncss/pdf/Procedures/NCSS/Logistic_Regression.pdf

Logistic Regression - Example

12

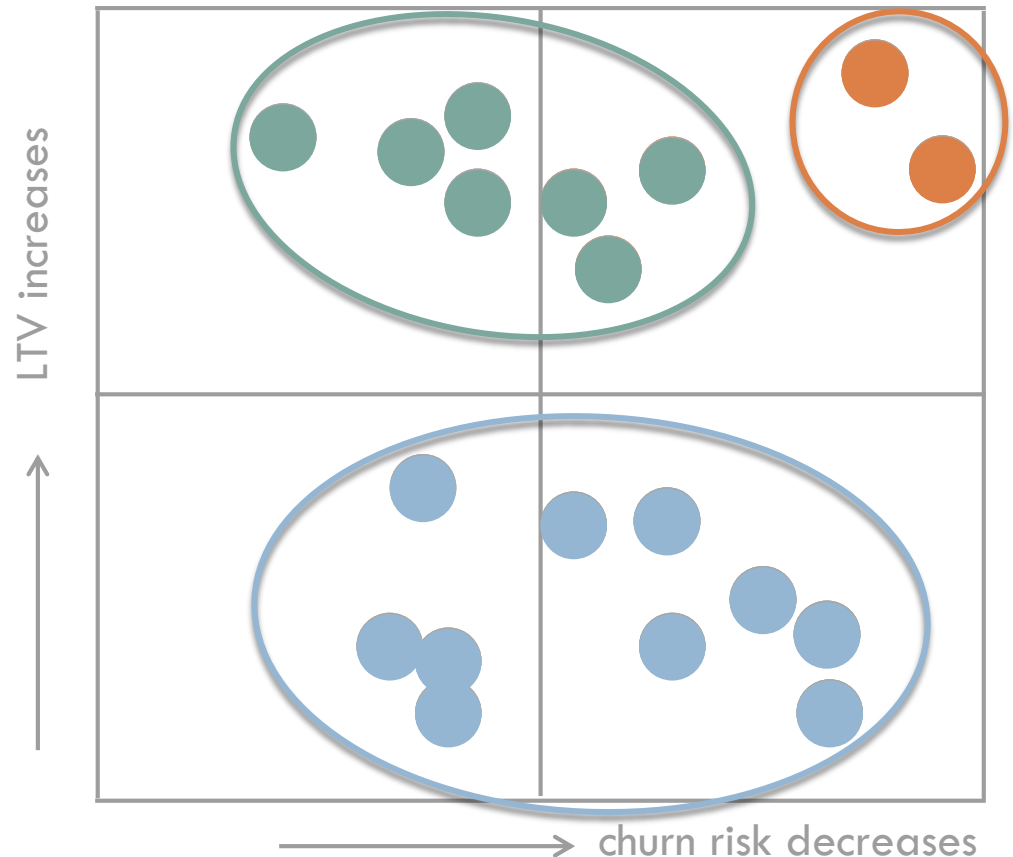
```
1.  from pyspark.mllib.classification import LogisticRegressionWithSGD, LogisticRegressionModel
2.  from pyspark.mllib.regression import LabeledPoint
3.  from pyspark import SparkConf, SparkContext
4.
   conf = SparkConf().setAppName("Testing Spark Commands")
5.  sc = SparkContext(conf = conf)
6.  data = [
       LabeledPoint(0.0, [0.0, 1.0]),
       LabeledPoint(1.0, [1.0, 0.0])
   ]
1.  lrm = LogisticRegressionWithSGD.train(sc.parallelize(data), iterations=10)
2.  lrm.predict([1.0, 0.0])
3.  #Output: 1.0
4.  lrm.predict([0.0, 1.0])
5.  #Output: 0.0
6.  # Save and load model
7.  lrm.save(sc, "lrsgd")
8.  sameModel = LogisticRegressionModel.load(sc, "lrsgd")
9.  sameModel.predict([1.0, 0.0])
10. sameModel.predict([0.0, 1.0])
```

[Code is same as logistic_regression_simple.py]

Clustering

13

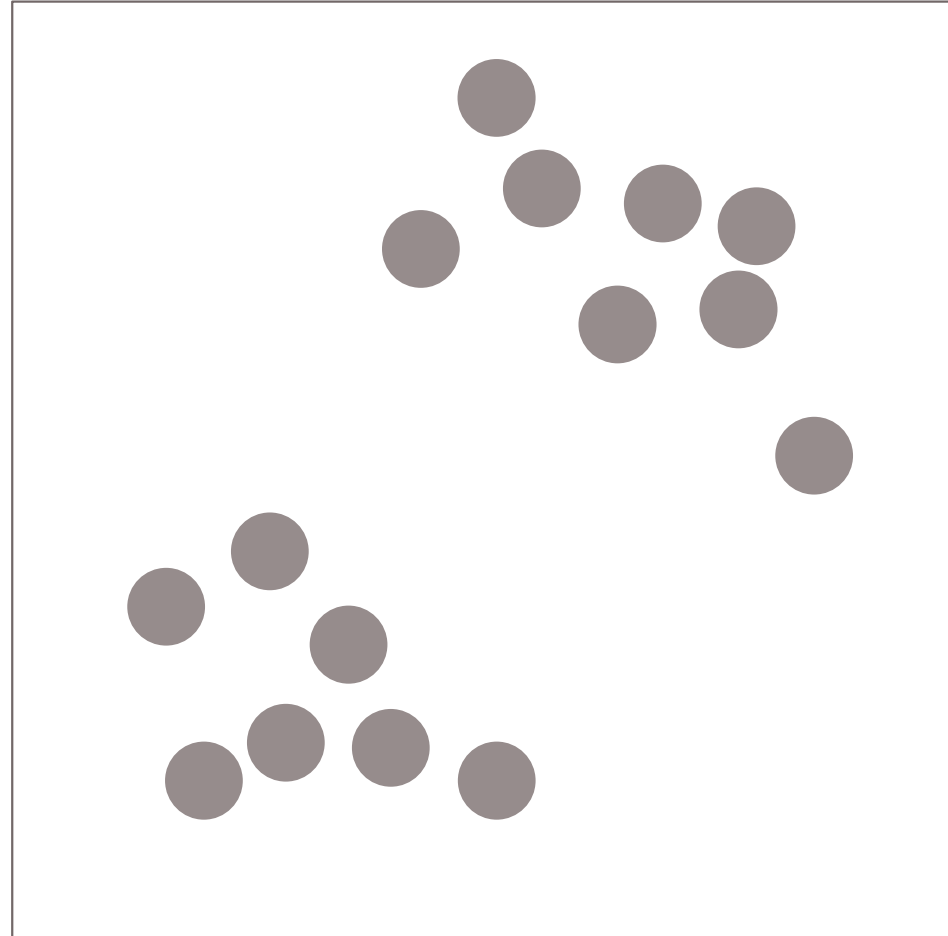
- Customer Value Segmentation
 - ▣ Value-based segmentation is useful for designing CRM strategies
 - ▣ Segmentation is easy to explain and useful for designing your campaigns/offers/messages
 - ▣ Clustering is more effective algorithmically but suffers explainability



K-Means – Single Node

14

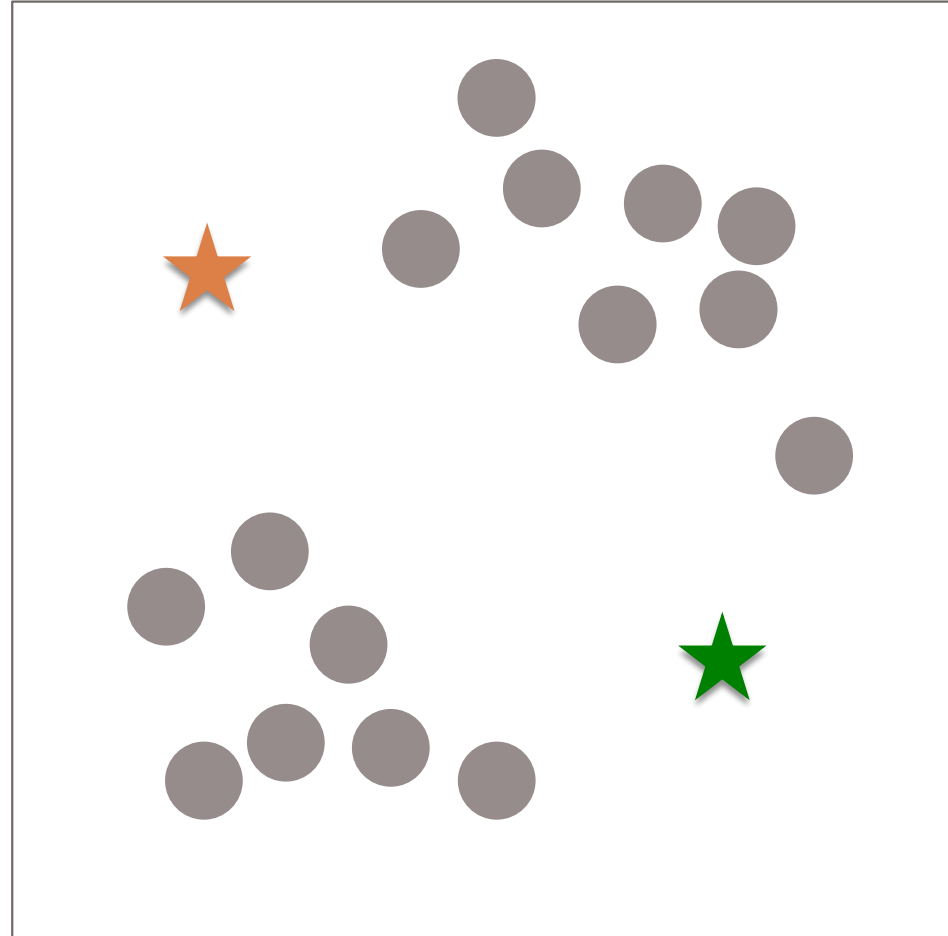
- Kmeans → Iterative algorithm until convergence
 1. Select K points at random as cluster centroids (centers)
 - Now we formed K clusters
 2. For each data point, assign it to the closest center
 3. For each cluster, re-compute the centers
 - E.g., in the case of 2D points →
 - X: average over all x-axis points in the cluster
 - Y: average over all y-axis points in the cluster
 4. If the new centers are different from the old centers (previous iteration) → Go to Step 2
 - Otherwise, stop



K-Means – Single Node

15

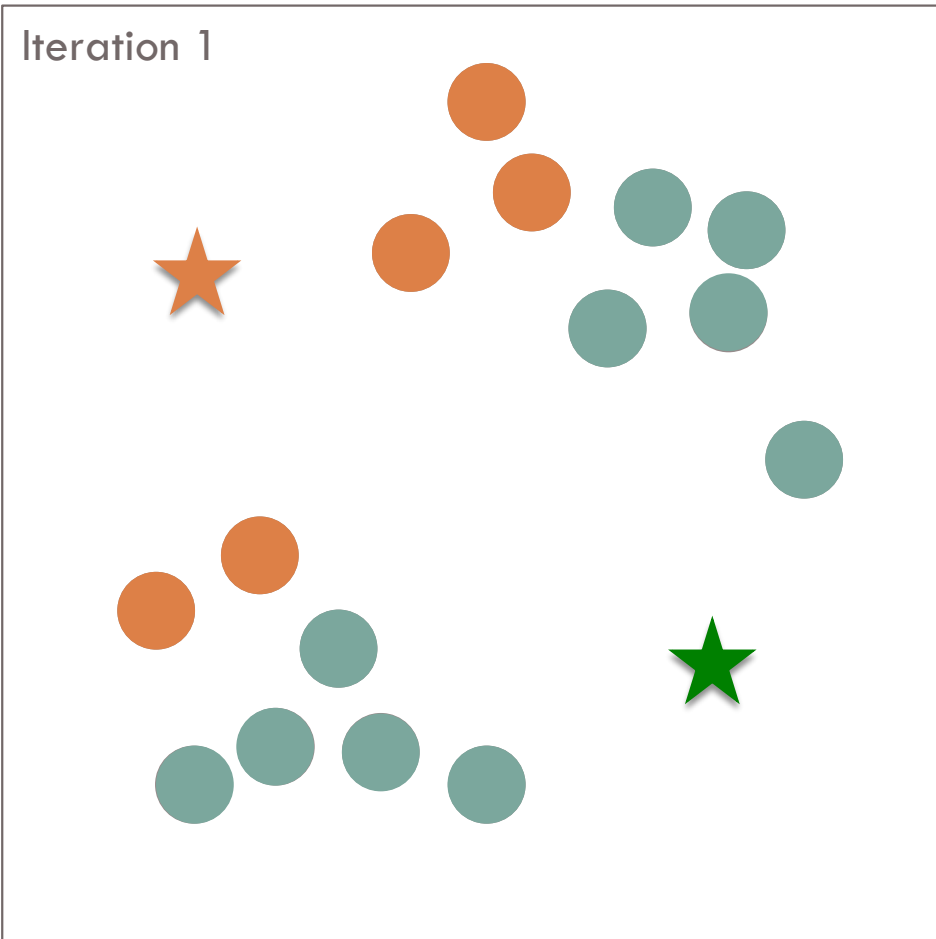
- Kmeans → Iterative algorithm until convergence
 1. Select K points at random as cluster centroids (centers)
 2. For each data point, assign it to the closest center
 - Now we formed K clusters
 3. For each cluster, re-compute the centers
 - E.g., in the case of 2D points →
 - X: average over all x-axis points in the cluster
 - Y: average over all y-axis points in the cluster
 4. If the new centers are different from the old centers (previous iteration) → Go to Step 2
 - Otherwise, stop



K-Means – Single Node

16

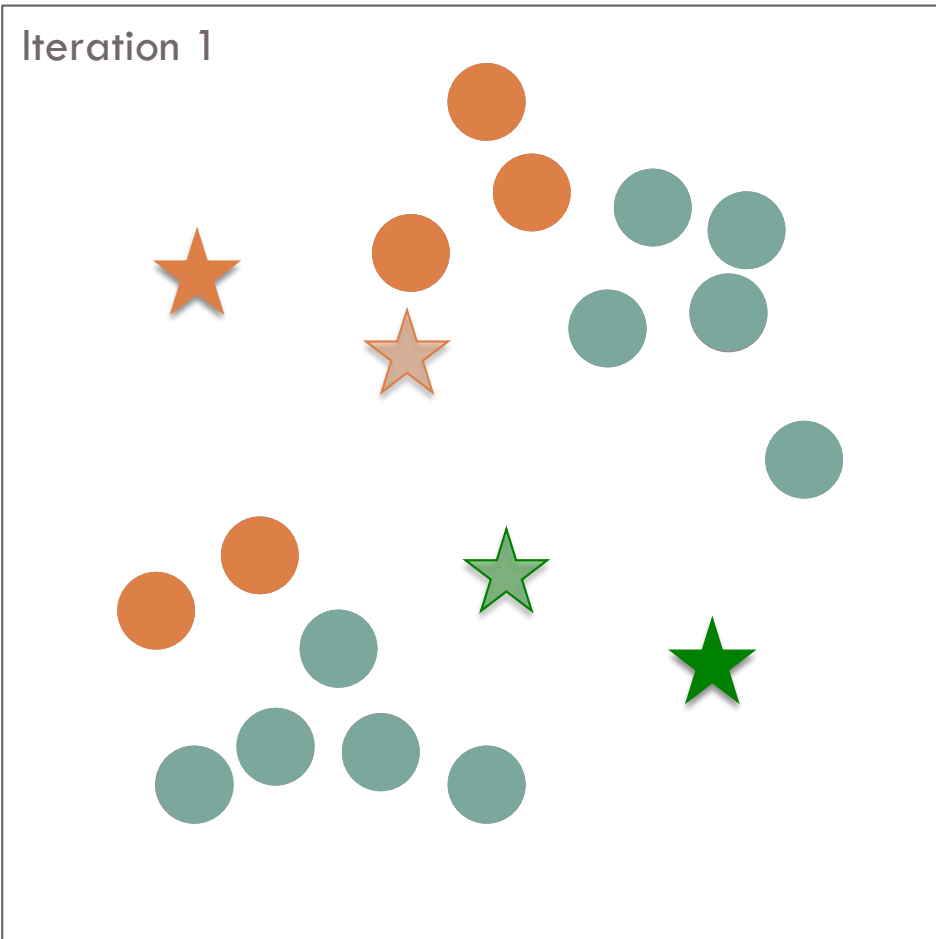
- Kmeans → Iterative algorithm until convergence
 1. Select K points at random as cluster centroids (centers)
 - Now we formed K clusters
 2. For each data point, assign it to the closest center
 - Now we formed K clusters
 3. For each cluster, re-compute the centers
 - E.g., in the case of 2D points →
 - X: average over all x-axis points in the cluster
 - Y: average over all y-axis points in the cluster
 4. If the new centers are different from the old centers (previous iteration) → Go to Step 2
 - Otherwise, stop



K-Means – Single Node

17

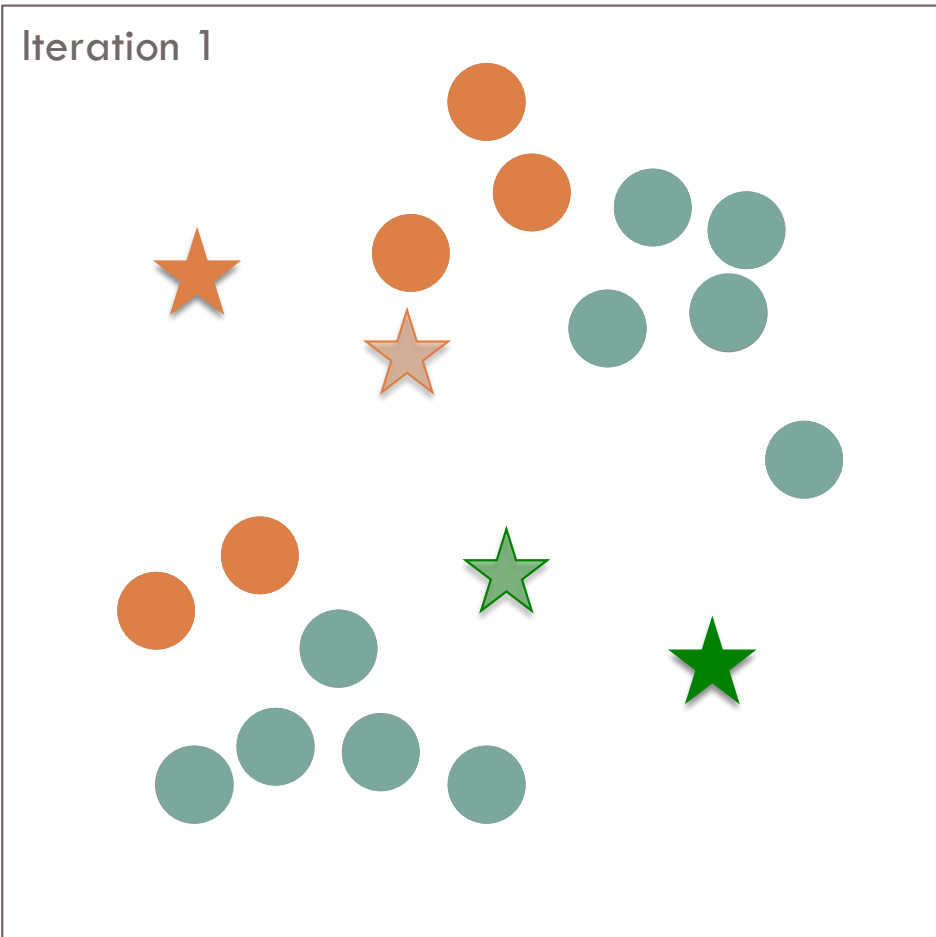
- Kmeans → Iterative algorithm until convergence
 1. Select K points at random as cluster centroids (centers)
 - Now we formed K clusters
 2. For each data point, assign it to the closest center
 3. For each cluster, re-compute the centers
 - E.g., in the case of 2D points →
 - X: average over all x-axis points in the cluster
 - Y: average over all y-axis points in the cluster
 4. If the new centers are different from the old centers (previous iteration) → Go to Step 2
 - Otherwise, stop



K-Means – Single Node

18

- Kmeans → Iterative algorithm until convergence
 1. Select K points at random as cluster centroids (centers)
 - 2. For each data point, assign it to the closest center
 - Now we formed K clusters
 3. For each cluster, re-compute the centers
 - E.g., in the case of 2D points →
 - X: average over all x-axis points in the cluster
 - Y: average over all y-axis points in the cluster
 4. If the new centers are significantly different from the old centers (previous iteration)
 - Set the new centers then Go to Step 2
 - Otherwise, stop



K-Means – Single Node

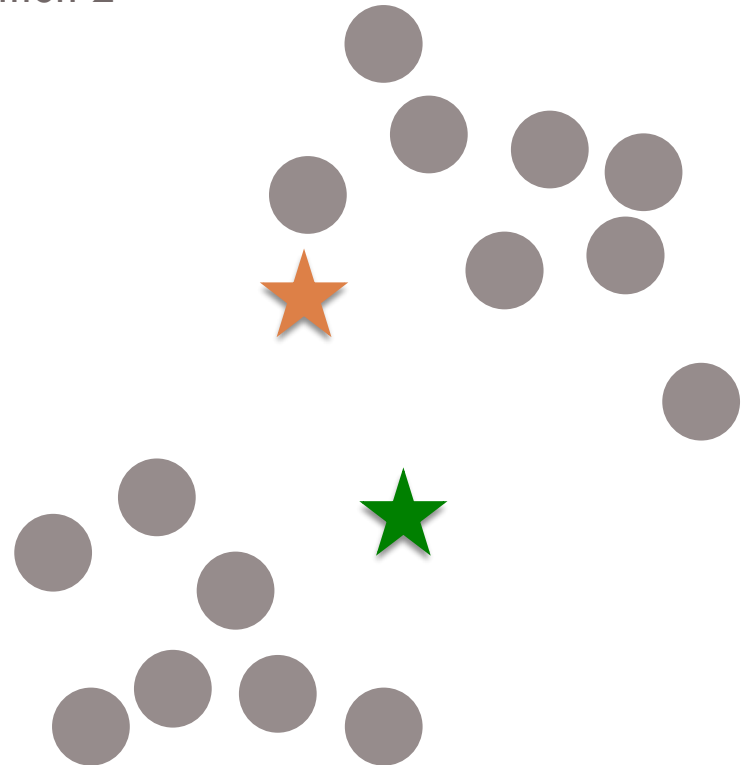
19

□ Kmeans → Iterative algorithm until convergence

1. Select K points at random as cluster centroids (centers)
2. For each data point, assign it to the closest cluster
 - Now we have formed K clusters
3. For each cluster, re-compute the centers
 - E.g., in the case of 2D points →
 - X: average over all x-axis points in the cluster
 - Y: average over all y-axis points in the cluster
4. If the new centers are significantly different from the old centers (previous iteration)
 - Set the new centers then Go to Step 2
 - Otherwise, stop

Iteration 2 Starts

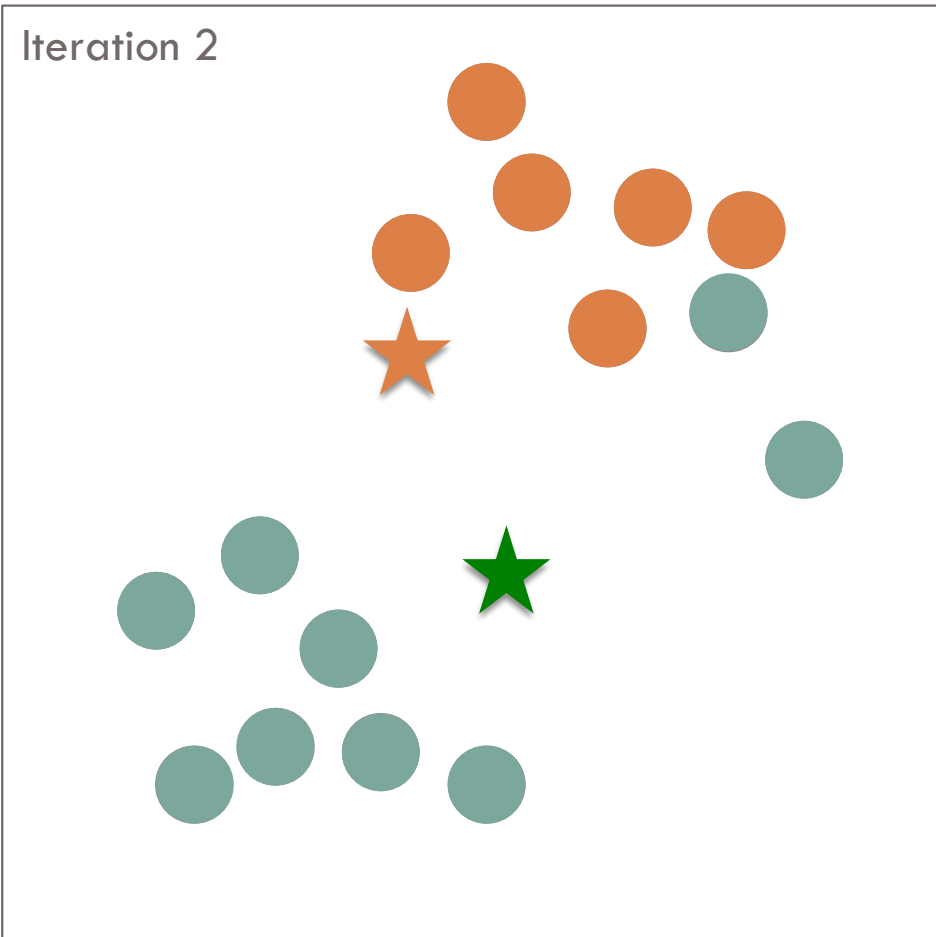
Iteration 2



K-Means – Single Node

20

- Kmeans → Iterative algorithm until convergence
 1. Select K points at random as cluster centroids (centers)
 2. For each data point, assign it to the closest center
 - Now we formed K clusters
 3. For each cluster, re-compute the centers
 - E.g., in the case of 2D points →
 - X: average over all x-axis points in the cluster
 - Y: average over all y-axis points in the cluster
 4. If the new centers are significantly different from the old centers (previous iteration)
 - Set the new centers then Go to Step 2
 - Otherwise, stop



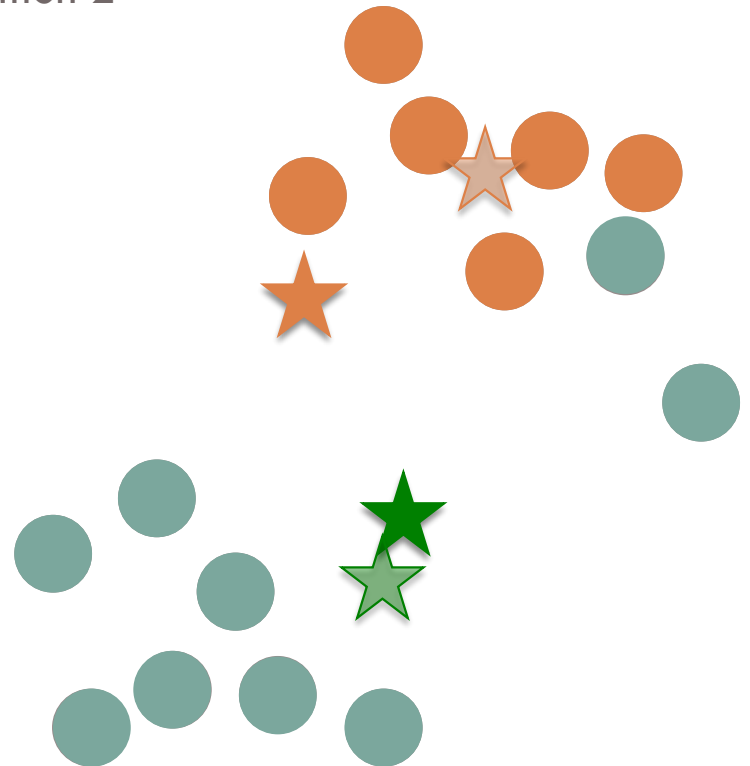
K-Means – Single Node

21

□ Kmeans → Iterative algorithm until convergence

1. Select K points at random as cluster centroids (centers)
2. For each data point, assign it to the closest center
 - Now we formed K clusters
3. For each cluster, re-compute the centers
 - E.g., in the case of 2D points →
 - X: average over all x-axis points in the cluster
 - Y: average over all y-axis points in the cluster
4. If the new centers are significantly different from the old centers (previous iteration)
 - Set the new centers then Go to Step 2
 - Otherwise, stop

Iteration 2

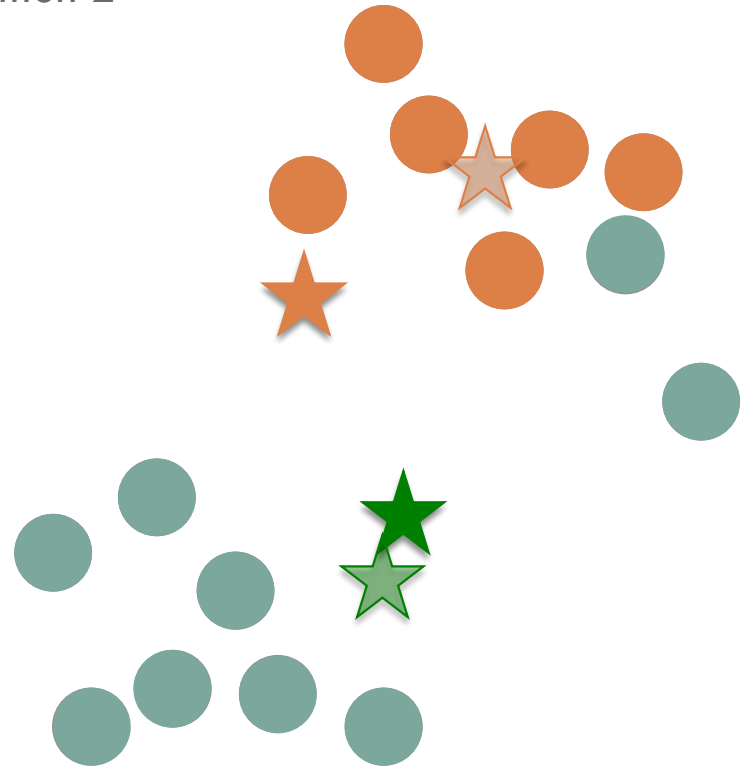


K-Means – Single Node

22

- Kmeans → Iterative algorithm until convergence
 1. Select K points at random as cluster centroids (centers)
 - 2. For each data point, assign it to the closest center
 - Now we formed K clusters
 3. For each cluster, re-compute the centers
 - E.g., in the case of 2D points →
 - X: average over all x-axis points in the cluster
 - Y: average over all y-axis points in the cluster
 4. If the new centers are significantly different from the old centers (previous iteration)
 - Set the new centers then Go to Step 2
 - Otherwise, stop

Iteration 2



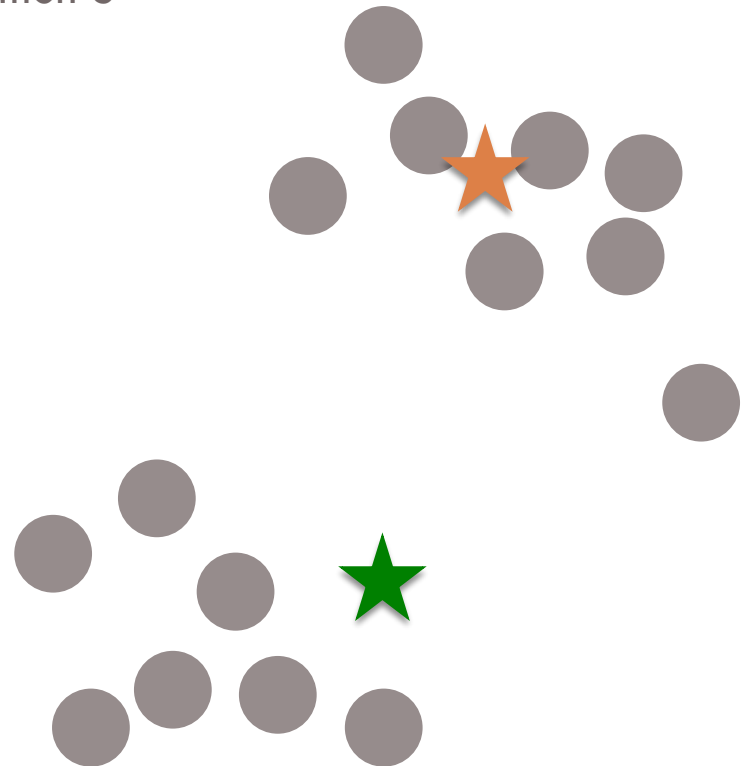
K-Means – Single Node

23

□ Kmeans → Iterative algorithm until convergence

1. Select K points at random as cluster centroids (centers)
2. For each data point, assign it to the closest cluster
 - Now we have K clusters
3. For each cluster, re-compute the centers
 - E.g., in the case of 2D points →
 - X: average over all x-axis points in the cluster
 - Y: average over all y-axis points in the cluster
4. If the new centers are significantly different from the old centers (previous iteration)
 - Set the new centers then Go to Step 2
 - Otherwise, stop

Iteration 3

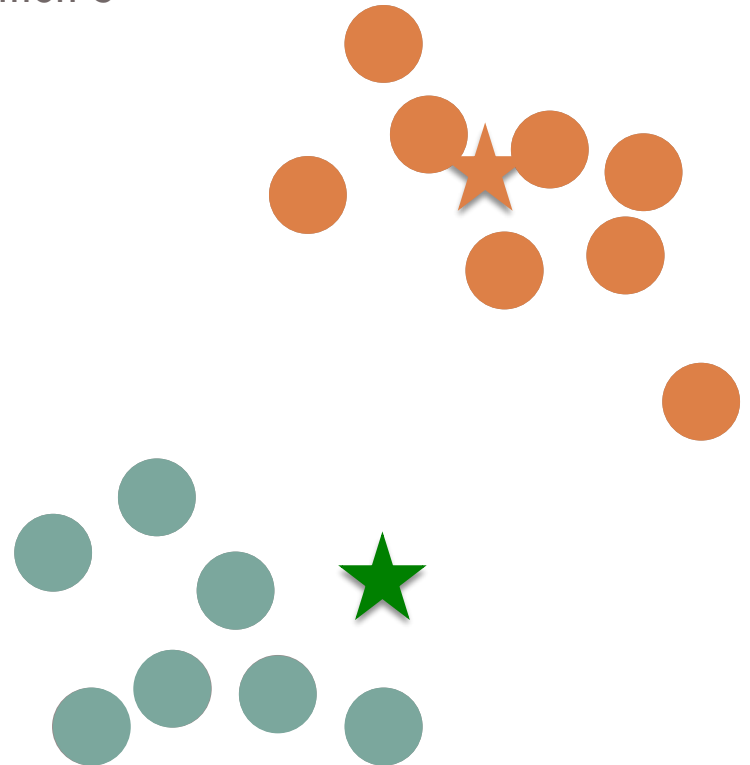


K-Means – Single Node

24

- Kmeans → Iterative algorithm until convergence
 1. Select K points at random as cluster centroids (centers)
 2. For each data point, assign it to the closest center
 - Now we formed K clusters
 3. For each cluster, re-compute the centers
 - E.g., in the case of 2D points →
 - X: average over all x-axis points in the cluster
 - Y: average over all y-axis points in the cluster
 4. If the new centers are significantly different from the old centers (previous iteration)
 - Set the new centers then Go to Step 2
 - Otherwise, stop

Iteration 3

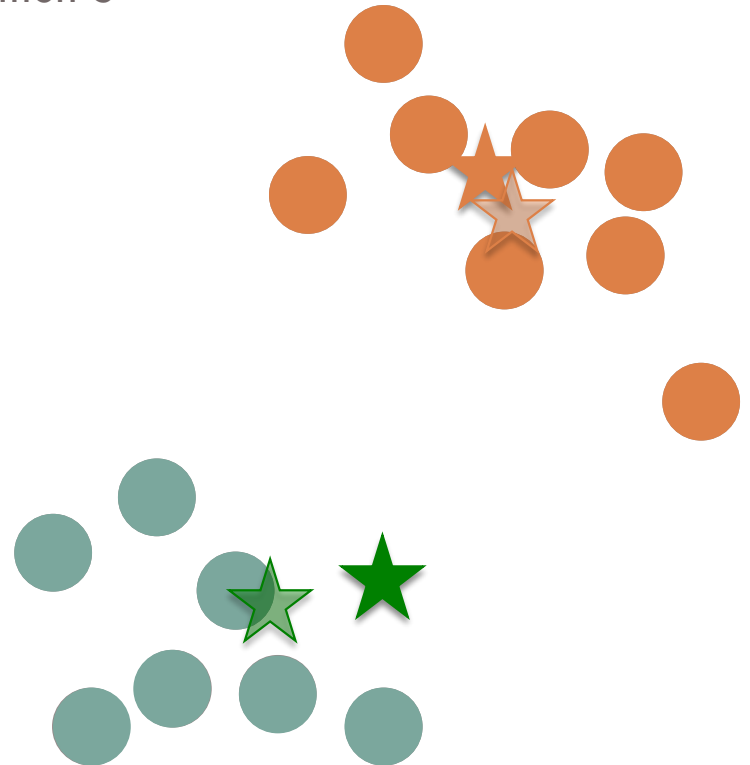


K-Means – Single Node

25

- Kmeans → Iterative algorithm until convergence
 1. Select K points at random as cluster centroids (centers)
 2. For each data point, assign it to the closest center
 - Now we formed K clusters
 3. For each cluster, re-compute the centers
 - E.g., in the case of 2D points →
 - X: average over all x-axis points in the cluster
 - Y: average over all y-axis points in the cluster
 4. If the new centers are significantly different from the old centers (previous iteration)
 - Set the new centers then Go to Step 2
 - Otherwise, stop

Iteration 3



K-Means – Single Node

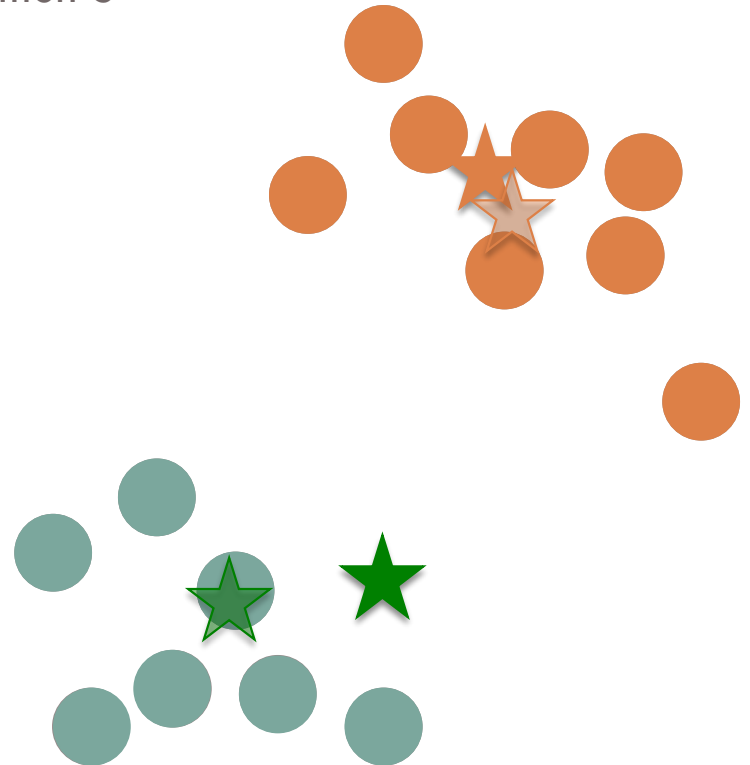
26

□ Kmeans → Iterative algorithm until convergence

1. Select K points at random as cluster centroids (centers)
2. For each data point, assign it to the closest center
 - Now we formed K clusters
3. For each cluster, re-compute the centers
 - E.g., in the case of 2D points →
 - X: average over all x-axis points in the cluster
 - Y: average over all y-axis points in the cluster
4. If the new centers are significantly different from the old centers (previous iteration)
 - Set the new centers then Go to Step 2
 - Otherwise, stop

Meets convergence criteria

Iteration 3

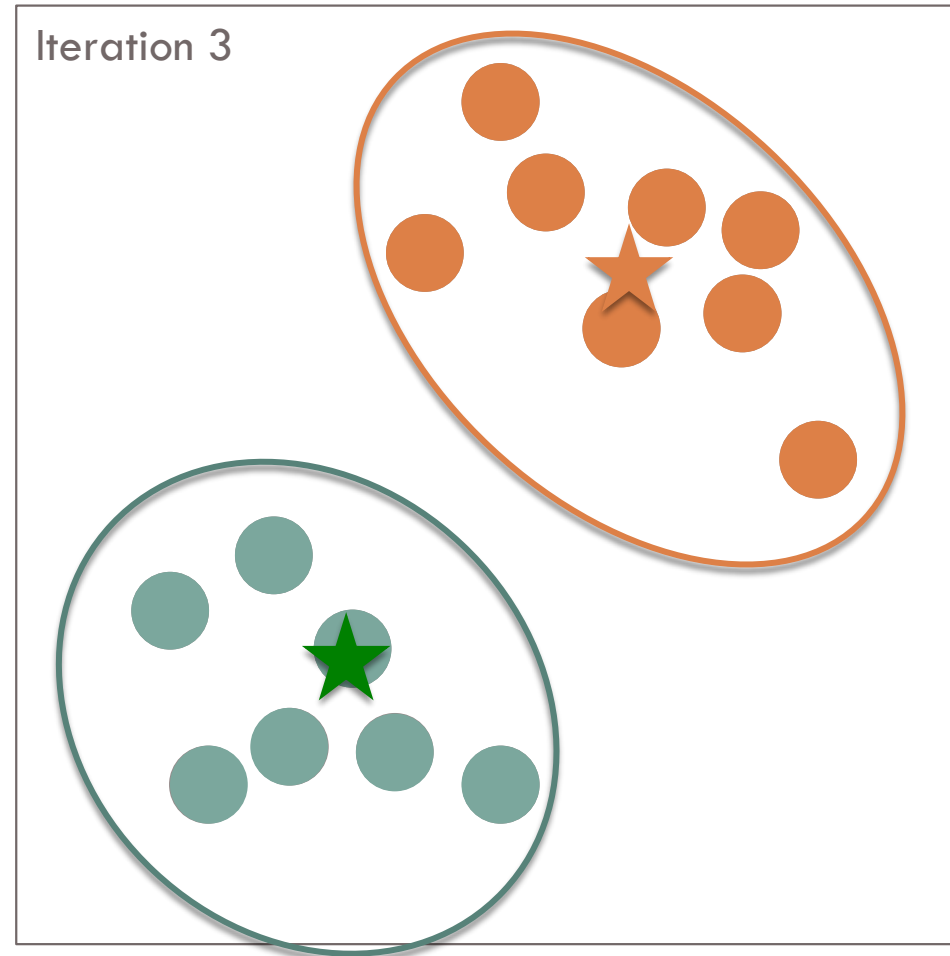


K-Means – Single Node

27

- Kmeans → Iterative algorithm until convergence
 1. Select K points at random as cluster centroids (centers)
 2. For each data point, assign it to the closest center
 - Now we formed K clusters
 3. For each cluster, re-compute the centers
 - E.g., in the case of 2D points →
 - X: average over all x-axis points in the cluster
 - Y: average over all y-axis points in the cluster
 4. If the new centers are significantly different from the old centers (previous iteration)
 - Set the new centers then Go to Step 2
 - Otherwise, stop

Meets convergence criteria



Distributed K-Means - MapReduce

28

□ Map-side

- Each map reads the K-centroids + one block from dataset
- Assign each point to the closest centroid
- Output <centroid, point>

□ Reduce-side

- Gets all points for a given centroid
- Re-compute a new centroid for this cluster
- Output: <new centroid>

□ Iteration Control

- Compare the old and new set of K-centroids
 - If similar → Stop
 - Else
 - If max iterations has reached → Stop
 - Else → Start another Map-Reduce Iteration

□ Use of combiners

- Similar to the reducer
- Computes for each centroid the local sums (and counts) of the assigned points
- Sends to the reducer <centroid, <partial sums>>

□ Use of single reducer

- Amount of data to reducers is very small
- Single reducer can tell whether any of the centers has changed or not
- Creates a single output file

Distributed K-Means - MapReduce

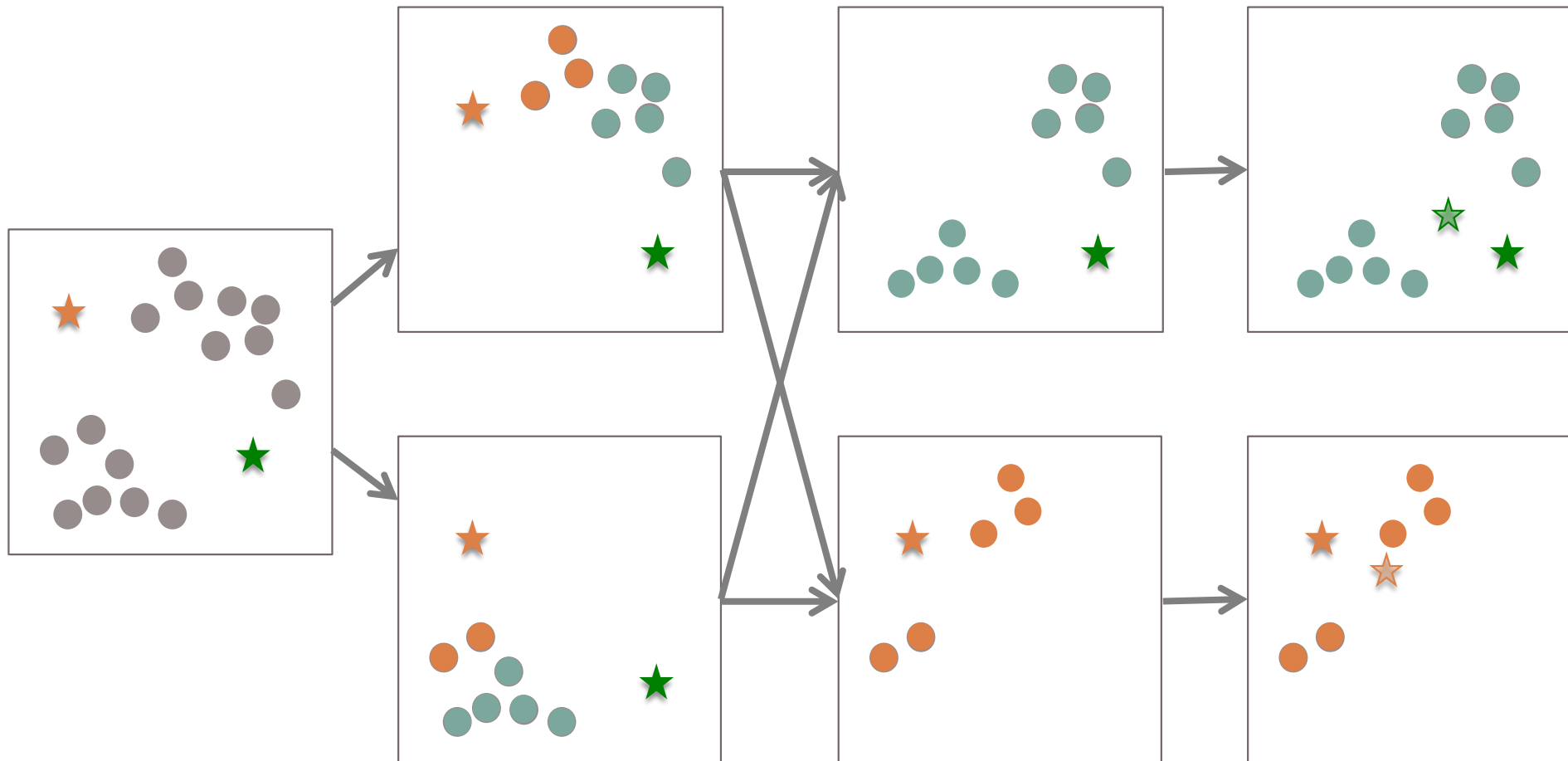
29

Splitting

Map

Shuffle/Sort
↑

Reduce



K-Means M/R One Iteration

KMeans - Example

30

```
1.  from pyspark import SparkConf, SparkContext
2.  from pyspark.mllib.clustering import KMeans, KMeansModel
3.  conf = SparkConf().setAppName("Testing Spark Commands")
4.  sc = SparkContext(conf = conf)
5.  data = [[1.0,1.0],[1.0,0.8],[-1.0,1.0],[-1.0,-1.0]]
6.  parsedData=sc.parallelize(data)
7.  kmeansModel = KMeans.train(parsedData, 2, maxIterations=10, runs=10, initializationMode="random")
8.  kmeansModel.predict([1.0, 1.0])
9.  #Output: 1
10. kmeansModel.predict([1.0, -2.0])
11. #Output: 0
12. # Save and load model
13. kmeansModel.save(sc, "KMeansModel")
14. model = KMeansModel.load(sc, "KmeansModel")
15. model.predict([1.0, 1.0])
16. model.predict([1.0, -2.0])
```

[Code same as kmeans_simple.py]

Collaborative Filtering: Recommendation

31

- Crowdsourced data
- Netflix Movie recommendation
- Amazon Product Recommendation

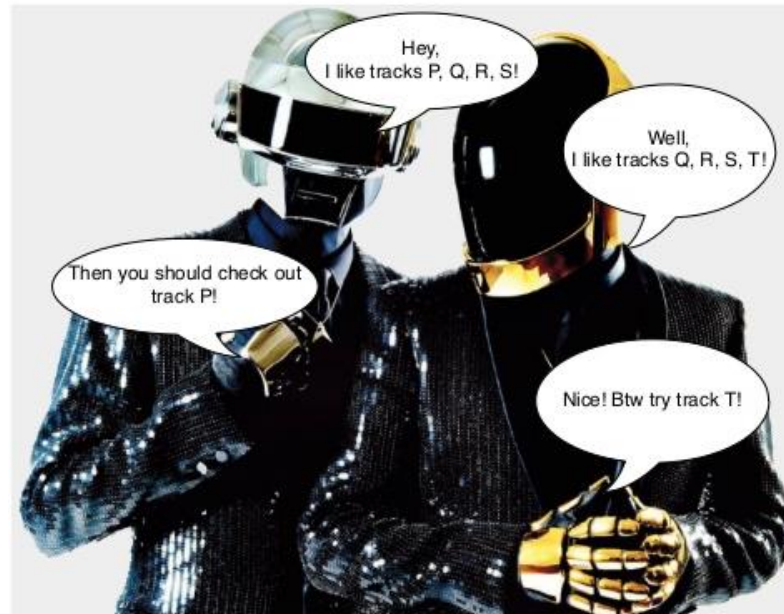



Image via Erik Bernhardsson







<http://www.slideshare.net/MrChrisJohnson/collaborative-filtering-with-spark>

Collaborative Filtering

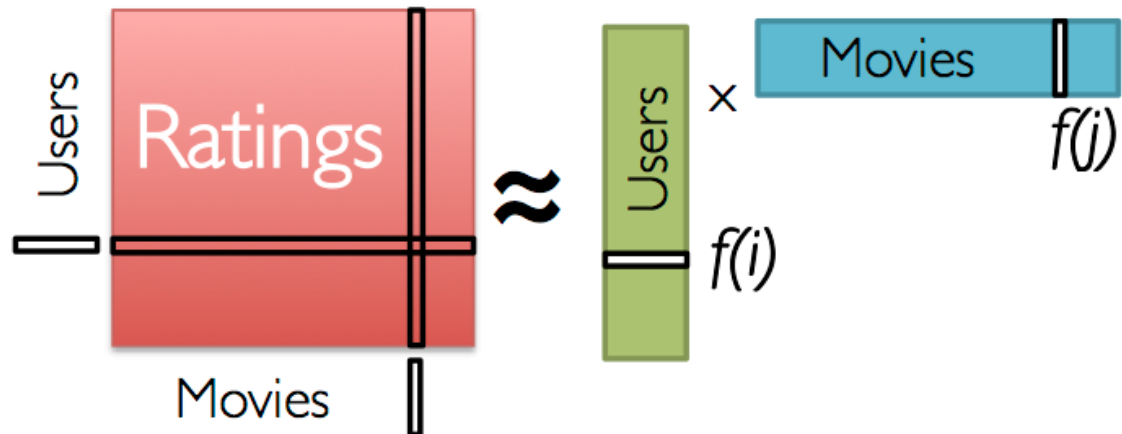
32

- Users rate some subset of items
- Goal: predict how users will rate new items



	★	★★★★★	
	★	★★★★	★★
	★★★★★		★
	★		★★
		★★★★	★★
	★★★★★	★★	

Low-Rank Matrix Factorization:



<https://databricks-training.s3.amazonaws.com/movie-recommendation-with-mllib.html>

- **ALS – Alternating Least Squares**

ALS - Example

33

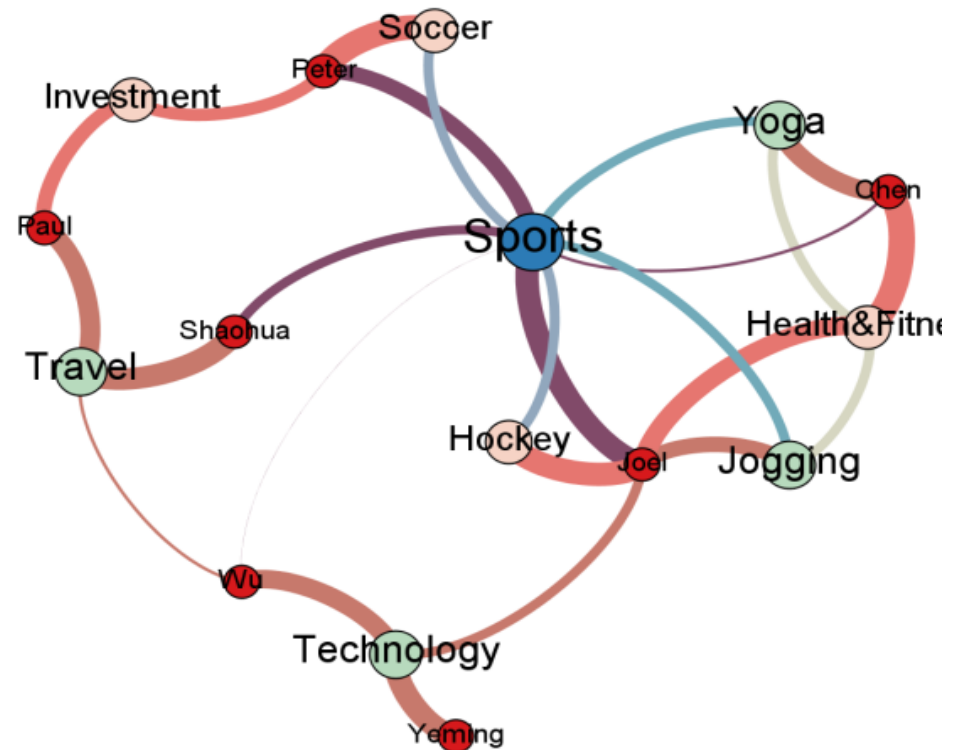
```
1.  from pyspark.mllib.recommendation import ALS, MatrixFactorizationModel, Rating
2.  from pyspark import SparkConf, SparkContext
3.  conf = SparkConf().setAppName("Testing Spark Commands")
4.  sc = SparkContext(conf = conf)
5.  r1 = (1, 1, 1.0)
6.  r2 = (1, 2, 2.0)
7.  r3 = (2, 1, 2.0)
8.  ratings = sc.parallelize([r1, r2, r3])
9.  model = ALS.train(ratings, 10, 10)
10. # Evaluate the model on training data
11. testdata = sc.parallelize([(2,2),(1,1)])
12. predictions = model.predictAll(testdata)
13. predictions.collect()
14. # Save and load model
15. model.save(sc, "myCollaborativeFilter")
16. sameModel = MatrixFactorizationModel.load(sc, "myCollaborativeFilter")
17. #Try using the loaded model
18. sameModel.predictAll(testdata).collect()
□ #Output: [Rating(user=1, product=1, rating=1.0000230136195141), Rating(user=2, product=2,
rating=0.96355230064639663)]
□ [these numbers can change]
```

Graph Processing

34

The Graph phenomenon!

- Web Graph (Google)
 - Social Graph (Facebook)
 - Follower Graph (Twitter)
 - Interest Graph (Pinterest!)
 - Music/Location/Food Graphs
-
- Available in Scala and Java
 - (Not yet in PySpark)



Graph Representation

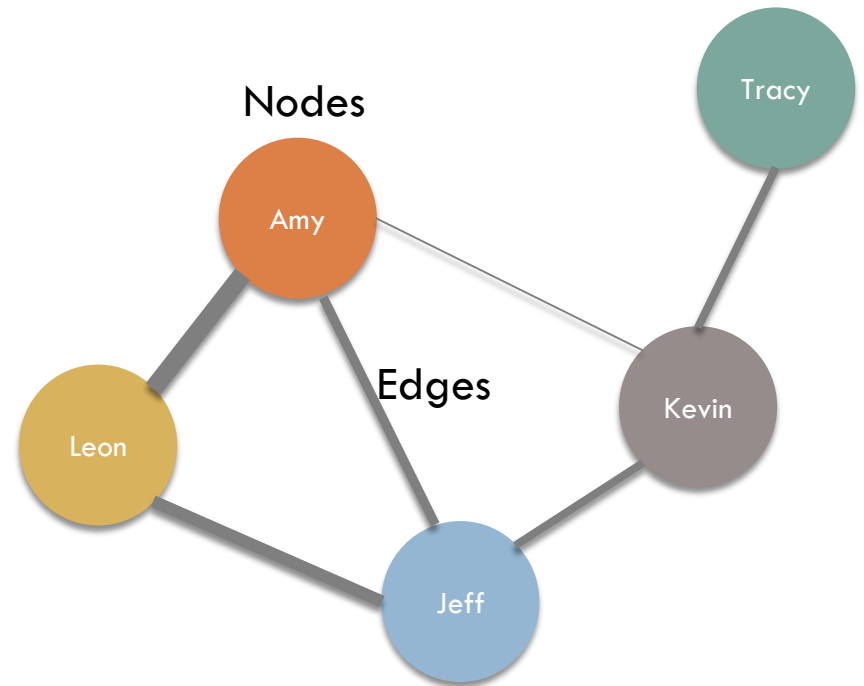
35

	Friend				
User	Amy	Kevin	Jeff	Tracy	Leon
Amy		0.8	1.7		4
Kevin	0.8		1.3	1	
Jeff	1.7	1.3			2.2
Tracy		1			
Leon	4		2.2		

Matrix

Key	Value
Amy	(Kevin,0.8), (Jeff, 1.7) (Leon,4)
Kevin	(Amy,0.8), (Jeff,1.3), (Tracy,1)
Jeff	(Amy, 1.7), (Kevin,1.3), (Leon,2.2)
Tracy	(Kevin,1)
Leon	(Amy,4), (Jeff,2.2)

Adjacency List

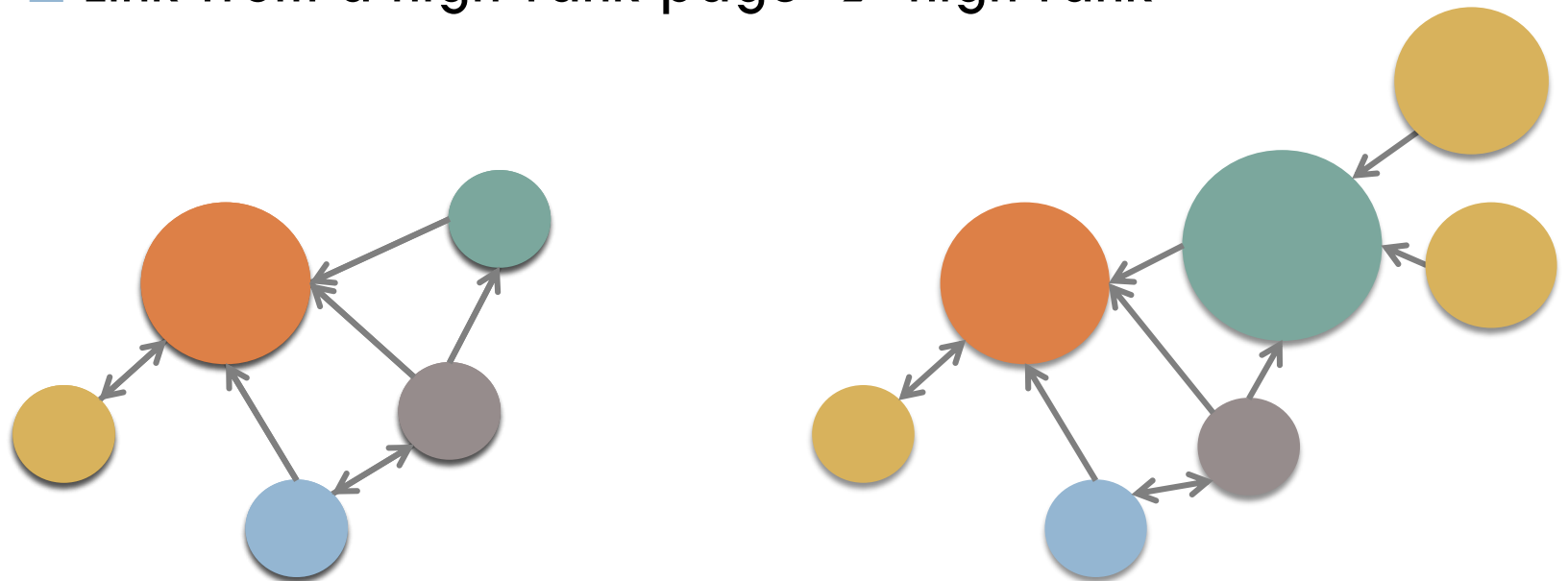


Graph (undirected)

PageRank Algorithm

36

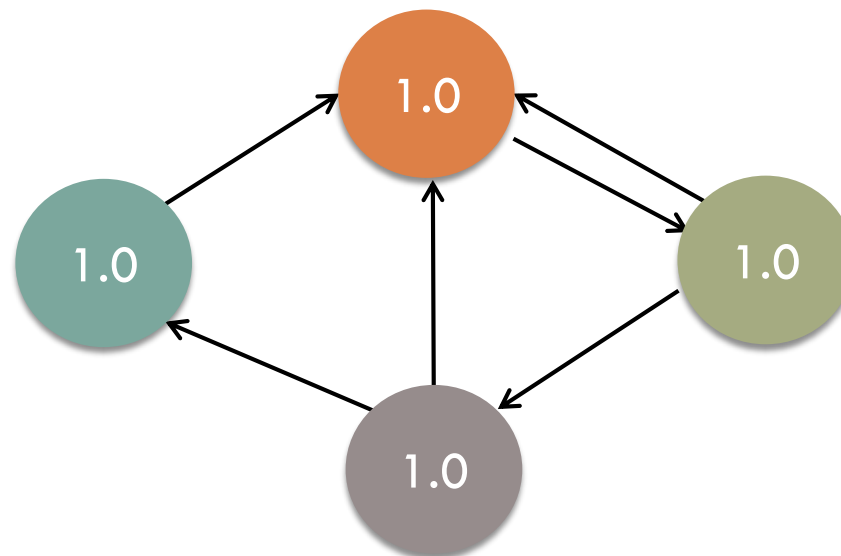
- PageRank gives web pages a ranking score on links from other pages
 - ▣ Links from many pages → high rank
 - ▣ Link from a high-rank page → high rank



Graph (directed)

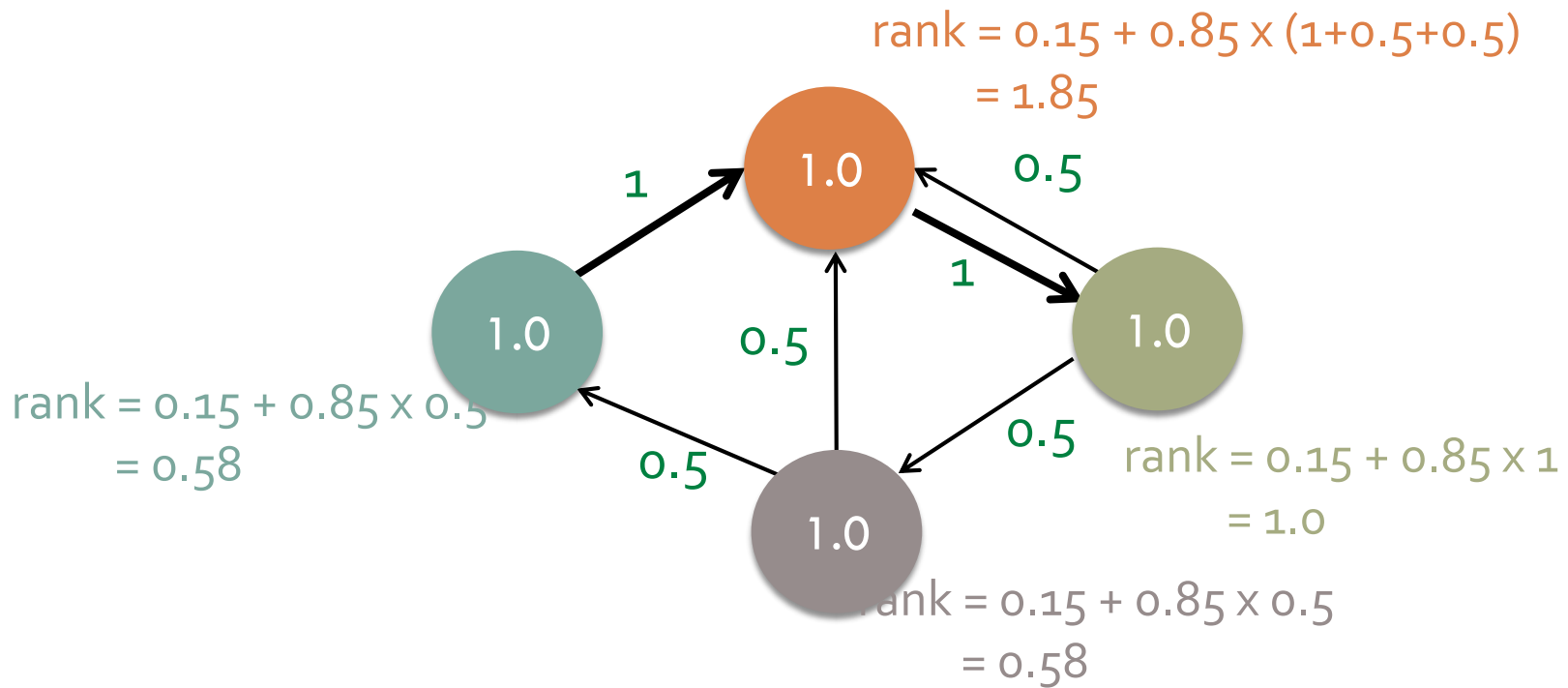
PageRank Explained – Initial State

1. Start each page at a rank of 1



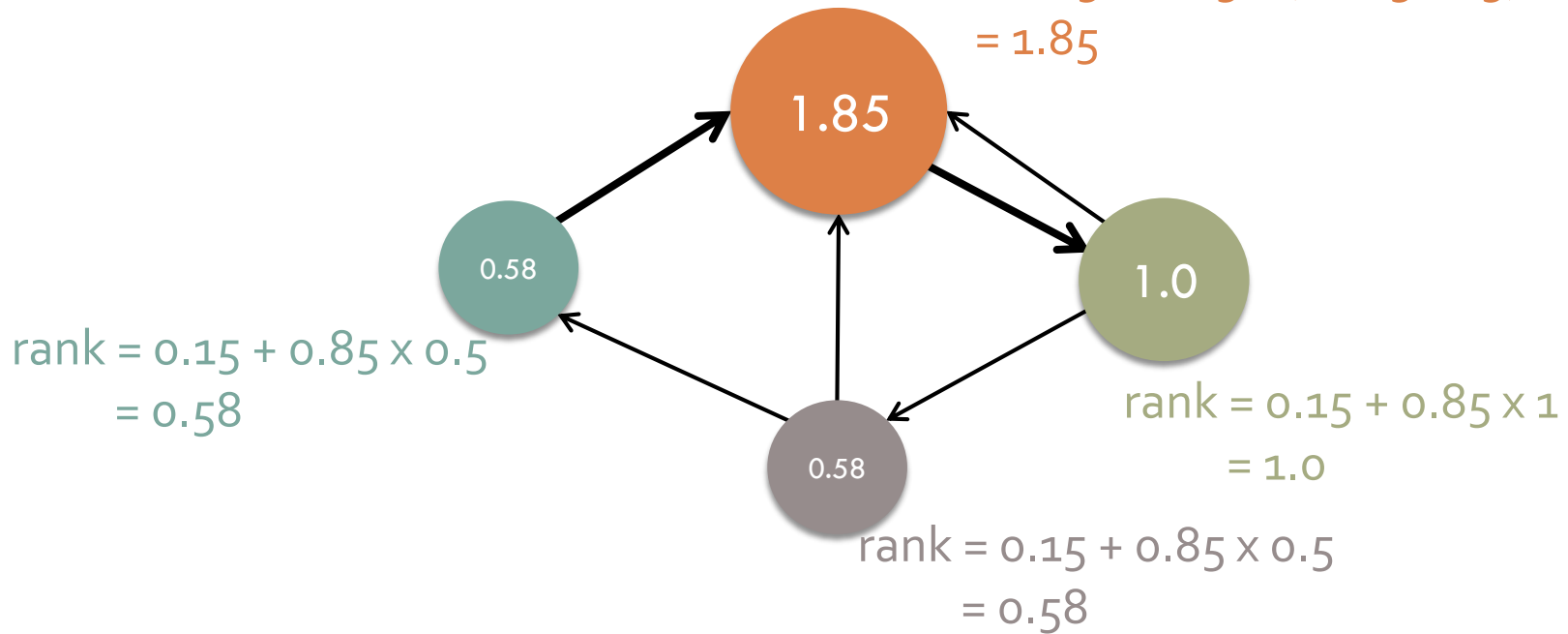
PageRank Explained – Iteration 1

1. Start each page at a rank of 1
2. On each iteration, have page **p** contribute $\text{contrib}_p = \text{rank}_p / \text{neighbors}_p$ to its neighbors



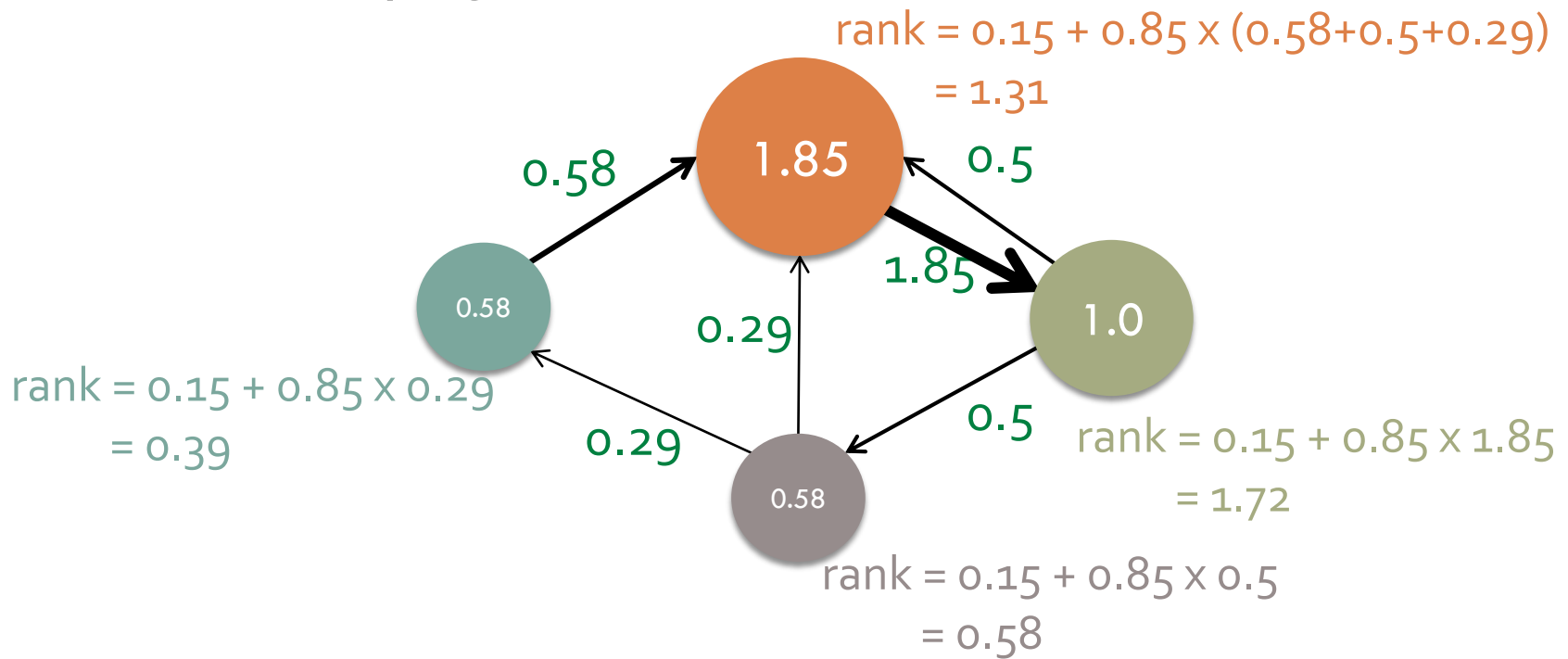
PageRank Explained – Iteration 1

1. Start each page at a rank of 1
2. On each iteration, have page p contribute $\text{contrib}_p = \text{rank}_p / \text{neighbors}_p$ to its neighbors
3. Set each page's rank to $0.15 + 0.85 \times \text{contribs}$



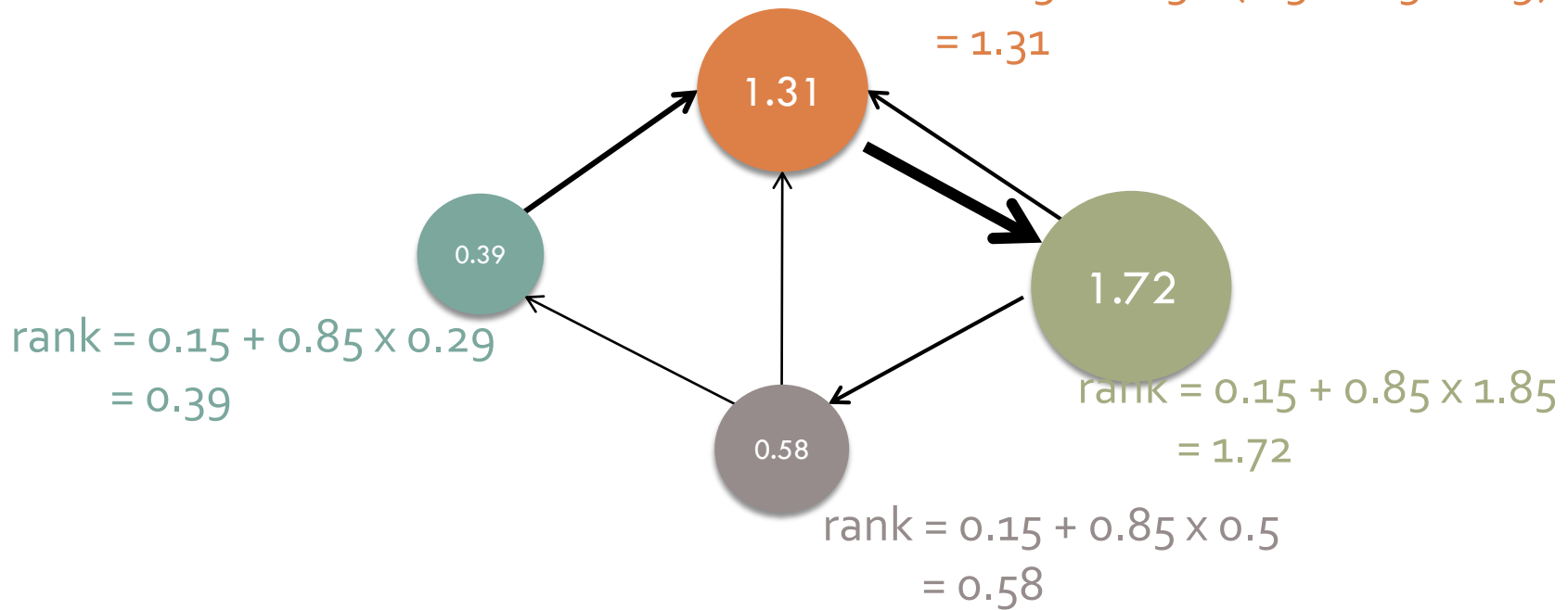
PageRank Explained – Iteration 2

1. Start each page at a rank of 1
2. On each iteration, have page **p** contribute $\text{contrib}_p = \text{rank}_p / \text{neighbors}_p$ to its neighbors
3. Set each page's rank to $0.15 + 0.85 \times \text{contribs}$



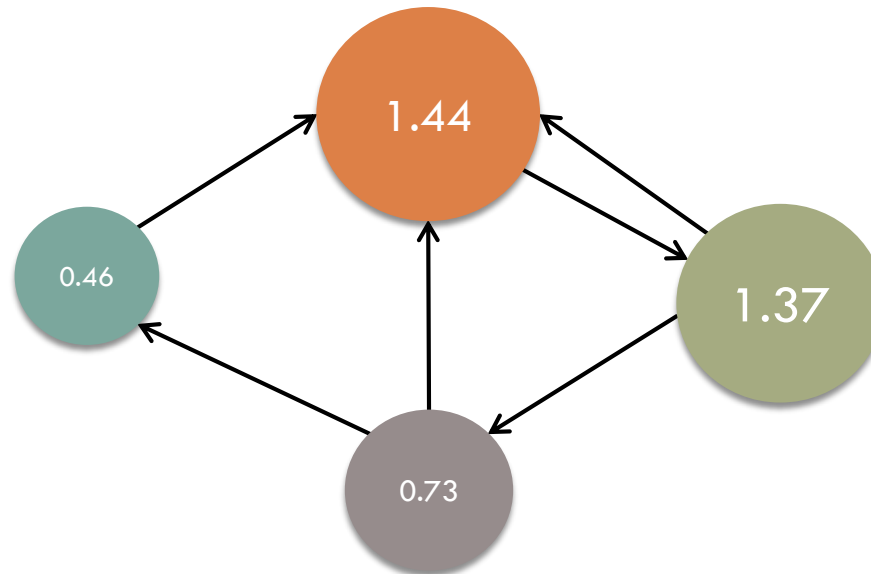
PageRank Explained – Iteration 2

1. Start each page at a rank of 1
2. On each iteration, have page p contribute $\text{contrib}_p = \text{rank}_p / \text{neighbors}_p$ to its neighbors
3. Set each page's rank to $0.15 + 0.85 \times \text{contribs}$
 $\text{rank} = 0.15 + 0.85 \times (0.58 + 0.5 + 0.29)$
 $= 1.31$



PageRank Explained – Final State

1. Start each page at a rank of 1
2. On each iteration, have page p contribute $\text{contrib}_p = \text{rank}_p / \text{neighbors}_p$ to its neighbors
3. Set each page's rank to $0.15 + 0.85 \times \text{contribs}$



PageRank using GraphFrames

43

- ❑ `from graphframes import *`
- ❑ `v = sqlContext.createDataFrame([
 ("a", "Alice", 34),
 ("b", "Bob", 36),
 ("c", "Charlie", 30),
], ["id", "name", "age"])`
- ❑ `# Create an Edge DataFrame with "src" and "dst" columns`
- ❑ `e = sqlContext.createDataFrame([
 ("a", "b", "friend"),
 ("b", "c", "follow"),
 ("c", "b", "follow"),
], ["src", "dst", "relationship"])`
- ❑ `# Create a GraphFrame`
- ❑ `g = GraphFrame(v, e)`
- ❑ `results = g.pageRank(resetProbability=0.15, tol=0.01)`
- ❑ `results.vertices().show()`

PageRank – Example in Scala using GraphX API

44

□ spark-shell [Command to Open Scala Spark shell]

```
1.      import org.apache.spark.graphx._
2.      val graph = GraphLoader.edgeListFile(sc, "/user/root/followers.txt")
3.      // Run PageRank
4.      val ranks = graph.pageRank(0.0001).vertices
5.      // Join the ranks with the usernames
6.      val users = sc.textFile("/user/root/users.txt").map { line =>
    val fields = line.split(",")
    (fields(0).toLong, fields(1))
  }
7.      val ranksByUsername = users.join(ranks).map {
    case (id, (username, rank)) => (username, rank)
  }
8.      // Print the result
9.      println(ranksByUsername.collect().mkString("\n"))
```

□ Output:

- (justinbieber,0.15)
- (matei_zaharia,0.7013599933629602)
- (ladygaga,1.390049198216498)
- (BarackObama,1.4588814096664682)
- (jeresig,0.9993442038507723)
- (odersky,1.2973176314422592)

□ [Copy and paste commands into “spark-shell” console; Copy followers.txt and users.txt from D2L into hdfs]

PageRank Implementation

□ PageRank with Pig

- ▣ Using PageRank to Detect Anomalies and Fraud in Healthcare (Hortonworks Blog Post)

- [\(PART1\) http://hortonworks.com/blog/using-pagerank-detect-anomalies-fraud-healthcare/](http://hortonworks.com/blog/using-pagerank-detect-anomalies-fraud-healthcare/)
- [\(PART2\) http://hortonworks.com/blog/using-pagerank-to-detect-anomalies-and-fraud-in-healthcare-part2/](http://hortonworks.com/blog/using-pagerank-to-detect-anomalies-and-fraud-in-healthcare-part2/)
- [\(PART3\) http://hortonworks.com/blog/using-pagerank-to-detect-anomalies-and-fraud-in-healthcare-part3/](http://hortonworks.com/blog/using-pagerank-to-detect-anomalies-and-fraud-in-healthcare-part3/)

□ PageRank with Spark

Getting Started with Spark

46

- Install Spark Standalone on Linux
- Cloudera/HDP Distributions
 - ▣ http://www.cloudera.com/content/cloudera/en/downloads/quickstart_vms/cdh-5-3-x.html

Spark Documentation

Spark 1.0.0 Documentation: <http://spark.apache.org/docs/latest/>

Spark Programming Guide (Scala,

Python): <http://spark.apache.org/docs/latest/programming-guide.html#overview>

Spark Cassandra Connector - DataStax ([github](#))

Spark HBase - lighting Spark with HBase ([link](#))

Spark MLlib Documentation ([link](#))

Spark GraphX Documentation ([link](#))

(Spark) Spark Cluster Mode Overview ([link](#))

(Spark) Running Spark on EC2 ([link](#))

(Cloudera) Pig is Flying: Apache Pig on Apache Spark ([link](#))