

# Determining the $k$ in $k$ -means with MapReduce

Thibault Debatty, Pietro Michiardi,  
Wim Mees & Olivier Thonnard



# Clustering & k-means

- Clustering

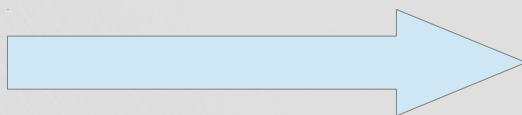
- K-means

[Stuart P. Lloyd. Least squares quantization in pcm. IEEE Transactions on Information Theory, 28:129–137, 1982.]

- 1982 (a great year!)
- But still largely used
- Drawbacks (amongst others):
  - Local minimum
  - K is a parameter!

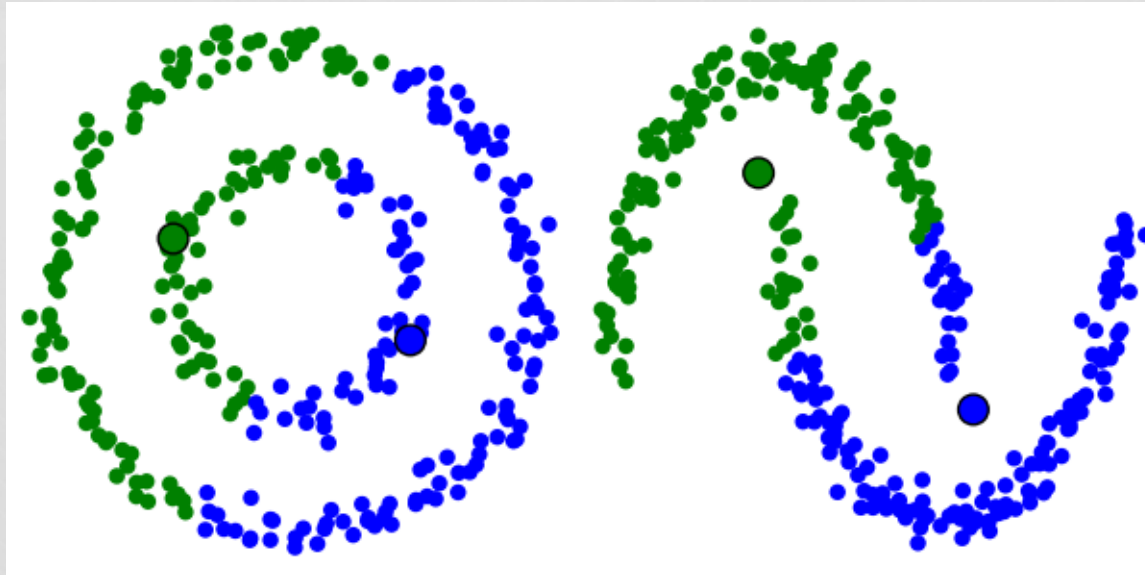
# Clustering & k-means

- Determine k:
  - VERY difficult  
[Anil K Jain. Data Clustering : 50 Years Beyond K-Means. Pattern Recognition Letters, 2009]
  - Using cluster evaluation metrics:  
Dunn's index, Elbow, Silhouette, “jump method” (based on information theory), “Gap statistic”, ...

  $O(k^2)$

# G-means

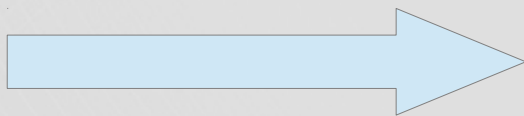
- G-means  
[Greg Hamerly and Charles Elkan. Learning the k in k-means. In Neural Information Processing Systems. MIT Press, 2003]
- K-means : points in each cluster are spherically distributed around the center



source: scikit-learn

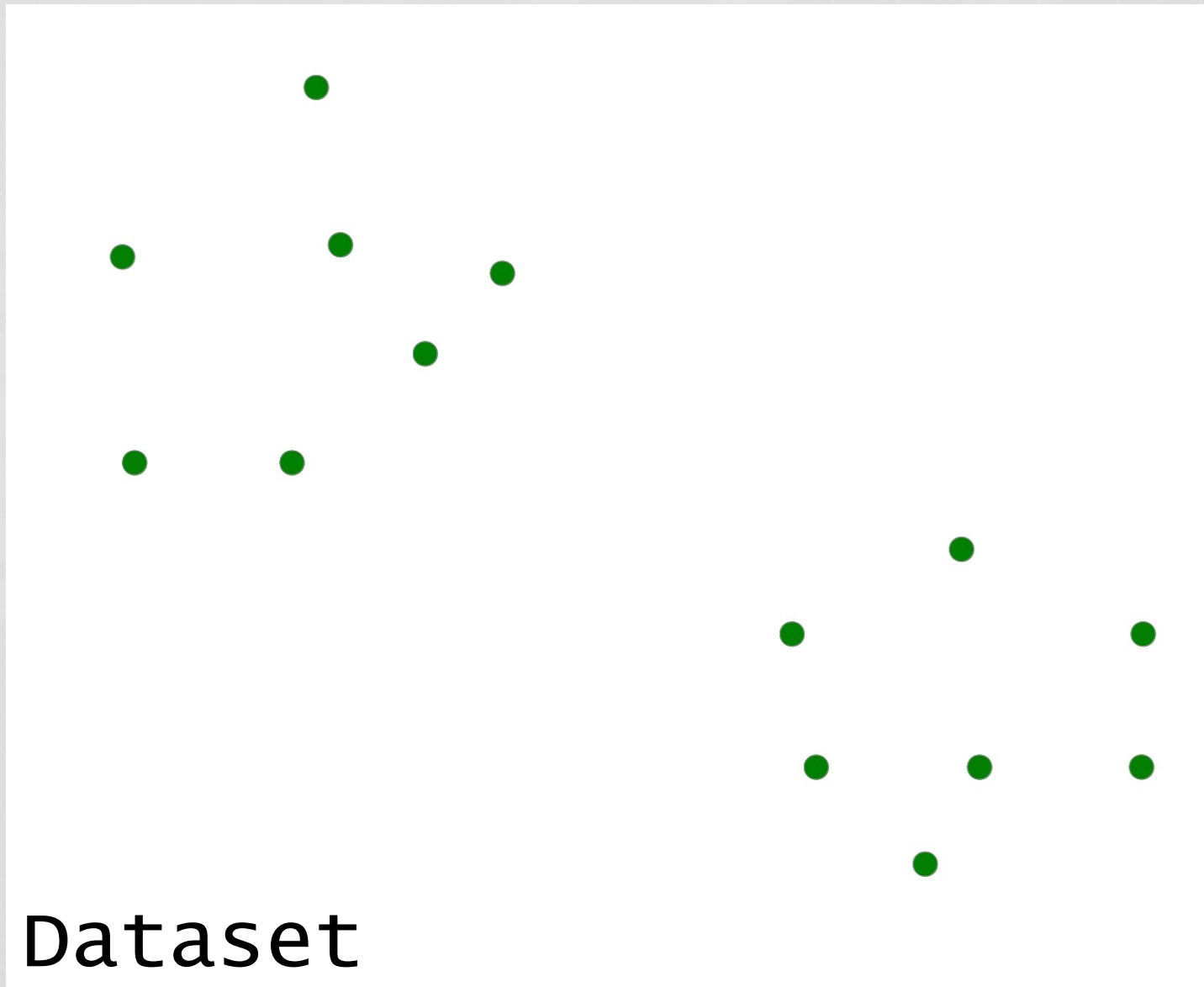
# G-means

- G-means  
[Greg Hamerly and Charles Elkan. Learning the k in k-means. In Neural Information Processing Systems. MIT Press, 2003]
- K-means : points in each cluster are spherically distributed around the center

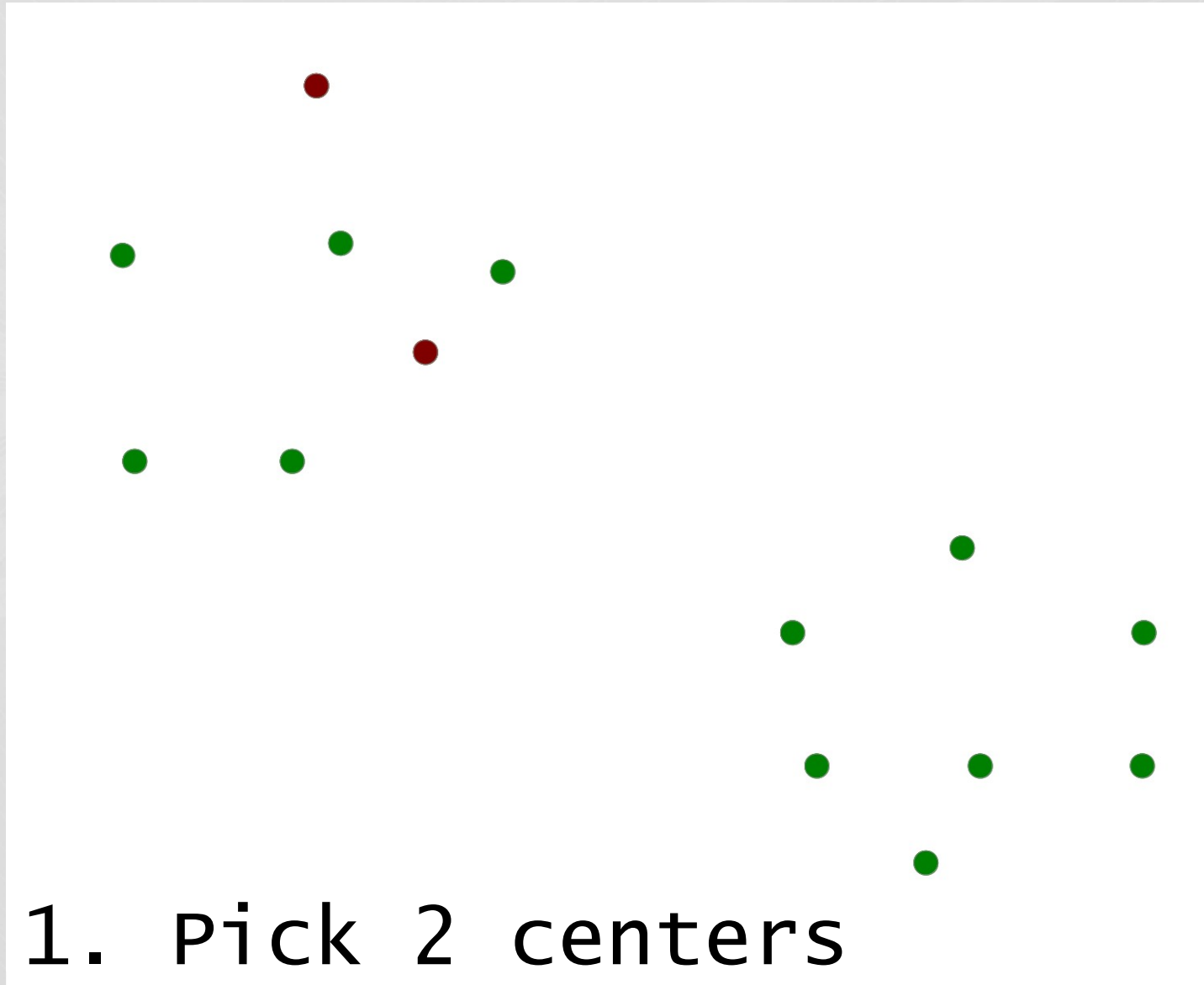


normality test &  
recursion

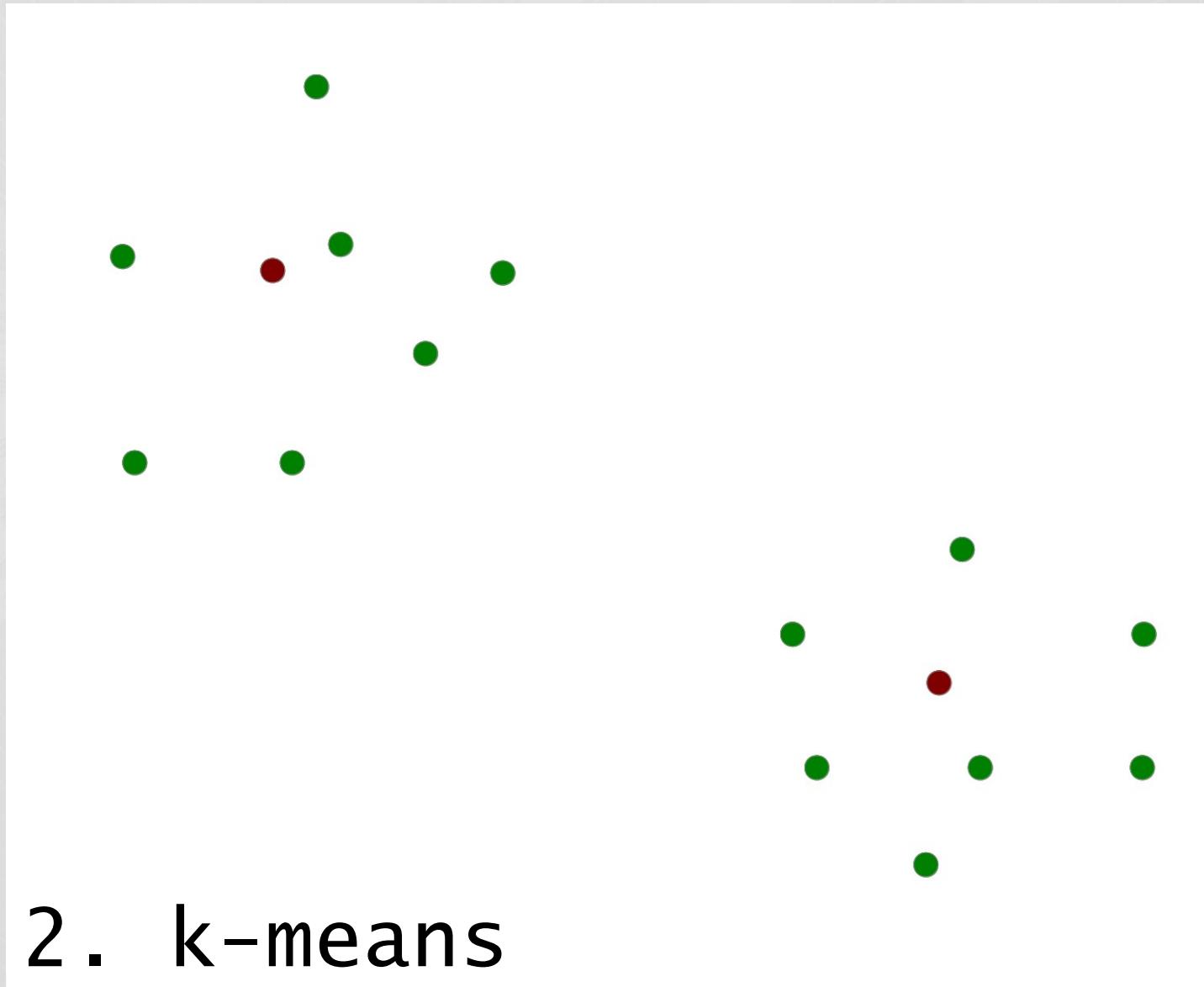
# G-means



# G-means

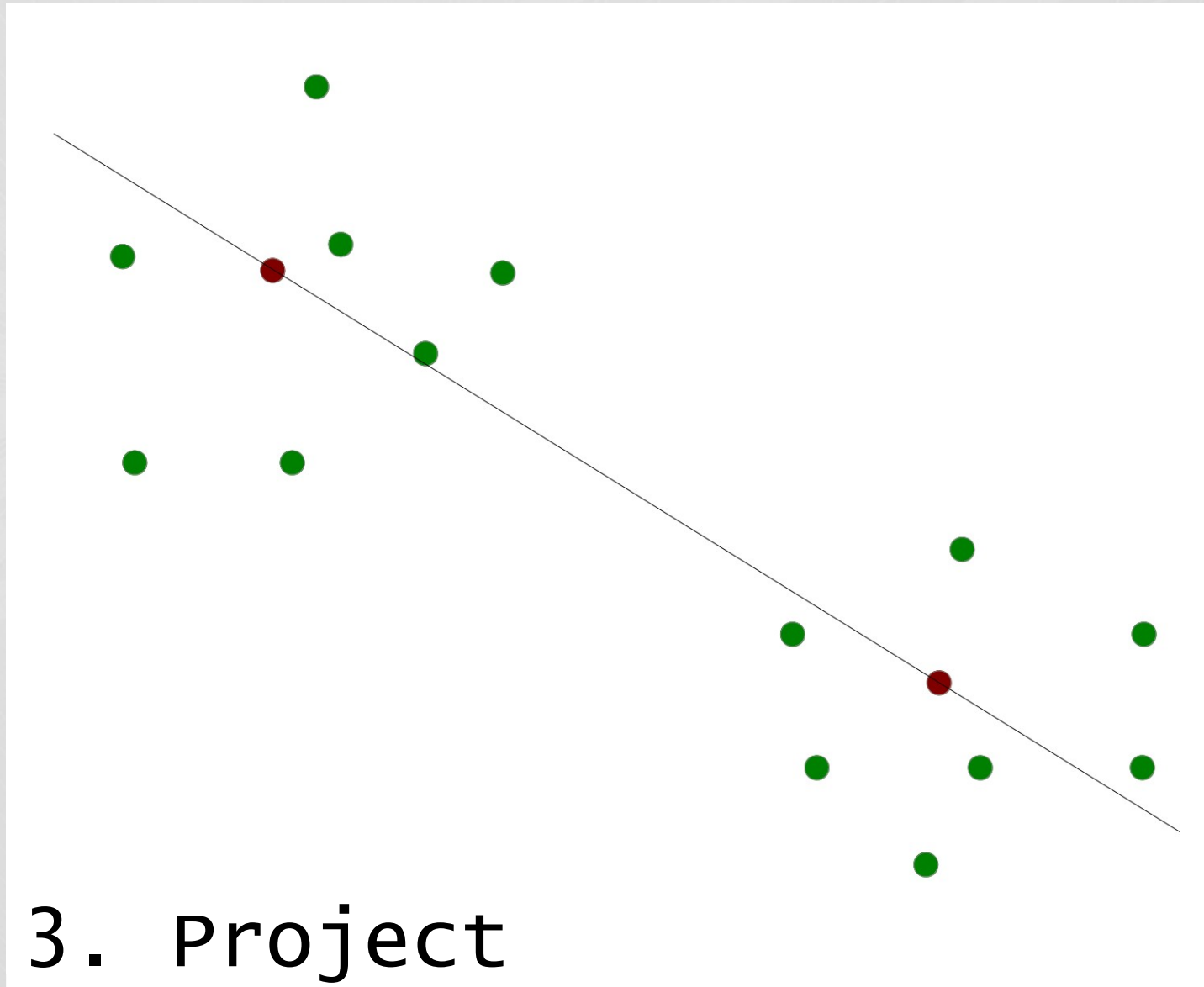


# G-means

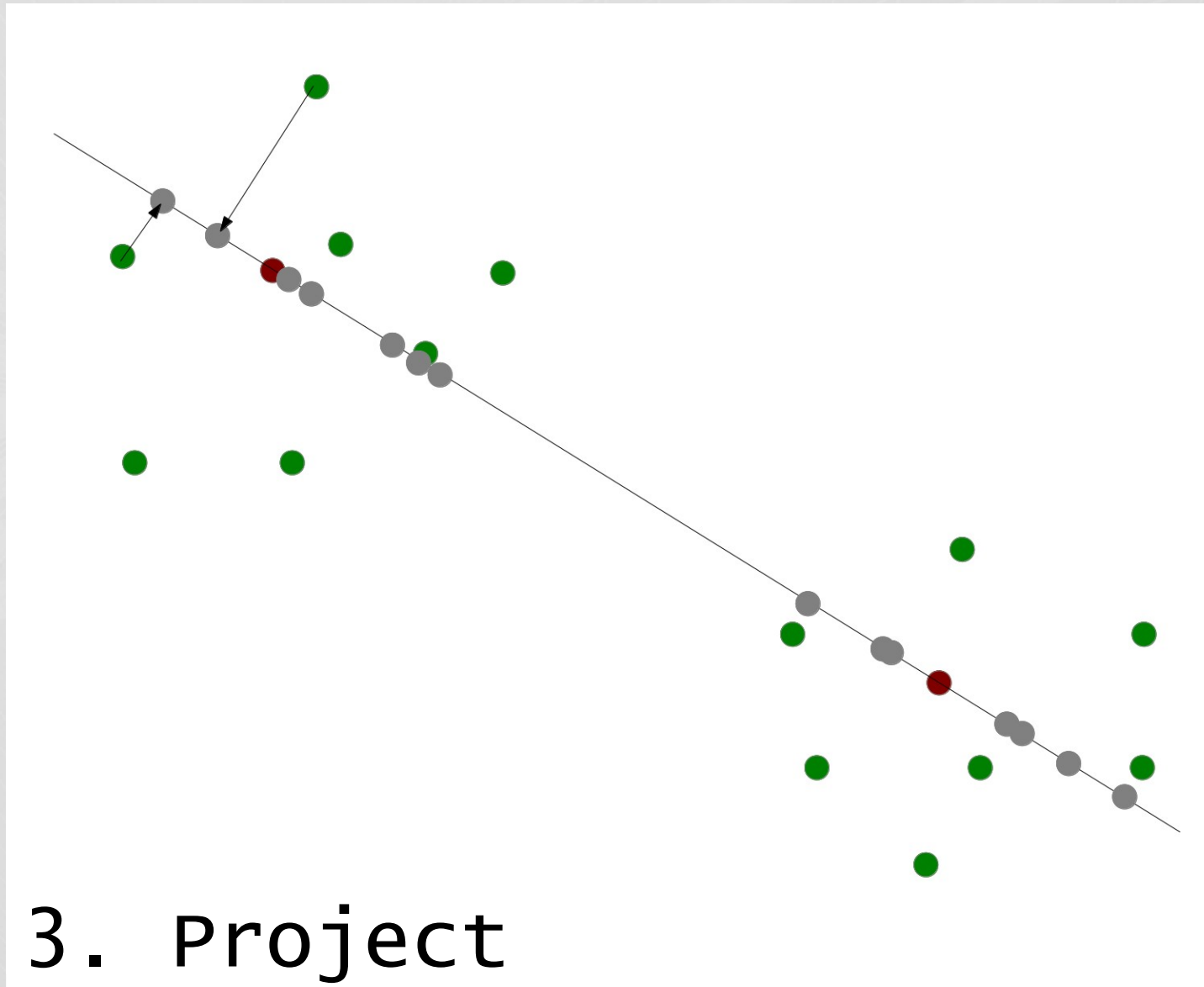


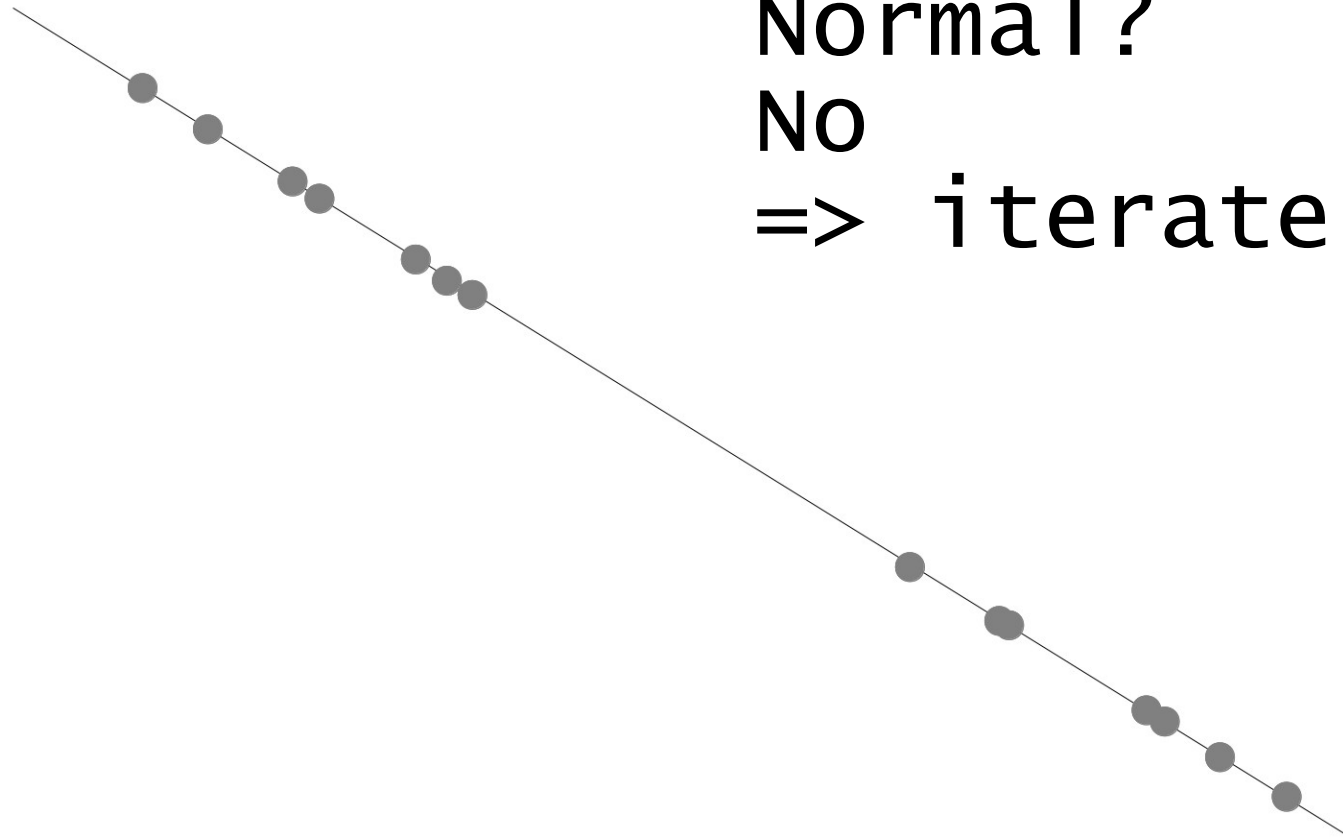


# G-means

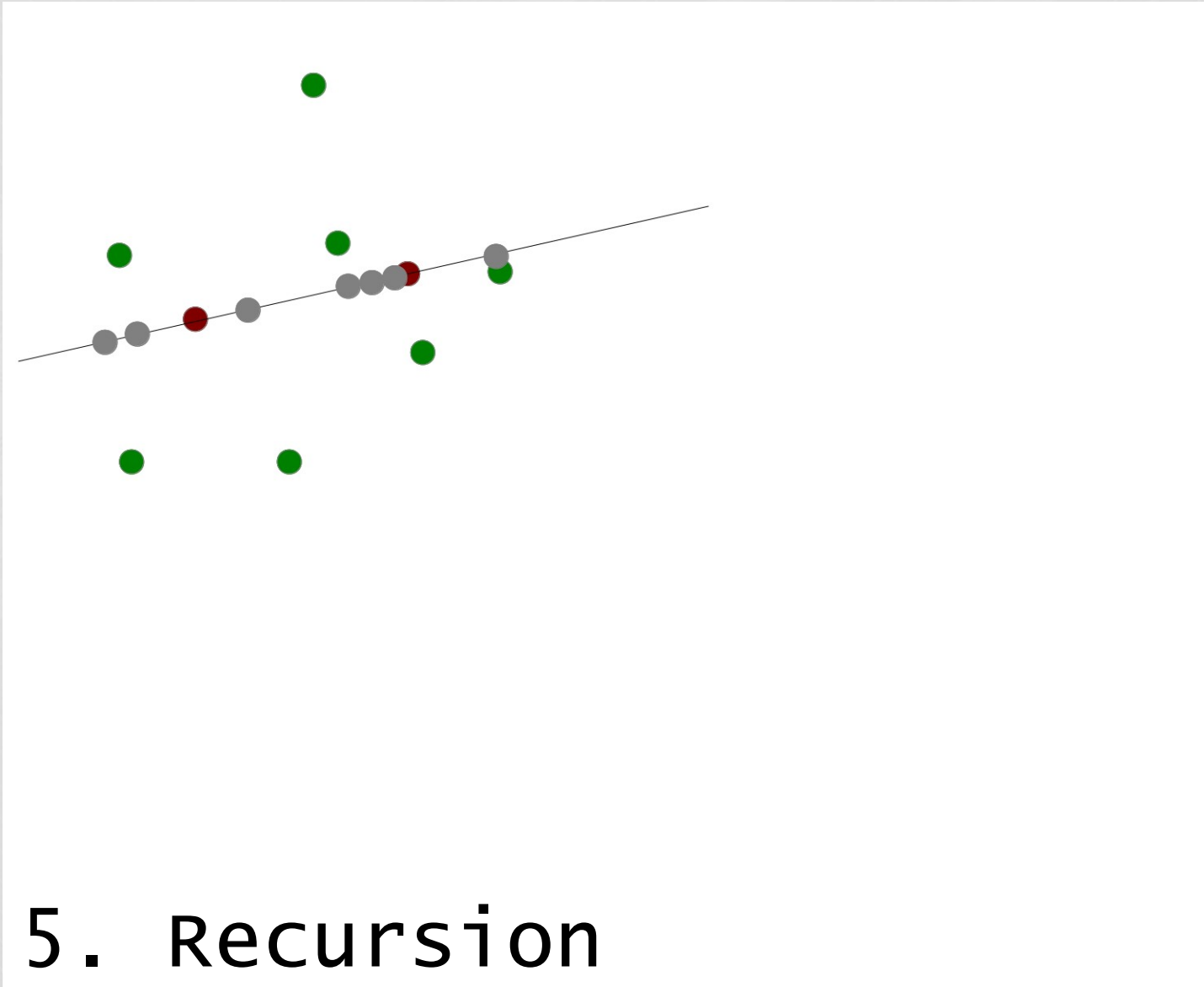


# G-means





## 4. Normality test



# MapReduce G-means

- Challenges:
  1. Reduce I/O operations
  2. Reduce number of jobs
  3. Maximize parallelism
  4. Limit memory usage

# MapReduce G-means

- Challenges:
  - 1. Reduce I/O operations**
  2. Reduce number of jobs
  3. Maximize parallelism
  4. Limit memory usage

# MapReduce G-means

## 2. Reduce number of jobs

```
PickInitialCenters  
while Not ClusteringCompleted do  
    KMeans  
    KMeansAndFindNewCenters  
    TestClusters  
end while
```

# MapReduce G-means

3. Maximize  
parallelism

4. Limit memory  
usage

```
TestClusters
```

```
Map(key, point)  
  Find cluster  
  Find vector  
  Project point on vector  
  Emit(cluster, projection)  
end procedure
```

```
Reduce(cluster, projections)  
  Build a vector  
  ADtest(vector)  
  if normal then  
    Mark cluster  
  end if  
end procedure
```



# MapReduce G-means

3. Maximize parallelism

4. Limit memory usage (risk of crash)

```
TestClusters
```

```
Map(key, point)
```

```
  Find cluster
```

```
  Find vector
```

```
  Project point on vector
```

```
  Emit(cluster, projection)
```

```
end procedure
```

```
Reduce(cluster, projections)
```

```
  Build a vector
```

```
  ADtest(vector)
```

```
  if normal then
```

```
    Mark cluster
```

```
  end if
```

```
end procedure
```

**Bottleneck**

# MapReduce G-means

## TestFewClusters

```
Map(key, point)  
  Find cluster  
  Find vector  
  Project point on vector  
  Add projection to list  
end procedure
```

```
Close()  
  For each list do  
    Build a vector  
     $A^2 = \text{ADtest}(\text{vector})$   
    Emit(cluster,  $A^2$ )  
  End for each  
end procedure
```

In memory combiner

## TestClusters

```
Map(key, point)  
  Find cluster  
  Find vector  
  Project point on vector  
  Emit(cluster, projection)  
end procedure
```

```
Reduce(cluster, projections)  
  Build a vector  
   $\text{ADtest}(\text{vector})$   
  if normal then  
    Mark cluster  
  end if  
end procedure
```

# MapReduce G-means

TestFewClusters

TestClusters

Map

r

FT

Pro

Ad

end

#clusters > #reducers

&

Estimated required memory < Java heap

Close()

For each *list* do

Build a *vector*

$A2 = ADtest(vector)$

Emit(*cluster*,  $A2$ )

End for each

end procedure

Reduce(*cluster*, *projections*)

Build a *vector*

$ADtest(vector)$

if normal then

Mark *cluster*

end if

end procedure

# MapReduce G-means

TestFewClusters

TestClusters

Map

#clusters > #reducers

FT

Pro

Ad

&

end

Estimated required memory < Java heap

Close()

For each

Build a vector

A2 = ADtest(vector)

Emit(cluster, A2)

End for each

end procedure

Experimentally:  
64 Bytes / point

reduce(cluster, projections)

Build a vector

ADtest(vector)

if normal then

Mark cluster

end if

end procedure

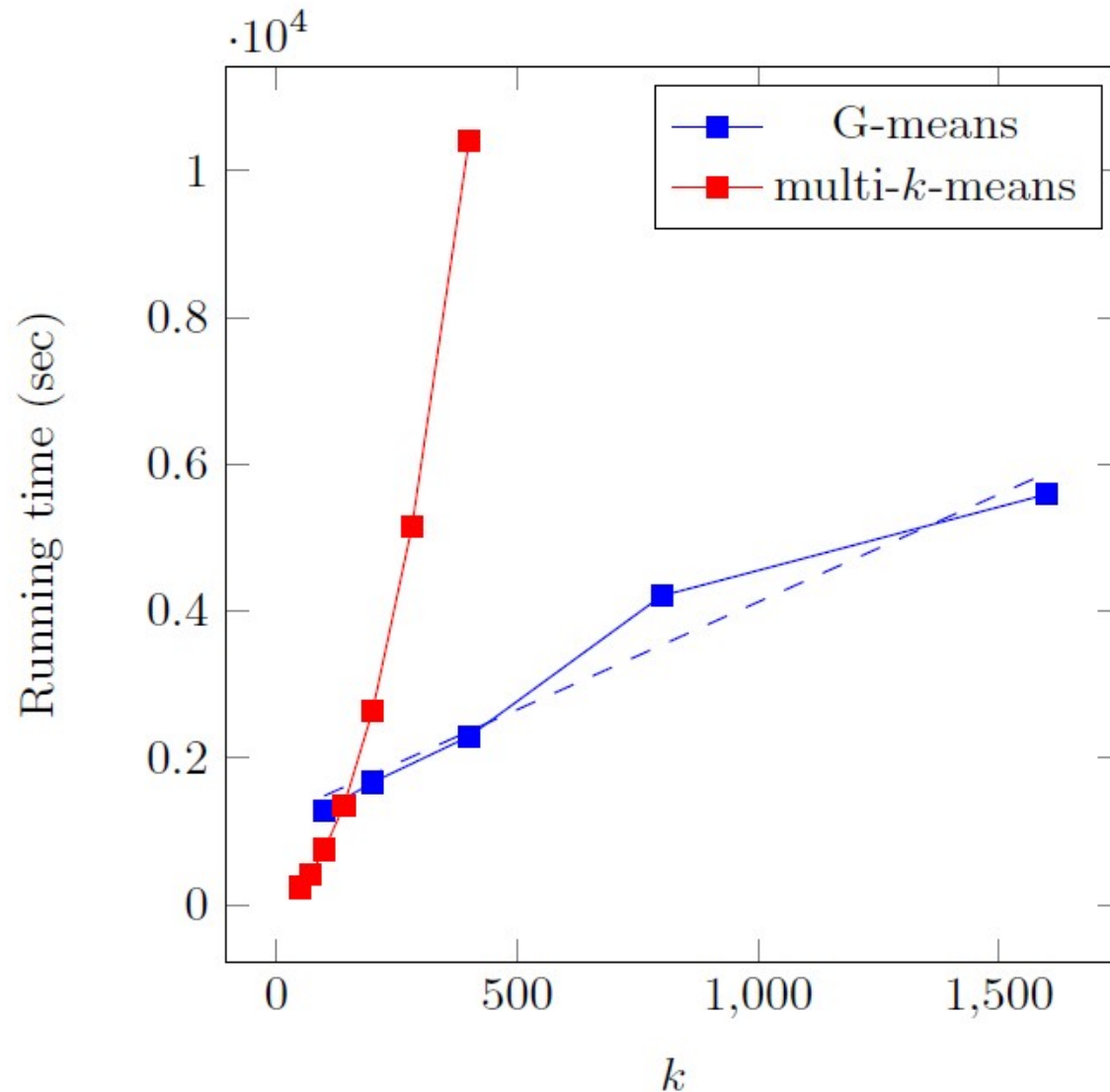
# Comparison

	MR multi-k-means	MR G-means
Speed		
	all possible values of k in a single job	
Quality		

# Comparison

	MR multi-k-means	MR G-means
Speed	$O(nk^2)$ computations	$O(nk)$ computations
		But: <ul style="list-style-type: none"><li>• more iterations</li><li>• more dataset reads</li><li>• <math>\log_2(k)</math></li></ul>
Quality		New centers added if and where needed
		But: tends to overestimate k!

# Experimental results : Speed




- Hadoop
- Synthetic dataset
- 10M points in  $R^{10}$
- Euclidean distance
- 8 machines

# Experimental results : Quality

k	100	200	400
$k_{\text{found}}$	150	279	639
<b>Within Cluster Sum of Square</b> (less is better)			
MR G-means	3.34	3.33	3.23
multi-k-means (with same k)	3.71	3.6	3.39

x ~1.5



- Hadoop
- Synthetic dataset
- 10M points in  $R^{10}$
- Euclidean distance
- 8 machines



# Conclusions & future work...

- MapReduce algorithm to determine  $k$
- Running time proportional to  $k$
- Future:
  - Overestimation of  $k$
  - Test on real data
  - Test scalability
  - Reduce I/O (using Spark)
  - Consider skewed data
  - Consider impact of machine failure

# Thank you!