

# Data Structures for Fast K-means

The Auton Lab  
Carnegie Mellon University

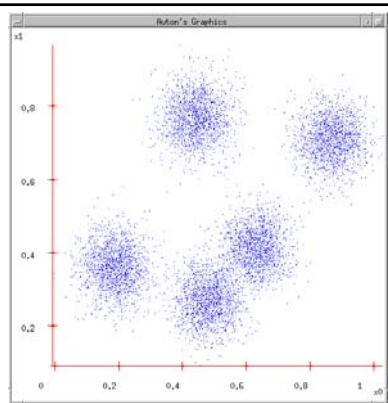


[www.autonlab.org](http://www.autonlab.org)

What does k-means do?

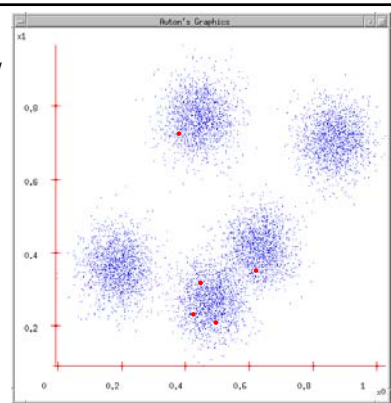
## K-means

1. Ask user how many clusters they'd like.  
(e.g.  $k=5$ )



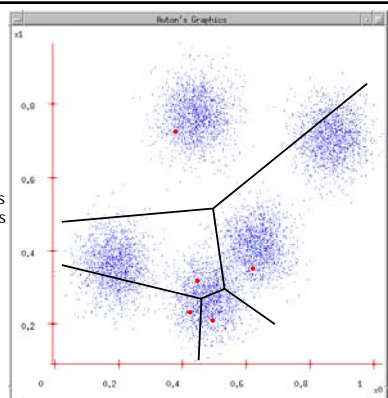
## K-means

1. Ask user how many clusters they'd like.  
(e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations



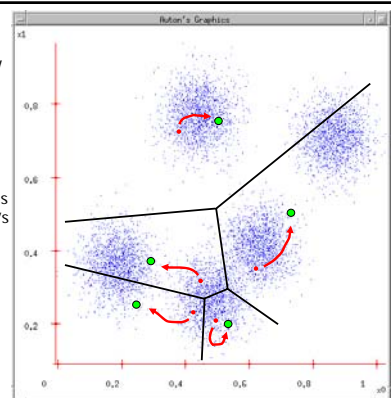
## K-means

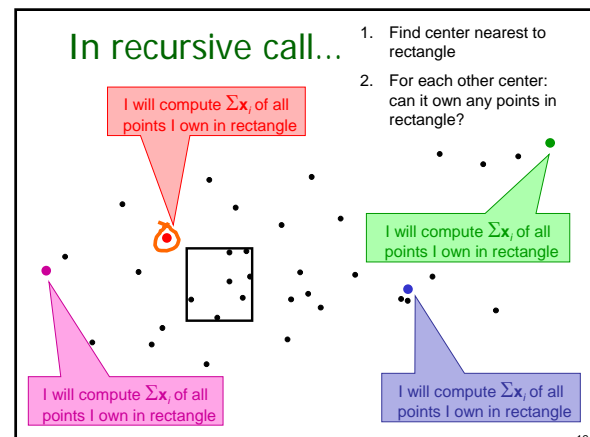
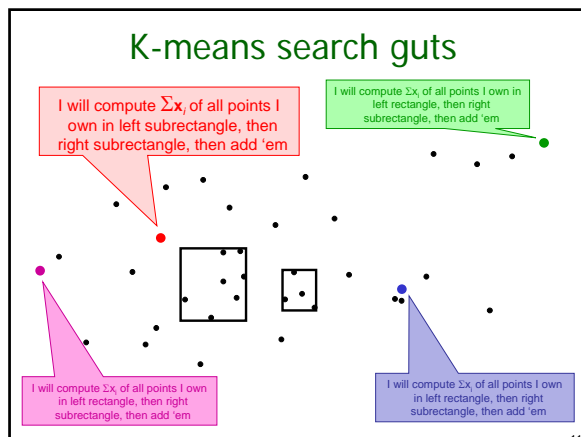
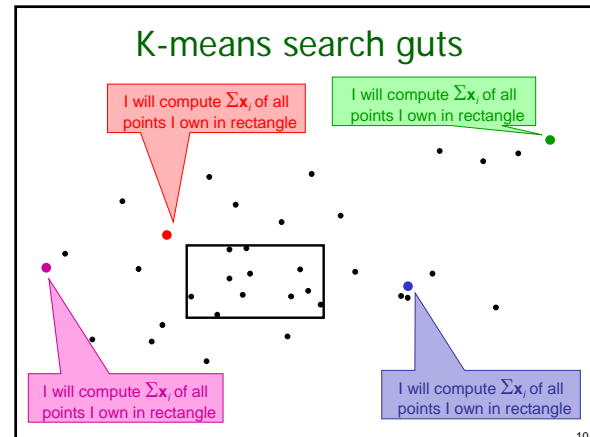
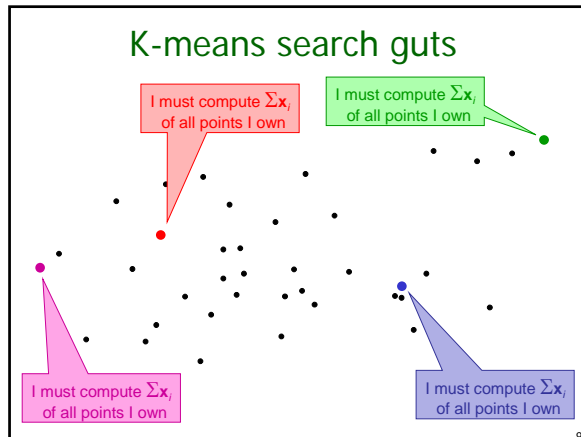
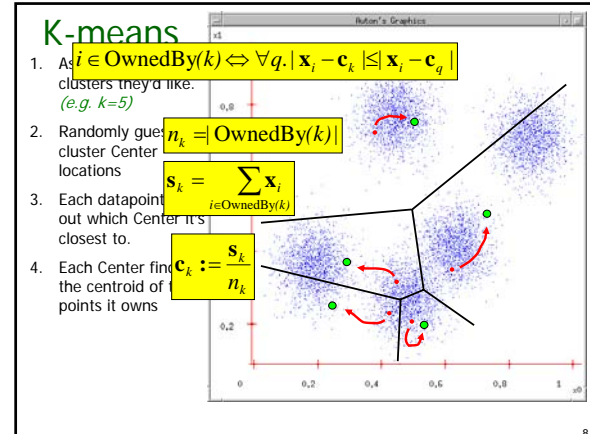
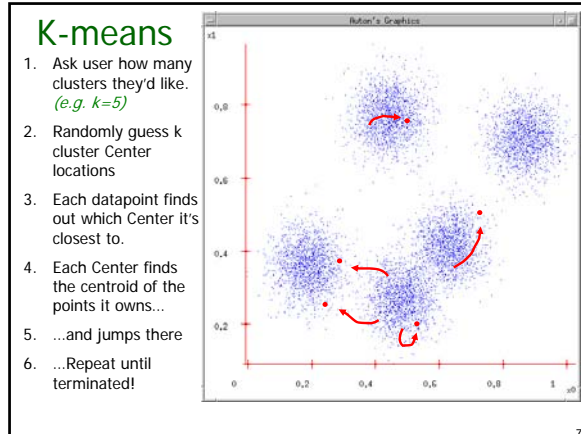
1. Ask user how many clusters they'd like.  
(e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations
3. Each datapoint finds out which Center it's closest to. (Thus each Center "owns" a set of datapoints)

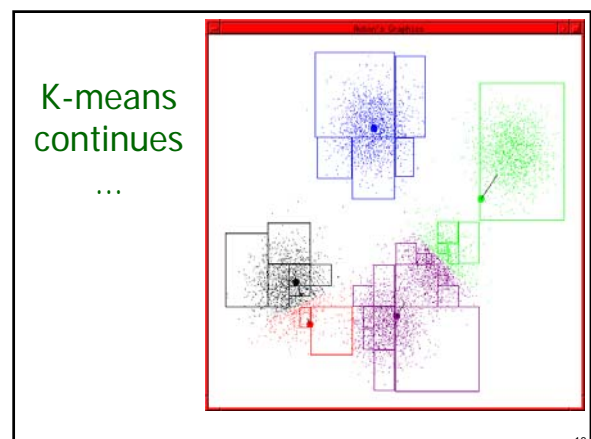
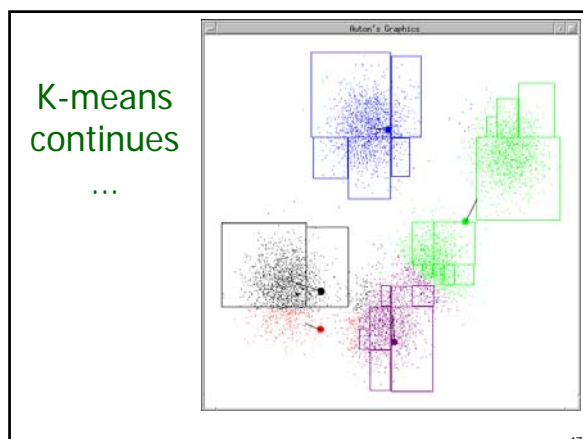
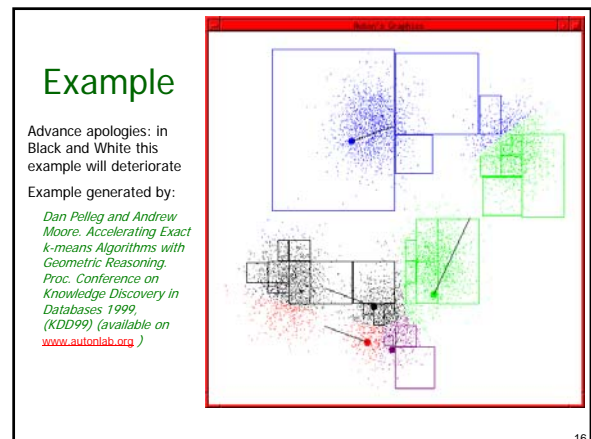
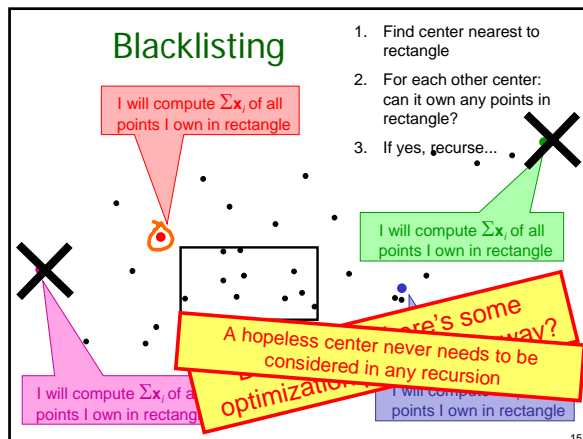
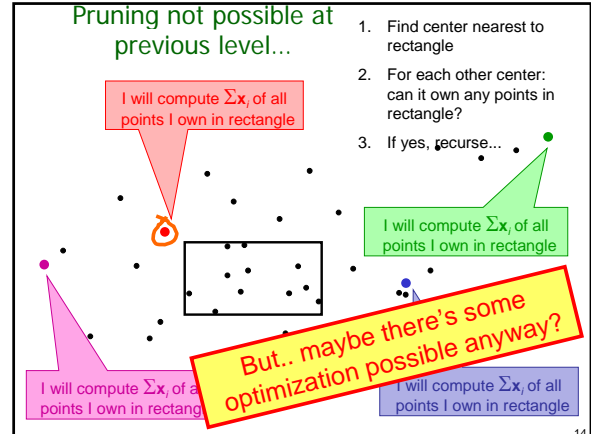
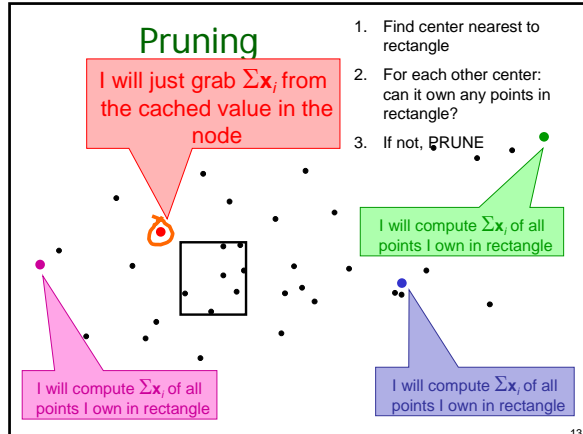


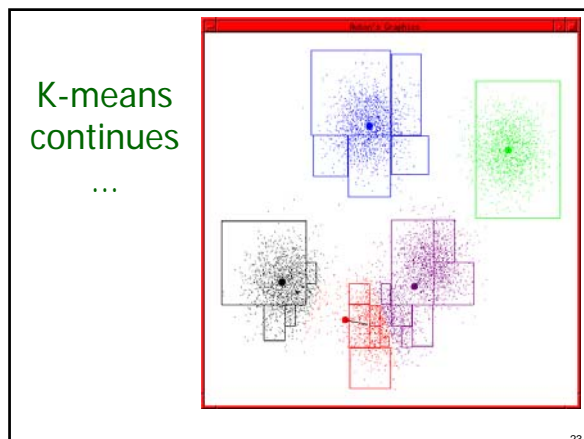
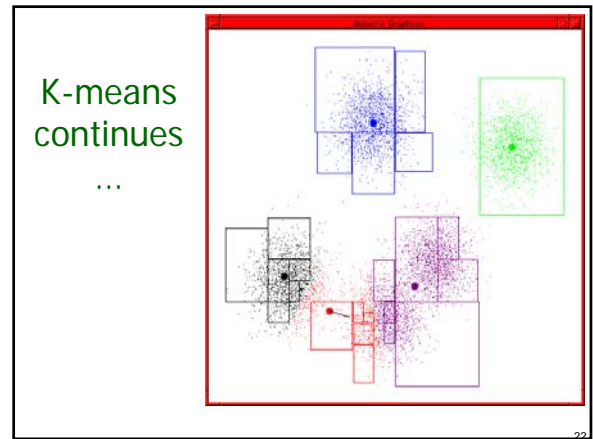
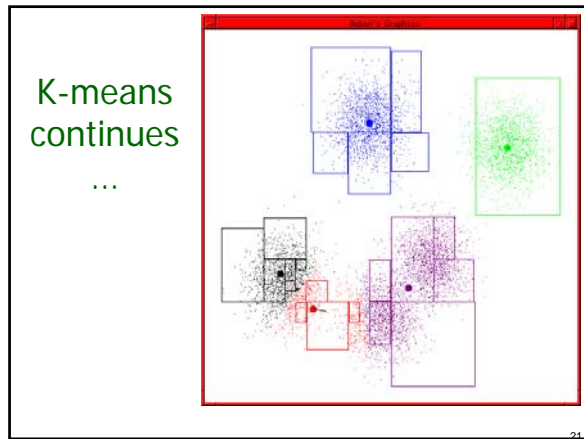
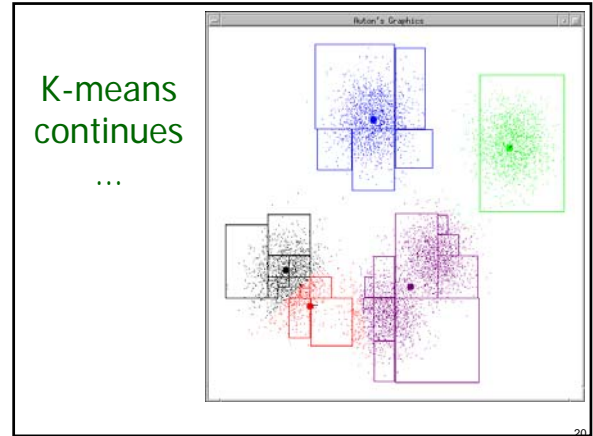
## K-means

1. Ask user how many clusters they'd like.  
(e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations
3. Each datapoint finds out which Center it's closest to.
4. Each Center finds the centroid of the points it owns

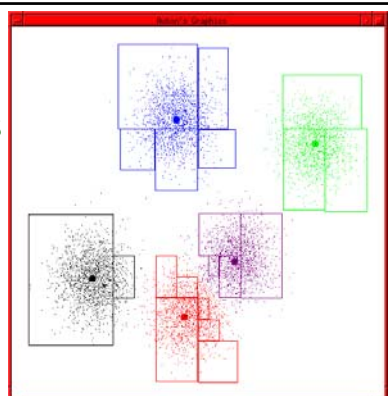








K-means  
terminates



25

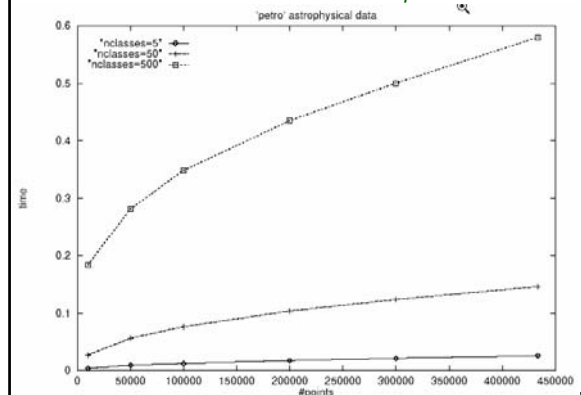
## Comparison to a linear algorithm

points	blacklisting	naive	speedup
50000	2.02	52.22	25.9
100000	2.16	134.82	62.3
200000	2.97	223.84	75.3
300000	1.87	328.80	176.3
433208	3.41	465.24	136.6

Astrophysics data (2-dimensions)

26

## Performance vs #records, #classes



## Performance vs #dimensions, #classes

