

# **COMP 551 - Applied Machine Learning**

## **Lecture 18 – Unsupervised learning**

---

William L. Hamilton

(with slides and content from Joelle Pineau)

\* Unless otherwise noted, all material posted for this course are copyright of the instructor, and cannot be reused or reposted without the instructor's written permission.

# What is unsupervised learning?

---

- Given only input data:  $D = \langle x_i \rangle, i=1:n$ , find some patterns or regularity in the data.
- Different classes of problems:
  1. Clustering
  2. Anomaly detection
  3. Dimensionality reduction

# What is unsupervised learning?

---

- Given only input data:  $D = \langle x_i \rangle, i=1:n$ , find some patterns or regularity in the data.
- Different classes of problems:
  1. Clustering
  2. Anomaly detection
  3. Dimensionality reduction

# A simple clustering example

---

- A fruit merchant approaches you, with a set of apples to classify according to their variety.
  - Tells you there are five varieties of apples in the dataset.
  - Tells you the weight and diameter of each apple in the dataset.
- Can you label each apple with the correct variety?
  - What would you need to know / assume?

Data =  $\langle x_1, ? \rangle, \langle x_2, ? \rangle, \dots, \langle x_n, ? \rangle$

# A simple clustering example

---

- You know there are 5 varieties.
- Assume each variety generates apples according to a (variety-specific) 2-dimensional Gaussian distribution.

# A simple clustering example

---

- You know there are 5 varieties.
- Assume each variety generates apples according to a (variety-specific) 2-dimensional Gaussian distribution.
- If you know  $\mu_i, \sigma_i^2$  for each class, it's easy to classify the apples.
- If you know the class of each apple, it's easy to estimate  $\mu_i, \sigma_i^2$ .

# A simple clustering example

---

- You know there are 5 varieties.
- Assume each variety generates apples according to a (variety-specific) 2-dimensional Gaussian distribution.
- If you know  $\mu_i, \sigma_i^2$  for each class, it's easy to classify the apples.
- If you know the class of each apple, it's easy to estimate  $\mu_i, \sigma_i^2$ .

But what if we know neither?

# A simple algorithm: K-means clustering

---

- **Objective:** Cluster  $n$  instances into  $K$  distinct classes.
- **Preliminaries:**
  - Step 1: Pick the desired number of clusters,  $K$ .
  - Step 2: Assume a parametric distribution for each class (e.g., Normal).
  - Step 3: Randomly estimate the parameters of the  $K$  distributions.
- **Iterate, until convergence:**
  - Step 4: Assign instances to the most likely classes.
  - Step 5: Estimate the parametric distribution of each class based on the latest assignment.

# K-means iterations: Simple K-means

- **Simple K-means:** just represent each cluster as a mean/centroid.
  - The likelihood that a point belongs to a cluster is simply proportional to the distance from the centroid.
  - Parametric cluster representation:  $\{ \mu_1, \dots, \mu_k \}$

Feature vector for point  $i$ .

Mean/centroid of cluster  $k$  at iteration  $t$ .

Cluster assignment of point  $i$  at iteration  $t$ .

$$c_i^{(t)} = \arg \min_k \| \mathbf{x}_i - \boldsymbol{\mu}_k^{(t)} \|$$

- Estimate the new parametric distribution of each class:

New mean/centroid for cluster  $k$ .

Set of points currently assigned to cluster  $k$ .

$$\boldsymbol{\mu}_k^{(t+1)} = \frac{1}{|\mathcal{C}_k^{(t)}|} \sum_{i \in \mathcal{C}_k^{(t)}} \mathbf{x}_i$$

# K-means iterations: Gaussian K-means

---

- Gaussian K-means: represent each cluster as a normal distribution.
  - The likelihood that a point belongs to a cluster is defined by a normal distribution.
  - Parametric cluster representation: means  $\{\mu_1, \dots, \mu_K\}$  and co-variances  $\{\Sigma_1, \dots, \Sigma_K\}$
- Assign instances to the most likely classes:  
Probability that point  $i$  was generated by a Gaussian with mean  $\mu_k^{(t)}$  and co-variance  $\Sigma_k^{(t)}$ .
$$c_i = \arg \max_k P_{\mathcal{N}(\mu_k^{(t)}, \Sigma_k^{(t)})}(\mathbf{x}_i)$$
- Estimate the new parametric distribution of each class:

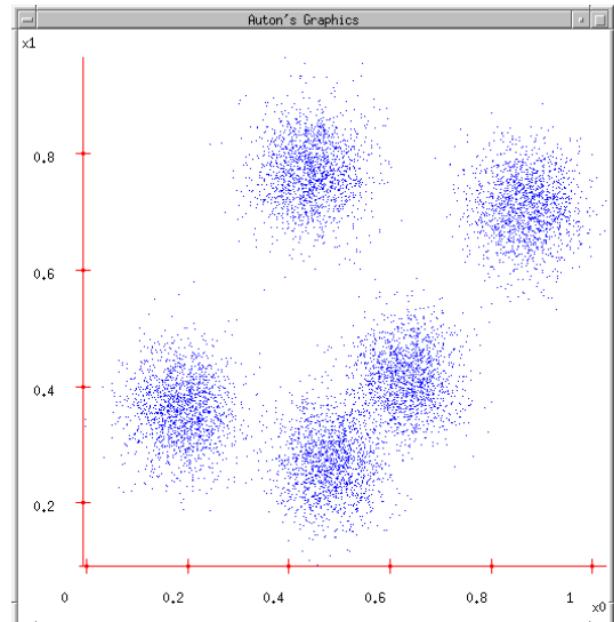
$$\mu_k^{(t+1)} = \frac{1}{|\mathcal{C}_k^{(t)}|} \sum_{i \in \mathcal{C}_k^{(t)}} \mathbf{x}_i \quad \Sigma_k^{(t+1)} = \frac{1}{|\mathcal{C}_k^{(t)}| - 1} \sum_{i \in \mathcal{C}_k^{(t)}} (\mathbf{x}_i - \mu_k^{(t)}) (\mathbf{x}_i - \mu_k^{(t)})^\top$$

# K-means algorithm

---

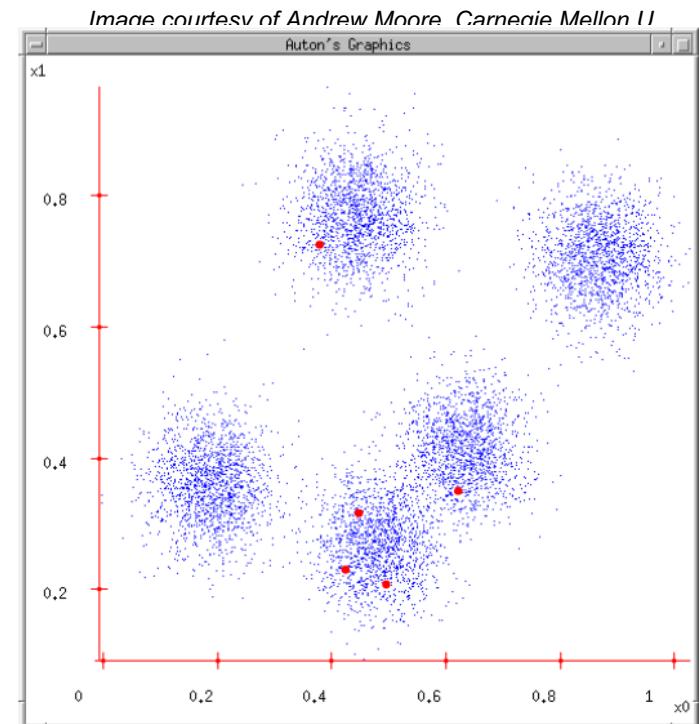
1. Ask user how many clusters.

*Image courtesy of Andrew Moore, Carnegie Mellon U.*



# K-means algorithm

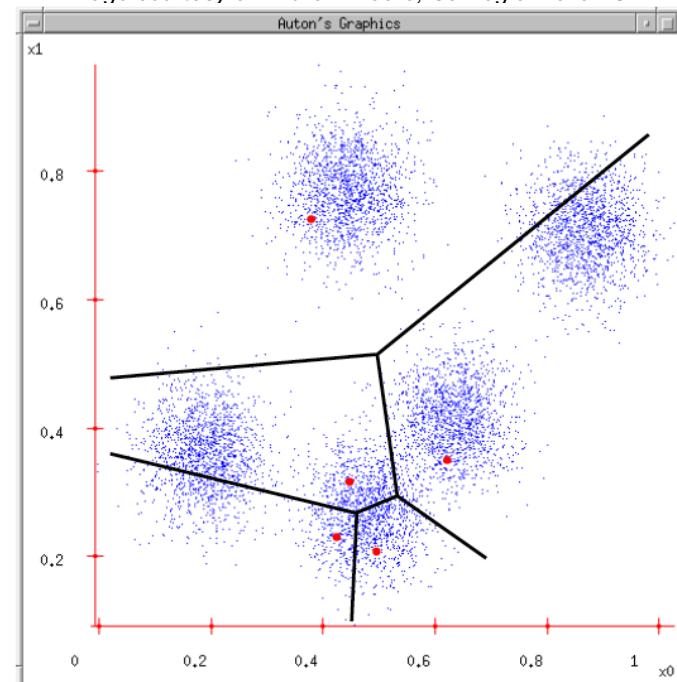
1. Ask user how many clusters.
2. Randomly guess k centers:  
 $\{ \mu_1, \dots, \mu_k \}$



# K-means algorithm

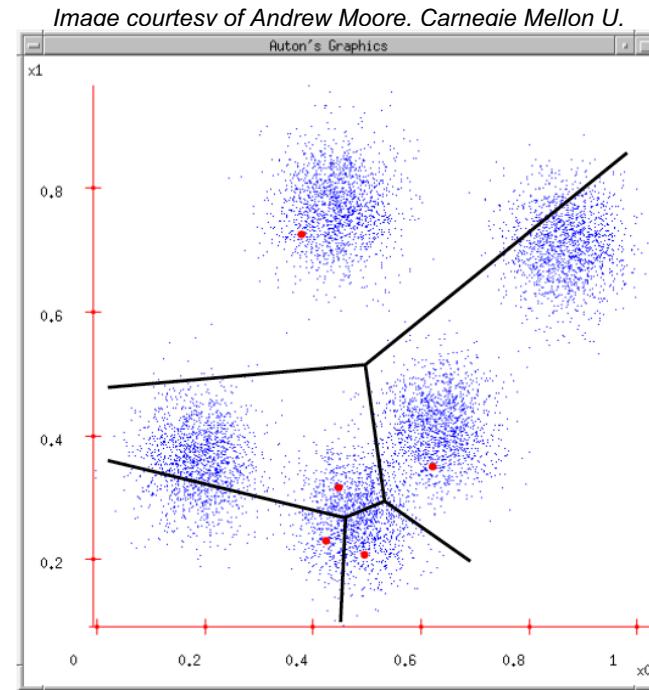
1. Ask user how many clusters.
2. Randomly guess k centers:  
 $\{ \mu_1, \dots, \mu_k \}$
3. Assign each data point to the closest center.

Image courtesy of Andrew Moore, Carnegie Mellon U.



# K-means algorithm

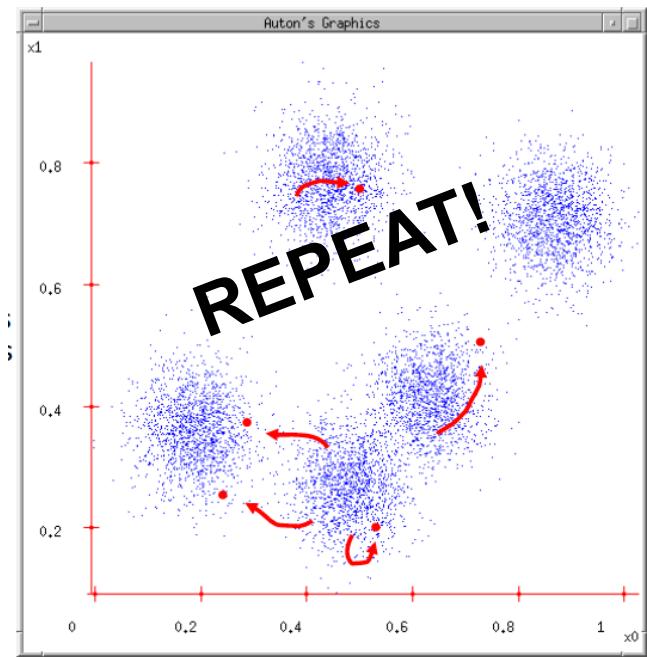
1. Ask user how many clusters.
2. Randomly guess k centers:  
 $\{ \mu_1, \dots, \mu_k \}$
3. Assign each data point to the closest center.
4. Each center finds the centroid of the points it owns... and jumps there.



# K-means algorithm

1. Ask user how many clusters.
2. Randomly guess k centers:  
 $\{ \mu_1, \dots, \mu_k \}$
3. Assign each data point to the closest center.
4. Each center finds the centroid of the points it owns... and jumps there.

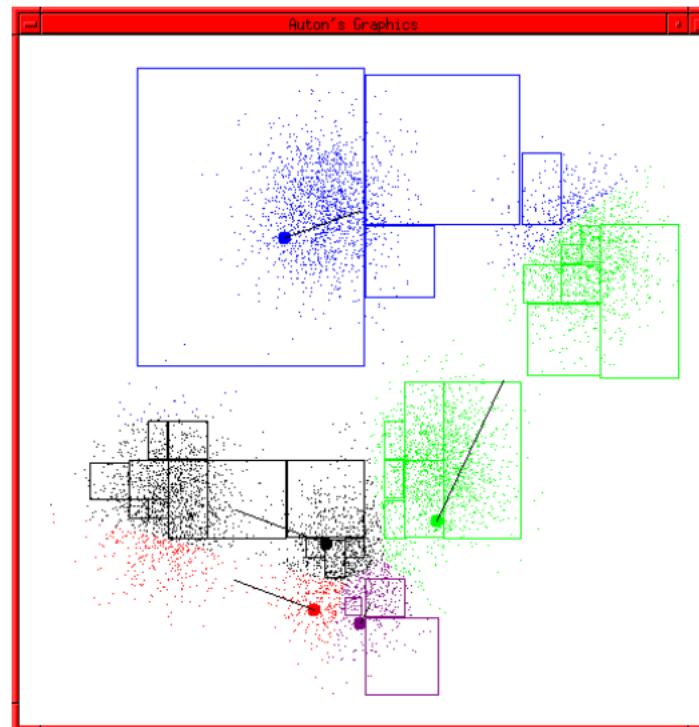
Image courtesy of Andrew Moore, Carnegie Mellon U.



# K-means algorithm starts

---

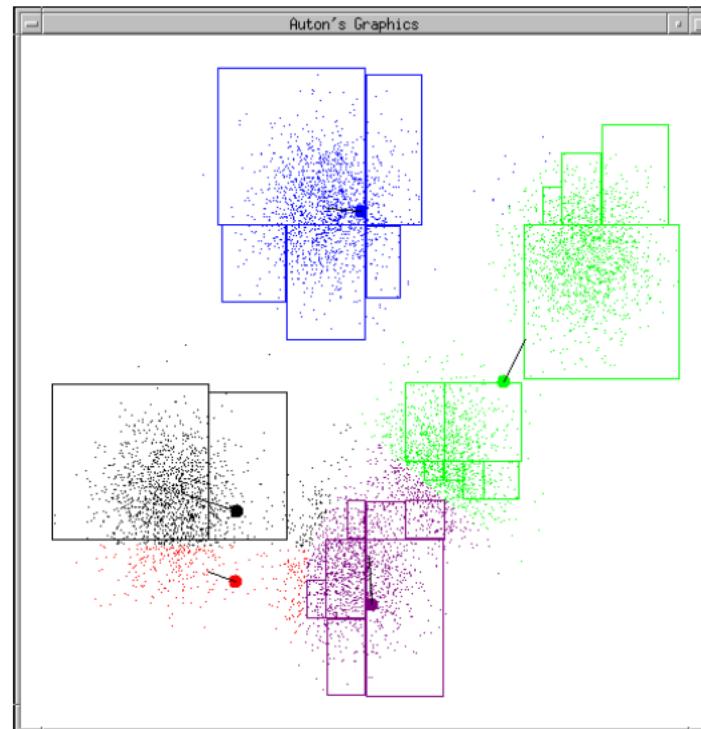
*Image courtesy of Andrew Moore, Carnegie Mellon U.*



# K-means algorithm continues (2)

---

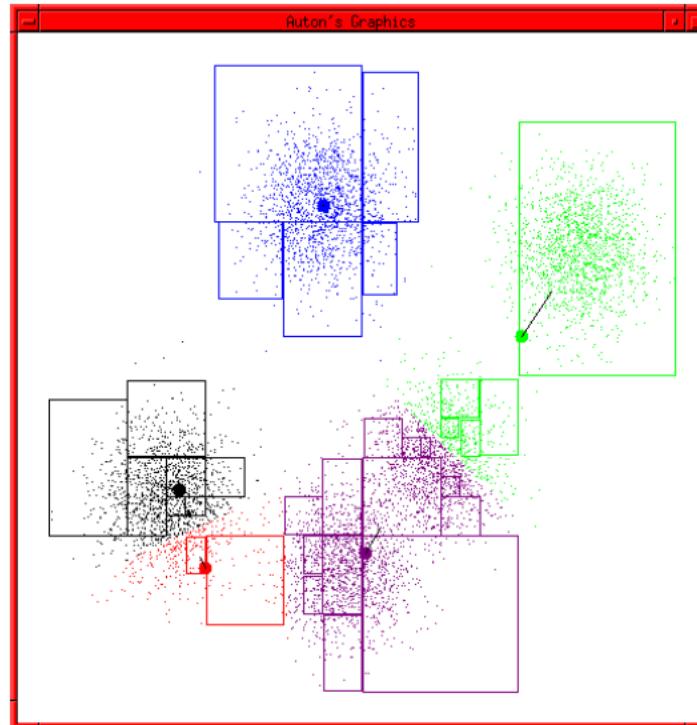
*Image courtesy of Andrew Moore, Carnegie Mellon U.*



# K-means algorithm continues (3)

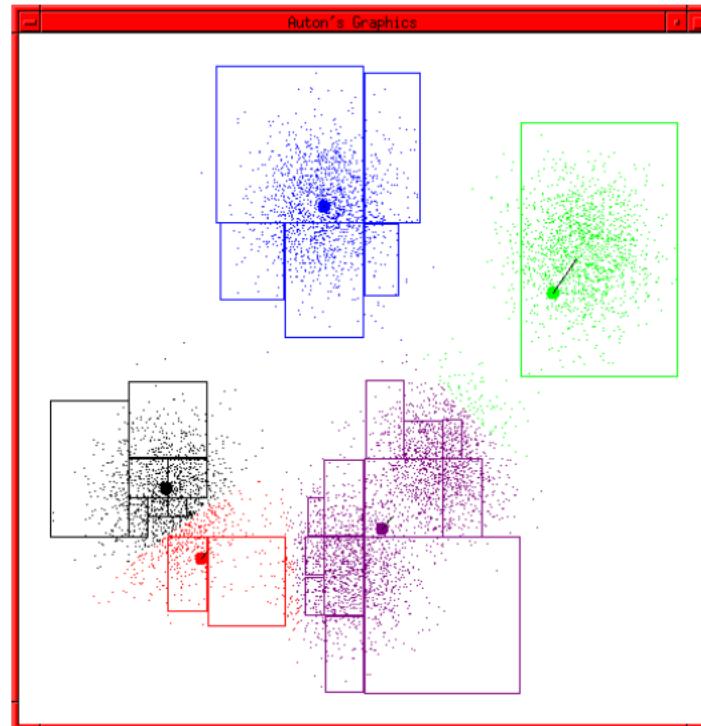
---

*Image courtesy of Andrew Moore, Carnegie Mellon U.*



# K-means algorithm continues (4)

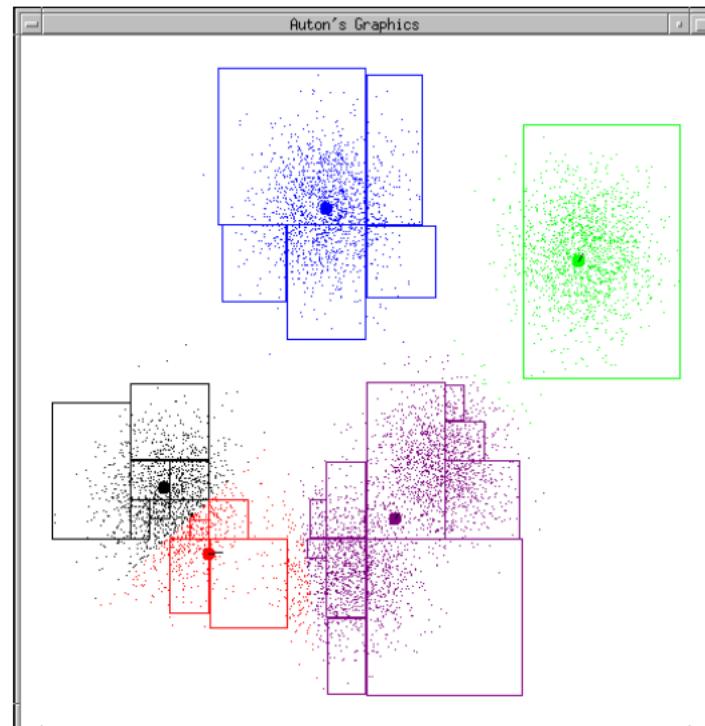
*Image courtesy of Andrew Moore, Carnegie Mellon U.*



# K-means algorithm continues (5)

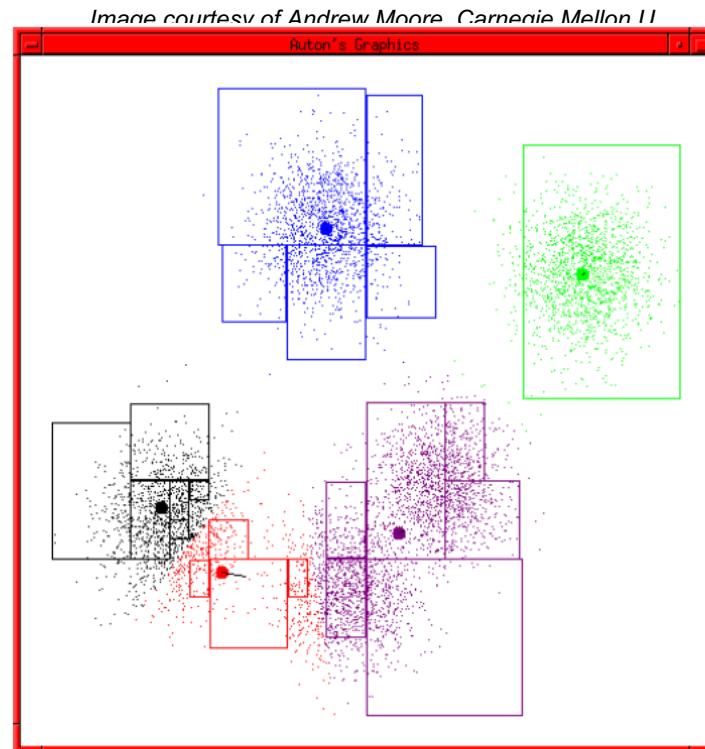
---

*Image courtesy of Andrew Moore, Carnegie Mellon U.*



# K-means algorithm continues (6)

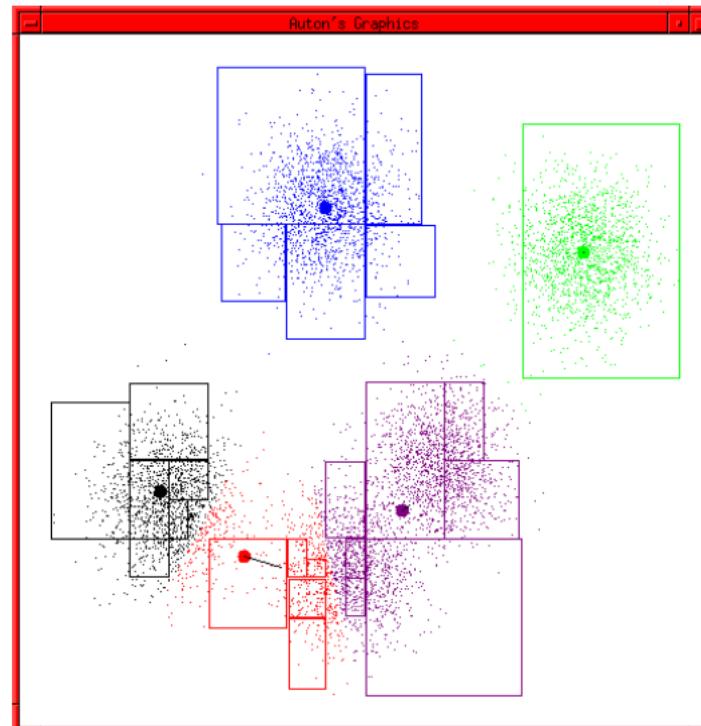
---



# K-means algorithm continues (7)

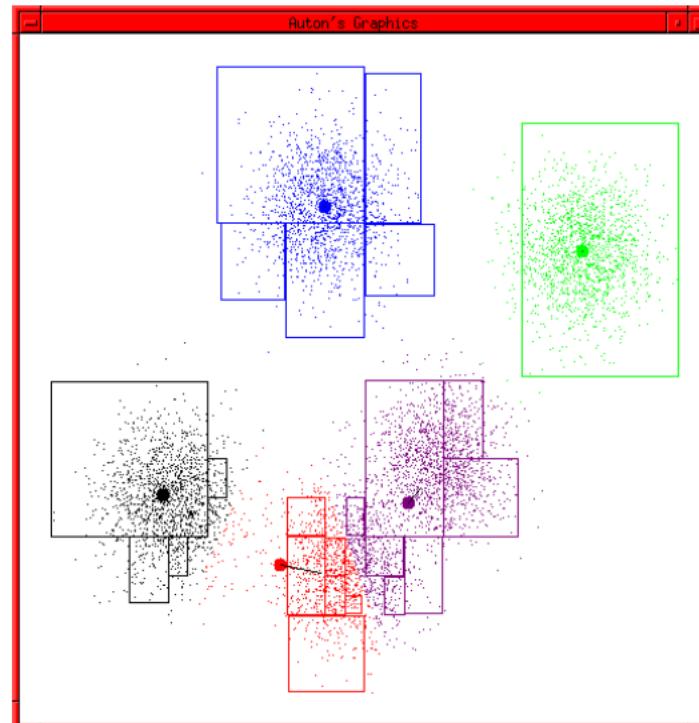
---

*Image courtesy of Andrew Moore, Carnegie Mellon U.*



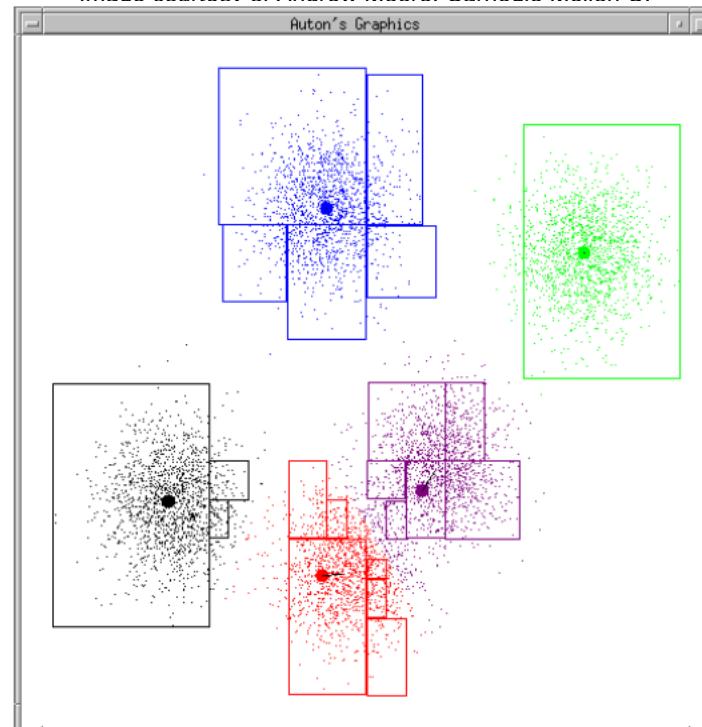
# K-means algorithm continues (8)

*Image courtesy of Andrew Moore, Carnegie Mellon U.*



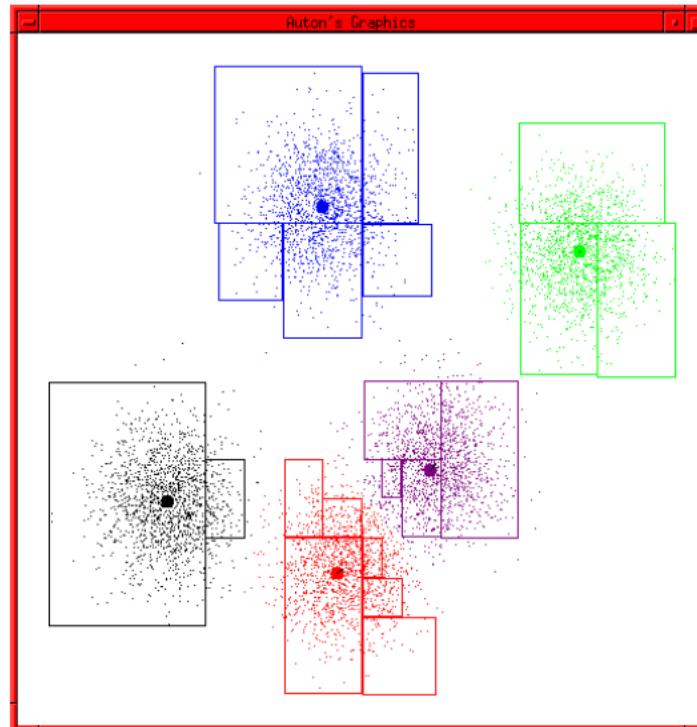
# K-means algorithm continues (9)

*Image courtesy of Andrew Moore, Carnegie Mellon U.*



# K-means algorithm terminates

*Image courtesy of Andrew Moore, Carnegie Mellon U.*



# Properties of K-means

---

- Optimality?
  - Converges to a local optimum.
  - Can use random re-starts to get better local optimum.
  - Alternately, can choose your initial centers carefully:
    - Place  $\mu_1$  on top of a randomly chosen datapoint.
    - Place  $\mu_2$  on top of datapoint that is furthest from  $\mu_1$ .
    - Place  $\mu_3$  on top of datapoint that is furthest from both  $\mu_1$  and  $\mu_2$ .

# Properties of K-means

---

- Optimality?
  - Converges to a local optimum.
  - Can use random re-starts to get better local optimum.
  - Alternately, can choose your initial centers carefully:
    - Place  $\mu_1$  on top of a randomly chosen datapoint.
    - Place  $\mu_2$  on top of datapoint that is furthest from  $\mu_1$ .
    - Place  $\mu_3$  on top of datapoint that is furthest from both  $\mu_1$  and  $\mu_2$ .
- Complexity?

# Properties of K-means

---

- Optimality?
  - Converges to a local optimum.
  - Can use random re-starts to get better local optimum.
  - Alternately, can choose your initial centers carefully:
    - Place  $\mu_1$  on top of a randomly chosen datapoint.
    - Place  $\mu_2$  on top of datapoint that is furthest from  $\mu_1$ .
    - Place  $\mu_3$  on top of datapoint that is furthest from both  $\mu_1$  and  $\mu_2$ .
- Complexity?  $O(knm)$ , where  $k=\# \text{centers}$ ,  $n=\# \text{datapoints}$ ,  $m=\text{dimensionality of data}$

# Properties of K-means

---

- Optimality?
  - Converges to a local optimum.
  - Can use random re-starts to get better local optimum.
  - Alternately, can choose your initial centers carefully:
    - Place  $\mu_1$  on top of a randomly chosen datapoint.
    - Place  $\mu_2$  on top of datapoint that is furthest from  $\mu_1$ .
    - Place  $\mu_3$  on top of datapoint that is furthest from both  $\mu_1$  and  $\mu_2$ .
- Complexity?  $O(knm)$ , where  $k=\# \text{centers}$ ,  $n=\# \text{datapoints}$ ,  $m=\text{dimensionality of data}$
- K-means is an instance of a family of learning algorithms called “Expectation-Maximization” (aka EM).

# A simple algorithm: K-means clustering

---

- **Objective:** Cluster  $n$  instances into  $K$  distinct classes.
- **Preliminaries:**
  - Step 1: Pick the desired number of clusters,  $K$ .
  - Step 2: Assume a parametric distribution for each class (e.g. Normal).
  - Step 3: Randomly estimate the parameters of the  $K$  distributions.
- Iterate, until convergence:
  - Step 4: Assign instances to the most likely classes based on the current parametric distributions.  
Maximization step
  - Step 5: Estimate the parametric distribution of each class based on the latest assignment.  
Expectation step

# Expectation Maximization (more generally)

---

- Iterative method for learning the maximum likelihood estimate of a probabilistic model, when the model contains unobservable variables.

# Expectation Maximization (more generally)

---

- Iterative method for learning the maximum likelihood estimate of a probabilistic model, when the model contains **unobservable variables**.
- Main idea:
  - If we had sufficient statistics for the data (e.g. counts of possible values), we could easily maximize the likelihood.
  - With missing data, we “fantasize” how the data should look based on the current parameter setting. I.e., compute expected sufficient statistics.
  - Then we maximize parameter setting, based on these statistics.

# Expectation Maximization (more generally)

---

- Start with some initial parameter setting.
- Repeat (as long as desired):
  - Expectation (E) step: Complete the data by assigning “values” to the missing items.
  - Maximization (M) step: Compute the maximum likelihood parameter setting based on the completed data.
- Once the data is completed (E-step), computing the log-likelihood and new parameters (M-step) is easy! **This is what we did for K-means.**

# Expectation Maximization: Properties

---

- Likelihood function is guaranteed to improve (or stay the same) with each iteration.
- Iterations can stop when no more improvements are achieved.
- Convergence to a local optimum of the likelihood function.
- Re-starts with different initial parameters are often necessary.
- Time complexity (per iteration) depends on model structure.

EM is very useful in practice!

# Aside: Hierarchical clustering

---

- K-means generates a “flat” cluster structure. However, it is also possible to learn “**hierarchical clusters**.”

# Aside: Hierarchical clustering

---

- K-means generates a “flat” cluster structure. However, it is also possible to learn “**hierarchical clusters**.”
- Two general approaches:
  - Top-down (agglomerative): Recursively merge a pair of clusters.
  - Bottom-up (divisive): Recursively split the existing clusters.

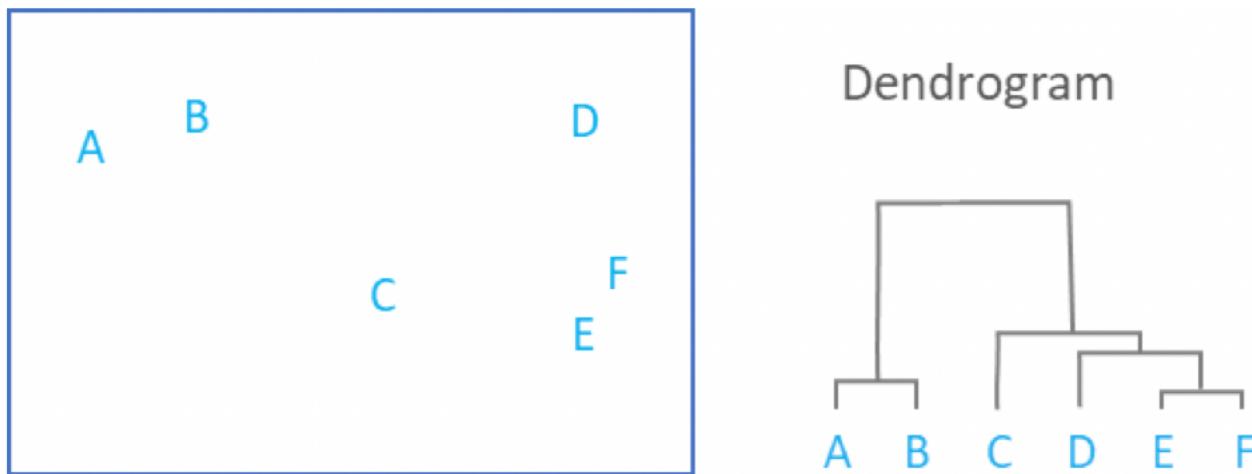
# Aside: Hierarchical clustering

---

- K-means generates a “flat” cluster structure. However, it is also possible to learn “**hierarchical clusters**”.
- Two general approaches:
  - Top-down (agglomerative): Recursively merge a pair of clusters.
  - Bottom-up (divisive): Recursively split the existing clusters.
- Use **dissimilarity measure** to recursively split/merge pairs:
  - Measure pairwise distance between any points in the 2 clusters.
    - E.g. Euclidean distance, Manhattan distance.
  - Measure distance over entire clusters using linkage criterion.
    - E.g. Min/Max/Mean over pairs of points.

# Aside: Hierarchical clustering

---



# Anomaly detection

---



<http://www.anomalydetectionresearch.com>

# Anomaly detection

---

- K-means (and other discriminative approaches) tend to be ineffective when one class/cluster is much more rare than the other.

# Anomaly detection

---

- K-means (and other discriminative approaches) tend to be ineffective when one class/cluster is much more rare than the other.
- A simple **generative** approach:
  - Fit a model,  $p(x)$  using the input data.
  - Set a decision threshold  $\varepsilon$  and predict  $Y = \{1 \text{ if } p(x) > \varepsilon, 0 \text{ otherwise}\}$ .
  - Use a validation set to measure performance (can use cross-validation to set  $\varepsilon$ ).

# Anomaly detection vs Supervised learning

---

## Anomaly detection

- Small number of positive examples (e.g. <10).
- Large number of negative examples (e.g. >100).

## Supervised learning

- Similar number of positive and negative examples.

# Anomaly detection vs Supervised learning

---

## Anomaly detection

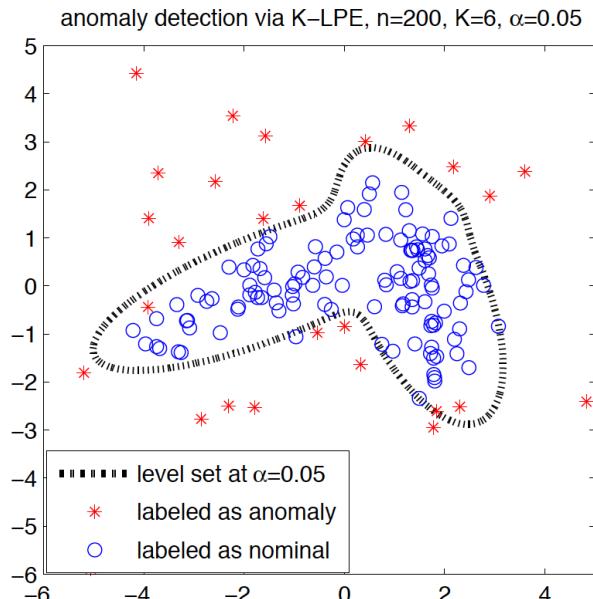
- Small number of positive examples (e.g. <10).
- Large number of negative examples (e.g. >100).
- Many different “types” of anomalies, so don’t want to fit a model for the positive class.

## Supervised learning

- Similar number of positive and negative examples.
- More homogeneity within classes, or enough data to sufficiently characterize each classes.

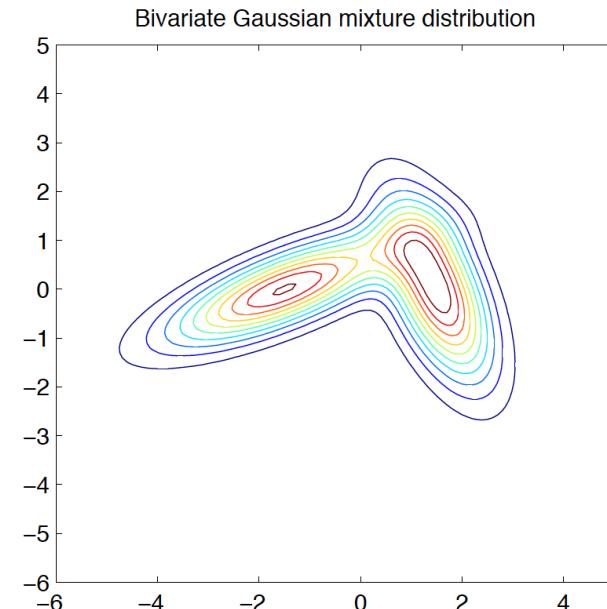
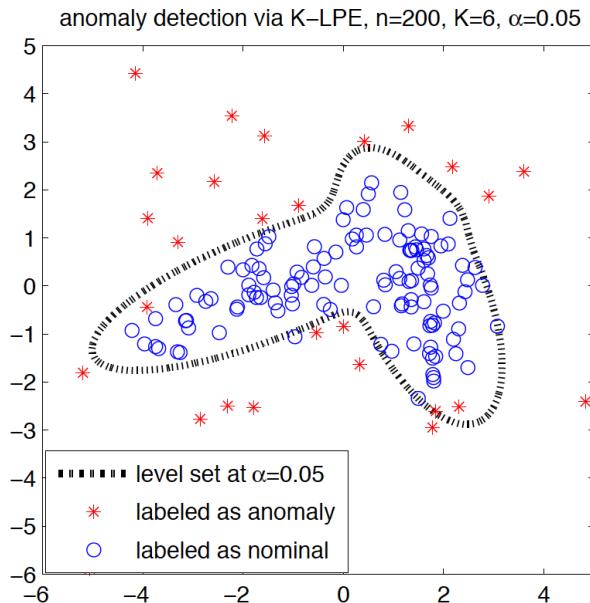
# A simple example

---



From: M. Zhao and V. Saligrama, "Anomaly Detection with Score functions based on Nearest Neighbor Graphs", Neural Information Processing Systems (NIPS) Conference, 2009

# A simple example



From: M. Zhao and V. Saligrama, "Anomaly Detection with Score functions based on Nearest Neighbor Graphs", Neural Information Processing Systems (NIPS) Conference, 2009

# Gaussian Mixture Model

---

- Idea: Fit data with a combination of Gaussian distributions.
- What defines a set of Gaussians?

# Gaussian Mixture Model

---

- Idea: Fit data with a combination of Gaussian distributions.
- Write  $p(x)$  as a linear combination of Gaussians:

$$p(x) = \sum_{k=1:K} p(z_k) p(x | z_k)$$

where  $p(z_k)$  is the probability of the  $k^{\text{th}}$  mixture component  
and  $p(x | z_k) = N(x | \mu_k, \sigma_k^2)$  is the prob. of  $x$  for the  $k^{\text{th}}$  mixture component.

- Estimate parameters  $p(z_k), \mu_k, \sigma_k^2$ , using Expectation-Maximization approach.

# Gaussian Mixture Model

---

- Idea: Fit data with a combination of Gaussian distributions.
- Write  $p(x)$  as a linear combination of Gaussians:

$$p(x) = \sum_{k=1:K} p(z_k) p(x | z_k)$$

where  $p(z_k)$  is the probability of the  $k^{\text{th}}$  mixture component  
and  $p(x | z_k) = N(x | \mu_k, \sigma_k^2)$  is the prob. of  $x$  for the  $k^{\text{th}}$  mixture component.

- Estimate parameters  $p(z_k), \mu_k, \sigma_k^2$ , using Expectation-Maximization approach.

# Generative modeling beyond anomalies

---

- Generative model  $p(x)$  is not only useful for anomaly detection.
- Can be used to:
  - Provide priors in generative classification models.
  - Generate artificial samples.
  - Provide likelihood model for reinforcement learning agents.
  - ....
- Gaussian mixture models are very common/standard, but we will cover recent advancements in deep generative models after the midterm.

# Practical issues

---

- In general, both clustering and anomaly detection are **very sensitive to the feature space!** Feature selection/construction is critical...

# Practical issues

---

- In general, both clustering and anomaly detection are **very sensitive to the feature space!** Feature selection/construction is critical...
- But how do we do good feature selection in the unsupervised setting?
  - Can't use cross-validation...
  - Can't use mutual information or correlation with the target label....

# Practical issues

---

- In general, both clustering and anomaly detection are **very sensitive to the feature space!** Feature selection/construction is critical...
- But how do we do good feature selection in the unsupervised setting?
  - Can't use cross-validation...
  - Can't use mutual information or correlation with the target label....
- **Idea:** Use unsupervised (deep) learning to infer high-quality low-dimensional representations (next lecture). Dimensionality reduction.