

Bayesian Networks (Structure) Learning

Machine Learning – 10701/15781
 Carlos Guestrin
 Carnegie Mellon University

November 12th, 2007

©2005-2007 Carlos Guestrin

1

~~Information-theoretic interpretation of maximum likelihood 1~~

$\log \prod_{i=1}^n P(x^{(i)} | \theta_G, G) = \sum_i \log P(x^{(i)} | \theta_G, G)$

$\log P(\mathcal{D} | \theta_G, G) = \sum_i \log P(f^{(i)}) P(a^{(i)}) P(s^{(i)} | f^{(i)}, a^{(i)}) P(h^{(i)} | s^{(i)}) P(n^{(i)} | s^{(i)})$

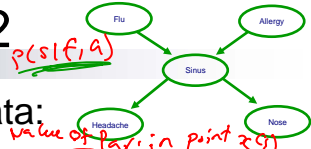
$= \underbrace{\left[\sum_i \log P(f^{(i)}) \right]}_{P(F)} + \underbrace{\left[\sum_i \log P(a^{(i)}) \right]}_{P(A)} + \underbrace{\left[\sum_i \log P(s^{(i)} | f^{(i)}, a^{(i)}) \right]}_{P(S|F,A)} + \underbrace{\left[\sum_i \log P(h^{(i)} | s^{(i)}) \right]}_{P(H|S)} + \underbrace{\left[\sum_i \log P(n^{(i)} | s^{(i)}) \right]}_{P(N|S)}$

one per CPT $\rightarrow P(x_i | \text{Pa}x_i)$

only this term changes about $P(A|F)$

$G:$

Information-theoretic interpretation of maximum likelihood 2



- Given structure, log likelihood of data:

$$\log P(\mathcal{D} | \theta_{\mathcal{G}}, \mathcal{G}) = \sum_{j=1}^m \sum_{i=1}^n \log P(X_i = x_i^{(j)} | \text{Pa}_{X_i} = x^{(j)}[\text{Pa}_{X_i}])$$

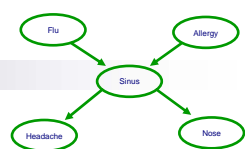
Handwritten notes: $\log P(\mathcal{D} | \theta_{\mathcal{G}}, \mathcal{G})$ is the log-likelihood. $\sum_{j=1}^m$ is the sum over points. $\sum_{i=1}^n$ is the sum over nodes. $\log P(X_i = x_i^{(j)} | \text{Pa}_{X_i} = x^{(j)}[\text{Pa}_{X_i}])$ is the log-probability for a point j where $X_i = x_i^{(j)}$ and its parents are $\text{Pa}_{X_i} = x^{(j)}[\text{Pa}_{X_i}]$. For point j where $F=t, A=f, S=t$, $\log P(S=t | F=t, A=f)$.

$$= \sum_{i=1}^n \sum_{j=1}^m \log P(X_i = x_i^{(j)} | \text{Pa}_{X_i} = x^{(j)}[\text{Pa}_{X_i}])$$

$$= \sum_{i=1}^n \sum_{x_i} \sum_u \underbrace{\text{count}(X_i = x_i, \text{Pa}_{X_i} = u)}_m \log P(X_i = x_i | \text{Pa}_{X_i} = u)$$

$\hat{P}(X_i = x_i, \text{Pa}_{X_i} = u)$
MLE estimate

Information-theoretic interpretation of maximum likelihood 3



- Given structure, log likelihood of data:

$$\log \hat{P}(\mathcal{D} | \theta, \mathcal{G}) = m \sum_i \sum_{x_i, \text{Pa}_{x_i, \mathcal{G}}} \hat{P}(x_i, \text{Pa}_{x_i, \mathcal{G}}) \log \hat{P}(x_i | \text{Pa}_{x_i, \mathcal{G}})$$

Handwritten notes: $\log \hat{P}(\mathcal{D} | \theta, \mathcal{G})$ is the log-likelihood. m is the number of points. \sum_i is the sum over points. $\sum_{x_i, \text{Pa}_{x_i, \mathcal{G}}}$ is the sum over nodes and their parents. $\hat{P}(x_i, \text{Pa}_{x_i, \mathcal{G}})$ is the MLE estimate of the joint probability. $\log \hat{P}(x_i | \text{Pa}_{x_i, \mathcal{G}})$ is the log-probability for a point i where $X_i = x_i$ and its parents are $\text{Pa}_{x_i, \mathcal{G}}$. $H(A|B)$ is the conditional entropy. $H(A|B) = -\sum_{a,b} P(a,b) \cdot \log P(a|b)$. $I(A, B)$ is the mutual information. $I(A, B) = -\sum_{a,b} P(a,b) \log \frac{P(a,b)}{P(a)P(b)}$. $I(A, B) = H(A) - H(A|B)$.

$$= m \sum_i H(X_i | \text{Pa}_{X_i})$$

Handwritten notes: $H(X_i | \text{Pa}_{X_i})$ is the conditional entropy. $H(X_i | \text{Pa}_{X_i})$ is small; little uncertainty about $X_i | \text{Pa}_{X_i}$.

$$= m \sum_{i=1}^n [I(X_i, \text{Pa}_{X_i}) - H(X_i)]$$

Handwritten notes: $I(X_i, \text{Pa}_{X_i})$ is the mutual information. $H(X_i)$ is the entropy. $I(X_i, \text{Pa}_{X_i}) - H(X_i)$ is the difference between mutual information and entropy. $I(X_i, \text{Pa}_{X_i})$ is the mutual information between X_i and its parents. $H(X_i)$ is the entropy of X_i . $I(X_i, \text{Pa}_{X_i}) - H(X_i)$ is the difference between mutual information and entropy.

Decomposable score

- Log data likelihood

↑ MLE

$$\log \hat{P}(D | \theta, G) = m \sum_i \hat{I}(X_i, \text{Pa}_{X_i, G}) - M \sum_i \hat{H}(X_i)$$

constants doesn't depend on G

- Decomposable score:

- Decomposes over families in BN (node and its parents)
- Will lead to significant computational efficiency!!!
- Score(G : D) = $\sum_{i=1}^n \text{FamScore}(X_i | \text{Pa}_{X_i} : D)$

graphs

$$\text{FamScore}(X_i | \text{Pa}_{X_i} : D) = m [I(X_i, \text{Pa}_{X_i}) - H(X_i)]$$

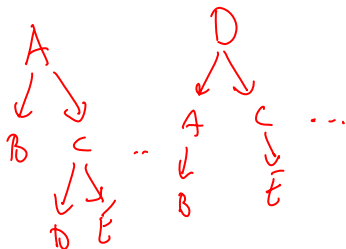
How many trees are there?

each node at most one parent
eg. NB, HNPAS, ...

Nonetheless - Efficient optimal algorithm finds best tree

$$X = \{A, B, C, \dots, Z\}$$

$$O(2^{\Theta(n \log n)})$$



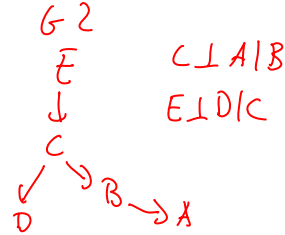
For trees only care about (no v. study) Symm: $I(A|B) = I(B|A)$
 edges not for directions

Scoring a tree 1: equivalent trees

$$\log \hat{P}(\mathcal{D} | \theta, \mathcal{G}) = M \sum_i \hat{I}(x_i, \text{Pa}_{x_i, \mathcal{G}}) - M \sum_i \hat{H}(X_i)$$



C ⊥ A | B
 E ⊥ D | C



C ⊥ A | B
 E ⊥ D | C

$$\text{Score}(G1) = I(A, B) + I(C, B) + I(C, D) + I(E, C)$$

$$\text{Score}(G2) = \text{Score}(G1)$$

Cool!!

GRAPHS MAKE SAME IND. ASSUMPTIONS
 ⇒ Same Score

Scoring a tree 2: similar trees

$$\log \hat{P}(\mathcal{D} | \theta, \mathcal{G}) = M \sum_i \hat{I}(x_i, \text{Pa}_{x_i, \mathcal{G}}) - M \sum_i \hat{H}(X_i)$$



$$\text{Score}(G1): I(A, B) + I(A, C) + I(C, D)$$



for trees T

$$\text{Score}(T) = \sum_{(i,j) \in T} I(X_i, X_j)$$

$$\begin{aligned} \text{Score}(G2): & I(A, B) + I(A, C) + I(B, D) \\ = & \text{Score}(G1) + I(B, D) - I(C, D) \end{aligned}$$

Chow-Liu tree learning algorithm 1

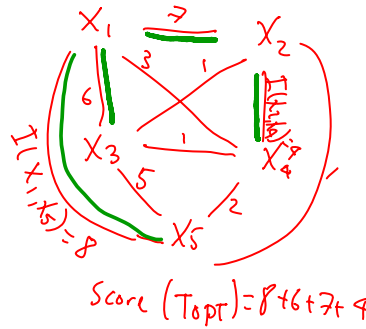
- For each pair of variables X_i, X_j
 - Compute empirical distribution:

$$\hat{P}(x_i, x_j) = \frac{\text{Count}(x_i, x_j)}{m}$$

- Compute mutual information:

$$\hat{I}(X_i, X_j) = \sum_{x_i, x_j} \hat{P}(x_i, x_j) \log \frac{\hat{P}(x_i, x_j)}{\hat{P}(x_i)\hat{P}(x_j)}$$

- Nodes X_1, \dots, X_n
- Edge (i, j) gets weight $\hat{I}(X_i, X_j)$



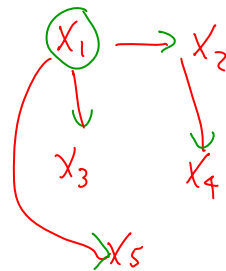
Maximum weight spanning
(Basic, very efficient alg)

Chow-Liu tree learning algorithm 2

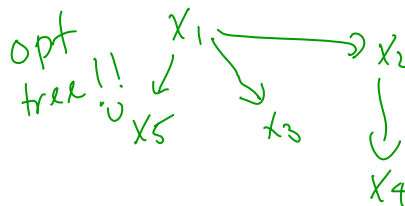
$$\log \hat{P}(\mathcal{D} | \theta, \mathcal{G}) = M \sum_i \hat{I}(x_i, \text{Pa}_{x_i, \mathcal{G}}) - M \sum_i \hat{H}(X_i)$$

Optimal tree BN

- Compute maximum weight spanning tree *don't get v-structure*
- Directions in BN: pick any node as root, breadth-first-search defines directions



$O(n^3)$



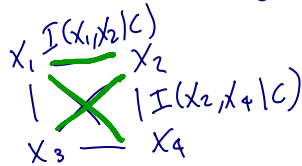
Can we extend Chow-Liu 1

- Tree augmented naïve Bayes (TAN) [Friedman et al. '97]

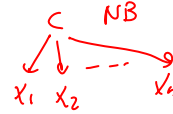
- Naïve Bayes model overcounts, because correlation between features not considered
- Same as Chow-Liu, but score edges with:

$$\hat{I}(X_i, X_j | C) = \sum_{c, x_i, x_j} P(c, x_i, x_j) \log \frac{P(x_i, x_j | c)}{P(x_i | c)P(x_j | c)}$$

$$\text{Score}(x_i | x_j, c) = I(x_i, x_j | c)$$

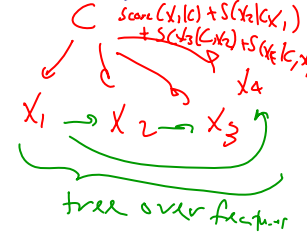


From HW1:



if features correlated doesn't do well

TAN: $\text{Score}(TAN) \neq \text{Score}(NB)$ always



Can we extend Chow-Liu 2

- (Approximately learning) models with tree-width up to k

(trees: treewidth=1)

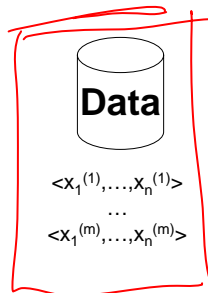
- [Checheta & Guestrin '07]
- But, $O(n^{2k+6})$.

What you need to know about learning BN structures so far

- Decomposable scores
 - Maximum likelihood
 - Information theoretic interpretation
- Best tree (Chow-Liu)
- Best TAN
- Nearly best k-treewidth (in $O(N^{2k+6})$)

Scoring general graphical models – Model selection problem

What's the best structure?



Score
→ 17



→ 19 > 17

fewer independencies
less bias, \Rightarrow fit data
fully connected!! (high variance) better (high bias)

The more edges, the fewer independence assumptions, the higher the likelihood of the data, but will overfit...

Maximum likelihood overfits!

$$\log \hat{P}(D | \theta, \mathcal{G}) = M \sum_i \hat{I}(x_i, \text{Pa}_{x_i, \mathcal{G}}) - M \sum_i \hat{H}(X_i)$$

- Information never hurts: $I(A; B, C) \geq I(A; B)$

$\Rightarrow I(x_i; \text{Pa}_{x_i})$ increases as $|\text{Pa}_{x_i}|$ increases

\Rightarrow fully connected graph ML

- Adding a parent always increases score!!!

Bayesian score avoids overfitting

(encodes "regularization")

$$\text{Score}(\theta) = \sum_i \text{Score}(x_i | \text{Pa}_{x_i})$$

- Given a structure, distribution over parameters

$$\log P(D | \mathcal{G}) = \log \int_{\theta_{\mathcal{G}}} P(D | \mathcal{G}, \theta_{\mathcal{G}}) P(\theta_{\mathcal{G}} | \mathcal{G}) d\theta_{\mathcal{G}}$$

- Difficult integral: use Bayes information criterion (BIC) approximation (equivalent as $M \rightarrow \infty$)

$$\log P(D | \mathcal{G}) \approx \log P(D | \mathcal{G}, \theta_{\mathcal{G}}) - \frac{\text{NumberParams}(\mathcal{G})}{2} (\log M) + \mathcal{O}(1)$$

Bayesian Score (regularized) = data Likelihood - num of params (# of parents)

as function of m	m	log m
Simple graph	lower	lower

- Note: regularize with MDL score

- Best BN under BIC still NP-hard

adding parents to x_i
NumParams: increase exponentially

Structure learning for general graphs

- In a tree, a node only has one parent

- **Theorem:**
 - The problem of learning a BN structure with at most d parents is NP-hard for any (fixed) $d \geq 2$

- Most structure learning approaches use heuristics
 - Exploit score decomposition
 - (Quickly) Describe ~~two~~ heuristics that exploit decomposition in different ways

Learn BN structure using local search

FamScore ~~is~~ includes regularization

Starting from Chow-Liu tree



can start from other points

Local search,

possible moves:

- Add edge
- Delete edge
- Invert edge

consider G , Adding edge $x_i \rightarrow x_j$ to generate G'

$$\text{Score}(G') - \text{Score}(G)$$

$$= \text{FamScore}(x_i | \text{Pa}_{x_i} \cup x_j)$$

$$- \text{FamScore}(x_i | \text{Pa}_{x_i})$$

local easy to compute

Score using BIC

17

19

20

Be Careful!

don't create a cycle!

What you need to know about learning BNs

■ Learning BNs

- Maximum likelihood or MAP learns parameters
- Decomposable score
- Best tree (Chow-Liu)
- Best TAN
- Other BNs, usually local search with BIC score

Unsupervised Learning Clustering K-means

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

November 12th, 2007

Some Data

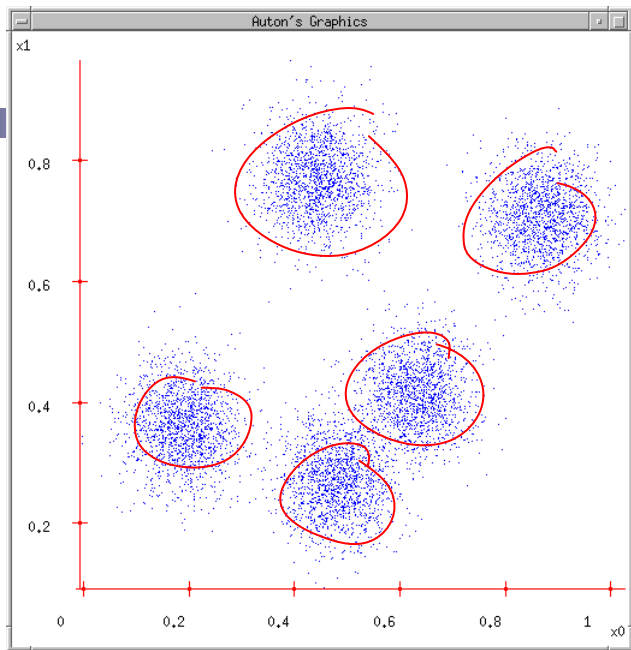
Supervised

x, y x, z
 $x, -$
:

Unsupervised

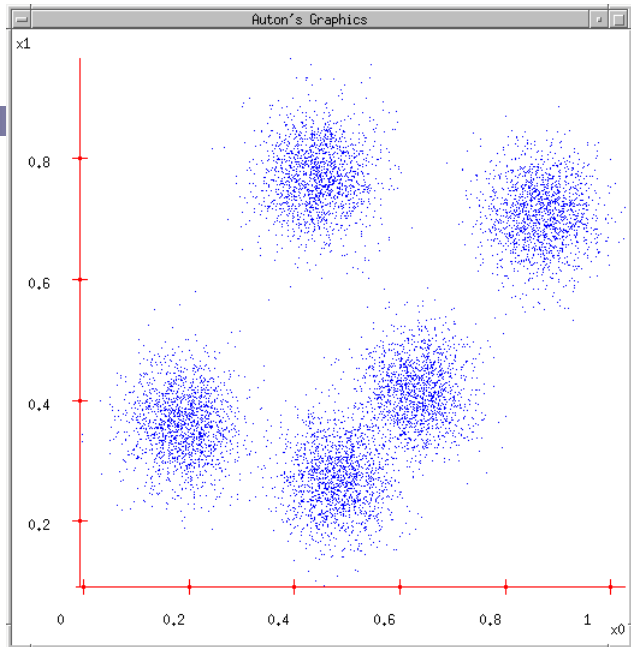
x ← "learn something"

Hand Cluster
split into
5 sets



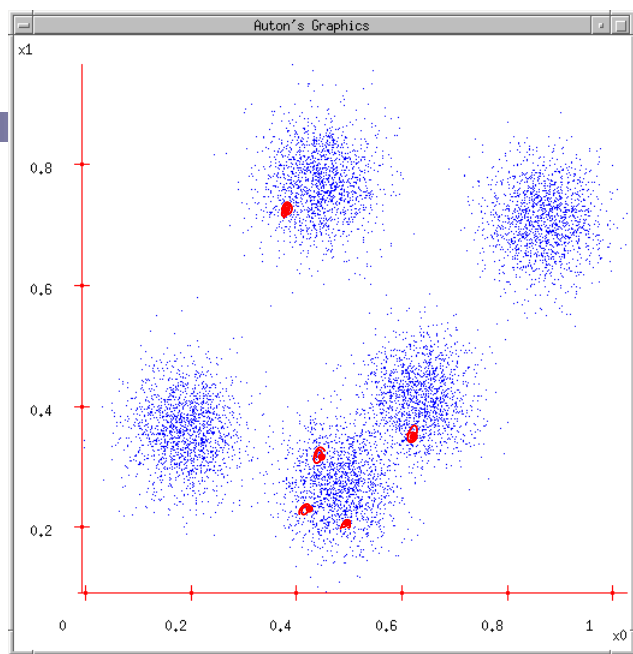
K-means

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



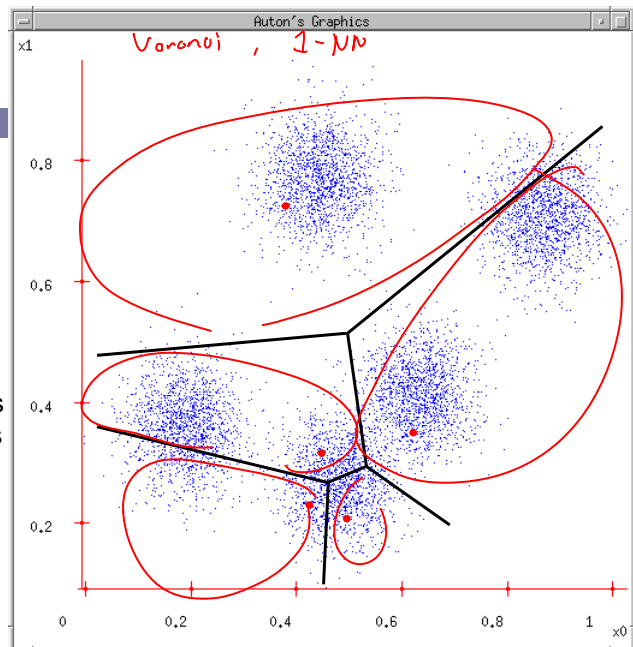
K-means

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



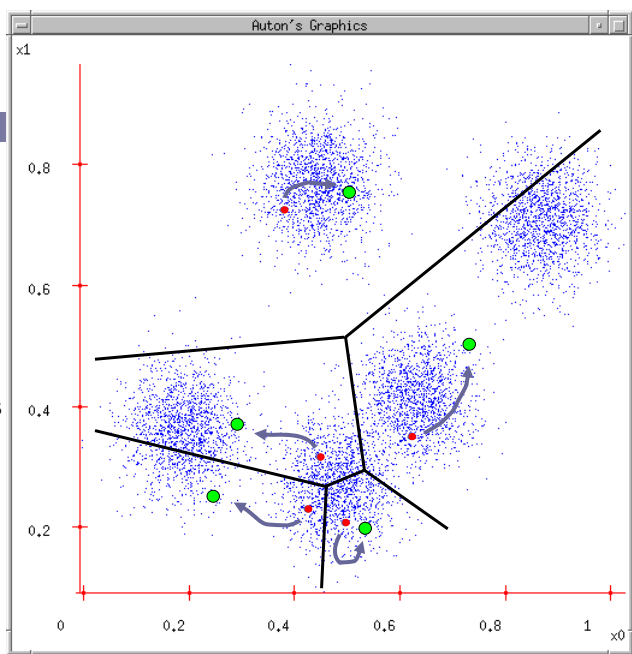
K-means

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)



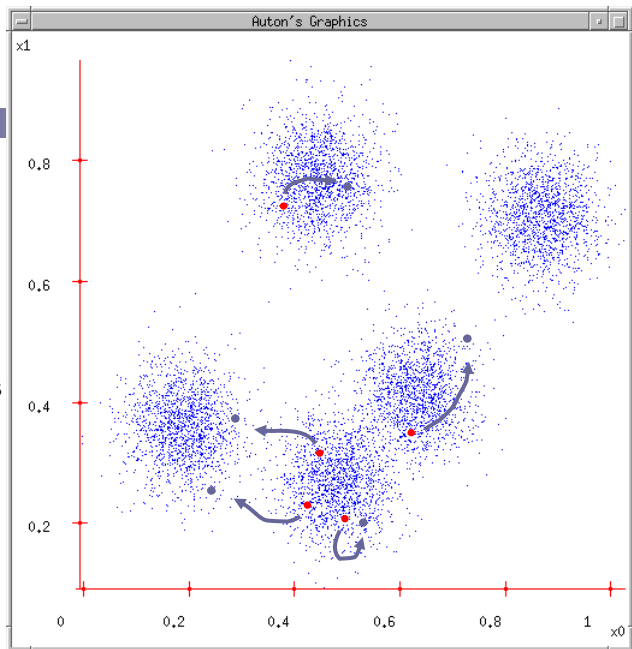
K-means

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

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...
5. ...and jumps there
6. ...Repeat until terminated!



K-means

- Randomly initialize k centers

- $\mu^{(0)} = \mu_1^{(0)}, \dots, \mu_k^{(0)}$

- Classify: Assign each point $j \in \{1, \dots, m\}$ to nearest center: *center of point j is closest to j*

- $C^{(t)}(j) \leftarrow \arg \min_i \|\mu_i - x_j\|^2$

- Recenter: μ_i becomes centroid of its point:

- $\mu_i^{(t+1)} \leftarrow \arg \min_{\mu} \sum_{j: C(j)=i} \|\mu - x_j\|^2$ *opt $\mu_i = \frac{\sum_{j: C(j)=i} x_j}{\sum_{j: C(j)=i} 1}$ is the mean!*

- Equivalent to $\mu_i \leftarrow$ average of its points!

What is K-means optimizing?

- Potential function $F(\mu, C)$ of centers μ and point allocations C :

- $F(\mu, C) = \sum_{j=1}^m \|\mu_{C(j)} - x_j\|^2$ *distance between x_j and its center $\mu_{C(j)}$*

- Optimal K-means:

- $\min_{\mu} \min_C F(\mu, C)$

Does K-means converge??? Part 1

- Optimize potential function:

$$\min_{\mu} \min_C F(\mu, C) = \min_{\mu} \min_C \sum_{i=1}^k \sum_{j: C(j)=i} \|\mu_i - x_j\|^2$$

distance

- Fix μ , optimize C $\mu = \bar{\mu}$

$$\begin{aligned} & \min_C \sum_{i=1}^k \sum_{j: C(j)=i} \|\bar{\mu}_i - x_j\|^2 \\ &= \min_C \sum_{j=1}^m \|\bar{\mu}_{C(j)} - x_j\|^2 \end{aligned}$$

optimize each center independently

$$= \sum_{j=1}^m \min_C \|\bar{\mu}_{C(j)} - x_j\|^2 \Rightarrow C(j) = \underset{i}{\operatorname{argmin}} \|\bar{\mu}_i - x_j\|^2$$

(classify step)

Does K-means converge??? Part 2

- Optimize potential function:

$$\min_{\mu} \min_C F(\mu, C) = \min_{\mu} \min_C \sum_{i=1}^k \sum_{j: C(j)=i} \|\mu_i - x_j\|^2$$

- Fix C, optimize μ

$$\begin{aligned} & \min_{\mu} \sum_{i=1}^k \sum_{j: C(j)=i} \|\mu_i - x_j\|^2 \\ &= \sum_{i=1}^k \min_{\mu_i} \sum_{j: C(j)=i} \|\mu_i - x_j\|^2 \end{aligned}$$

μ_i is mean of points in cluster i

\Rightarrow *re center in k-means.*

Coordinate descent algorithms

$$\min_{\mu} \min_C F(\mu, C) = \min_{\mu} \min_C \sum_{i=1}^k \sum_{j:C(j)=i} \|\mu_i - x_j\|^2$$

- Want: $\min_a \min_b F(a,b)$
- Coordinate descent:
 - fix a, minimize b
 - fix b, minimize a
 - repeat
- Converges!!!
 - if F is bounded
 - to a (often good) local optimum
 - as we saw in applet (play with it!)
- K-means is a coordinate descent algorithm!

$F(a,b)$



~~$F(a,b)$~~

