

Only covered supervised learning $X \rightarrow \mathbb{R}$ regression
 $X \rightarrow \{0, 1, \dots, k\}$ classification

Training data included labels

Clustering K-means

Machine Learning – CSE546
 Carlos Guestrin
 University of Washington

November 4, 2014
 ©Carlos Guestrin 2005-2014

Clustering images

Set of Images

given no labels

organize data into themes

beaches

flowers

C_1

C_2

C_3

C_4

C_5

©Carlos Guestrin 2005-2014 [Goldberger et al.] 2

K-means

d-dim vectors

- Randomly initialize k centers *or "smartly"*

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

Repeat until convergence: no point changes, cluster membership

- **Classify:** Assign each point $j \in \{1, \dots, N\}$ to nearest center:

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

fix μ , OPT C

- **Recenter:** $\mu_i^{(t+1)}$ becomes centroid of its point:

- $\mu_i^{(t+1)} \leftarrow \arg \min_{\mu} \sum_{j: C^{(t)}(j)=i} \|\mu - x_j\|^2$

sum of points in cluster i

fix C , OPT μ

$$\mu_i^{(t+1)} = \frac{\sum_{j: C^{(t)}(j)=i} x_j}{|\{j: C^{(t)}(j)=i\}|}$$

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

©Carlos Guestrin 2005-2014

3

Mixtures of Gaussians

Machine Learning – CSE546

Carlos Guestrin

University of Washington

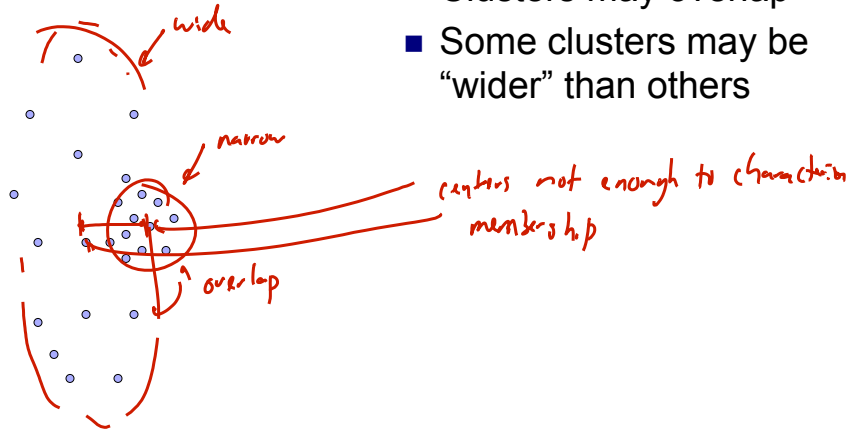
November 4, 2014

©Carlos Guestrin 2005-2014

4

(One) bad case for k-means

- Clusters may overlap
- Some clusters may be "wider" than others

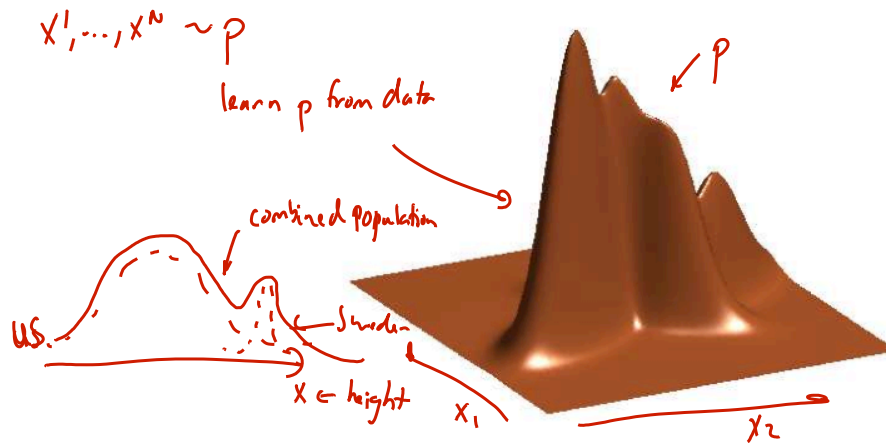


©Carlos Guestrin 2005-2014

5

Density Estimation

- Estimate a density based on x^1, \dots, x^N



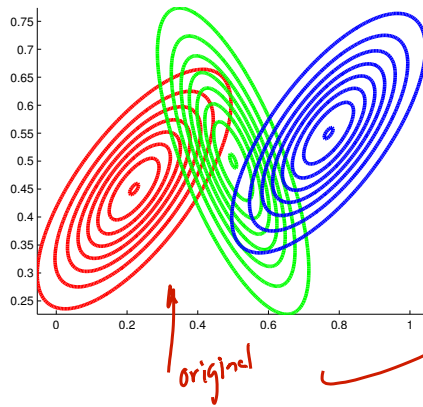
©Carlos Guestrin 2005-2014

6

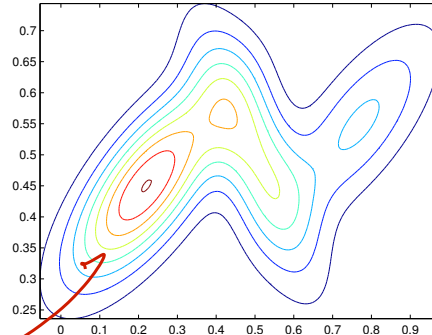
Density as Mixture of Gaussians

- Approximate density with a mixture of Gaussians

Mixture of 3 Gaussians



Contour Plot of Joint Density



p sum with weights π_i of each Gaussian

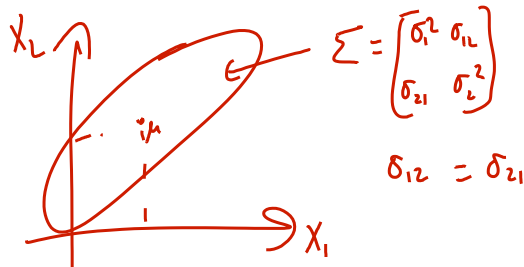
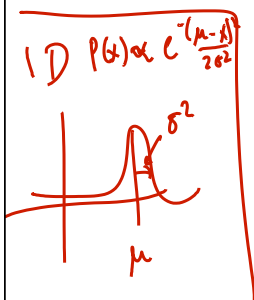
©Carlos Guestrin 2005-2014

7

Gaussians in d Dimensions

$$P(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} \|\Sigma\|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)\right]$$

mean vector
covariance matrix



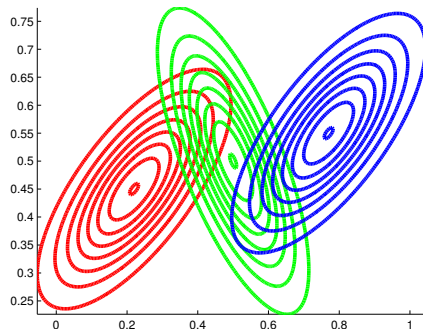
©Carlos Guestrin 2005-2014

8

Density as Mixture of Gaussians

- Approximate density with a mixture of Gaussians $\pi_1, \pi_2, \dots, \pi_k$

Mixture of 3 Gaussians



$$p(x^i | \pi, \mu, \Sigma) = \sum_{i=1}^k \pi_i N(x^i | \mu_i, \Sigma_i)$$

Handwritten notes:
 - $\pi_1, \pi_2, \dots, \pi_k$ are weights, $\sum_{i=1}^k \pi_i = 1$
 - $N(x^i | \mu_i, \Sigma_i)$ is the target density
 - $\ln 10$ is written near the plot

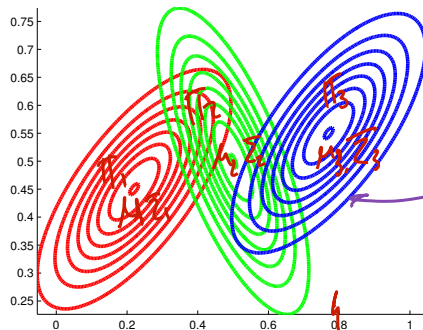
©Carlos Guestrin 2005-2014

9

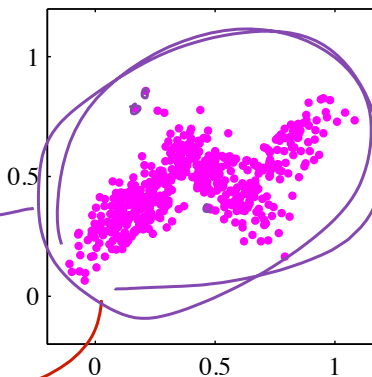
Density as Mixture of Gaussians

- Approximate with density with a mixture of Gaussians

Mixture of 3 Gaussians



Our actual observations

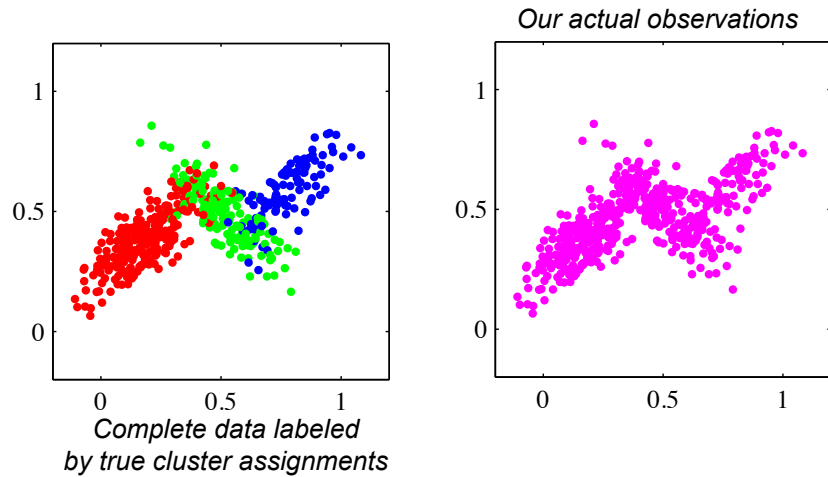


Handwritten notes:
 - μ_1, μ_2, μ_3 are written near the contours in the left plot.
 - "recon original densities How??" is written below the scatter plot.

C. Bishop, Pattern Recognition & Machine Learning

Clustering our Observations

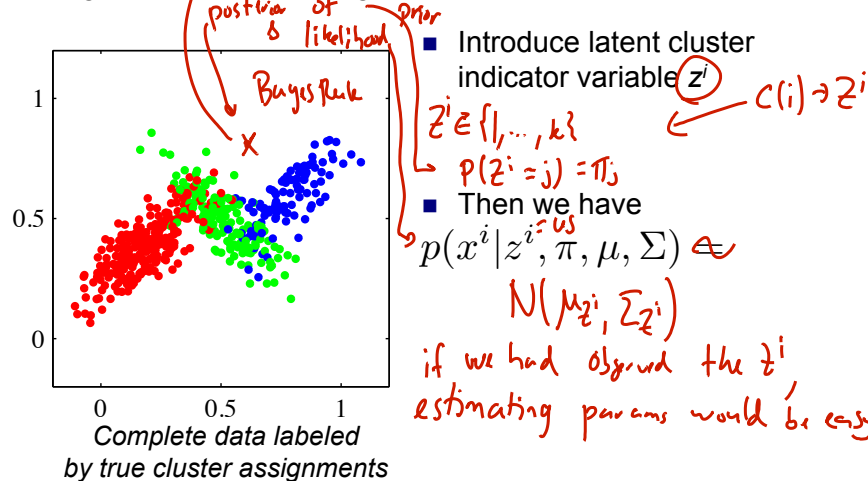
- Imagine we have an assignment of each x^i to a Gaussian



C. Bishop, Pattern Recognition & Machine Learning

Clustering our Observations

- Imagine we have an assignment of each x^i to a Gaussian



$p(z^i | x^i, \pi, \mu, \Sigma) = \frac{p(x^i | z^i, \pi, \mu, \Sigma) p(z^i | \pi)}{p(x^i | \pi, \mu, \Sigma)}$
 $z^i \in \{1, \dots, k\}$
 $p(z^i = j) = \pi_j$
 $p(x^i | z^i = j, \pi, \mu, \Sigma) \propto N(\mu_j^i, \Sigma_j^i)$
 if we had observed the z^i , estimating params would be easy

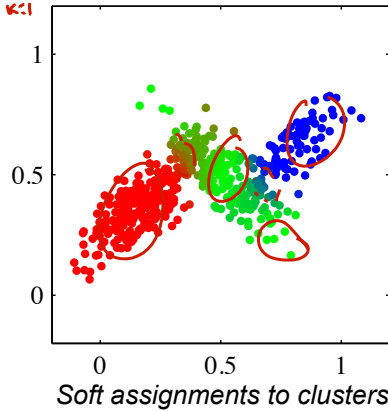
C. Bishop, Pattern Recognition & Machine Learning

Clustering our Observations

$$P(z|x) = \frac{P(x|z)P(z)}{P(x)}$$

- We must infer the cluster assignments from the observations

$$\sum_{k=1}^K v_{ik} = 1$$



Soft assignments to clusters

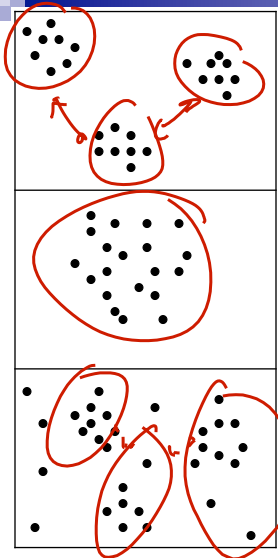
- Posterior probabilities of assignments to each cluster *given* model parameters:

$$r_{ik} = p(z^i = k | x^i, \pi, \mu, \Sigma) = \frac{\pi_k N(x^i | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x^i | \mu_j, \Sigma_j)}$$

if I had params, the "classify" step is simple.

C. Bishop, Pattern Recognition & Machine Learning

Unsupervised Learning: not as hard as it looks



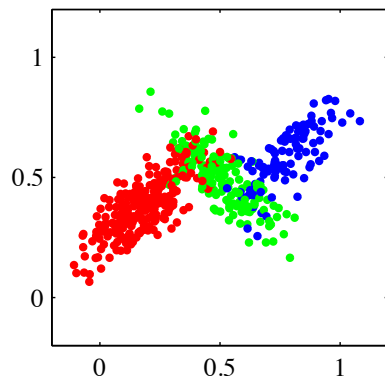
Sometimes easy

Sometimes impossible

and sometimes in between

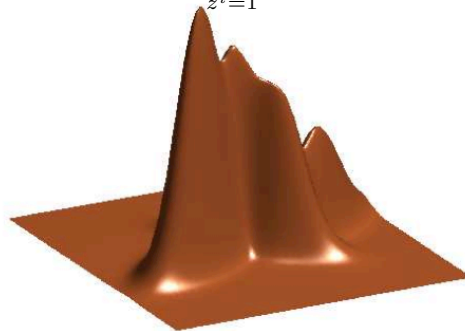
Summary of GMM Concept

- Estimate a density based on x^1, \dots, x^N



Complete data labeled by true cluster assignments

$$p(x^i | \pi, \mu, \Sigma) = \sum_{z^i=1}^K \pi_{z^i} \mathcal{N}(x^i | \mu_{z^i}, \Sigma_{z^i})$$



Surface Plot of Joint Density, Marginalizing Cluster Assignments

©Carlos Guestrin 2005-2014

15

Summary of GMM Components

- Observations $x^i \in \mathbb{R}^d, \quad i = 1, 2, \dots, N$
- Hidden cluster labels $z_i \in \{1, 2, \dots, K\}, \quad i = 1, 2, \dots, N$
- Hidden mixture means $\mu_k \in \mathbb{R}^d, \quad k = 1, 2, \dots, K$
- Hidden mixture covariances $\Sigma_k \in \mathbb{R}^{d \times d}, \quad k = 1, 2, \dots, K$
- Hidden mixture probabilities $\pi_k, \quad \sum_{k=1}^K \pi_k = 1$

Gaussian mixture marginal and conditional likelihood :

$$p(x^i | \pi, \mu, \Sigma) = \sum_{z^i=1}^K \pi_{z^i} p(x^i | z^i, \mu, \Sigma)$$

$$p(x^i | z^i, \mu, \Sigma) = \mathcal{N}(x^i | \mu_{z^i}, \Sigma_{z^i})$$

©Carlos Guestrin 2005-2014

16

*- coord. ascent alg.
- generalizes k-means*

Expectation Maximization

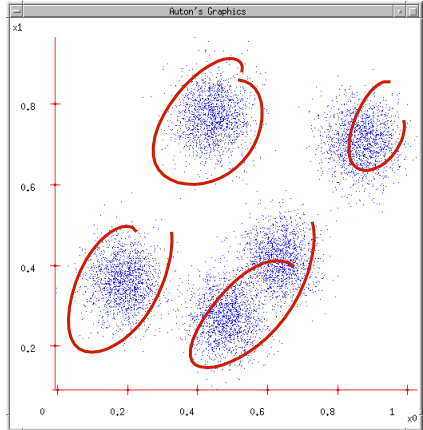
Machine Learning – CSE546
Carlos Guestrin
University of Washington

November 6, 2014
©Carlos Guestrin 2005-2014

17

Next... back to Density Estimation

What if we want to do density estimation with multimodal or clumpy data?

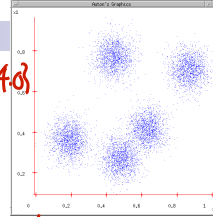


©Carlos Guestrin 2005-2014

18

But we don't see class labels!!!

- MLE: $\text{argmax}_{\pi, \mu, \Sigma} \prod_i P(z^i, x^i)$
 - in classification: $x^i = \{GPA=3.9, NLGrade=4.0\}$
 - $z^i = \{role = VP\}$
 - simplex convex
 - "latent" or "nuisance" variable
- But we don't know z^i
- Maximize marginal likelihood: $\text{argmax}_{\pi, \mu, \Sigma} \prod_i P(x^i) = \text{argmax}_{\pi, \mu, \Sigma} \prod_i \sum_{k=1}^K P(z^i=k, x^i)$
 - only over observables
 - Sum out latent variables
 - Sum over role: $\{VP, DS, Barista, \dots\}$
 - weigh by prob of each role
 - non-convex & hard \approx coordinate or gradient descent



©Carlos Guestrin 2005-2014

19

Special case: spherical Gaussians and hard assignments

$$P(z^i = k, \mathbf{x}^i) = \frac{1}{(2\pi)^{m/2} \|\Sigma_k\|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x}^i - \mu_k)^T \Sigma_k^{-1} (\mathbf{x}^i - \mu_k)\right] P(z^i = k)$$

- If $P(X|z=k)$ is spherical, with same σ for all classes:

$$P(\mathbf{x}^i | z^i = k) \propto \exp\left[-\frac{1}{2\sigma^2} \|\mathbf{x}^i - \mu_k\|^2\right]$$

- If each x^i belongs to one class $C(i)$ (hard assignment), marginal likelihood:

$$P(x^i) = \prod_{i=1}^N \sum_{k=1}^K P(x^i, z^i = k) \propto \prod_{i=1}^N \exp\left[-\frac{1}{2\sigma^2} \|\mathbf{x}^i - \mu_{C(i)}\|^2\right]$$

- Same as K-means!!!

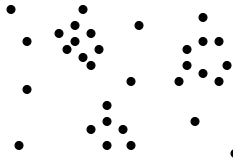
$$\max_{C, \mu} \sum_{i=1}^N -\frac{1}{2\sigma^2} \|\mathbf{x}^i - \mu_{C(i)}\|^2 \equiv \min_{C, \mu} \sum_{i=1}^N \|\mathbf{x}^i - \mu_{C(i)}\|^2$$

©Carlos Guestrin 2005-2014

20

EM: "Reducing" Unsupervised Learning to Supervised Learning

- If we knew assignment of points to classes → Supervised Learning!



- Expectation-Maximization (EM)

- **E** Guess assignment of points to classes
 - In standard ("soft") EM: each point associated with prob. of being in each class
 - **M** Recompute model parameters
 - Iterate
- coordinate descent like k-means*

©Carlos Guestrin 2005-2014

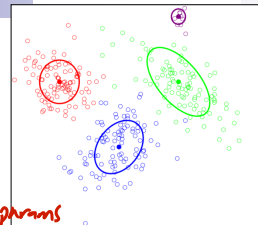
21

Generic Mixture Models

MoG Example:

- Observations: x^1, \dots, x^N with $x^i \in \mathbb{R}^d$ *eg. $x^i \in \mathbb{R}^d$*

- Parameters: $\pi = \{\pi_1, \dots, \pi_k\}$ *mix weights*
 $\phi = \{\phi_1, \dots, \phi_k\}$ *mix component params*
 e.g. $\phi_k = \{\mu_k, \Sigma_k\}$



- Likelihood:

$$P(x^i | \theta) = \sum_{k=1}^k \pi_k P(x^i | \phi_k) \leftarrow N(x^i | \mu_k, \Sigma_k)$$

- Ex. z^i = country of origin, x^i = height of i^{th} person
 - k^{th} mixture component = distribution of heights in country k

©Carlos Guestrin 2005-2014

22

ML Estimate of Mixture Model Params

- Log likelihood

$$L_x(\theta) \triangleq \log p(\{x^i\} | \theta) = \sum_{i=1}^N \log \sum_{z^i} p(x^i, z^i | \theta)$$

- Want ML estimate

$$\hat{\theta}^{ML} = \underset{\theta}{\operatorname{argmax}} L_x(\theta)$$

- Neither convex nor concave and local optima

If “complete” data were observed...

- Assume class labels z^i were observed in addition to x^i

$$\begin{aligned} \max_{\theta} L_{x,z}(\theta) &= \sum_{i=1}^N \log p(x^i, z^i | \theta) = \sum_{i=1}^N \log P(z^i | \theta) P(x^i | z^i, \theta) \\ \text{Simply (count)} &= \sum_{i=1}^N \log P(z^i | \theta) + \sum_{i=1}^N \log P(x^i | z^i, \theta) \end{aligned}$$

Params for each component $\approx N_k$

$$\hat{\pi}_k = \frac{\text{Count}(z^i=k)}{N}$$

- Compute ML estimates
 - Separates over clusters k !

- Example: mixture of Gaussians (MoG) $\theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$

$$\hat{\mu}_k = \frac{1}{N_k} \sum_{i: z^i=k} x^i, \text{ similarly for } \Sigma_k$$

Iterative Algorithm

- Motivates a coordinate ascent-like algorithm:

1. Infer missing values z^i given estimate of parameters $\hat{\theta}$
2. Optimize parameters to produce new $\hat{\theta}$ given "filled in" data z^i
3. Repeat

- Example: MoG (derivation soon...)

1. Infer "responsibilities"

$$r_{ik} = p(z^i = k | x^i, \hat{\theta}^{(t-1)}) = \frac{\pi_k P(x^i | \phi_k^{(t-1)})}{\sum_j \pi_j P(x^i | \phi_j^{(t-1)})}$$

weighted data

2. Optimize parameters

max w.r.t. π_k : $\hat{\pi}_k = \frac{1}{N} \sum_{i=1}^N r_{ik}$

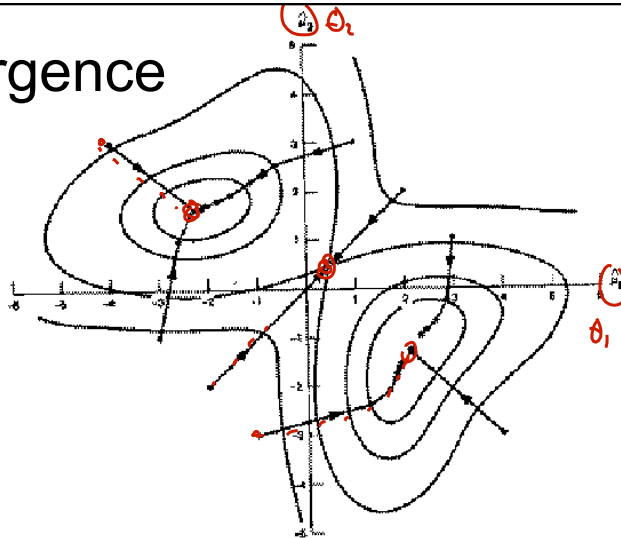
- max w.r.t. μ_k, Σ_k :

$$\hat{\mu}_k^{(t)} = \frac{1}{N} \sum_{i=1}^N r_{ik} x^i \quad \text{similarly for } \Sigma_k$$

weighted mean

E.M. Convergence

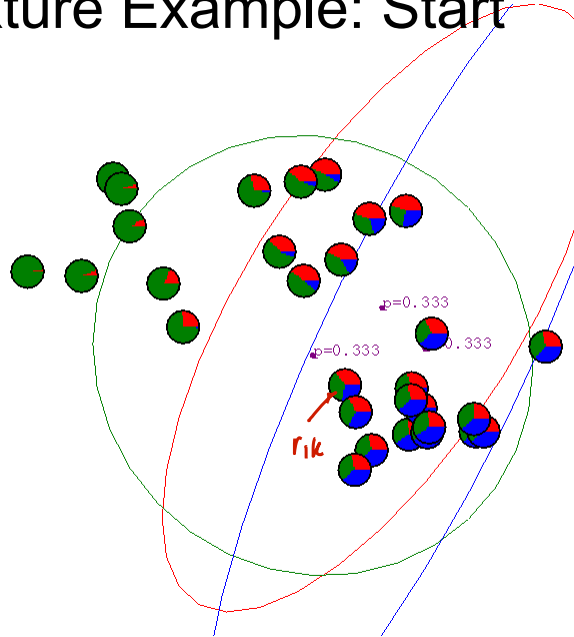
- EM is coordinate ascent on an interesting potential function
- Coord. ascent for bounded pot. func. → convergence to a local optimum guaranteed



- This algorithm is REALLY USED. And in high dimensional state spaces, too. E.G. Vector Quantization for Speech Data

Gaussian Mixture Example: Start

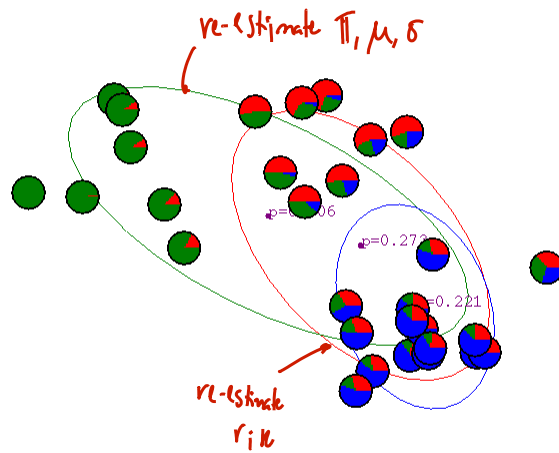
guess some $\theta^{(0)}$



©Carlos Guestrin 2005-2014

27

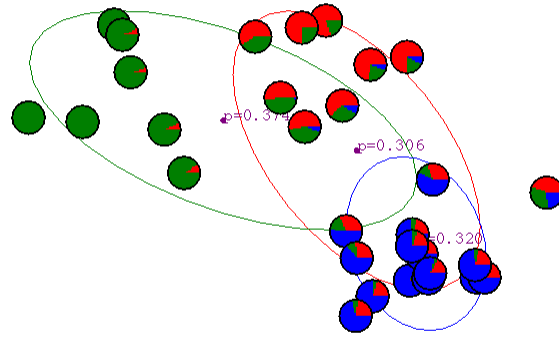
After first iteration



©Carlos Guestrin 2005-2014

28

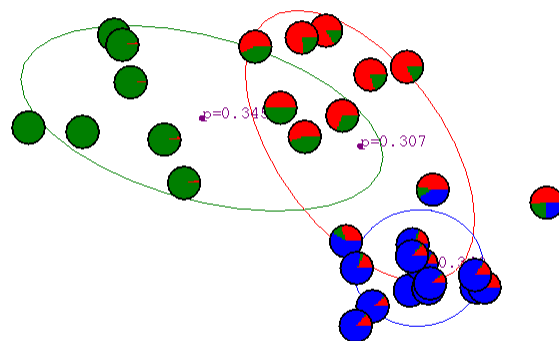
After 2nd iteration



©Carlos Guestrin 2005-2014

29

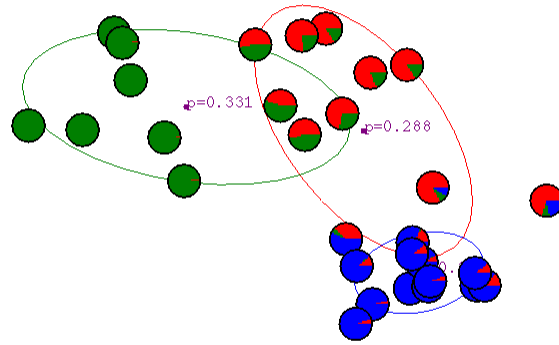
After 3rd iteration



©Carlos Guestrin 2005-2014

30

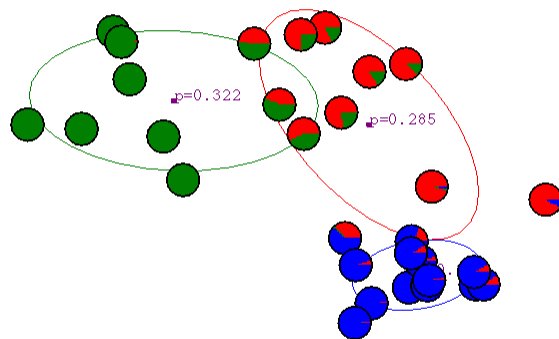
After 4th iteration



©Carlos Guestrin 2005-2014

31

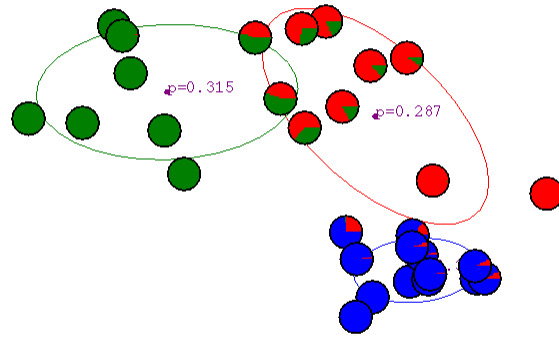
After 5th iteration



©Carlos Guestrin 2005-2014

32

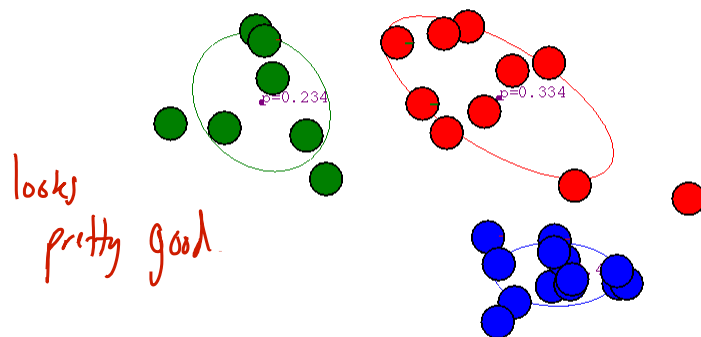
After 6th iteration



©Carlos Guestrin 2005-2014

33

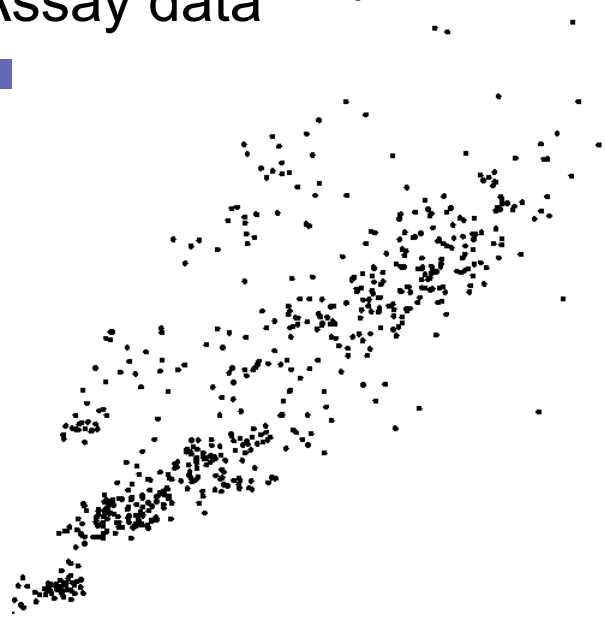
After 20th iteration



©Carlos Guestrin 2005-2014

34

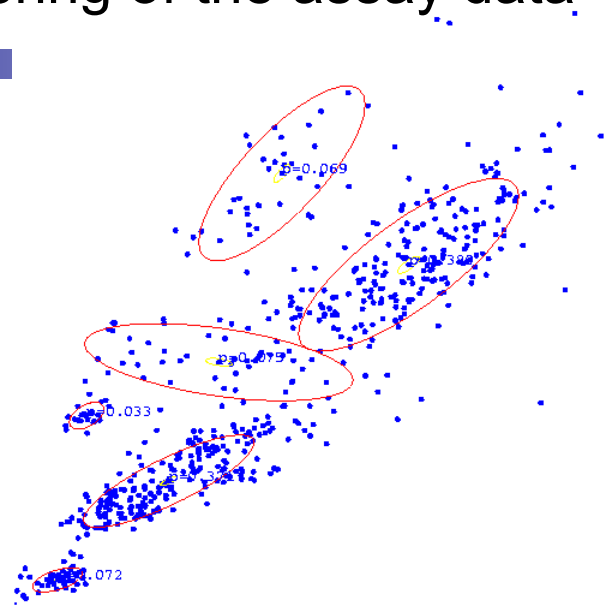
Some Bio Assay data



©Carlos Guestrin 2005-2014

35

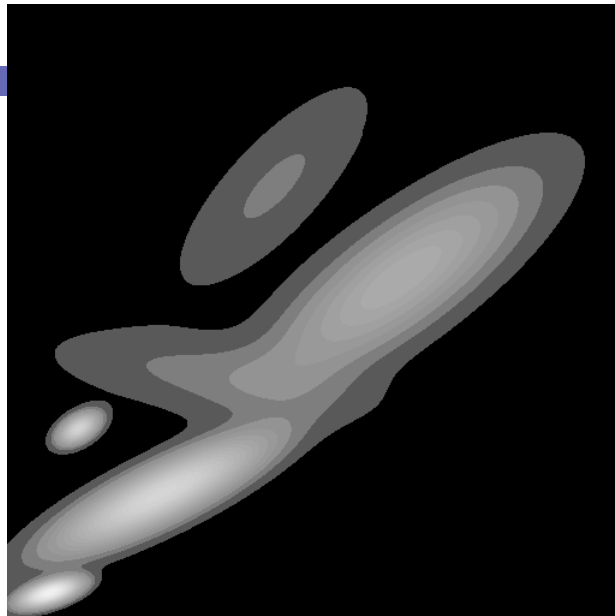
GMM clustering of the assay data



©Carlos Guestrin 2005-2014

36

Resulting Density Estimator



©Carlos Guestrin 2005-2014

37

E.M.: The General Case

- E.M. widely used beyond mixtures of Gaussians
 - The recipe is the same...
- Expectation Step: Fill in missing data, given current values of parameters, $\theta^{(t)}$
 - If variable y is missing (could be many variables)
 - Compute, for each data point \mathbf{x}^i , for each value i of y :
 - $P(y=i|\mathbf{x}^i, \theta^{(t)})$
- Maximization step: Find maximum likelihood parameters for (weighted) “completed data”:
 - For each data point \mathbf{x}^i , create k weighted data points
 -
 - Set $\theta^{(t+1)}$ as the maximum likelihood parameter estimate for this weighted data
- Repeat

©Carlos Guestrin 2005-2013

38

Initialization

- In mixture model case where $y^i = \{z^i, x^i\}$ there are many ways to initialize the EM algorithm
- Examples:
 - Choose K observations at random to define each cluster. Assign other observations to the nearest “centroid” to form initial parameter estimates
 - Pick the centers sequentially to provide good coverage of data
 - Grow mixture model by splitting (and sometimes removing) clusters until K clusters are formed
- Can be quite important to quality of solution in practice

What you should know

- K-means for clustering:
 - algorithm
 - converges because it's coordinate ascent
- EM for mixture of Gaussians:
 - How to “learn” maximum likelihood parameters (locally max. like.) in the case of unlabeled data
- Remember, E.M. can get stuck in local minima, and empirically it DOES
- EM is coordinate ascent

Expectation Maximization (EM) – Setup

- More broadly applicable than just to mixture models considered so far

- Model: x observable – “incomplete” data
 y not (fully) observable – “complete” data
 θ parameters

- Interested in maximizing (wrt θ):

$$\max_{\theta} p(x | \theta) = \sum_y p(x, y | \theta)$$

- Special case:

$$x = g(y)$$

$y = \begin{bmatrix} z \\ x \end{bmatrix}$ hidden cluster labels
 observation

Expectation Maximization (EM) – Derivation

- Step 1
 - Rewrite desired likelihood in terms of complete data terms

$$p(y | \theta) = p(y | x, \theta) p(x | \theta)$$

$$\Rightarrow \log P(x | \theta) = \log P(y | \theta) - \log P(y | x, \theta)$$

what we want to max *quantity we want to maximize*

- Step 2
 - Assume estimate of parameters $\hat{\theta}$
 - Take expectation with respect to $p(y | x, \hat{\theta})$

don't observe y , so

average over y
 $E[\cdot | x, \hat{\theta}]$

$$L_x(\theta) = \underbrace{E[\log P(y | \theta) | x, \hat{\theta}]}_{V(\theta, \hat{\theta})} + \underbrace{E[-\log P(y | x, \theta) | x, \hat{\theta}]}_{V(\theta, \hat{\theta})}$$

Expectation Maximization (EM) – Derivation

Jensen's inequality: f is concave
 $E[f(x)] \leq f(E[x])$

Step 3

new θ we are trying to maximize

θ from previous iteration

Consider log likelihood of data at any θ relative to log likelihood at $\hat{\theta}$

$$L_x(\theta) - L_x(\hat{\theta}) = [U(\theta, \hat{\theta}) - U(\hat{\theta}, \hat{\theta})] + [V(\theta, \hat{\theta}) - V(\hat{\theta}, \hat{\theta})]$$

focus on max this ≥ 0

Aside: **Gibbs Inequality** $E_p[\log p(x)] \geq E_p[\log q(x)]$

Proof:

?

↓

Expectation Maximization (EM) – Derivation

want to max

$$L_x(\theta) - L_x(\hat{\theta}) = [U(\theta, \hat{\theta}) - U(\hat{\theta}, \hat{\theta})] - [V(\theta, \hat{\theta}) - V(\hat{\theta}, \hat{\theta})]$$

focus ≥ 0

Step 4

Determine conditions under which log likelihood at θ exceeds that at $\hat{\theta}$

Using Gibbs inequality:

$$V(\theta, \hat{\theta}) = E[-\log P(y|x, \theta) | x, \hat{\theta}] \geq E[-\log P(y|x, \hat{\theta}) | x, \hat{\theta}] = V(\hat{\theta}, \hat{\theta})$$

If $U(\theta, \hat{\theta}) \geq U(\hat{\theta}, \hat{\theta})$

Then

$$L_x(\theta) \geq L_x(\hat{\theta})$$

making some progress

choose θ such that U increases

Motivates EM Algorithm

- Initial guess: $\hat{\theta}^{(0)}$
- Estimate at iteration t : $\hat{\theta}^{(t)}$
- **E-Step**
Compute $U(\theta, \hat{\theta}^{(t)}) = E[\log p(y|\theta) | x, \hat{\theta}^{(t)}]$
- **M-Step**
Compute

©Carlos Guestrin 2005-2014

45

Example – Mixture Models

- **E-Step** Compute $U(\theta, \hat{\theta}^{(t)}) = E[\log p(y | \theta) | x, \hat{\theta}^{(t)}]$
- **M-Step** Compute $\hat{\theta}^{(t+1)} = \arg \max_{\theta} U(\theta, \hat{\theta}^{(t)})$

- Consider $y^i = \{z^i, x^i\}$ i.i.d.

$$p(x^i, z^i | \theta) = \pi_{z^i} p(x^i | \phi_{z^i}) =$$

$$E_{q_t}[\log p(y | \theta)] = \sum_i E_{q_t}[\log p(x^i, z^i | \theta)] =$$

©Carlos Guestrin 2005-2014

46

Coordinate Ascent Behavior

- Bound log likelihood:

$$\begin{aligned} L_x(\theta) &= U(\theta, \hat{\theta}^{(t)}) + V(\theta, \hat{\theta}^{(t)}) \\ &\geq \\ L_x(\hat{\theta}^{(t)}) &= U(\hat{\theta}^{(t)}, \hat{\theta}^{(t)}) + V(\hat{\theta}^{(t)}, \hat{\theta}^{(t)}) \end{aligned}$$

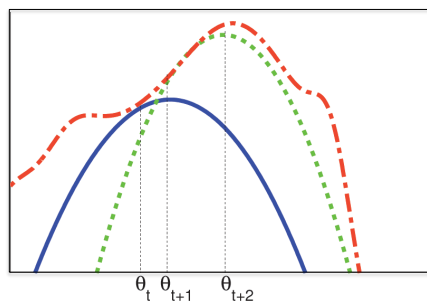


Figure from
KM textbook

©Carlos Guestrin 2005-2014

47

Comments on EM

- Since Gibbs inequality is satisfied with equality only if $p=q$, any step that changes θ should strictly **increase likelihood**
- In practice, can replace the **M-Step** with increasing U instead of maximizing it (**Generalized EM**)
- Under certain conditions (e.g., in exponential family), can show that EM **converges to a stationary point** of $L_x(\theta)$
- Often there is a **natural choice for y** ... has physical meaning
- If you want to choose any y , not necessarily $x=g(y)$, replace $p(y | \theta)$ in U with $p(y, x | \theta)$

©Carlos Guestrin 2005-2014

48

Initialization

- In mixture model case where $y^i = \{z^i, x^i\}$ there are many ways to initialize the EM algorithm
- Examples:
 - Choose K observations at random to define each cluster. Assign other observations to the nearest “centroid” to form initial parameter estimates
 - Pick the centers sequentially to provide good coverage of data
 - Grow mixture model by splitting (and sometimes removing) clusters until K clusters are formed
- Can be quite important to convergence rates in practice

©Carlos Guestrin 2005-2014

49

What you should know

- K-means for clustering:
 - algorithm
 - converges because it's coordinate ascent
- EM for mixture of Gaussians:
 - How to “learn” maximum likelihood parameters (locally max. like.) in the case of unlabeled data
- Be happy with this kind of probabilistic analysis
- Remember, E.M. can get stuck in local minima, and empirically it DOES
- EM is coordinate ascent

©Carlos Guestrin 2005-2014

50