

Mixtures of Gaussians

Machine Learning – CSE546

Sham Kakade

University of Washington

November 15, 2016

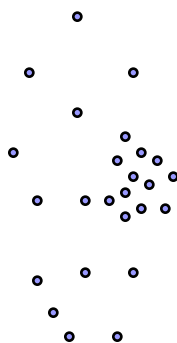
©Sham Kakade 2016

©Sham Kakade 2016

1

(One) bad case for k-means

- Clusters may overlap
- Some clusters may be “wider” than others



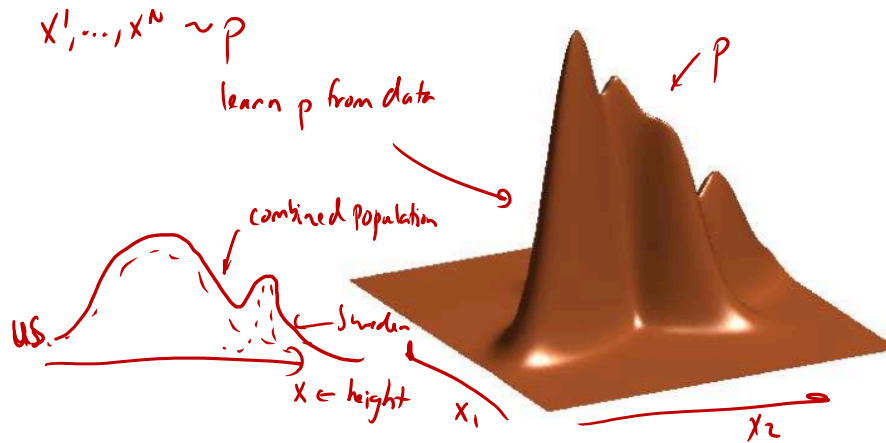
©Sham Kakade 2016

©Sham Kakade 2016

2

Density Estimation

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



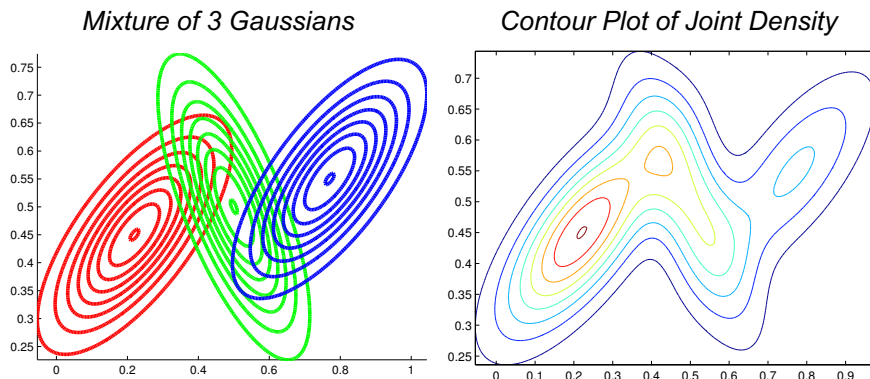
©Sham Kakade 2016

©Sham Kakade 2016

3

Density as Mixture of Gaussians

- Approximate density with a mixture of Gaussians



©Sham Kakade 2016

©Sham Kakade 2016

4

Gaussians in d Dimensions

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

©Sham Kakade 2016

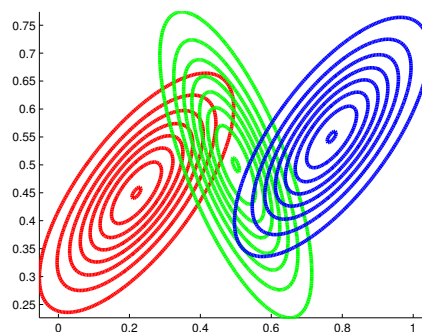
©Sham Kakade 2016

5

Density as Mixture of Gaussians

- Approximate density with a mixture of Gaussians

Mixture of 3 Gaussians



$$p(x^i | \pi, \mu, \Sigma) =$$

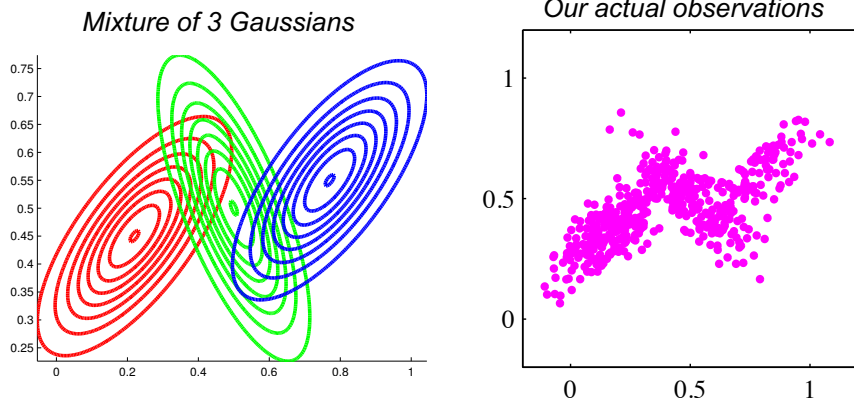
©Sham Kakade 2016

©Sham Kakade 2016

6

Density as Mixture of Gaussians

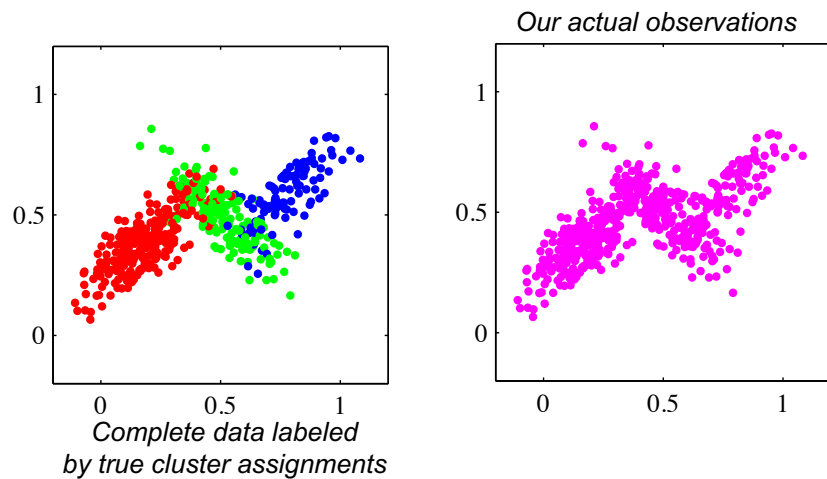
- Approximate with density with a mixture of Gaussians



C. Bishop, *Pattern Recognition & Machine Learning*

Clustering our Observations

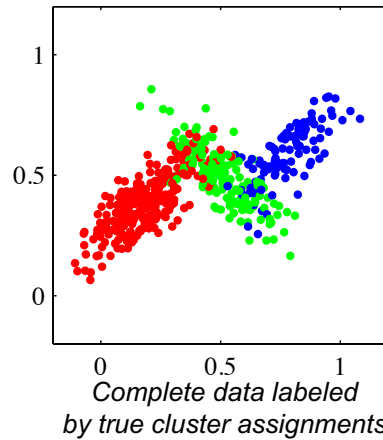
- Imagine we have an assignment of each x^i to a Gaussian



C. Bishop, *Pattern Recognition & Machine Learning*

Clustering our Observations

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



- Introduce latent cluster indicator variable z^i

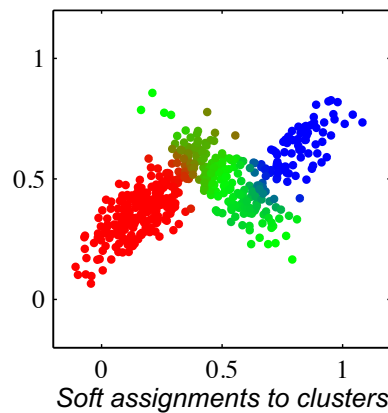
- Then we have
$$p(x^i | z^i, \pi, \mu, \Sigma) =$$

©Sham Kakade 2016

C. Bishop, *Pattern Recognition & Machine Learning*

Clustering our Observations

- We must infer the cluster assignments from the observations

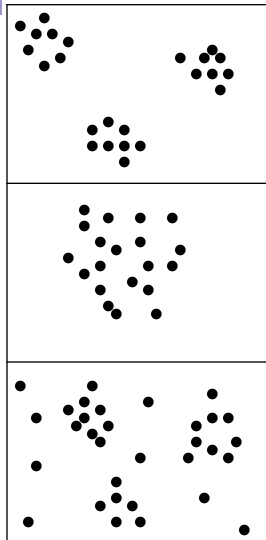


- Posterior probabilities of assignments to each cluster *given* model parameters:
$$r_{ik} = p(z^i = k | x^i, \pi, \mu, \Sigma) =$$

©Sham Kakade 2016

C. Bishop, *Pattern Recognition & Machine Learning*

Unsupervised Learning: not as hard as it looks



Sometimes easy

Sometimes impossible

and sometimes in between

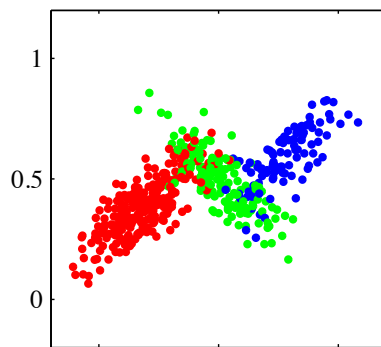
©Sham Kakade 2016

©Sham Kakade 2016

11

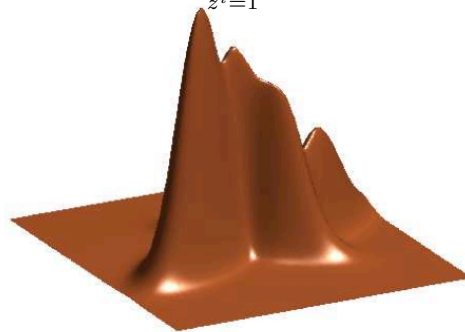
Summary of GMM Concept

- 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

©Sham Kakade 2016

©Sham Kakade 2016

12

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 :

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

©Sham Kakade 2016

©Sham Kakade 2016

13

Expectation Maximization

Machine Learning – CSE546

Sham Kakade

University of Washington

November 15, 2016

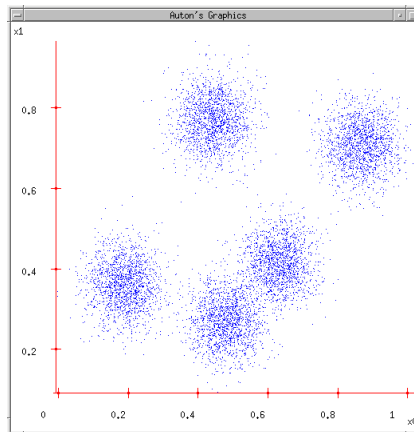
©Sham Kakade 2016

©Sham Kakade 2016

14

Next... back to Density Estimation

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



©Sham Kakade 2016

©Sham Kakade 2016

15

But we don't see class labels!!!

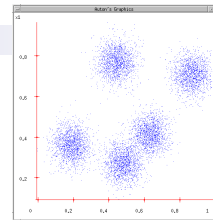
■ MLE:

□ $\operatorname{argmax} \prod_i P(z^i, x^i)$

■ But we don't know z^i

■ Maximize marginal likelihood:

□ $\operatorname{argmax} \prod_i P(x^i) = \operatorname{argmax} \prod_i \sum_{k=1}^K P(z^i=k, x^i)$



©Sham Kakade 2016

©Sham Kakade 2016

16

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 $P(\mathbf{X}|z=k)$ 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 \mathbf{x}^i belongs to one class $C(i)$ (hard assignment), marginal likelihood:

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

- Same as K-means!!!

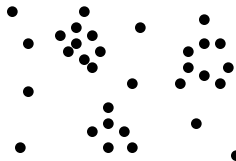
©Sham Kakade 2016

©Sham Kakade 2016

17

EM: “Reducing” Unsupervised Learning to Supervised Learning

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



- Expectation-Maximization (EM)

- Guess assignment of points to classes
 - In standard (“soft”) EM: each point associated with prob. of being in each class
- Recompute model parameters
- Iterate

©Sham Kakade 2016

©Sham Kakade 2016

18

Generic Mixture Models

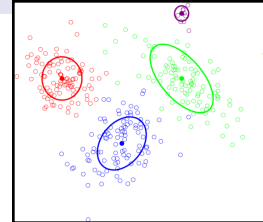
- Observations:

- Parameters:

- Likelihood:

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

MoG Example:



©Sham Kakade 2016

©Sham Kakade 2016

19

ML Estimate of Mixture Model Params

- Log likelihood

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

- Want ML estimate

$$\hat{\theta}^{ML} =$$

- Neither convex nor concave and local optima

©Sham Kakade 2016

©Sham Kakade 2016

20

If “complete” data were observed...

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

$$L_{x,z}(\theta) = \sum_i \log p(x^i, z^i | \theta)$$

- Compute ML estimates
 - Separates over clusters k !

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

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”

$$r_{ik} = p(z^i = k | x^i, \hat{\theta}^{(t-1)}) =$$

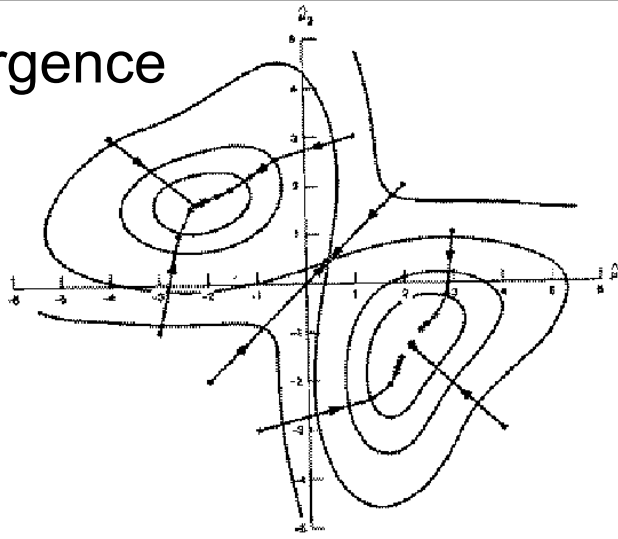
2. Optimize parameters

max w.r.t. π_k :

max w.r.t. μ_k, Σ_k :

E.M. Convergence

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



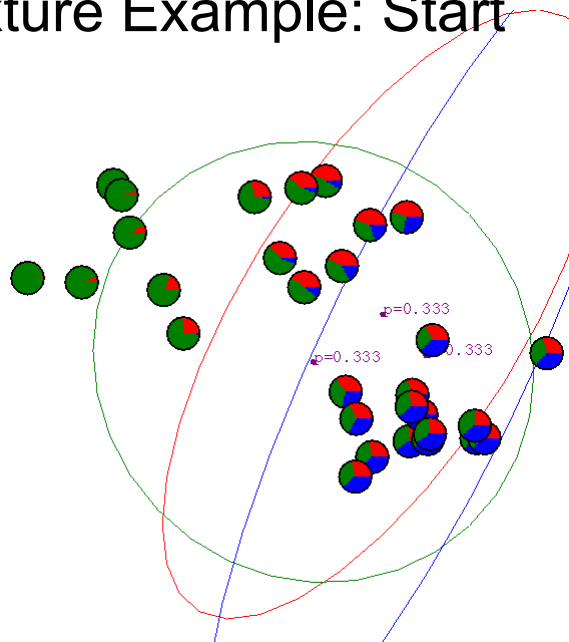
- This algorithm is REALLY USED. And in high dimensional state spaces, too. E.G. Vector Quantization for Speech Data

©Sham Kakade 2016

©Sham Kakade 2016

23

Gaussian Mixture Example: Start

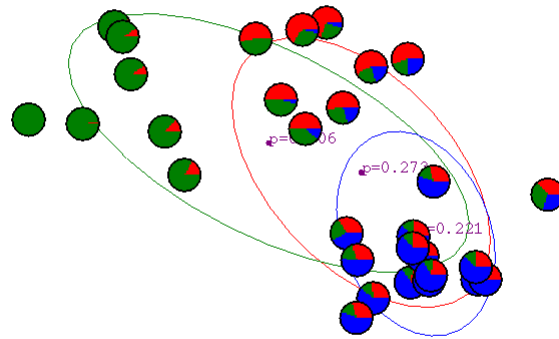


©Sham Kakade 2016

©Sham Kakade 2016

24

After first iteration

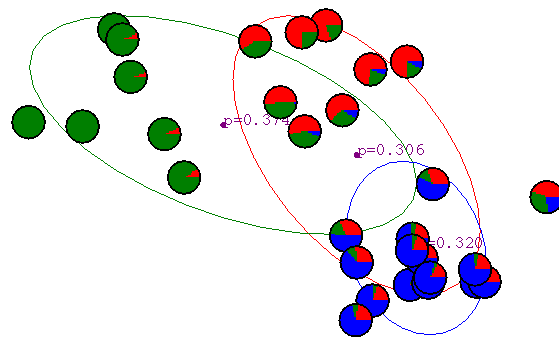


©Sham Kakade 2016

©Sham Kakade 2016

25

After 2nd iteration

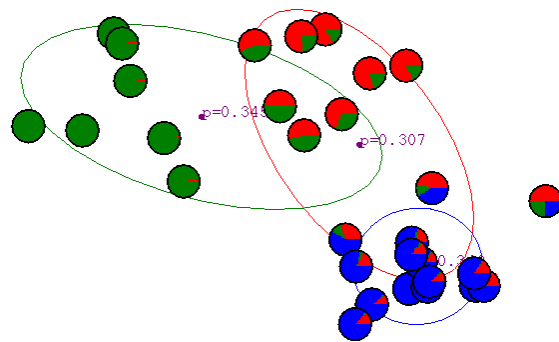


©Sham Kakade 2016

©Sham Kakade 2016

26

After 3rd iteration

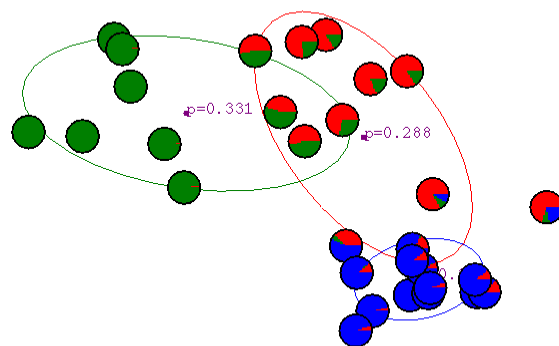


©Sham Kakade 2016

©Sham Kakade 2016

27

After 4th iteration

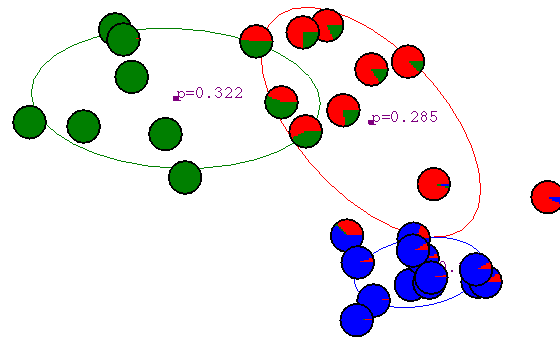


©Sham Kakade 2016

©Sham Kakade 2016

28

After 5th iteration

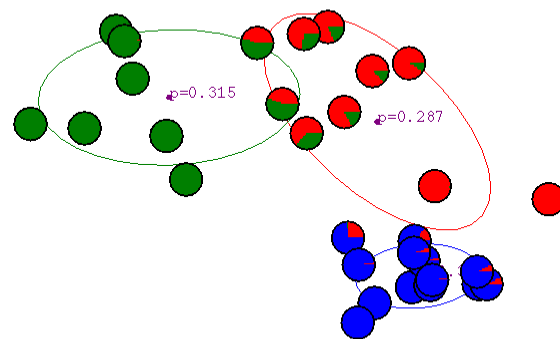


©Sham Kakade 2016

©Sham Kakade 2016

29

After 6th iteration

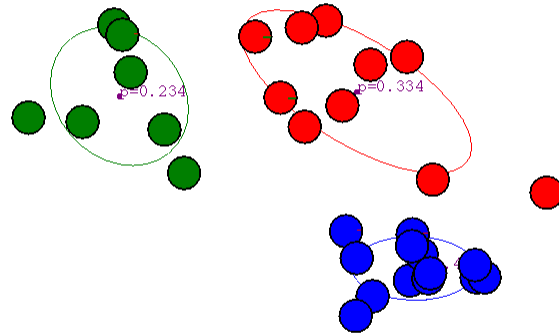


©Sham Kakade 2016

©Sham Kakade 2016

30

After 20th iteration

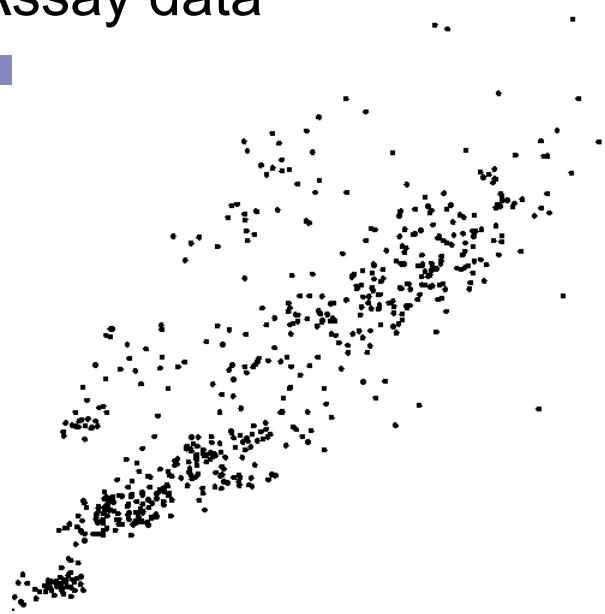


©Sham Kakade 2016

©Sham Kakade 2016

31

Some Bio Assay data

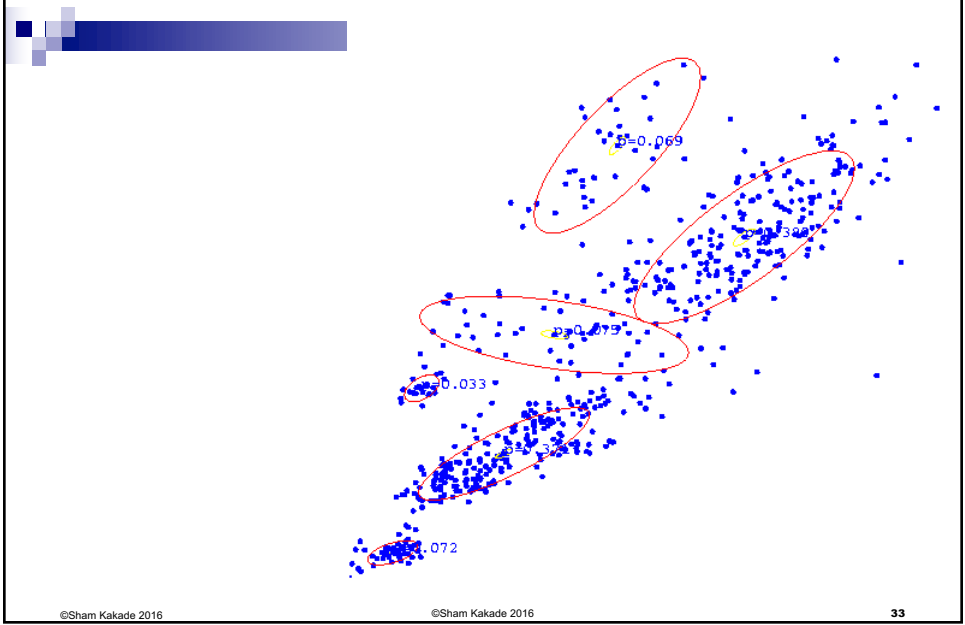


©Sham Kakade 2016

©Sham Kakade 2016

32

GMM clustering of the assay data



Resulting Density Estimator



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

©Sham Kakade 2016

©Sham Kakade 2016

35

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

©Sham Kakade 2016

©Sham Kakade 2016

36

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

©Sham Kakade 2016

©Sham Kakade 2016

37

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):
$$p(x | \theta) = \sum_y p(x, y | \theta)$$
- Special case:
$$x = g(y)$$

©Sham Kakade 2016

©Sham Kakade 2016

38

Expectation Maximization (EM) – Derivation

■ Step 1

- Rewrite desired likelihood in terms of complete data terms

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

■ Step 2

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

©Sham Kakade 2016

©Sham Kakade 2016

39

Expectation Maximization (EM) – Derivation

■ Step 3

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

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

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

Proof:

©Sham Kakade 2016

©Sham Kakade 2016

40

Expectation Maximization (EM) – Derivation

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

- Step 4
 - Determine conditions under which log likelihood at θ exceeds that at $\hat{\theta}$
Using Gibbs inequality:

If

Then

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

Motivates EM Algorithm

- Initial guess:
- Estimate at iteration t :

- **E-Step**

Compute

- **M-Step**

Compute

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 $p(y, x | \theta)$

©Sham Kakade 2016

©Sham Kakade 2016

45

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

©Sham Kakade 2016

©Sham Kakade 2016

46

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