# Clustering
# K-means

Machine Learning – CSE546

Emily Fox
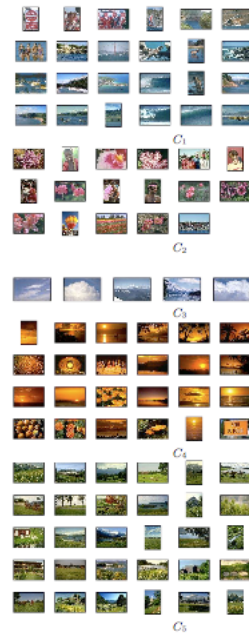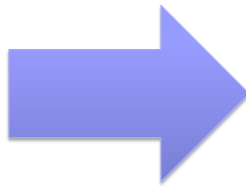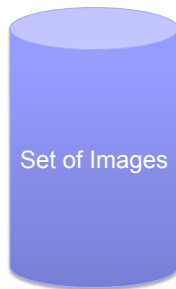
University of Washington

November 4, 2013

1

---

# Clustering images

[Goldberger et al.] 2

---

1

# Clustering web search results

3

# Some Data

4

# K-means

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

5

---

# K-means

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

2. Randomly guess k cluster Center locations

6

3

# 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)

7

# 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

8

4

# 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!

9

---

# K-means

- Randomly initialize *k* centers
  - □ $\mu^{(0)} = \mu_1^{(0)}, \ldots, \mu_k^{(0)}$

- **Classify**: Assign each point j∈{1,…N} to nearest center:
  - □ $C^{(t)}(j) \leftarrow \arg\min_i ||\mu_i - x_j||^2$

- **Recenter**: $\mu_i$ becomes centroid of its point:
  - □ $\mu_i^{(t+1)} \leftarrow \arg\min_\mu \sum_{j:C(j)=i} ||\mu - x_j||^2$

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

10

---

5

# What is K-means optimizing?

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

  □ $F(\mu, C) = \sum_{j=1}^{N} ||\mu_{C(j)} - x_j||^2$

- Optimal K-means:

  □ $\min_{\mu} \min_{C} F(\mu, C)$

11

# Does K-means converge??? Part 1

- Optimize potential function:

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

- Fix $\mu$, optimize C

12

# 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 $\mu$

13

# 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!)
      - □ (For LASSO it converged to the global optimum, because of convexity)

- K-means is a coordinate descent algorithm!

14

7

# Mixtures of Gaussians

Machine Learning – CSE546

Emily Fox

University of Washington

November 4, 2013

15

---

# (One) bad case for k-means

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

16

# Density Estimation

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

# Density Estimation

*Contour Plot of Joint Density*

9

# Density as Mixture of Gaussians

- Approximate density with a mixture of Gaussians

*Mixture of 3 Gaussians*                    *Contour Plot of Joint Density*

**19**

---

# Gaussians in *d* Dimensions

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

**20**

10

# Density as Mixture of Gaussians

- Approximate density with a mixture of Gaussians

*Mixture of 3 Gaussians*

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

# Density as Mixture of Gaussians

- Approximate with density with a mixture of Gaussians

*Mixture of 3 Gaussians*

*Our actual observations*

*C. Bishop, Pattern Recognition & Machine Learning*

# Clustering our Observations

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

*Our actual observations*



*Complete data labeled by true cluster assignments*
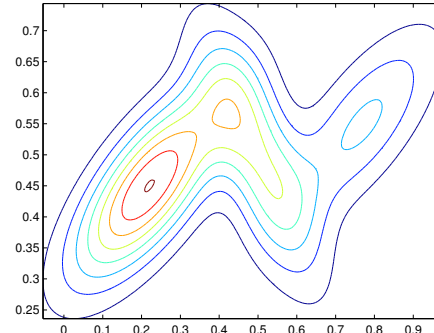
C. Bishop, Pattern Recognition & Machine Learning

---

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

*Complete data labeled by true cluster assignments*

C. Bishop, Pattern Recognition & Machine Learning

# Clustering our Observations

- We must infer the cluster assignments from the observations



0 0.5 1

*Soft assignments to clusters*

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

$$r_{ik} = p(z^i = k | x^i, \pi, \mu, \Sigma) =$$

---

# Unsupervised Learning:
# not as hard as it looks



Sometimes easy

Sometimes impossible

and sometimes in between

26

# Summary of GMM Concept

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

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



*Complete data labeled*
*by true cluster assignments*

*Surface Plot of Joint Density,*
*Marginalizing Cluster Assignments*

©Emily Fox 2013                    27

# Summary of GMM Components

- Observations $\qquad\qquad x^i \in \mathbb{R}^d, \quad i = 1, 2, \ldots, N$

- Hidden cluster labels $\quad z_i \in \{1, 2, \ldots, K\}, \quad i = 1, 2, \ldots, N$

- Hidden mixture means $\qquad\qquad \mu_k \in \mathbb{R}^d, \quad k = 1, 2, \ldots, K$

- Hidden mixture covariances $\quad \Sigma_k \in \mathbb{R}^{d \times d}, \quad k = 1, 2, \ldots, K$

- Hidden mixture probabilities $\qquad \pi_k, \quad \sum_{k=1}^{K} \pi_k = 1$

***Gaussian mixture marginal and conditional likelihood :***

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

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

©Emily Fox 2013                    28

# Expectation Maximization

Machine Learning – CSE546

Emily Fox

University of Washington

November 6, 2013

29

---

Next…  back to Density Estimation

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

30

# But we don't see class labels!!!

- MLE:
  - argmax $\prod_i P(z^i, x^i)$

- But we don't know $z^i$
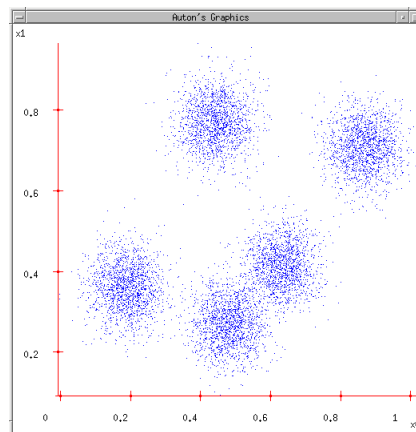- Maximize marginal likelihood:
  - argmax $\prod_i P(x^i)$ = argmax $\prod_i \sum_{k=1}^{K} P(z^i=k, x^i)$

---

# Special case: spherical Gaussians and hard assignments

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

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

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

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

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

- Same as K-means!!!

## Supervised Learning of Mixtures of Gaussians

- Mixtures of Gaussians:
  - Prior class probabilities: $P(z=k)$
  - Likelihood function per class: $P(\mathbf{x}|z=k)$

- Suppose, for each data point, we know location $\mathbf{x}$ and class $z$
  - Learning is easy… ☺

  - For prior $P(z)$

  - For likelihood function:

33

## 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

34

17

# Form of Likelihood

- Conditioned on class of point **x**$^i$...

$$p(x^i \mid z^i, \mu, \Sigma) =$$

- Marginalizing class assignment:

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

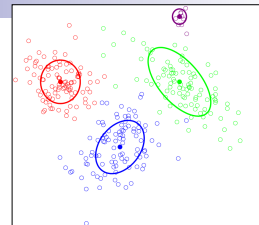# Gaussian Mixture Model

- Most commonly used mixture model
- 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 $k$

# Example

(Taken from Kevin Murphy's ML textbook)

- Data: gene expression levels
- Goal: cluster genes with similar expression trajectories



yeast microarray data

---

# Mixture models are useful for…

- Density estimation
    - □ Allows for multimodal density
- Clustering
    - □ Want membership information for each observation
        - e.g., topic of current document
    - □ Soft clustering:

$$p(z^i = k \mid x^i, \theta) =$$

    - □ Hard clustering:

$$z^{i*} = \arg \max_k p(z^i = k \mid x^i, \theta) =$$

# Issues

- Label switching
  - □ Color = label does not matter
  - □ Can switch labels and likelihood is unchanged



- Log likelihood is not convex in the parameters
  - □ Problem is simpler for "complete data likelihood"

---

# ML Estimate of Mixture Model Params

- Log likelihood

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

- Want ML estimate

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

- Neither convex nor concave and local optima

# If "complete" data were observed...

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

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

- Compute ML estimates
  - □ Separates over clusters *k*!

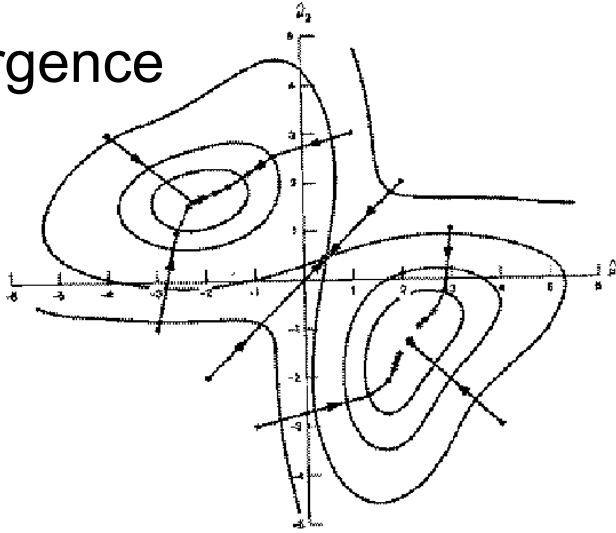- Example: mixture of Gaussians (MoG)  $\theta = \left\{ \pi_k, \mu_k, \Sigma_k \right\}_{k=1}^K$

41

---

# 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... + HW)
  1. Infer "responsibilities"

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

  2. Optimize parameters

$$\max \text{ w.r.t. } \pi_k :$$

$$\max \text{ w.r.t. } \mu_k, \Sigma_k :$$

42

# E.M. Convergence

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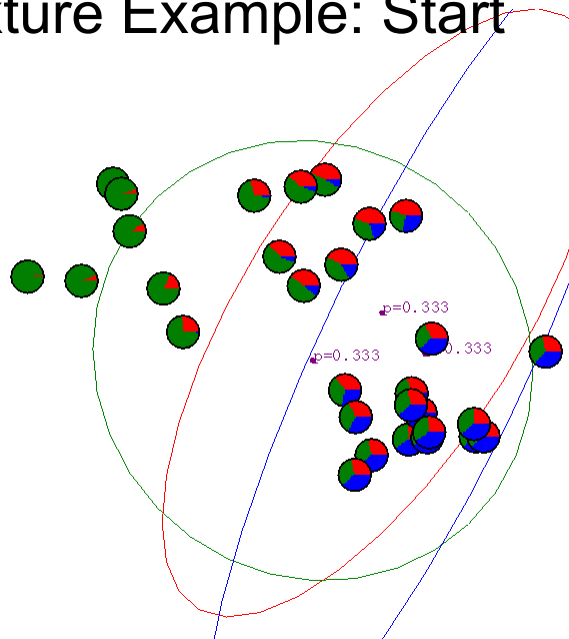- Coord. ascent for bounded pot. func. ➔ convergence to a local optimum guaranteed



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

43

# Gaussian Mixture Example: Start



p=0.333

p=0.333          0.333

44

22

# After first iteration

45

# After 2nd iteration

46

23

# After 3rd iteration

47

# After 4th iteration

48

# After 5th iteration



p=0.322

p=0.285

49

# After 6th iteration



p=0.315

p=0.287

50

# After 20th iteration

# Some Bio Assay data

# GMM clustering of the assay data



p=0.069

p=0.075

p=0.033

p=0.072

53

# Resulting Density Estimator

54

# 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*

    $\theta$  parameters

- Interested in maximizing (wrt  $\theta$ ):

$$p(x \mid \theta) = \sum_y p(x, y \mid \theta)$$

- Special case:

$$x = g(y)$$

# Expectation Maximization (EM) – Derivation

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

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

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

# Expectation Maximization (EM) – Derivation

- Step 3
  - Consider log likelihood of data at any $\theta$ relative to log likelihood at $\hat{\theta}$

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

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

57

---

# 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 $\theta$ exceeds that at $\hat{\theta}$
  Using Gibbs inequality:

  If

  Then

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

58

29

# Motivates EM Algorithm

- Initial guess:
- Estimate at iteration *t*:

- **E-Step**

  Compute

- **M-Step**

  Compute

---

# Example – Mixture Models

- **E-Step**  Compute     $U(\theta, \hat{\theta}^{(t)}) = E[\log p(y \mid \theta) \mid x, \hat{\theta}^{(t)}]$
- **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.

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

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

# Coordinate Ascent Behavior

- Bound log likelihood:

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

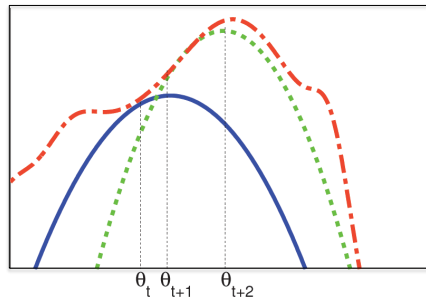

Figure from KM textbook

# Comments on EM

- Since Gibbs inequality is satisfied with equality only if *p=q*, any step that changes $\theta$ 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 $p(y \mid \theta)$ in *U* with $p(y, x \mid \theta)$

# 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 "centriod" 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

©Emily Fox 2013                                                                    63

# 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 2005-2013                                                       64