Unsupervised Learning with Mixtures of Gaussians

# EM Algorithm - continued

Machine Learning – CSE446

Carlos Guestrin

University of Washington

May 20, 2013

1

---

# Supervised Learning of Mixtures of Gaussians

K Components $(y^j, x^j)$

■ Mixtures of Gaussians:
  □ Prior class probabilities: $P(y)$ ← K-1 params (multinomial)
  □ Likelihood function per class: $P(x|y=i)$ ← $N(\mu_i, \Sigma_i)$

total:
K-1 +
$K\left(m^2 + \frac{m}{2}\right)$
+ km

$x \to m$ dims

↑ symmetric $m \times m$ matrix
$\frac{m^2}{2} + \frac{m}{2}$ params

■ Suppose, for each data point, we know location **x** and class y
  □ Learning is easy… ☺

  □ For prior $P(y)$ ↔ $P(y=i) = \dfrac{\text{count}(y=i) \text{ in data}}{N}$

  □ For likelihood function:

  $P(x|y=i)$ ⎡ $\mu_i$ is average of $x^j$ for points in class i
            ⎣ $\Sigma_i$ ← $\sigma_{uv} = \dfrac{\sum_{j \text{ in cluster } i}(x_u^j - \mu_{iu})(x_v^j - \mu_{iv})}{\text{num points in Cluster } i}$

2

1

# Unsupervised Learning: not as hard as it looks

*we don't have $y^j$*

Sometimes easy

*well separated*

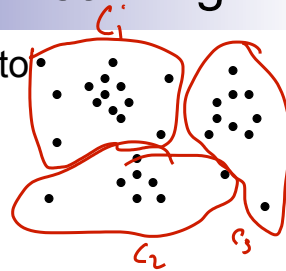*no good mixture of Gaussians*

Sometimes impossible

*IN CASE YOU'RE WONDERING WHAT THESE DIAGRAMS ARE, THEY SHOW 2-d UNLABELED DATA (X VECTORS) DISTRIBUTED IN 2-d SPACE. THE TOP ONE HAS THREE VERY CLEAR GAUSSIAN CENTERS*

and sometimes in between

*overlapping but meaningful clusters*

3

---

# EM: "Reducing" Unsupervised Learning to Supervised Learning

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

  $C_1$ $C_2$ $C_3$

- Expectation-Maximization (EM)
  - ☐ Guess assignment of points to classes *or clusters*
  - ☐ Recompute model parameters
  - ☐ Iterate

4

2

# Back to Unsupervised Learning of Mixtures of Gaussians – a simple version

*spherical*

A simple case:

We have unlabeled data $x^1 \, x^2 \, \ldots \, x^N$

We know there are k classes

We know $P(y_1) \, P(y_2) \, P(y_3) \, \ldots \, P(y_k)$ ← *prior*

We *don't* know $\mu_1 \, \mu_2 \, .. \, \mu_k$

*also know $\sigma^2$ and same for all classes*

We can write $P(\text{data} \mid \mu_1 \ldots \mu_k)$

*want to max $\mu$*

$$= p\left(x^1 \ldots x^N \mid \mu_1 \ldots \mu_k\right)$$

*iid*
*don't know $y_j$, so avg.*

$$= \prod_{j=1}^{N} p\left(x_j \mid \mu_1 \ldots \mu_k\right)$$

*prob $y = i$* , *prior prob of cluster*

$$= \prod_{j=1}^{N} \sum_{i=1}^{k} p\left(x^j \mid \mu_i\right) P(y = i)$$

*optimize objective wrt. $\mu$*

$$\propto \prod_{j=1}^{N} \sum_{i=1}^{k} \exp\left(-\frac{1}{2\sigma^2}\left\|x^j - \mu_i\right\|^2\right) P(y = i)$$

*plug in spherical gaussian*

# EM for simple version of Mixtures of Gaussians: The E-step

■ If we know $\mu_1, \ldots, \mu_k$  → easily compute prob. point $x^j$ belongs to class y=i

$x^j = (GPA: 3.99, \, 446 \, \text{(mcat)}: 3.95)$

$$p\left(y = i \mid x^j, \mu_1 \ldots \mu_k\right) \propto \exp\left(-\frac{1}{2\sigma^2}\left\|x^j - \mu_i\right\|^2\right) P(y = i)$$

$P(y=0 \mid x^j, \mu) \propto 3.7$

$P(y=1 \mid x^j, \mu) \propto 3.2$  $\Big\}$ ⇒

$P(y=0 \mid x^j, \mu) = \dfrac{3.7}{3.7 + 3.2} \approx 0.6$

$P(y=1 \mid x^j, \mu) = \dfrac{3.2}{3.7 + 3.2} \approx 0.4$

*it's like 2 data points* : $(Y=0, x_j)$ weight $0.6$

$(Y=1, x_j)$ weight $0.4$

3

# EM for simple version of Mixtures of Gaussians: The M-step

- If we know prob. point $x^j$ belongs to class $y=i$

    $\rightarrow$ MLE for $\mu_i$ is weighted average

  □ imagine k copies of each $x^j$, each with weight $P(y=i|x^j)$:

$$\mu_i = \frac{\sum_{j=1}^{N} P(y=i|x^j) x^j}{\sum_{j=1}^{N} P(y=i|x^j)}$$

|  | $x^0$ | $x^1$ | $x^2$ ... |
|---|---|---|---|
| $P(y=0|x^j)$ | 0.7 | 0.8 | 0.1 |
| $P(y=1|x^j)$ | 0.3 | 0.2 | 0.9 |

weighted average

---

# E.M. for Simple version of Mixtures of Gaussians

**E-step**

Compute "expected" classes of all datapoints for each class

$$p\left(y=i|x^j,\mu_1...\mu_k\right) \propto \exp\left(-\frac{1}{2\sigma^2}\left\|x^j-\mu_i\right\|^2\right) P(y=i)$$

*Just evaluate a Gaussian at $x^j$*

**M-step**

Compute Max. like **μ** given our data's class membership distributions

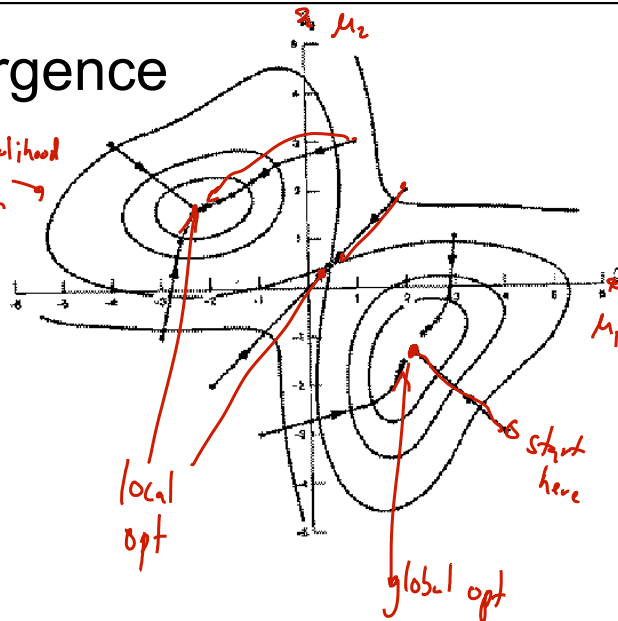$$\mu_i = \frac{\sum_{j=1}^{m} P(y=i|x^j) x^j}{\sum_{j=1}^{m} P(y=i|x^j)}$$

# E.M. Convergence



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

*(handwritten: not a global optima, necessarily)*

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

9

---

# E.M. for axis-aligned GMM

$$\Sigma = \begin{pmatrix} \sigma_1^2 & 0 & 0 & \cdots & 0 & 0 \\ 0 & \sigma_2^2 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \sigma_3^2 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & \sigma_{m-1}^2 & 0 \\ 0 & 0 & 0 & \cdots & 0 & \sigma_m^2 \end{pmatrix}$$

Iterate. On the $t$'th iteration let our estimates be

$$\lambda_t = \{ \mu_1^{(t)}, \mu_2^{(t)} \ldots \mu_k^{(t)}, \Sigma_1^{(t)}, \Sigma_2^{(t)} \ldots \Sigma_k^{(t)}, p_1^{(t)}, p_2^{(t)} \ldots p_k^{(t)} \}$$

*(handwritten labels: means of Gaussian per class / Cov. gaussian per Class / Priors)*

$p_i^{(t)}$ is shorthand for estimate of prior $P(y=i)$ on $t$'th iteration

**E-step**

Compute "expected" classes of all datapoints for each class

$$P\left(y = i \mid x^j, \lambda_t\right) \propto p_i^{(t)} p\left(x^j \mid \mu_i^{(t)}, \Sigma_i^{(t)}\right)$$

*Just evaluate a Gaussian at $x^j$*

M-step

Compute Max. like $\mu$, given our data's class membership distributions

$$\mu_i^{(t+1)} = \frac{\sum_j P\left(y = i \mid x^j, \lambda_t\right) x^j}{\sum_j P\left(y = i \mid x^j, \lambda_t\right)}$$

*(handwritten: weighted average)*

$$p_i^{(t+1)} = \frac{\sum_j P\left(y = i \mid x^j, \lambda_t\right)}{N}$$

*(handwritten: weighted count)*

$m$ = #records

*(handwritten: covariance: same as usual, but with weighted data)*

10

---

5

# E.M. for General GMMs

Iterate. On the *t*'th iteration let our estimates be

$$\lambda_t = \{ \mu_1^{(t)}, \mu_2^{(t)} \ldots \mu_k^{(t)}, \Sigma_1^{(t)}, \Sigma_2^{(t)} \ldots \Sigma_k^{(t)}, p_1^{(t)}, p_2^{(t)} \ldots p_k^{(t)} \}$$

**E-step**

Compute "expected" classes of all datapoints for each class

$$P\left(y = i \middle| x^j, \lambda_t\right) \propto p_i^{(t)} p\left(x^j \middle| \mu_i^{(t)}, \Sigma_i^{(t)}\right)$$

*Just evaluate a Gaussian at $x^j$*

*Compact version of earlier slide*

M-step

Compute Max. like **μ** given our data's class membership distributions

*general law-min of cov matrices*

$$\mu_i^{(t+1)} = \frac{\sum_j P\left(y = i \middle| x^j, \lambda_t\right) x^j}{\sum_j P\left(y = i \middle| x^j, \lambda_t\right)}$$

$$\Sigma_i^{(t+1)} = \frac{\sum_j P\left(y = i \middle| x^j, \lambda_t\right) \left[x^j - \mu_i^{(t+1)}\right]\left[x^j - \mu_i^{(t+1)}\right]^T}{\sum_j P\left(y = i \middle| x^j, \lambda_t\right)}$$

*weighted data*

$$p_i^{(t+1)} = \frac{\sum_j P\left(y = i \middle| x^j, \lambda_t\right)}{n}$$

*m = #records*

11

---

# Gaussian Mixture Example: Start

*Started by guessing $M, \Sigma, P(Y)$*

*mostly green, a little red*

p=0.333
p=0.333   0.333

*final prob*



12

# After first iteration



more red

more green

p=0.2
p=0.273
=0.221

more blue

# After 2nd iteration



p=0.37
p=0.306
=0.320

7

# After 3rd iteration



p=0.34

p=0.307

**15**

# After 4th iteration



p=0.331

p=0.288

**16**

8

# After 5th iteration



p=0.322

p=0.285

# After 6th iteration



p=0.315

p=0.287

# After 20th iteration



p=0.234    p=0.334

still a
little
a green,
but that's
OK

19
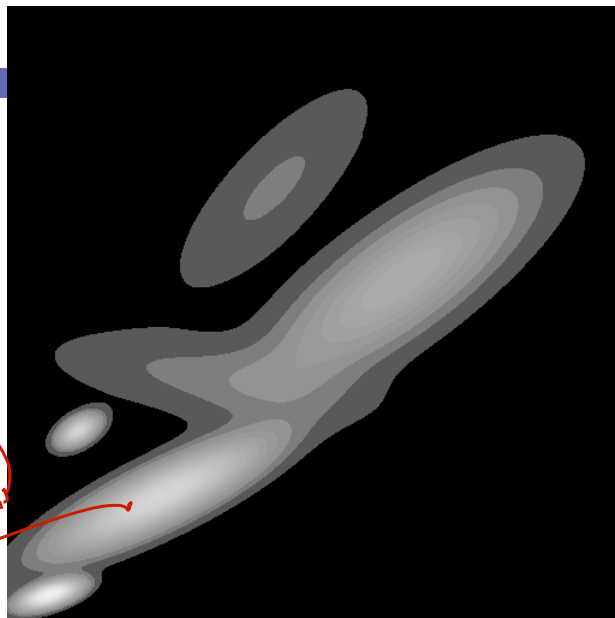
# Some Bio Assay data



20

# GMM clustering of the assay data

# Resulting Density Estimator



$P(x \mid \text{mixture of Gaussians})$

Three classes of assay
(each learned with it's own mixture model)



Resulting Bayes Classifier

Resulting Bayes Classifier, using posterior probabilities to alert about ambiguity and anomalousness

**Yellow means anomalous**

**Cyan means ambiguous**

©Carlos Guestrin 2005-2013       25

---

# 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 $\mathbf{x}^j$, for each value $i$ of $y$:
    - $P(y=i|\mathbf{x}^j,\theta^{(t)})$

- Maximization step: Find maximum likelihood parameters for (weighted) "completed data":
  - For each data point $\mathbf{x}^j$, create $k$ weighted data points
    - $(y=i, x_j)$ weight $P(y=i|x^j, \theta^{(t)})$
  - Set $\theta^{(t+1)}$ as the maximum likelihood parameter estimate for this weighted data

- Repeat

©Carlos Guestrin 2005-2013       26

# 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 <u>DOES</u>
- EM is coordinate ascent

27

# Acknowledgements

- K-means & Gaussian mixture models presentation contains material from excellent tutorial by Andrew Moore:
  - □ http://www.autonlab.org/tutorials/
- K-means Applet:
  - □ http://www.elet.polimi.it/upload/matteucc/Clustering/tutorial_html/AppletKM.html
- Gaussian mixture models Applet:
  - □ http://www.neurosci.aist.go.jp/%7Eakaho/MixtureEM.html

28

# Dimensionality Reduction PCA

Machine Learning – CSE446

Carlos Guestrin

University of Washington

May 20, 2013

29

---

# Dimensionality reduction

- Input data may have thousands or millions of dimensions!

  *x with 10 000 — 10 000 000 dims*

  - e.g., text data has

- **Dimensionality reduction**: represent data with fewer dimensions
  - easier learning – fewer parameters
  - visualization – hard to visualize more than 3D or 4D
  - discover "intrinsic dimensionality" of data
    - high dimensional data that is truly lower dimensional

30

# Lower dimensional projections

- Rather than picking a <u>subset of the features</u>, we can new <u>features</u> that are combinations of existing features
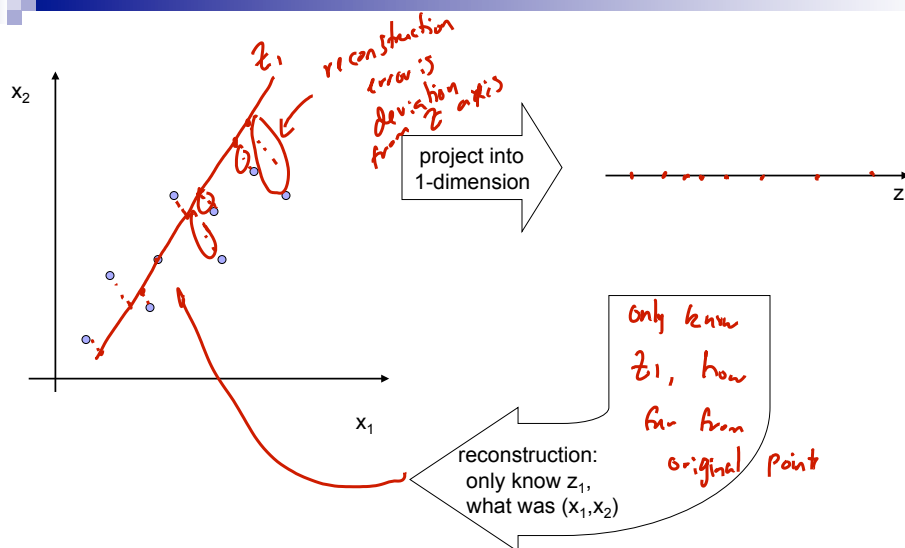
$$z_7 = 2.5 x_1 + 2.9 x_2 - 3.7 x_3 \cdots$$

model   linear

$$z = Ax$$

learn   loss/accuracy   reconstruction error

- Let's see this in the unsupervised setting
  - just **X**, but no Y

# Linear projection and reconstruction

$x_2$

$z_1$   reconstruction error is deviation from z axis

project into 1-dimension

$z_1$

$x_1$

reconstruction: only know $z_1$, what was $(x_1, x_2)$

Only know $z_1$, how far from original point