

EM (cont.)

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

November 26th, 2007

©2005-2007 Carlos Guestrin

1

Silly Example

Let events be “grades in a class”

w_1 = Gets an <u>A</u>	$P(A) = \frac{1}{2}$
w_2 = Gets a <u>B</u>	$P(B) = \mu$
w_3 = Gets a <u>C</u>	$P(C) = 2\mu$
w_4 = Gets a <u>D</u>	$P(D) = \frac{1}{2} - 3\mu$

(Note $0 \leq \mu \leq 1/6$)

Assume we want to estimate μ from data. In a given class there were

- a A's
- b B's
- c C's
- d D's

What's the maximum likelihood estimate of μ given a,b,c,d ?

$\hat{\mu}_{MLE}$

©2005-2007 Carlos Guestrin

2

Trivial Statistics

$$P(A) = \frac{1}{2} \quad P(B) = \mu \quad P(C) = 2\mu \quad P(D) = \frac{1}{2} - 3\mu$$

$$P(a, b, c, d | \mu) = K \left(\frac{1}{2}\right)^a (\mu)^b (2\mu)^c \left(\frac{1}{2} - 3\mu\right)^d$$

$$\log P(a, b, c, d | \mu) = \log K + a \log \frac{1}{2} + b \log \mu + c \log 2\mu + d \log \left(\frac{1}{2} - 3\mu\right)$$

FOR MAX LIKE μ , SET $\frac{\partial \text{Log} P}{\partial \mu} = 0$

$$\frac{\partial \text{Log} P}{\partial \mu} = \frac{b}{\mu} + \frac{2c}{2\mu} - \frac{3d}{1/2 - 3\mu} = 0$$

Gives max like $\mu = \frac{b + c}{6(b + c + d)}$

So if class got

A	B	C	D
14	6	9	10

Max like $\mu = \frac{1}{10}$

Boring, but true!

Same Problem with Hidden Information

Someone tells us that

Number of High grades (A's + B's) = h

Number of C's = c

Number of D's = d

What is the max. like estimate of μ now?

We can answer this question circularly:

EXPECTATION

If we know the value of μ we could compute the expected value of a and b

iterate

Since the ratio $a:b$ should be the same as the ratio $\frac{1}{2} : \mu$

$$\bar{a} = \frac{\frac{1}{2}}{\frac{1}{2} + \mu} h \quad \bar{b} = \frac{\mu}{\frac{1}{2} + \mu} h$$

MAXIMIZATION

If we know the expected values of \bar{a} and \bar{b} we could compute the maximum likelihood value of μ

$$\mu = \frac{\bar{b} + c}{6(\bar{b} + c + d)}$$

REMEMBER

$$P(A) = \frac{1}{2}$$

$$P(B) = \mu$$

$$P(C) = 2\mu$$

$$P(D) = \frac{1}{2} - 3\mu$$

E.M. for our Trivial Problem

REMEMBER

$$P(A) = \frac{1}{2}$$

$$P(B) = \mu$$

$$P(C) = 2\mu$$

$$P(D) = \frac{1}{2} - 3\mu$$

We begin with a guess for μ

We iterate between EXPECTATION and MAXIMALIZATION to improve our estimates of μ and a and b .

Define $\mu^{(t)}$ the estimate of μ on the t 'th iteration

$b^{(t)}$ the estimate of b on t 'th iteration

$\mu^{(0)}$ = initial guess

$$b^{(t)} = \frac{\mu^{(t)} h}{\frac{1}{2} + \mu^{(t)}} = E[b | \mu^{(t)}]$$

$$\mu^{(t+1)} = \frac{b^{(t)} + c}{6(b^{(t)} + c + d)}$$

= max like est. of μ given $b^{(t)}$

E-step

M-step

Continue iterating until converged.

Good news: Converging to local optimum is assured.

Bad news: I said "local" optimum.

©2005-2007 Carlos Guestrin

5

E.M. Convergence

- Convergence proof based on fact that $\text{Prob}(\text{data} | \mu)$ must increase or remain same between each iteration [NOT OBVIOUS]
 - But it can never exceed 1 [OBVIOUS]
- So it must therefore converge [OBVIOUS]

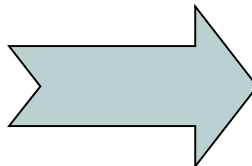
In our example, suppose we had

$$h = 20$$

$$c = 10$$

$$d = 10$$

$$\mu^{(0)} = 0$$



Convergence is generally linear: error decreases by a constant factor each time step.

t	$\mu^{(t)}$	$b^{(t)}$
0	0	0
1	0.0833	2.857
2	0.0937	3.158
3	0.0947	3.185
4	0.0948	3.187
5	0.0948	3.187
6	0.0948	3.187

act, 20

©2005-2007 Carlos Guestrin

6

Back to Unsupervised Learning of GMMs – a simple case

A simple case:

We have unlabeled data $\mathbf{x}_1 \mathbf{x}_2 \dots \mathbf{x}_m$

We know there are k classes

We know $P(y_1) P(y_2) P(y_3) \dots P(y_k)$

We don't know $\mu_1 \mu_2 \dots \mu_k$

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

$$= p(\mathbf{x}_1 \dots \mathbf{x}_m \mid \mu_1 \dots \mu_k)$$

$$= \prod_{j=1}^m p(\mathbf{x}_j \mid \mu_1 \dots \mu_k)$$

$$= \prod_{j=1}^m \sum_{i=1}^k p(\mathbf{x}_j \mid \mu_i) P(y = i)$$

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

EM for simple case of GMMs: The E-step

- If we know $\mu_1, \dots, \mu_k \rightarrow$ easily compute prob. point \mathbf{x}_j belongs to class $y=i$

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

EM for simple case of GMMs: The M-step

- If we know prob. point x_j belongs to class $y=i$
 - MLE for μ_i is weighted average
- imagine k copies of each x_j , each with weight $P(y=i|x_j)$:

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

E.M. for GMMs

E-step

Compute “expected” classes of all datapoints for each class

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

Just evaluate
a Gaussian at
 x_j

M-step

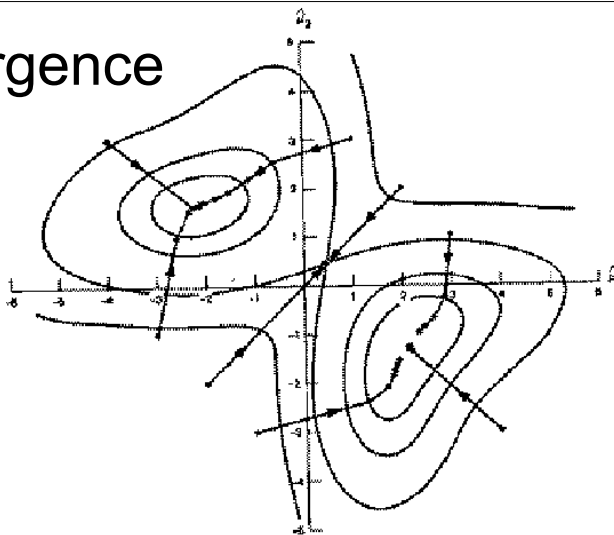
Compute Max. like μ given our data's class membership distributions

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

E.M. Convergence



- EM is coordinate ascent on an interesting potential function
- Coord. ascent for bounded pot. func. ! convergence to a local optimum guaranteed
- See Neal & Hinton reading on class webpage



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

©2005-2007 Carlos Guestrin

11

E.M. for axis-aligned GMN

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

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

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

E-step

Compute "expected" classes of all datapoints for each class

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

$$P(y=i|x_j, \lambda_t) \propto p_i^{(t)} p(x_j | \mu_i^{(t)}, \Sigma_i^{(t)})$$

Just evaluate a Gaussian at x_j

M-step

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

$$\hat{i}_j^{(t+1)} = \frac{\sum_j P(y=i|x_j, \lambda_t) x_j}{\sum_j P(y=i|x_j, \lambda_t)}$$

$$p_i^{(t+1)} = \frac{\sum_j P(y=i|x_j, \lambda_t)}{m}$$

$m = \text{\#records}$

©2005-2007 Carlos Guestrin

12

E.M. for General GMMs

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

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

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

E-step

Compute "expected" classes of all datapoints for each class

$$P(y=i|x_j, \lambda_t) \propto p_i^{(t)} p(x_j|\mu_i^{(t)}, \Sigma_i^{(t)})$$

Just evaluate a Gaussian at x_j

M-step

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

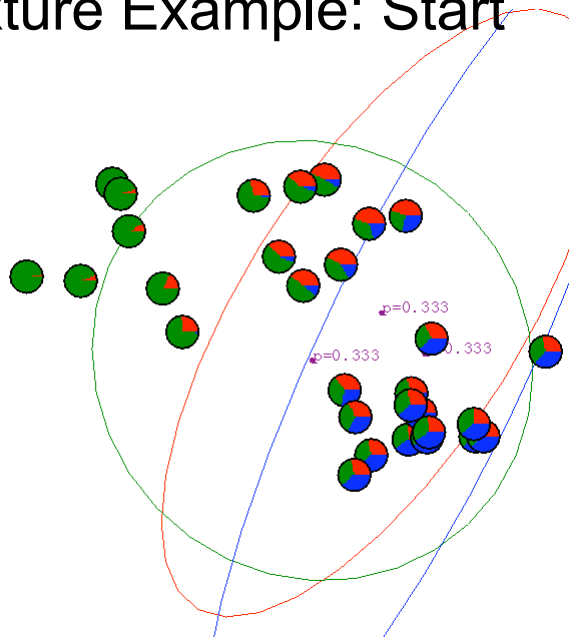
$$\hat{\mu}_i^{(t+1)} = \frac{\sum_j P(y=i|x_j, \lambda_t) x_j}{\sum_j P(y=i|x_j, \lambda_t)}$$

$$\hat{\Sigma}_i^{(t+1)} = \frac{\sum_j P(y=i|x_j, \lambda_t) [x_j - \mu_i^{(t+1)}][x_j - \mu_i^{(t+1)}]^T}{\sum_j P(y=i|x_j, \lambda_t)}$$

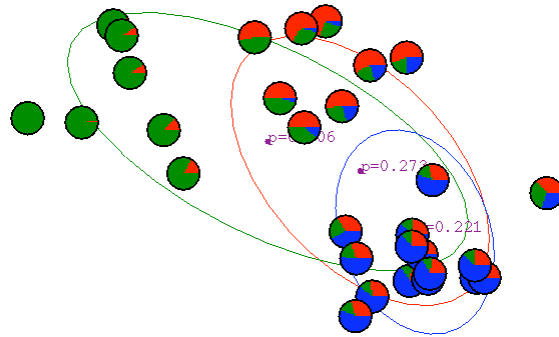
$$p_i^{(t+1)} = \frac{\sum_j P(y=i|x_j, \lambda_t)}{m}$$

$m = \text{\#records}$

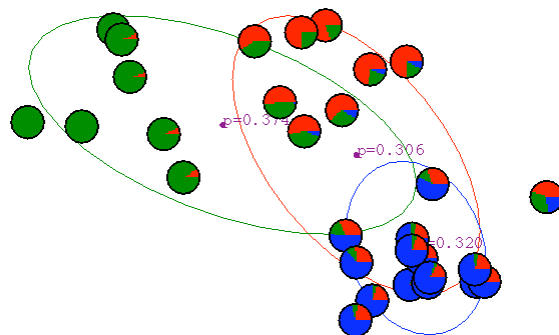
Gaussian Mixture Example: Start



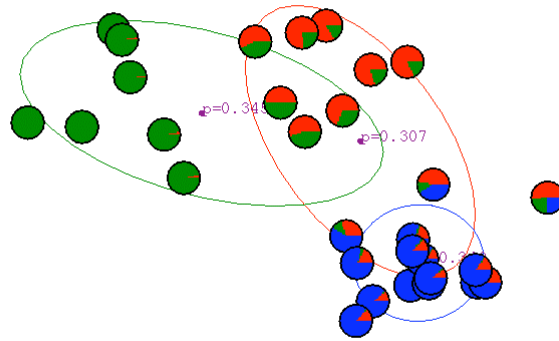
After first iteration



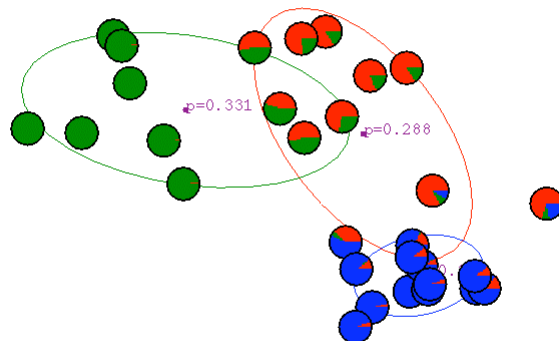
After 2nd iteration



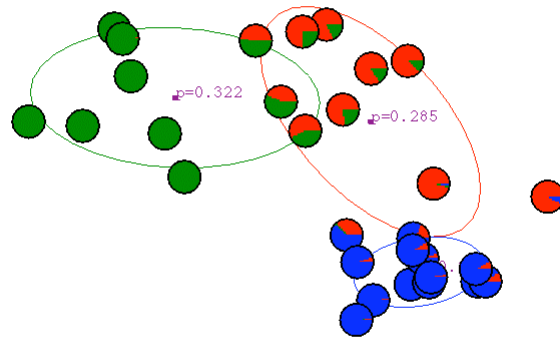
After 3rd iteration



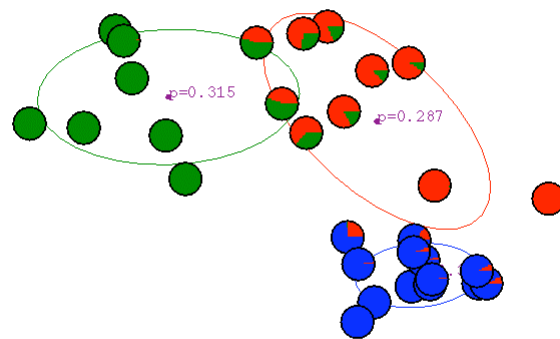
After 4th iteration



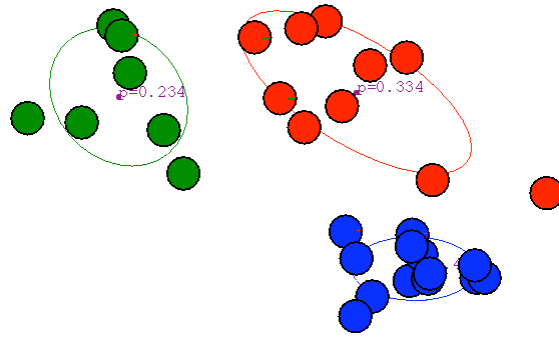
After 5th iteration



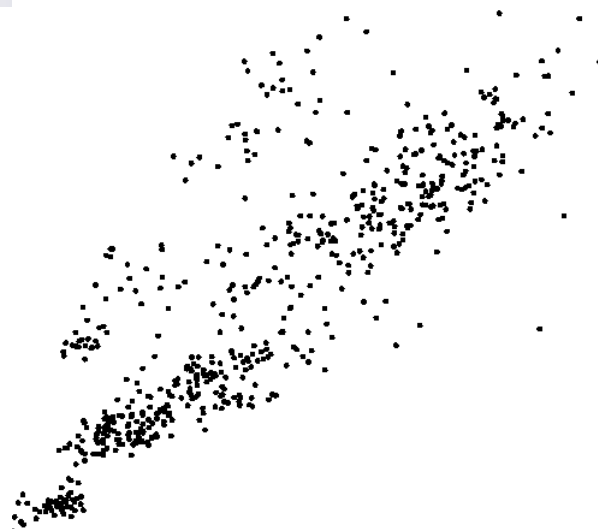
After 6th iteration



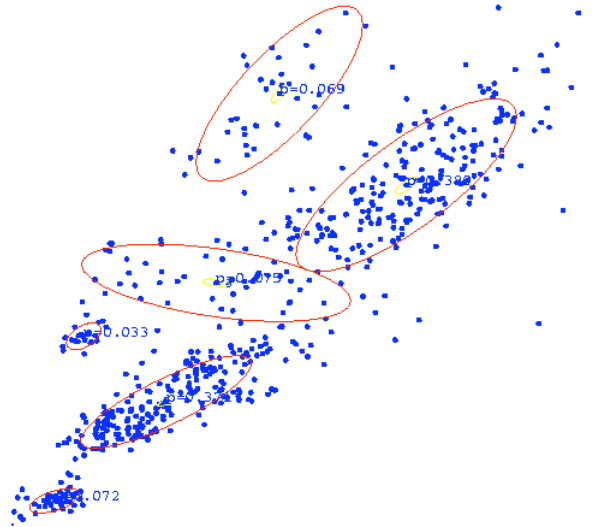
After 20th iteration



Some Bio Assay data



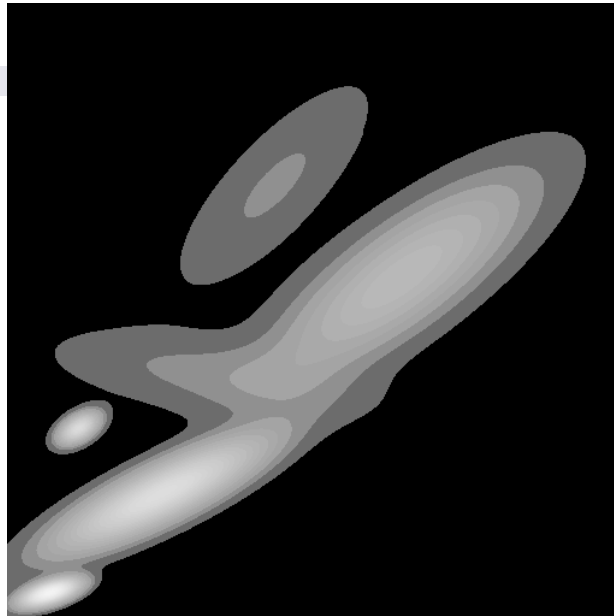
GMM clustering of the assay data



©2005-2007 Carlos Guestrin

23

Resulting
Density
Estimator



©2005-2007 Carlos Guestrin

24



Three classes of assay

(each learned with its own mixture model)

Compound =

IL-1

TNF

none

n
u
c
l
e
u
s

ing

©2005-2007 Carlos Guestrin

25



Resulting Bayes Classifier

Compound =

IL-1

TNF

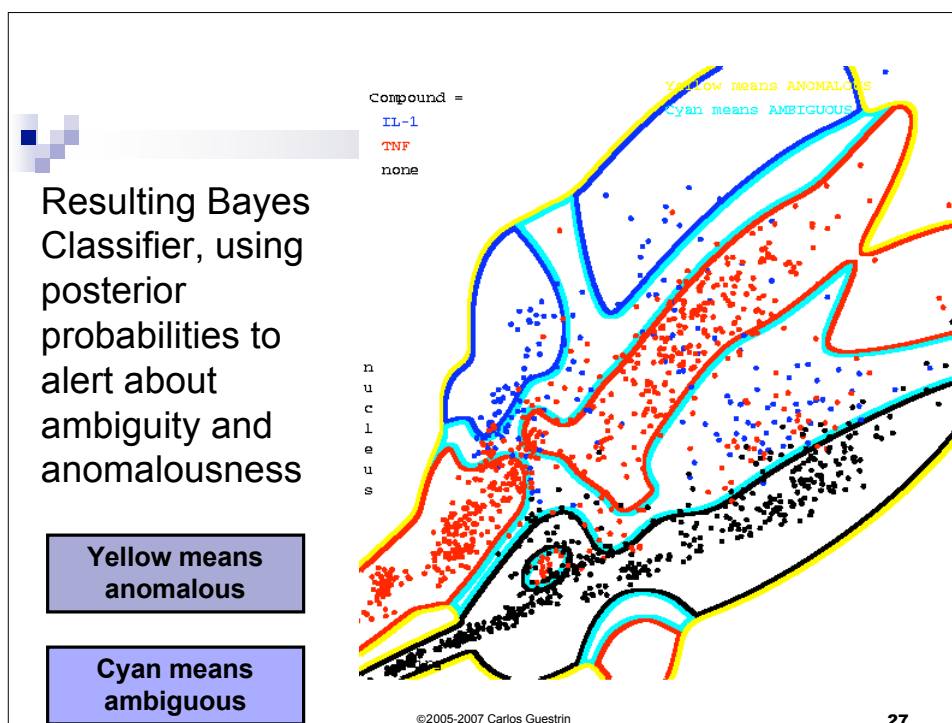
none

n
u
c
l
e
u
s

ing

©2005-2007 Carlos Guestrin

26



Announcements

■ Project:

- Poster session: NSH Atrium, Friday 11/30, 2-5pm
 - Print your poster early!!!
 - SCS facilities has a poster printer, ask helpdesk
 - Students from outside SCS should check with their departments
 - It's OK to print separate pages
 - We'll provide pins, posterboard and an easel
 - Poster size: 32x40 inches
 - Invite your friends, there will be a prize for best poster, by popular vote

■ Last lecture:

- Thursday, 11/29, 5-6:20pm, Wean 7500

The general learning problem with missing data

- Marginal likelihood – \mathbf{x} is observed, \mathbf{z} is missing:

$$\begin{aligned}\ell(\theta : \mathcal{D}) &= \log \prod_{j=1}^m P(\mathbf{x}_j | \theta) \\ &= \sum_{j=1}^m \log P(\mathbf{x}_j | \theta) \\ &= \sum_{j=1}^m \log \sum_{\mathbf{z}} P(\mathbf{x}_j, \mathbf{z} | \theta)\end{aligned}$$

E-step

- \mathbf{x} is observed, \mathbf{z} is missing
- Compute probability of missing data given current choice of θ
 - $Q(\mathbf{z} | \mathbf{x}_j)$ for each \mathbf{x}_j
 - e.g., probability computed during classification step
 - corresponds to “classification step” in K-means

$$Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) = P(\mathbf{z} | \mathbf{x}_j, \theta^{(t)})$$

Jensen's inequality

$$\ell(\theta : \mathcal{D}) = \sum_{j=1}^m \log \sum_{\mathbf{z}} P(\mathbf{z} | \mathbf{x}_j) P(\mathbf{x}_j | \theta)$$

■ **Theorem:** $\log \sum_{\mathbf{z}} P(\mathbf{z}) f(\mathbf{z}) \geq \sum_{\mathbf{z}} P(\mathbf{z}) \log f(\mathbf{z})$

Applying Jensen's inequality

■ Use: $\log \sum_{\mathbf{z}} P(\mathbf{z}) f(\mathbf{z}) \geq \sum_{\mathbf{z}} P(\mathbf{z}) \log f(\mathbf{z})$

$$\ell(\theta^{(t)} : \mathcal{D}) = \sum_{j=1}^m \log \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \frac{P(\mathbf{z}, \mathbf{x}_j | \theta^{(t)})}{Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j)}$$

The M-step maximizes lower bound on weighted data

- Lower bound from Jensen's:

$$\ell(\theta^{(t)} : \mathcal{D}) \geq \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta^{(t)}) + m.H(Q^{(t+1)})$$

- Corresponds to weighted dataset:

- ☐ $\langle \mathbf{x}_1, \mathbf{z}=1 \rangle$ with weight $Q^{(t+1)}(\mathbf{z}=1 | \mathbf{x}_1)$
- ☐ $\langle \mathbf{x}_1, \mathbf{z}=2 \rangle$ with weight $Q^{(t+1)}(\mathbf{z}=2 | \mathbf{x}_1)$
- ☐ $\langle \mathbf{x}_1, \mathbf{z}=3 \rangle$ with weight $Q^{(t+1)}(\mathbf{z}=3 | \mathbf{x}_1)$
- ☐ $\langle \mathbf{x}_2, \mathbf{z}=1 \rangle$ with weight $Q^{(t+1)}(\mathbf{z}=1 | \mathbf{x}_2)$
- ☐ $\langle \mathbf{x}_2, \mathbf{z}=2 \rangle$ with weight $Q^{(t+1)}(\mathbf{z}=2 | \mathbf{x}_2)$
- ☐ $\langle \mathbf{x}_2, \mathbf{z}=3 \rangle$ with weight $Q^{(t+1)}(\mathbf{z}=3 | \mathbf{x}_2)$
- ☐ ...

The M-step

$$\ell(\theta^{(t)} : \mathcal{D}) \geq \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta^{(t)}) + m.H(Q^{(t+1)})$$

- Maximization step:

$$\theta^{(t+1)} \leftarrow \arg \max_{\theta} \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta)$$

- Use expected counts instead of counts:

- ☐ If learning requires $\text{Count}(\mathbf{x}, \mathbf{z})$
- ☐ Use $E_{Q^{(t+1)}}[\text{Count}(\mathbf{x}, \mathbf{z})]$

Convergence of EM

- Define potential function $F(\theta, Q)$:

$$\ell(\theta : \mathcal{D}) \geq F(\theta, Q) = \sum_{j=1}^m \sum_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j | \theta)}{Q(\mathbf{z} | \mathbf{x}_j)}$$

- EM corresponds to coordinate ascent on F
 - Thus, maximizes lower bound on marginal log likelihood

M-step is easy

$$\theta^{(t+1)} \leftarrow \arg \max_{\theta} \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta)$$

- Using potential function

$$F(\theta, Q^{(t+1)}) = \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta) + m \cdot H(Q^{(t+1)})$$

E-step also doesn't decrease potential function 1

- Fixing θ to $\theta^{(t)}$:

$$\ell(\theta^{(t)} : \mathcal{D}) \geq F(\theta^{(t)}, Q) = \sum_{j=1}^m \sum_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j | \theta^{(t)})}{Q(\mathbf{z} | \mathbf{x}_j)}$$

KL-divergence

- Measures distance between distributions

$$KL(Q||P) = \sum_z Q(z) \log \frac{Q(z)}{P(z)}$$

- KL=zero if and only if $Q=P$

E-step also doesn't decrease potential function 2

- Fixing θ to $\theta^{(t)}$:

$$\begin{aligned}\ell(\theta^{(t)} : \mathcal{D}) \geq F(\theta^{(t)}, Q) &= \ell(\theta^{(t)} : \mathcal{D}) + \sum_{j=1}^m \sum_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}_j) \log \frac{P(\mathbf{z} | \mathbf{x}_j, \theta^{(t)})}{Q(\mathbf{z} | \mathbf{x}_j)} \\ &= \ell(\theta^{(t)} : \mathcal{D}) - m \sum_{j=1}^m KL(Q(\mathbf{z} | \mathbf{x}_j) || P(\mathbf{z} | \mathbf{x}_j, \theta^{(t)}))\end{aligned}$$

E-step also doesn't decrease potential function 3

$$\ell(\theta^{(t)} : \mathcal{D}) \geq F(\theta^{(t)}, Q) = \ell(\theta^{(t)} : \mathcal{D}) - m \sum_{j=1}^m KL(Q(\mathbf{z} | \mathbf{x}_j) || P(\mathbf{z} | \mathbf{x}_j, \theta^{(t)}))$$

- Fixing θ to $\theta^{(t)}$
- Maximizing $F(\theta^{(t)}, Q)$ over $Q \rightarrow$ set Q to posterior probability:

$$Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \leftarrow P(\mathbf{z} | \mathbf{x}_j, \theta^{(t)})$$

- Note that

$$F(\theta^{(t)}, Q^{(t+1)}) = \ell(\theta^{(t)} : \mathcal{D})$$

EM is coordinate ascent



$$\ell(\theta : \mathcal{D}) \geq F(\theta, Q) = \sum_{j=1}^m \sum_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j | \theta)}{Q(\mathbf{z} | \mathbf{x}_j)}$$

- **M-step:** Fix Q , maximize F over θ (a lower bound on $\ell(\theta : \mathcal{D})$):

$$\ell(\theta : \mathcal{D}) \geq F(\theta, Q^{(t)}) = \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t)}(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta) + m \cdot H(Q^{(t)})$$

- **E-step:** Fix θ , maximize F over Q :

$$\ell(\theta^{(t)} : \mathcal{D}) \geq F(\theta^{(t)}, Q) = \ell(\theta^{(t)} : \mathcal{D}) - m \sum_{j=1}^m KL(Q(\mathbf{z} | \mathbf{x}_j) || P(\mathbf{z} | \mathbf{x}_j, \theta^{(t)}))$$

- “Realigns” F with likelihood:

$$F(\theta^{(t)}, Q^{(t+1)}) = \ell(\theta^{(t)} : \mathcal{D})$$

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
- General case for EM

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>

Dimensionality Reduction

Machine Learning – 10701/15781
Carlos Guestrin
Carnegie Mellon University

November 26th, 2007

Dimensionality reduction

- Input data may have thousands or millions of dimensions!
 - 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

Feature selection

- Want to learn $f: \mathbf{X} \rightarrow Y$
 - $\mathbf{X} = \langle X_1, \dots, X_n \rangle$
 - but some features are more important than others
- **Approach:** select subset of features to be used by learning algorithm
 - **Score** each feature (or sets of features)
 - **Select** set of features with best score

Simple greedy **forward** feature selection algorithm

- Pick a dictionary of features
 - e.g., polynomials for linear regression
- Greedy heuristic:
 - Start from empty (or simple) set of features $F_0 = \emptyset$
 - Run learning algorithm for current set of features F_t
 - Obtain h_t
 - Select **next best feature** X_i
 - e.g., X_j that results in lowest cross-validation error learner when learning with $F_t \cup \{X_j\}$
 - $F_{t+1} \leftarrow F_t \cup \{X_i\}$
 - Recurse

Simple greedy **backward** feature selection algorithm

- Pick a dictionary of features
 - e.g., polynomials for linear regression
- Greedy heuristic:
 - Start from all features $F_0 = F$
 - Run learning algorithm for current set of features F_t
 - Obtain h_t
 - Select **next worst feature** X_i
 - e.g., X_j that results in lowest cross-validation error learner when learning with $F_t - \{X_j\}$
 - $F_{t+1} \leftarrow F_t - \{X_i\}$
 - Recurse

Impact of feature selection on classification of fMRI data [Pereira et al. '05]

Accuracy classifying
category of word read
by subject

#voxels	mean	subjects							
		233B	329B	332B	424B	474B	496B	77B	86B
50	0.735	0.783	0.817	0.55	0.783	0.75	0.8	0.65	0.75
100	0.742	0.767	0.8	0.533	0.817	0.85	0.783	0.6	0.783
200	0.737	0.783	0.783	0.517	0.817	0.883	0.75	0.583	0.783
300	0.75	0.8	0.817	0.567	0.833	0.883	0.75	0.583	0.767
400	0.742	0.8	0.783	0.583	0.85	0.833	0.75	0.583	0.75
800	0.735	0.833	0.817	0.567	0.833	0.833	0.7	0.55	0.75
1600	0.698	0.8	0.817	0.45	0.783	0.833	0.633	0.5	0.75
all (~2500)	0.638	0.767	0.767	0.25	0.75	0.833	0.567	0.433	0.733

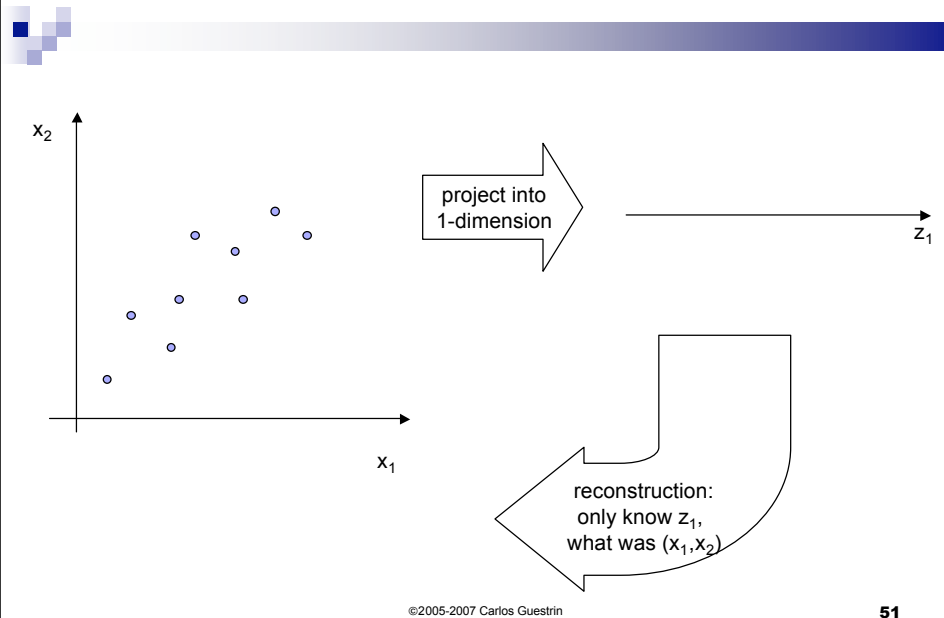
Table 1: Average accuracy across all pairs of categories, restricting the procedure to use a certain number of voxels for each subject. The highlighted line corresponds to the best mean accuracy, obtained using 300 voxels.

Voxels scored by p-value of regression to predict voxel value from the task

Lower dimensional projections

- Rather than picking a subset of the features, we can new features that are combinations of existing features
- Let's see this in the unsupervised setting
 - just **X**, but no **Y**

Linear projection and reconstruction



51

Principal component analysis – basic idea

- Project n -dimensional data into k -dimensional space while preserving information:
 - e.g., project space of 10000 words into 3-dimensions
 - e.g., project 3-d into 2-d
- Choose projection with minimum reconstruction error

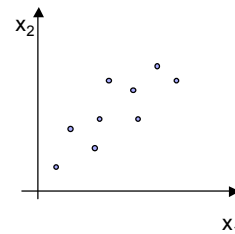
Linear projections, a review

- Project a point into a (lower dimensional) space:
 - **point**: $\mathbf{x} = (x_1, \dots, x_n)$
 - **select a basis** – set of basis vectors – $(\mathbf{u}_1, \dots, \mathbf{u}_k)$
 - we consider orthonormal basis:
 - $\mathbf{u}_i \cdot \mathbf{u}_i = 1$, and $\mathbf{u}_i \cdot \mathbf{u}_j = 0$ for $i \neq j$
 - **select a center** – $\bar{\mathbf{x}}$, defines offset of space
 - **best coordinates** in lower dimensional space defined by dot-products: (z_1, \dots, z_k) , $z_i = (\mathbf{x} - \bar{\mathbf{x}}) \cdot \mathbf{u}_i$
 - minimum squared error

PCA finds projection that minimizes reconstruction error

- Given m data points: $\mathbf{x}^i = (x_1^i, \dots, x_n^i)$, $i=1 \dots m$
- Will represent each point as a projection:
 - $\hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^k z_j^i \mathbf{u}_j$ where: $\bar{\mathbf{x}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}^i$ and $z_j^i = (\mathbf{x}^i - \bar{\mathbf{x}}) \cdot \mathbf{u}_j$
- PCA:
 - Given $k \leq n$, find $(\mathbf{u}_1, \dots, \mathbf{u}_k)$ minimizing reconstruction error:

$$error_k = \sum_{i=1}^m (\mathbf{x}^i - \hat{\mathbf{x}}^i)^2$$



Understanding the reconstruction error



- Note that \mathbf{x}^i can be represented exactly by n -dimensional projection:

$$\mathbf{x}^i = \bar{\mathbf{x}} + \sum_{j=1}^n z_j^i \mathbf{u}_j$$

$$\hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^k z_j^i \mathbf{u}_j$$

$$z_j^i = (\mathbf{x}^i - \bar{\mathbf{x}}) \cdot \mathbf{u}_j$$

- Given $k \leq n$, find $(\mathbf{u}_1, \dots, \mathbf{u}_k)$ minimizing reconstruction error:

$$error_k = \sum_{i=1}^m (\mathbf{x}^i - \hat{\mathbf{x}}^i)^2$$

- Rewriting error:

Reconstruction error and covariance matrix



$$error_k = \sum_{i=1}^m \sum_{j=k+1}^n [\mathbf{u}_j \cdot (\mathbf{x}^i - \bar{\mathbf{x}})]^2$$

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}^i - \bar{\mathbf{x}})(\mathbf{x}^i - \bar{\mathbf{x}})^T$$

Minimizing reconstruction error and eigen vectors

- Minimizing reconstruction error equivalent to picking orthonormal basis ($\mathbf{u}_1, \dots, \mathbf{u}_n$) minimizing:

$$error_k = \sum_{j=k+1}^n \mathbf{u}_j^T \Sigma \mathbf{u}_j$$

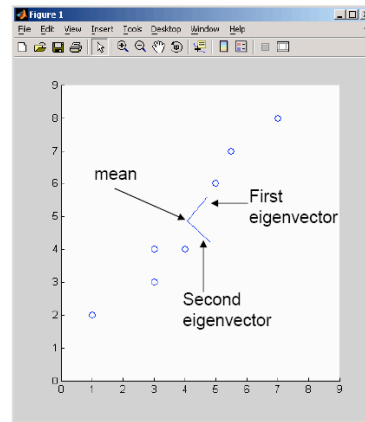
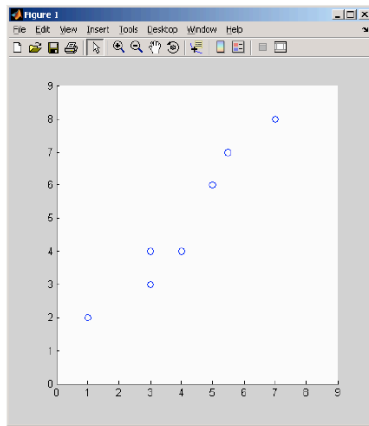
- Eigen vector:
- Minimizing reconstruction error equivalent to picking ($\mathbf{u}_{k+1}, \dots, \mathbf{u}_n$) to be eigen vectors with smallest eigen values

Basic PCA algorithm

- Start from m by n data matrix \mathbf{X}
- **Recenter**: subtract mean from each row of \mathbf{X}
 - $\mathbf{X}_c \leftarrow \mathbf{X} - \bar{\mathbf{X}}$
- **Compute covariance matrix**:
 - $\Sigma \leftarrow 1/m \mathbf{X}_c^T \mathbf{X}_c$
- Find **eigen vectors and values** of Σ
- **Principal components**: k eigen vectors with highest eigen values

PCA example

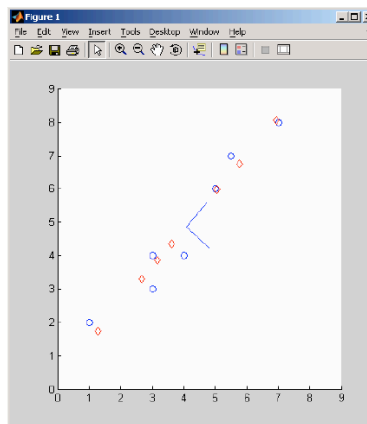
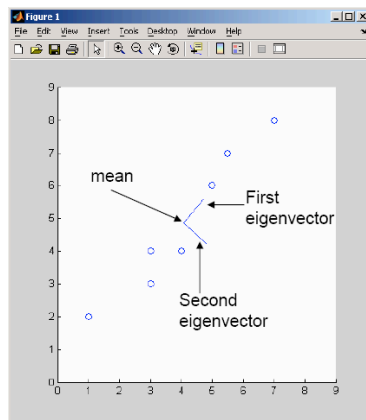
$$\hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^k z_j^i \mathbf{u}_j$$



PCA example – reconstruction

$$\hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^k z_j^i \mathbf{u}_j$$

only used first principal component



Eigenfaces [Turk, Pentland '91]



■ Input images:



■ Principal components:



Eigenfaces reconstruction



■ Each image corresponds to adding 8 principal components:



Relationship to Gaussians

- PCA assumes data is Gaussian

- $\mathbf{x} \sim N(\bar{\mathbf{x}}; \Sigma)$

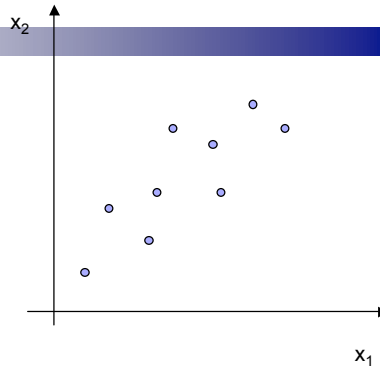
- Equivalent to weighted sum of simple Gaussians:

$$\mathbf{x} = \bar{\mathbf{x}} + \sum_{j=1}^n z_j \mathbf{u}_j; \quad z_j \sim N(0; \sigma_j^2)$$

- Selecting top k principal components equivalent to lower dimensional Gaussian approximation:

$$\mathbf{x} \approx \bar{\mathbf{x}} + \sum_{j=1}^k z_j \mathbf{u}_j + \varepsilon; \quad z_j \sim N(0; \sigma_j^2)$$

- $\varepsilon \sim N(0; \sigma^2)$, where σ^2 is defined by error_k



Scaling up

- Covariance matrix can be really big!

- Σ is n by n
 - 10000 features $\rightarrow |\Sigma|$
 - finding eigenvectors is very slow...

- Use singular value decomposition (SVD)

- finds to k eigenvectors
 - great implementations available, e.g., Matlab svd

SVD

- Write $\mathbf{X} = \mathbf{W} \mathbf{S} \mathbf{V}^T$
 - $\mathbf{X} \leftarrow$ data matrix, one row per datapoint
 - $\mathbf{W} \leftarrow$ weight matrix, one row per datapoint – coordinate of \mathbf{x}^i in eigenspace
 - $\mathbf{S} \leftarrow$ singular value matrix, diagonal matrix
 - in our setting each entry is eigenvalue λ_j
 - $\mathbf{V}^T \leftarrow$ singular vector matrix
 - in our setting each row is eigenvector \mathbf{v}_j

PCA using SVD algorithm

- Start from m by n data matrix \mathbf{X}
- **Recenter**: subtract mean from each row of \mathbf{X}
 - $\mathbf{X}_c \leftarrow \mathbf{X} - \bar{\mathbf{X}}$
- Call SVD algorithm on \mathbf{X}_c – ask for k singular vectors
- **Principal components**: k singular vectors with highest singular values (rows of \mathbf{V}^T)
 - **Coefficients** become:

Using PCA for dimensionality reduction in classification

- Want to learn $f: \mathbf{X} \rightarrow Y$
 - $\mathbf{X} = \langle X_1, \dots, X_n \rangle$
 - but some features are more important than others
- **Approach:** Use PCA on \mathbf{X} to select a few important features

PCA for classification can lead to problems...

- Direction of maximum variation may be unrelated to “discriminative” directions:
- PCA often works very well, but sometimes must use more advanced methods
 - e.g., Fisher linear discriminant

What you need to know



- Dimensionality reduction
 - why and when it's important
- Simple feature selection
- Principal component analysis
 - minimizing reconstruction error
 - relationship to covariance matrix and eigenvectors
 - using SVD
 - problems with PCA