

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 A

w_2 = Gets a B

w_3 = Gets a C

w_4 = Gets a D

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

$P(B) = \mu$

$P(C) = 2\mu$

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

(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

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

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

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

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

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

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

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

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

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
eg., $\mu^{(0)} = 0$

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

if true param was $\mu^{(t)}$

E-step

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

= max like est. of μ given $b^{(t)}$

M-step

Continue iterating until converged.

Good news: Converging to local optimum is assured.

Bad news: I said "local" optimum.

E.M. Convergence

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

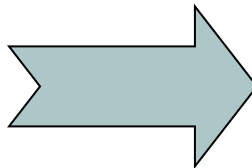
In our example, suppose we had

$h = 20$

$c = 10$

$d = 10$

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



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

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

Back to Unsupervised Learning of GMMs – a simple case

A simple case:

We have unlabeled data $x_1 x_2 \dots x_m$

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 $\mu_1 \mu_2 \dots \mu_k$

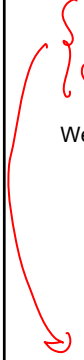
We can write $P(\text{data} | \mu_1 \dots \mu_k)$

$$= p(x_1 \dots x_m | \mu_1 \dots \mu_k)$$

$$\stackrel{\text{iid}}{=} \prod_{j=1}^m p(x_j | \mu_1 \dots \mu_k)$$

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

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



summing over possible classes

EM for simple case of GMMs: The E-step

expected value of hidden vars

- If we know $\mu_1, \dots, \mu_k \rightarrow$ easily compute prob.
point x_j belongs to class $y=i$

Bayes Rule

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

for each point j

$$P(Y=1 | x_j, \mu_1, \dots, \mu_k) = 0.7$$

$$P(Y=2 | x_j, \mu_1, \dots, \mu_k) = 0.2$$

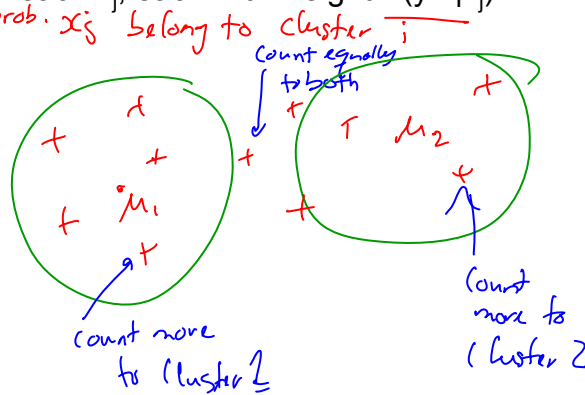
$$P(Y=3 | x_j, \mu_1, \dots, \mu_k) = 0.1$$

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 $P(y=i|x_j)$:

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



E.M. for GMMs

E-step

Compute "expected" classes of ~~all~~ ^{each} 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

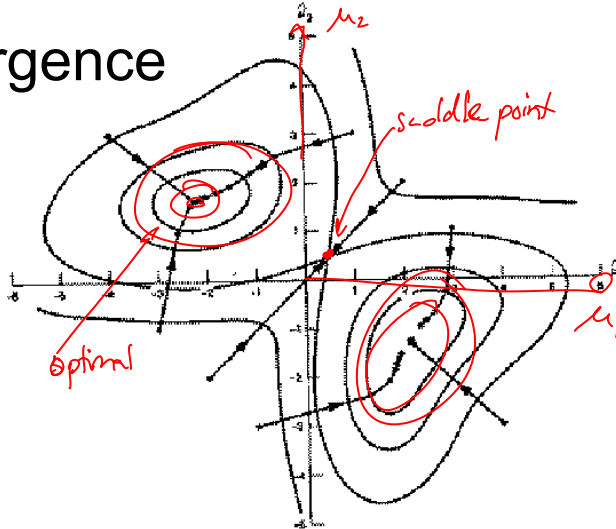
M-step

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



E.M. for axis-aligned GMN

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

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

$$\Sigma_i = \begin{pmatrix} \sigma_{i,1}^2 & 0 & 0 & \dots & 0 & 0 \\ 0 & \sigma_{i,2}^2 & 0 & \dots & 0 & 0 \\ 0 & 0 & \sigma_{i,3}^2 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & \sigma_{i,m-1}^2 & 0 \\ 0 & 0 & 0 & \dots & 0 & \sigma_{i,m}^2 \end{pmatrix}$$

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

with obs. data

$$P(y=i) = \frac{\sum_j \mathbb{1}(y=i)}{m}$$

M-step

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

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

$$\sigma_{i,k}^{2(t+1)} = \frac{\sum_j (x_{j,k} - \mu_{i,k}^{(t+1)})^2 P(y=i | x_j)}{\sum_j P(y=i | x_j)}$$

i 4th class k th feature

Same

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

$m = \# \text{records}$

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

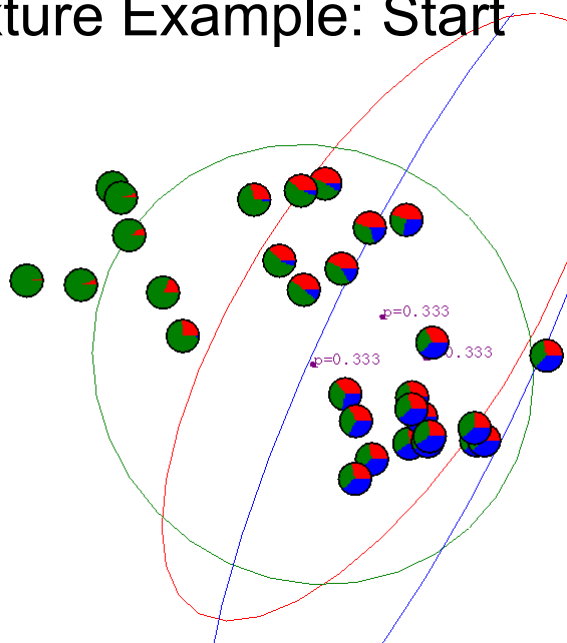
Compute Max. like μ given our data's class membership distributions

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

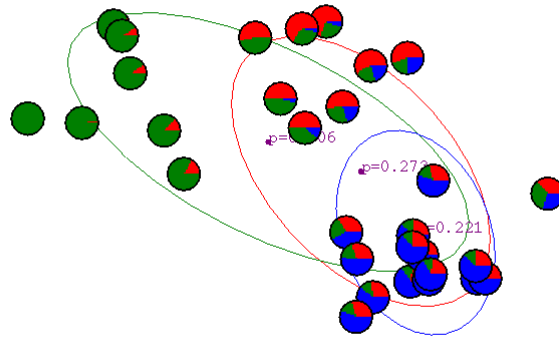
$$p_i^{(t+1)} = \frac{\sum_j P(y = i | x_j, \lambda_t)}{m}$$

$m = \text{\#records}$

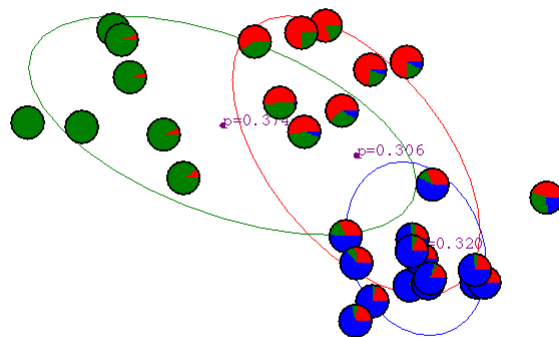
Gaussian Mixture Example: Start



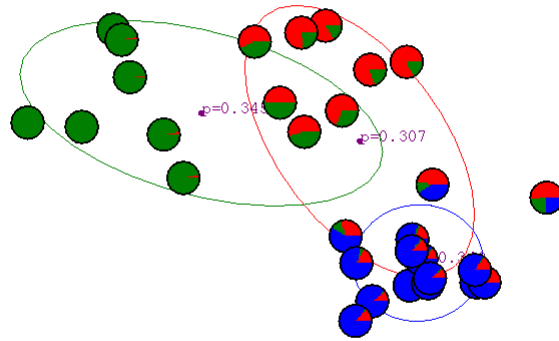
After first iteration



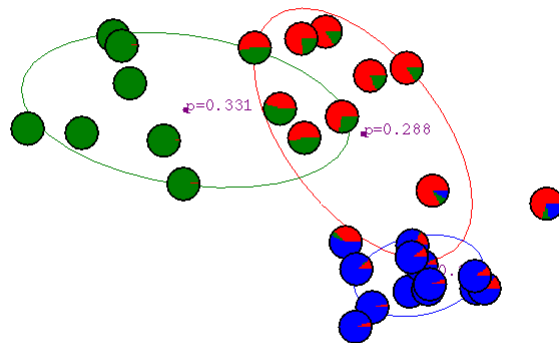
After 2nd iteration



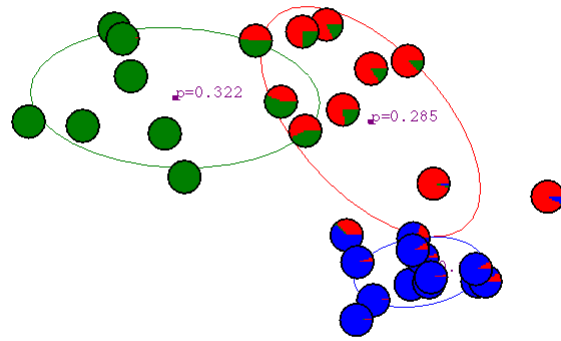
After 3rd iteration



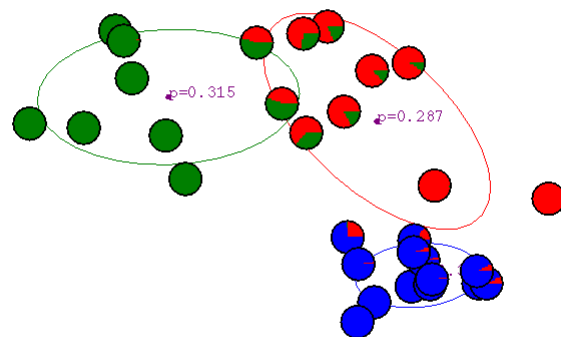
After 4th iteration



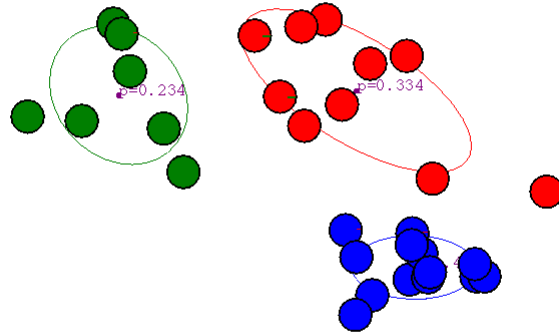
After 5th iteration



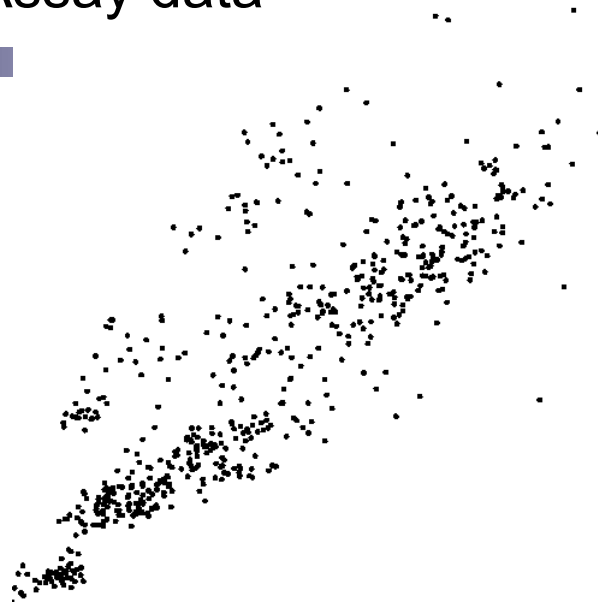
After 6th iteration



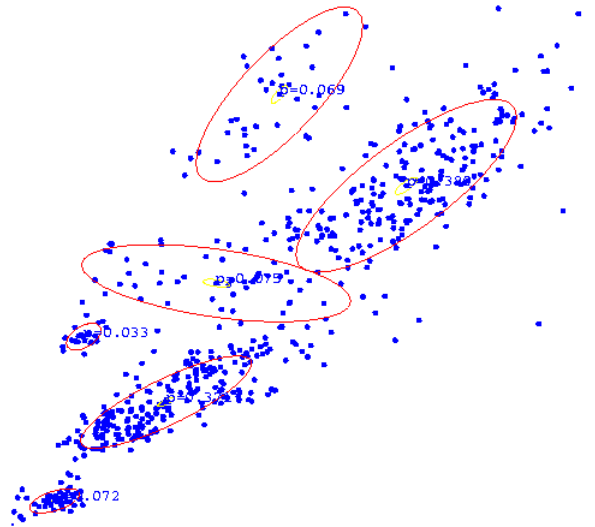
After 20th iteration



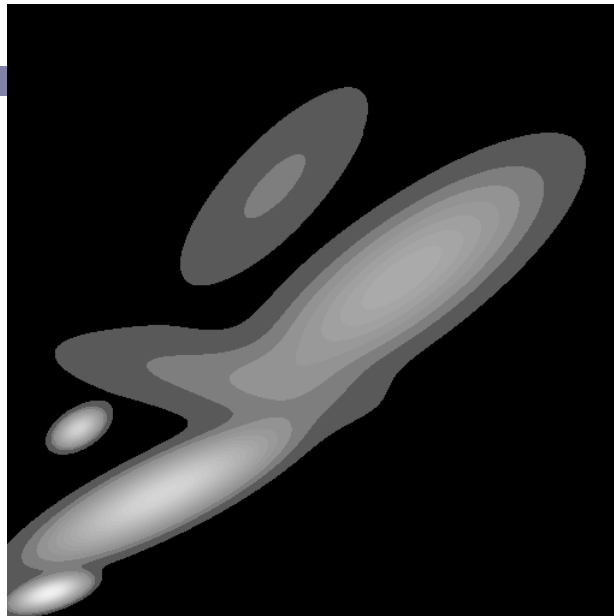
Some Bio Assay data

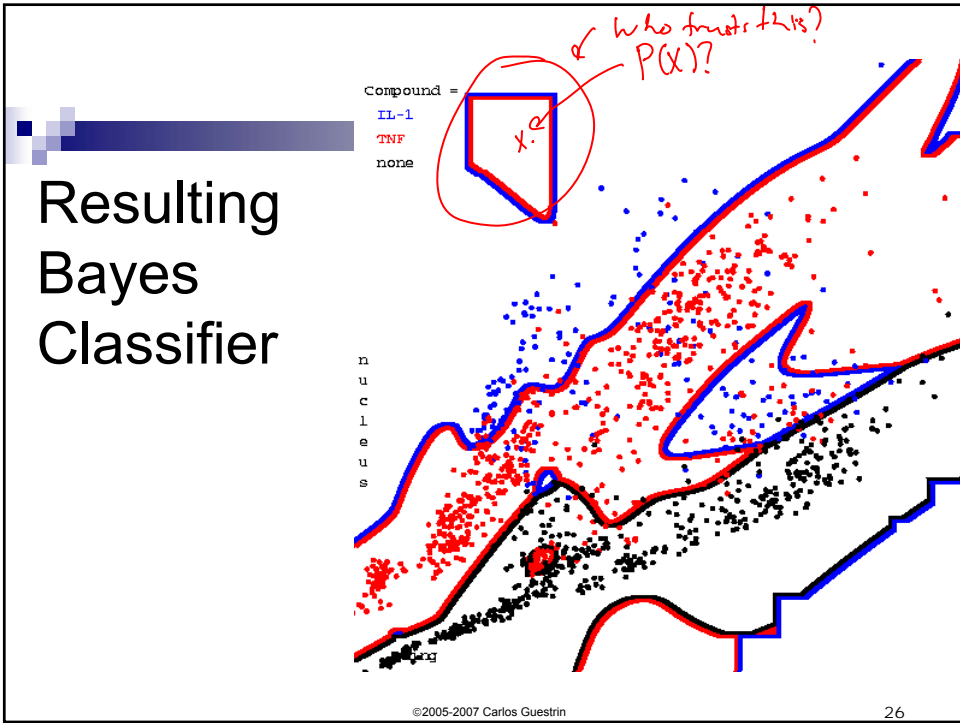
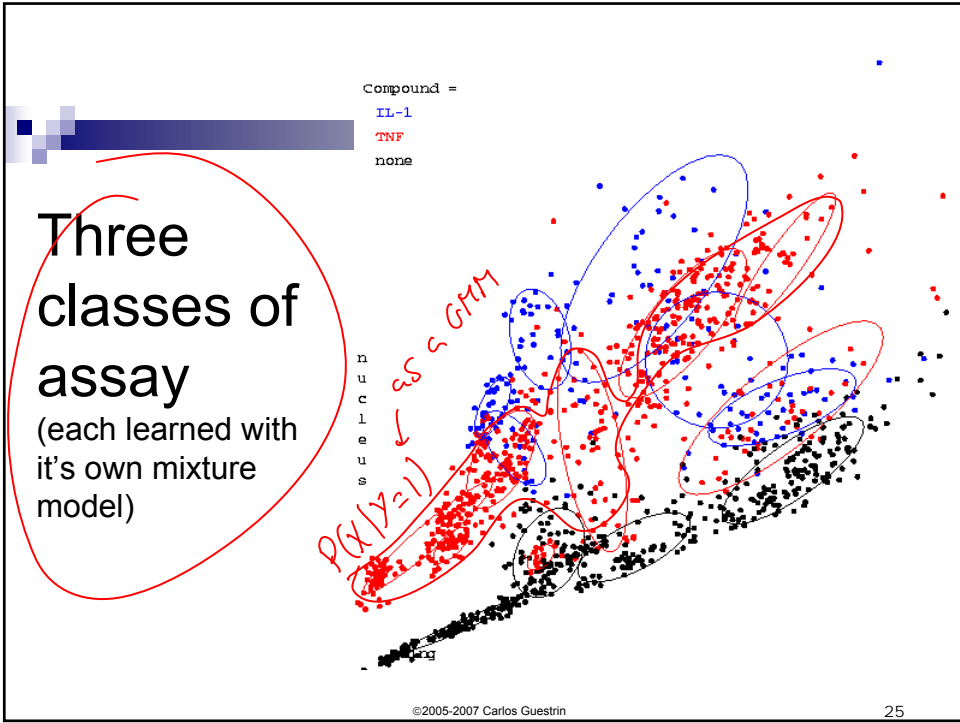


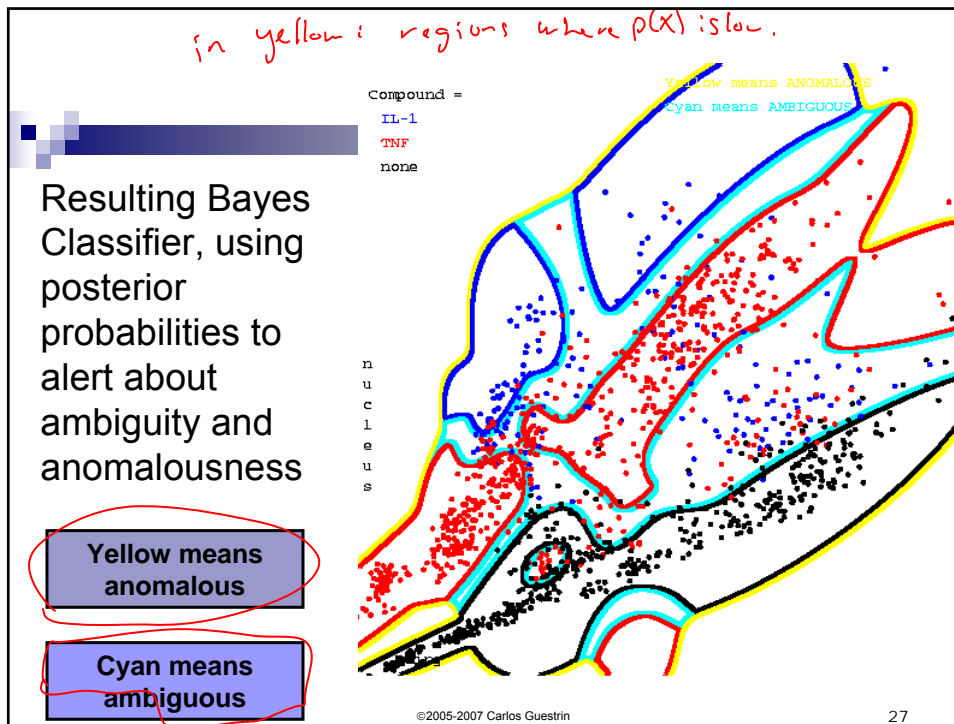
GMM clustering of the assay data



Resulting Density Estimator







Homeworks to Monica Please

Announcements

Please: fill out your University Course Assessment

rest of away next week

Monday/Tuesday 12/10-11 many extra office hours arrive 15 mins early

- **Project:**
 - Poster session: NSH Atrium, Friday 11/30, 2-5pm *4:45pm*
 - 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

©2005-2007 Carlos Guestrin 28

The general learning problem with missing data

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

want to next $\ell(\theta : \mathcal{D})$ ^{iid}

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

summing out hidden vars \mathbf{z}

E-step

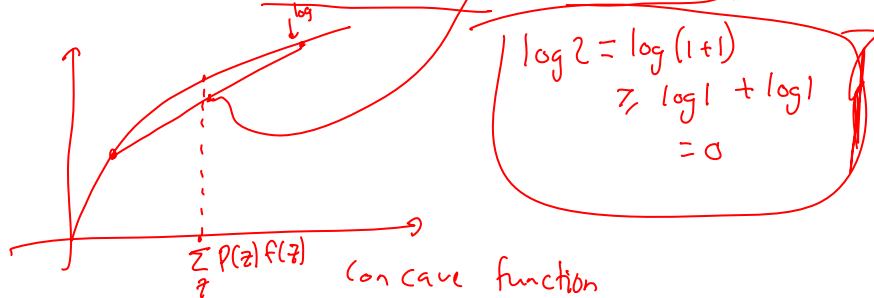
- \mathbf{x} is observed, \mathbf{z} is missing
- Compute probability of missing data given current choice of θ
 - $Q(\mathbf{z} | \mathbf{x}_j)$ for each \mathbf{x}_j
 - e.g., probability computed during classification step
 - corresponds to “classification step” in K-means

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

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 $\log \frac{a}{b} = \log a - \log b$

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

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

The M-step maximizes lower bound on weighted data

- Lower bound from Jensen's:

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

if $\ell(\theta) = \sum_j \log P(z_j, x_j | \theta)$ fully obs data will maximize lower bound
wait to max
entropy
I don't know z , but I introduce $|z|$ data points with weight $Q^{(t+1)}(z|x_j)$

- 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

$Q^{(t+1)}$ is fixed !!
maximize lower bound over θ

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

constant wrt θ

- Maximization step:

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

\equiv MLE w. weighted data

- Use expected counts instead of counts:

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

Convergence of EM

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

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

fix θ , max Q
 fix Q , max θ

- 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

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

MLE using weighted data

constant

E-step also doesn't decrease potential function 1

- Fixing θ to $\theta^{(t)}$:

$$\begin{aligned}
 \ell(\theta^{(t)} : \mathcal{D}) &\geq F(\theta^{(t)}, Q) = \sum_{j=1}^m \sum_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j | \theta^{(t)})}{Q(\mathbf{z} | \mathbf{x}_j)} \\
 &\stackrel{\text{chain rule}}{=} \sum_{j=1}^m \sum_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}_j) \log \underbrace{P(\mathbf{z} | \mathbf{x}_j; \theta^{(t)}) \cdot P(\mathbf{x}_j | \theta^{(t)})}_{Q(\mathbf{z} | \mathbf{x}_j)} \\
 &= \underbrace{\sum_{j=1}^m \sum_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}_j) \log \frac{P(\mathbf{z} | \mathbf{x}_j; \theta^{(t)})}{Q(\mathbf{z} | \mathbf{x}_j)}}_{- \text{KL divergence}} + \sum_{j=1}^m \underbrace{\sum_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}_j)}_{=1} \log P(\mathbf{x}_j | \theta^{(t)}) \\
 &\qquad\qquad\qquad \underbrace{\sum_{j=1}^m \log P(\mathbf{x}_j | \theta^{(t)})}_{\ell(\theta^{(t)} : \mathcal{D})}
 \end{aligned}$$

$\log a \cdot b = \log a + \log b$
 chain rule $P(\mathbf{z}, \mathbf{x}_j | \theta) = P(\mathbf{z} | \mathbf{x}_j; \theta) \cdot P(\mathbf{x}_j | \theta)$
 doesn't depend on \mathbf{z}

KL-divergence

- Measures distance between distributions

$$\underline{KL(Q||P)} = \sum_{\mathbf{z}} Q(\mathbf{z}) \log \frac{Q(\mathbf{z})}{P(\mathbf{z})}$$

- KL=zero if and only if Q=P

E-step also doesn't decrease potential function 2

- Fixing θ to $\theta^{(t)}$:

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

maximize the right side, we know that $KL(Q || P) \geq 0$
 $= 0 \Leftrightarrow P = Q$

as small as possible to max. $\ell(\theta : \mathcal{D})$
 \hookrightarrow by setting $Q(\mathbf{z} | \mathbf{x}_j) = P(\mathbf{z} | \mathbf{x}_j, \theta^{(t)})$

E-step also doesn't decrease potential function 3

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

- 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

$$KL = 0 \iff \begin{matrix} \downarrow \\ F(\theta^{(t)}, Q^{(t+1)}) = \ell(\theta^{(t)} : \mathcal{D}) \end{matrix}$$

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 :

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

- “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 (PCA)

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 *40K words*
- 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 $f: X \rightarrow Y$
 - $X = \langle X_1, \dots, X_n \rangle$ *40K*
 - 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	0.507	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

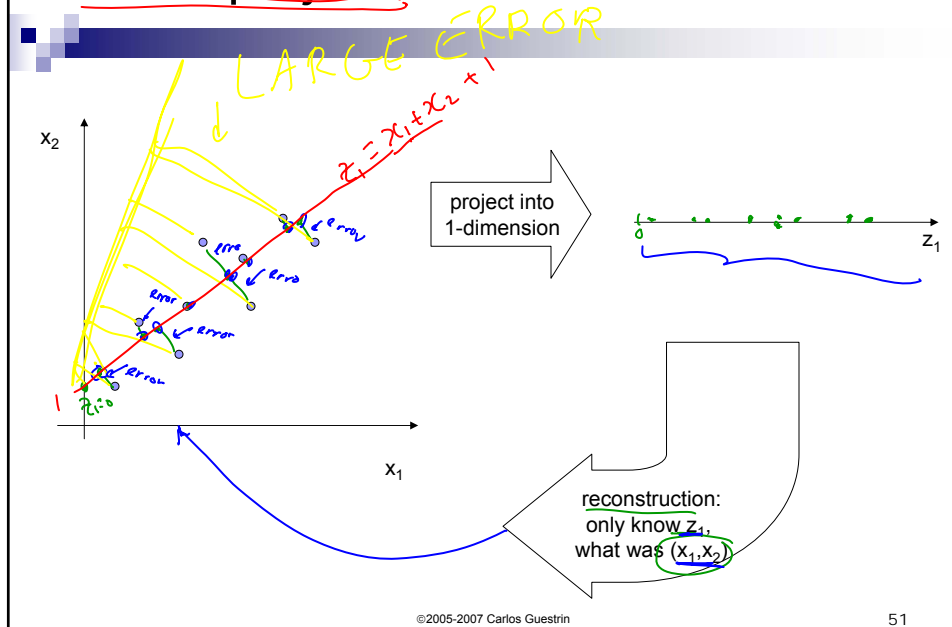
- Rather than picking a subset of the features, we can new features that are combinations of existing features

e.g., feature selection:
use x_1, x_7, x_{11}

low. dim. proj.
 $\tilde{x} = 0.1x_1 + 0.7x_2 - 0.35x_3 \dots$

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

Linear projection and reconstruction



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

$u \cdot v \leftarrow$ dot product

Linear projections, a review

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

