



EM!

Expectation Maximization

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

April 10th, 2006

Announcements



- Reminder: Project milestone due Wednesday beginning of class

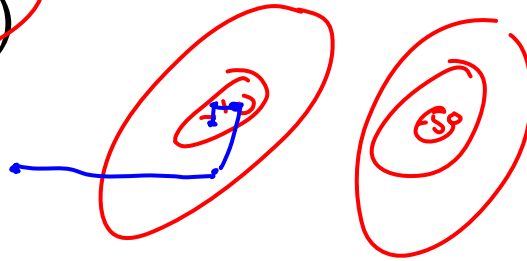
Coordinate descent algorithms

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

Want: $\min_a \min_b F(a,b)$

Coordinate descent:

- fix a, minimize b
- fix b, minimize a
- repeat



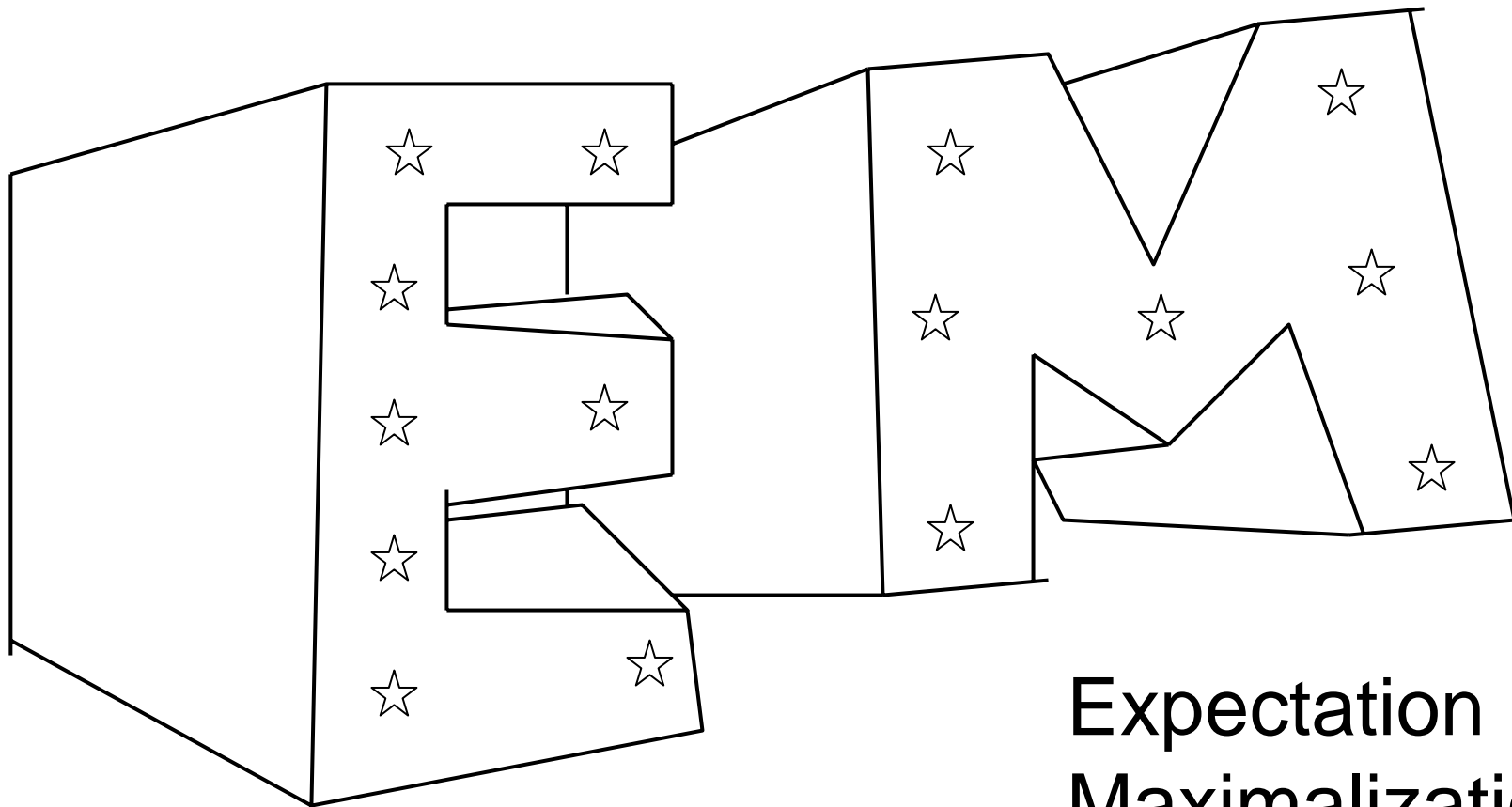
Converges!!!

$$\exists K: \forall a, b \quad f(a,b) \geq K$$

- if F is bounded
- to a (often good) local optimum
 - as we saw in applet (play with it!)



K-means is a coordinate descent algorithm!



Expectation
Maximalization

Back to Unsupervised Learning of GMMs – a simple case

Remember:

We have unlabeled data $x_1 x_2 \dots x_m$

We know there are k classes/clusters

We know $P(y_1) P(y_2) P(y_3) \dots P(y_k)$ ← prob. X_j comes from class $i = P(y_i)$

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

Spherical Gaussians with var. σ^2 (known)

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

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

$$= \prod_{j=1}^m p(x_j | \mu_1 \dots \mu_k) \quad (\text{i.i.d.})$$

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

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

$\underbrace{\hspace{10em}}_{\mathcal{N}(\mu_i, \sigma^2)} \quad \leftarrow \text{prior}$

EM for simple case of GMMs: The E-step

- If we know μ_1, \dots, μ_k → easily compute prob.

prob. x_j comes from $y=i$ point x_j belongs to class $y=i$

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

normalize

$$\left. \begin{array}{l} P(y=1|x_j, \mu) \propto 0.3 \\ P(y=0|x_j, \mu) \propto 0.2 \end{array} \right\} \begin{array}{l} = 0.6 \\ = 0.4 \end{array}$$

repeat for each x_j .

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 $P(y=i|x_j)$:

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

$$\mu_0 = \frac{0.8 \times 1 + 0.2 \times 2 + 0.1 \times 3}{0.8 + .2 + .1}$$

weighted version

x_j	$P(y=0 x_j)$
1	0.8
2	0.2
3	0.1

if we knew x_j comes from $C(j)$

$$\mu_i = \sum_{j:C(j)=i} x_j$$

$$\sum_{j:C(j)=i} 1$$

E.M. for GMMs

E-step

Compute “expected” classes of all datapoints for each class

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

∝ prior

*Just evaluate
a Gaussian at
 x_j*

M-step

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

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

*} weighted data
using $p(y|x)$*

E.M. for General GMMs

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

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

$$\lambda_t = \{ \underbrace{\mu_1^{(t)}, \mu_2^{(t)} \dots \mu_k^{(t)}}_{\text{means}}, \underbrace{\Sigma_1^{(t)}, \Sigma_2^{(t)} \dots \Sigma_k^{(t)}}_{\text{covariance}}, \underbrace{p_1^{(t)}, p_2^{(t)} \dots p_k^{(t)}}_{\text{priors}} \}$$

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

solution depends on starting λ_0
use random restarts

M-step

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

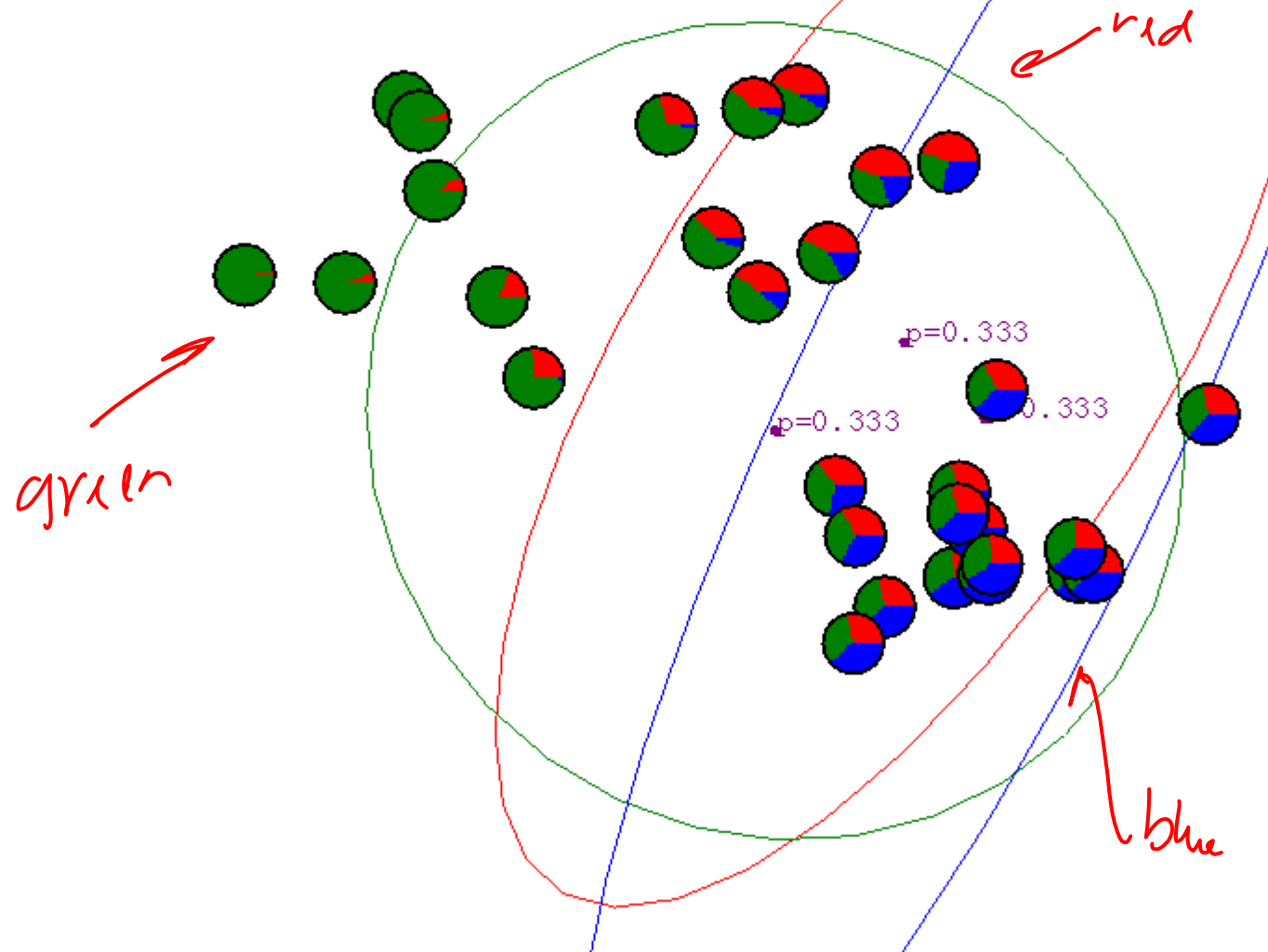
$$\mu_i^{(t+1)} = \frac{\sum_j P(y = i | x_j, \lambda_t) x_j}{\sum_j P(y = i | x_j, \lambda_t)} \quad \Sigma_i^{(t+1)} = \frac{\sum_j P(y = i | x_j, \lambda_t) [x_j - \mu_i^{(t+1)}][x_j - \mu_i^{(t+1)}]^T}{\sum_j P(y = i | x_j, \lambda_t)}$$

weighted data

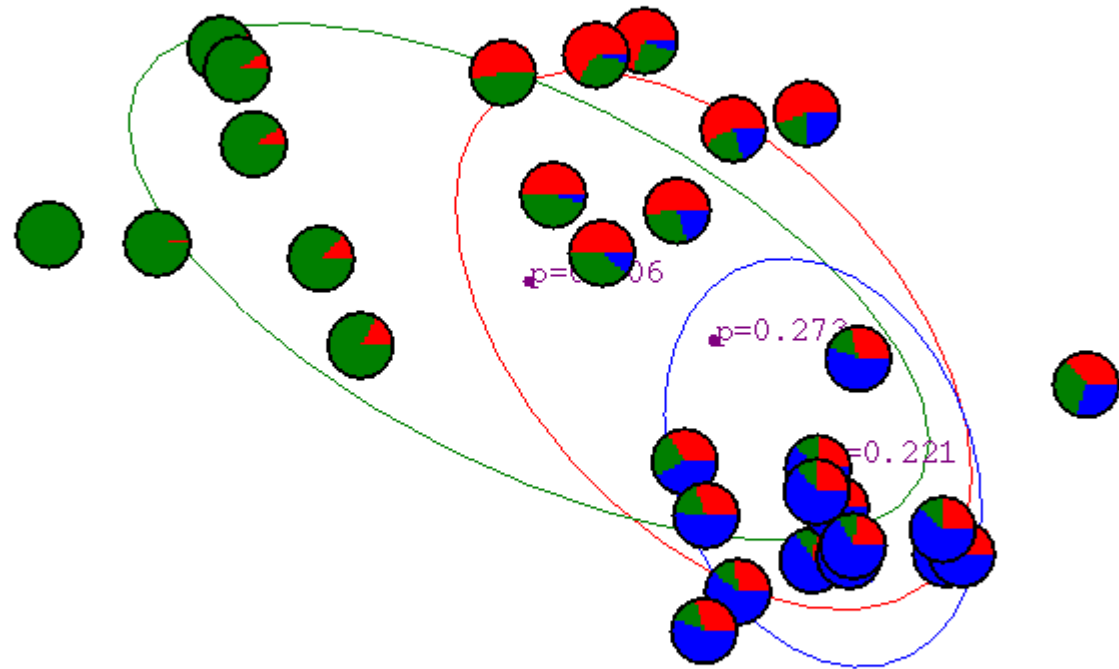
$$p_i^{(t+1)} = \frac{\sum_j P(y = i | x_j, \lambda_t)}{m}$$

$m = \#$ records

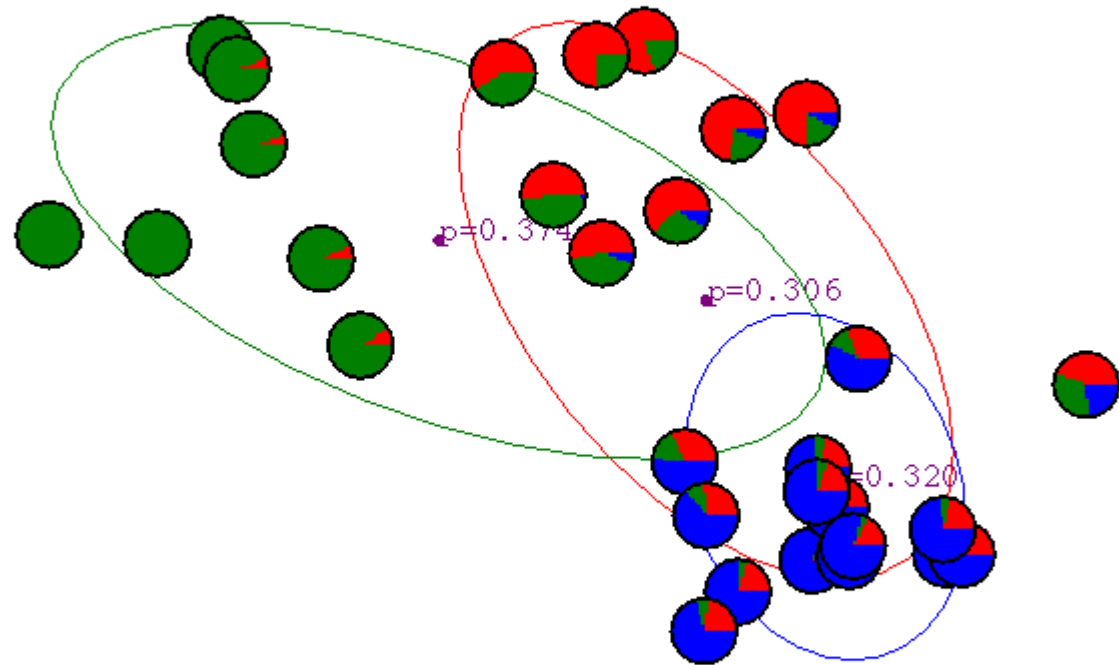
Gaussian Mixture Example: Start



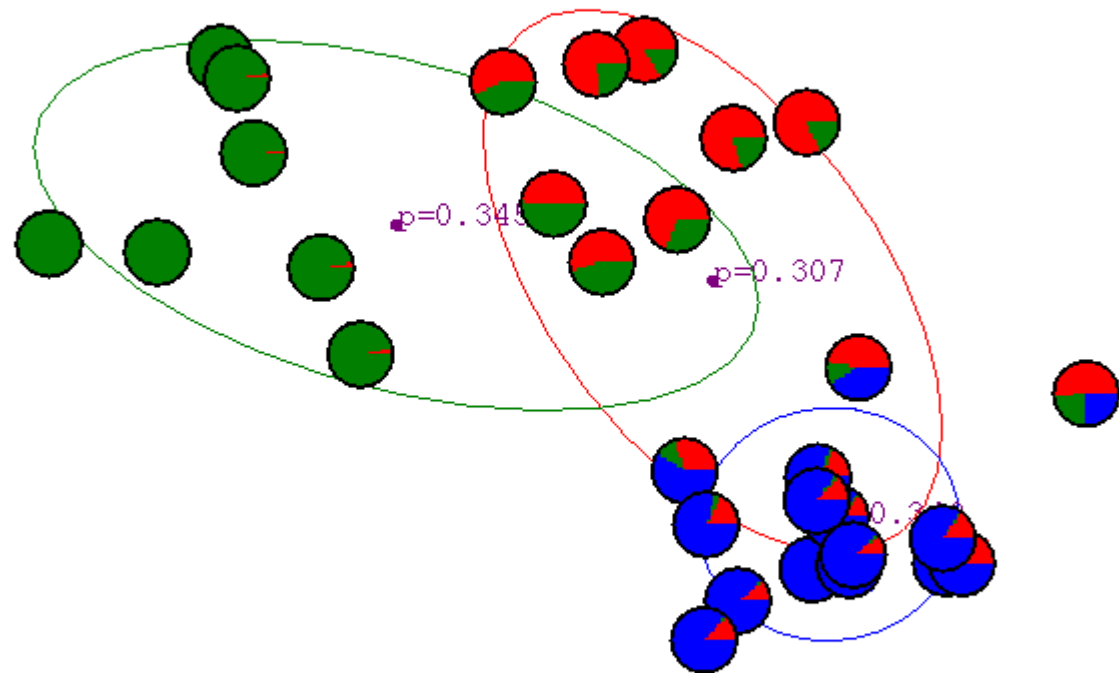
After first iteration



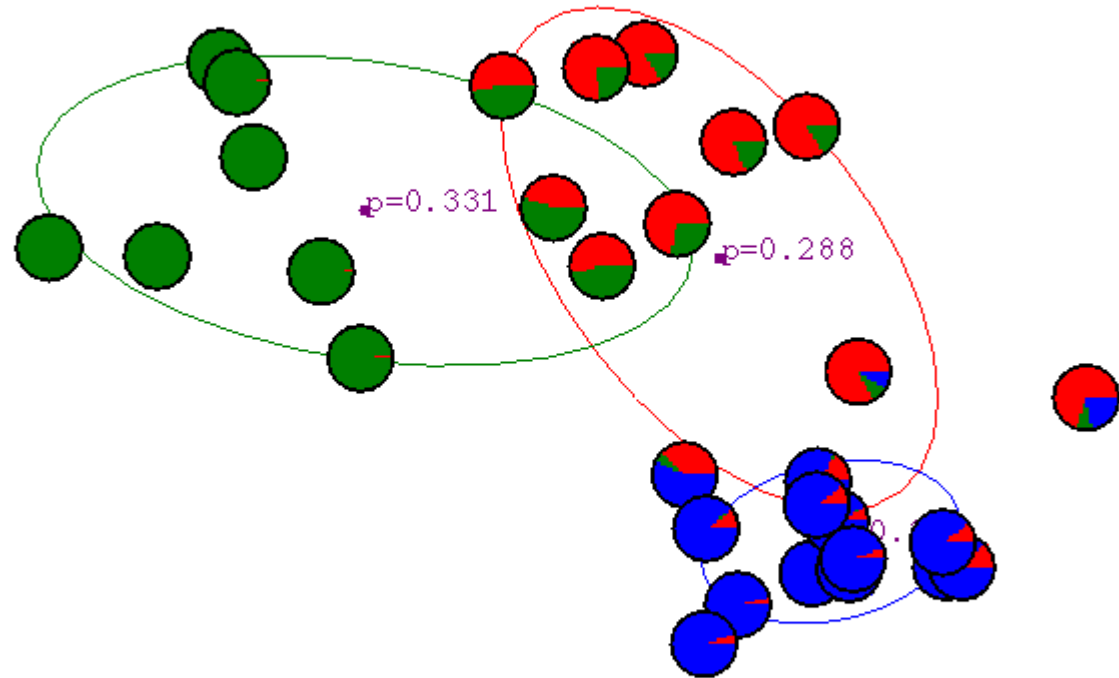
After 2nd iteration



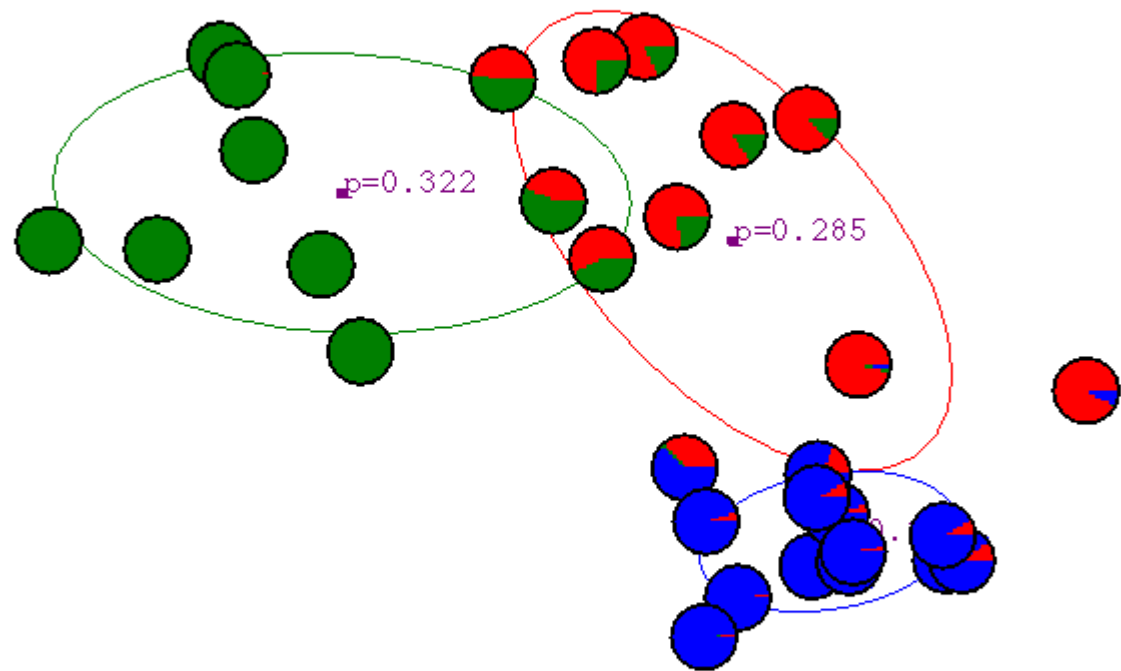
After 3rd iteration



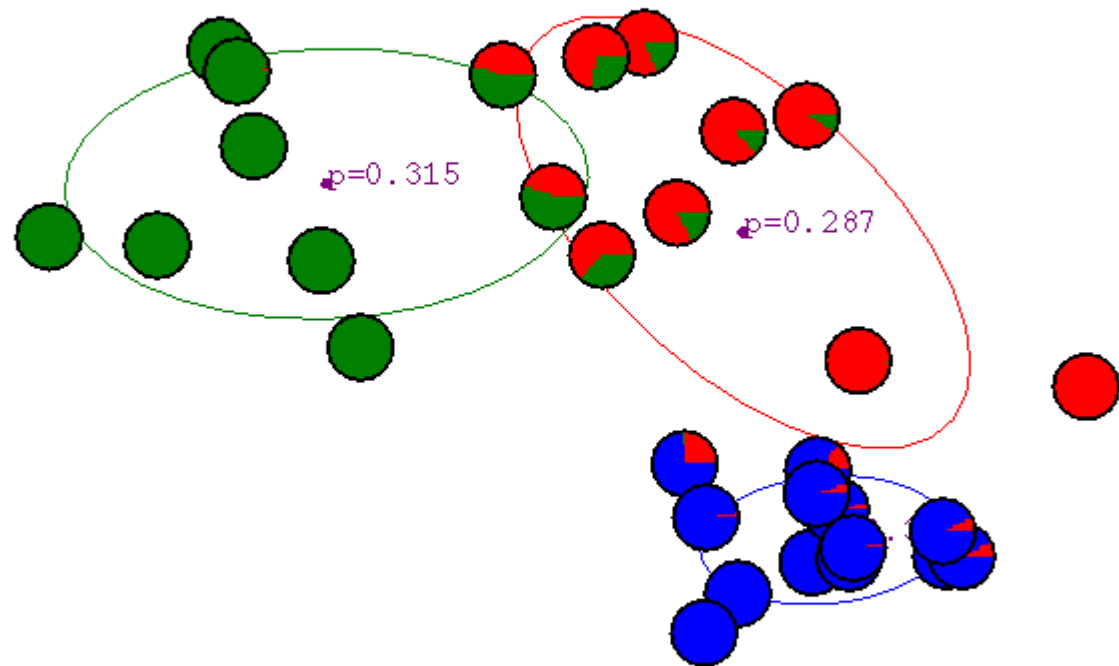
After 4th iteration



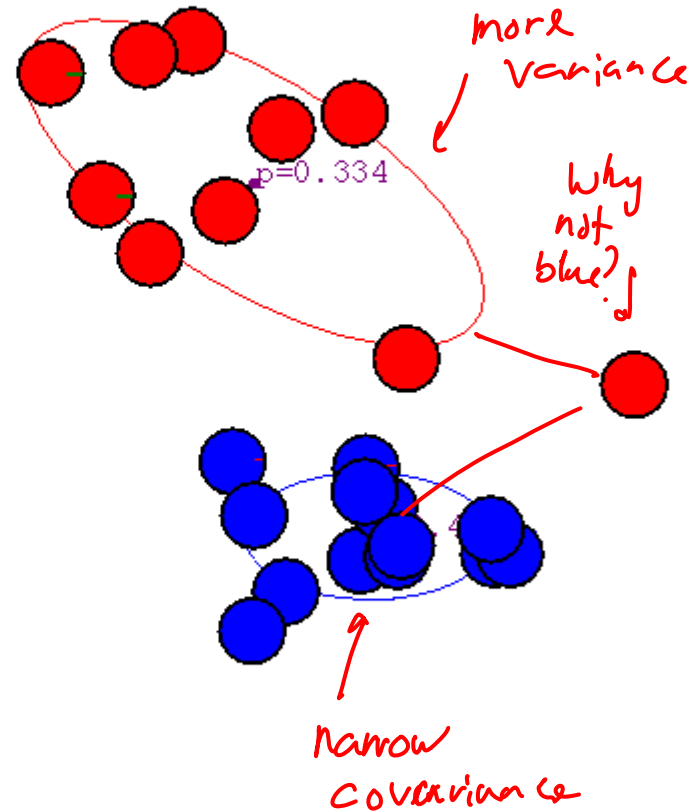
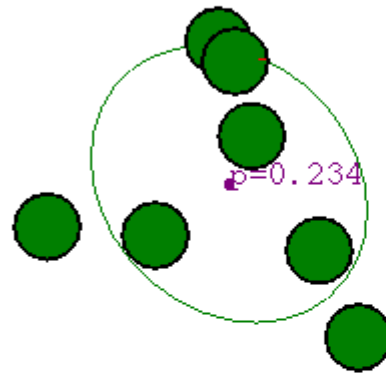
After 5th iteration



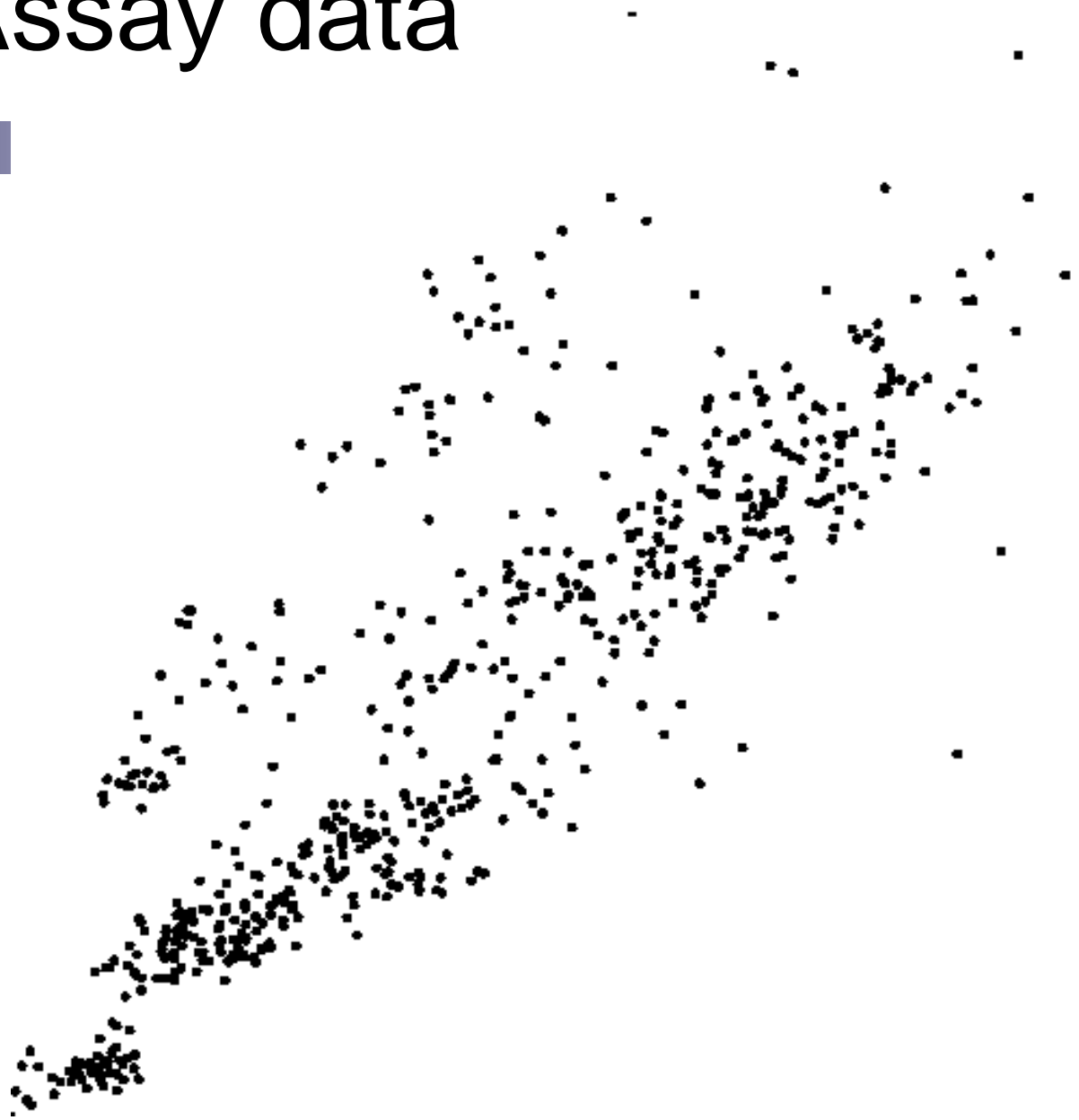
After 6th iteration



After 20th iteration



Some Bio Assay data

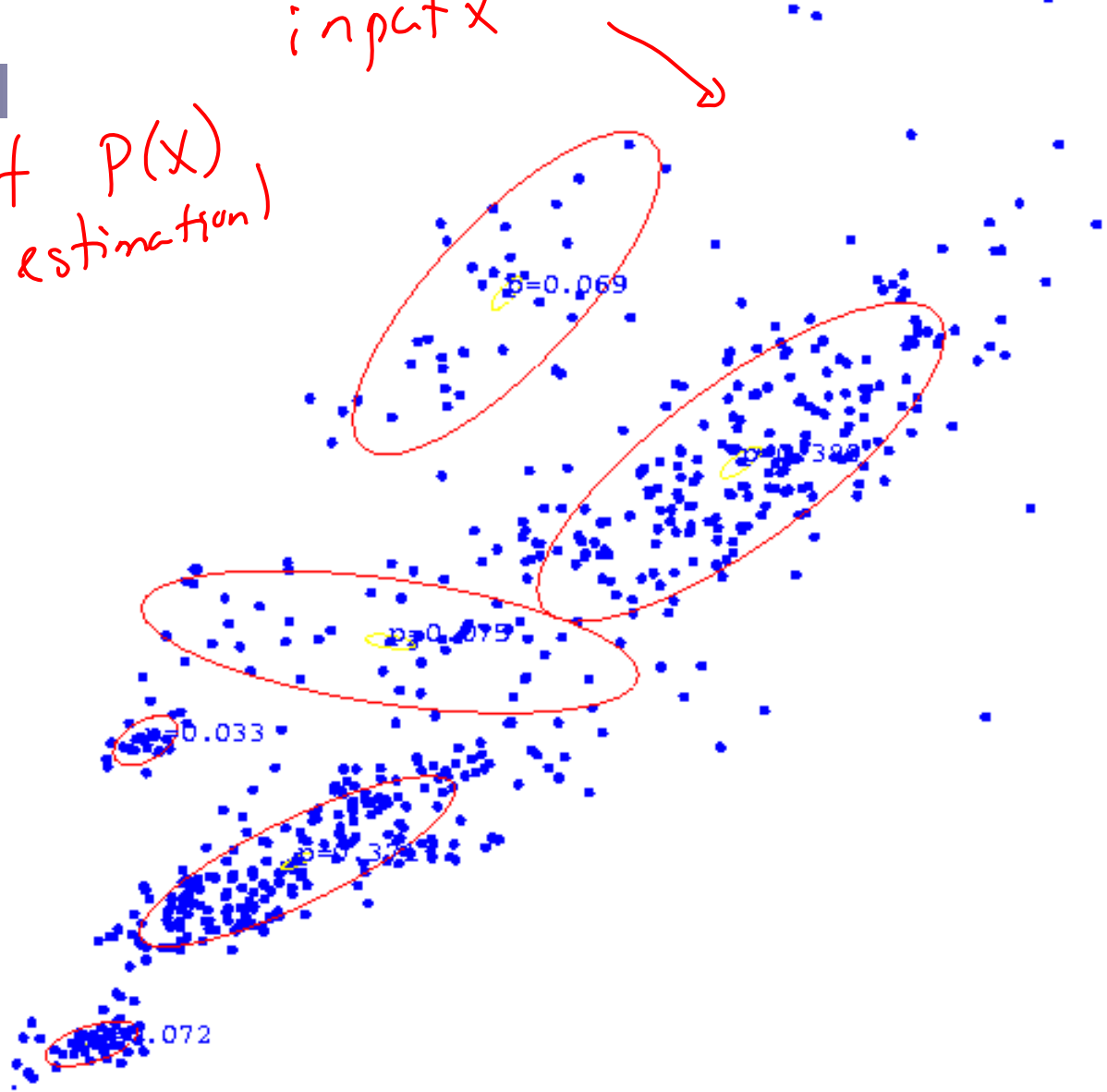


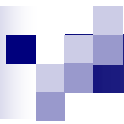
GMM clustering of the assay data



represent $P(x)$
(density estimation)

input x





Resulting Density Estimator

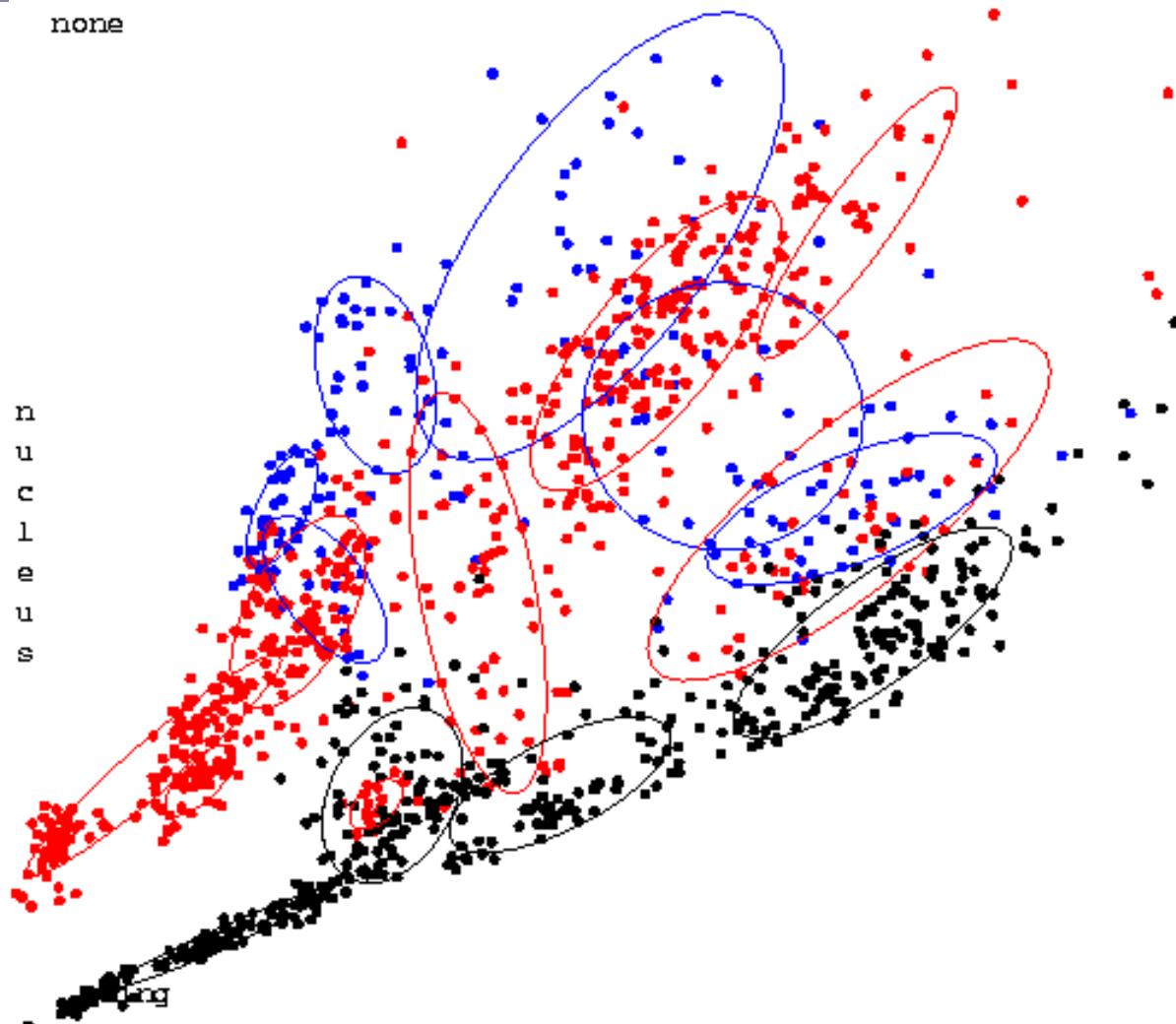


Three classes of assay

(each learned with its own mixture model)

Compound =
IL-1
TNF
none

n
u
c
l
e
u
s

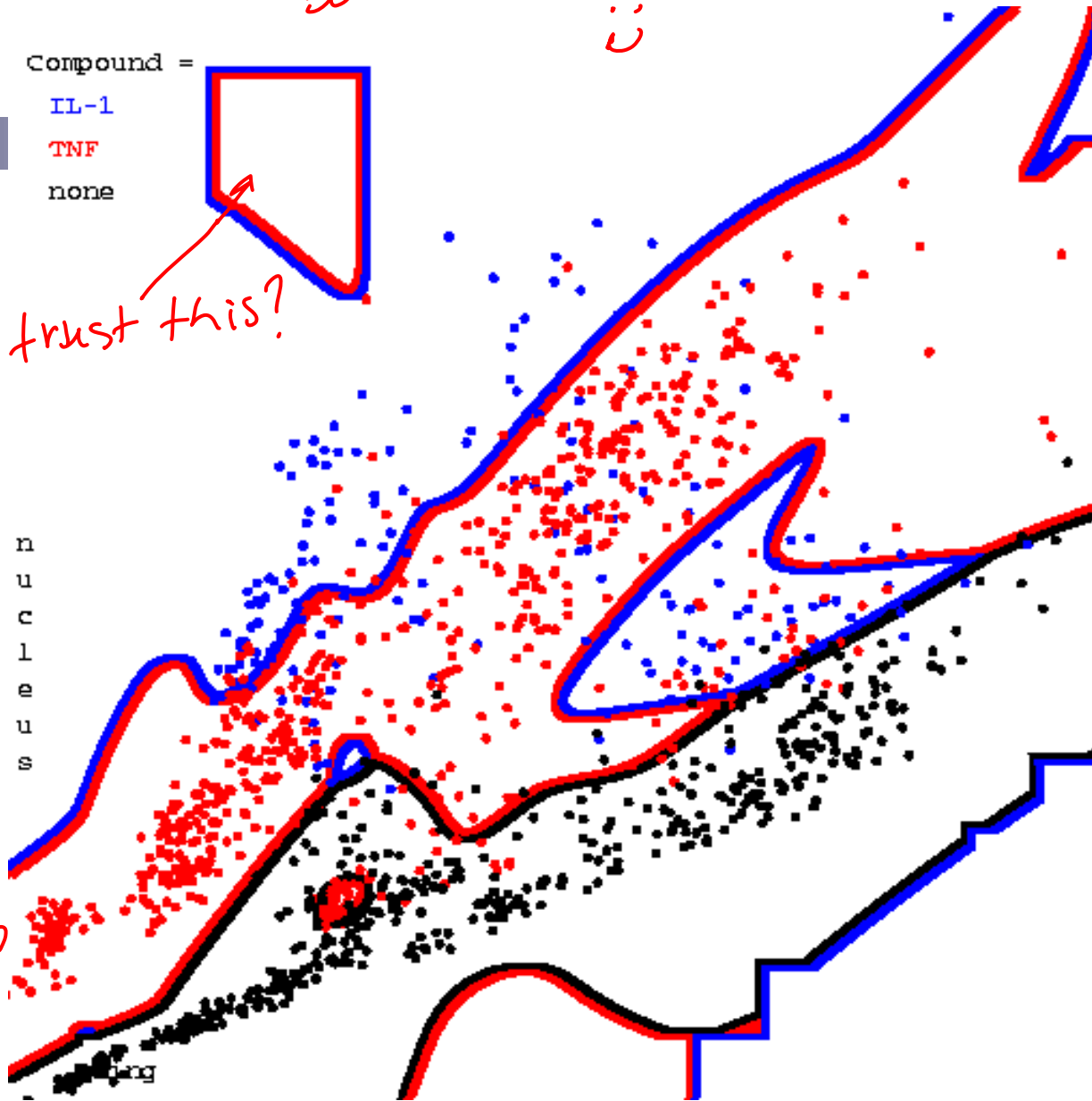


Resulting Bayes Classifier

Can learn complicated boundaries!

Compound =
IL-1
TNF
none

n
u
c
l
e
u
s



$P(\text{red} | x)$
 $\rightarrow P(\text{blue} | x)$

can't do this easily with discriminative classifier.

Classifier. →

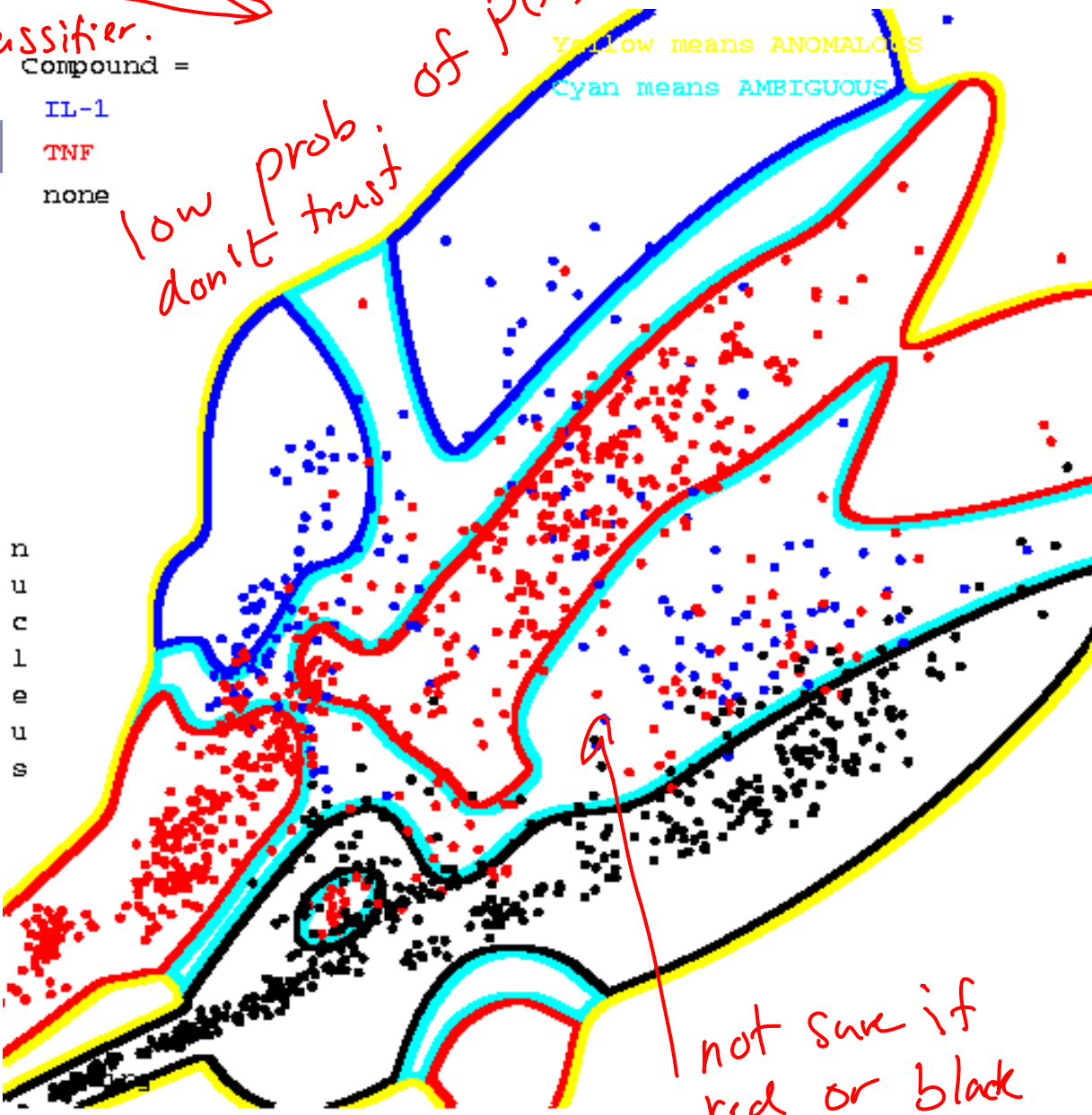
Compound =
IL-1
TNF
none

Yellow means ANOMALOUS
Cyan means AMBIGUOUS

low prob. of $P(x)$
don't trust

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

n
u
c
l
e
s



Yellow means anomalous

Cyan means ambiguous

not sure if red or black

The general learning problem with missing data

more general than GMMs

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

$$\underline{\ell(\theta : \mathcal{D})} = \log \prod_{j=1}^m P(\mathbf{x}_j | \theta) \quad (\text{i.i.d.})$$

↑
marginal
likelihood
of data

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

want to $\operatorname{argmax}_{\theta} \ell(\theta : \mathcal{D})$

E-step

- \mathbf{x} is observed, \mathbf{z} is missing

$\underbrace{\langle 1, 0, true, \dots \rangle}_{\mathbf{x}} \quad \underbrace{\langle ?, ??, \dots \rangle}_{\mathbf{z}}$

- 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

iteration t of EM

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

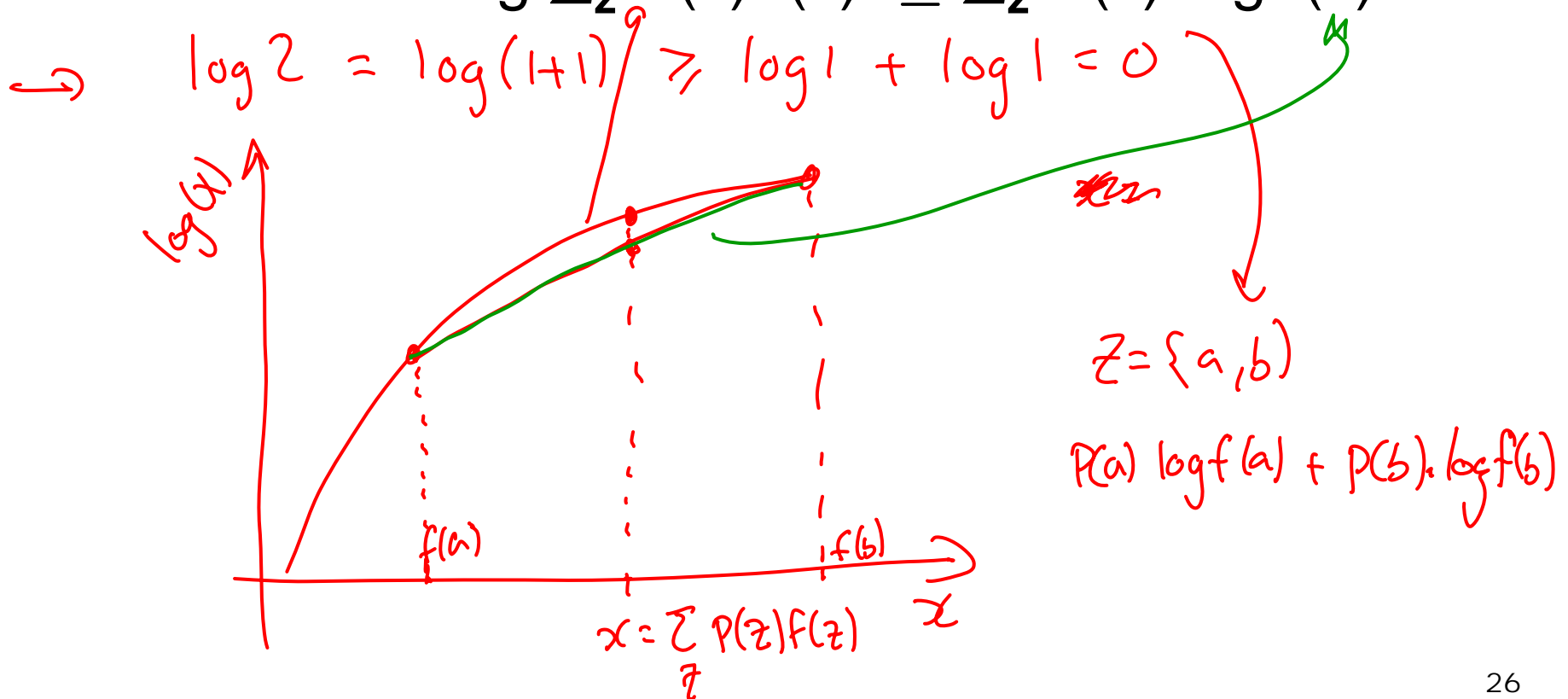
\mathbf{z} is an assignment to hidden vars.
→ example: $\mathbf{z} = (y=1)$

Jensen's inequality

$$\log(0.6f(a) + 0.4f(b)) \geq 0.6 \log f(a) + 0.4 \log f(b)$$

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

■ **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

$$H(x) = -\sum_x p(x) \log p(x)$$

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

$$\ell(\theta; \mathcal{D}) =$$

$$\sum_{j=1}^m \log \sum_{\mathbf{z}} P(\mathbf{z}, x_j | \theta)$$

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

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

Jensen's inequality: $\geq \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \cdot \log \frac{P(\mathbf{z}, \mathbf{x}_j | \theta^{(t)})}{Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j)}$

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

$$- \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \log Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j)$$

$$m \cdot \hat{H}_{Q^{(t+1)}}(\mathbf{z} | \mathbf{x})$$

$$H(Q^{(t+1)})$$

The M-step maximizes lower bound on weighted data

- Lower bound from Jensen's:

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

likelihood
all data

doesn't
depend on
 θ

fix Q , max over θ \uparrow MLE for weighted data \rightarrow irrelevant

- 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

doesn't play
in M-step

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

■ Maximization step:

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


same as fully observed data
weighted by $Q^{(t+1)}$

■ Use expected counts instead of counts:

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

Convergence of EM

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

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


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

fix $Q \rightarrow \max$ over θ

fix $\theta \rightarrow \max$ over Q
will see soon

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

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

$\theta^{(t+1)}$ maximizes this part

$$\Rightarrow F(\theta^{(t+1)}, Q^{(t+1)}) \geq F(\theta^{(t)}, Q^{(t+1)})$$

Fixed $Q^{(t+1)}$, improved over θ

doesn't change
in M step

E-step also doesn't decrease potential function 1

$$\log a \cdot b = \log a + \log b$$

$$P(z, x_j | \theta^{(t)}) = P(z | x_j, \theta^{(t)}) \cdot P(x_j | \theta^{(t)})$$

- Fixing θ to $\theta^{(t)}$:

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

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

doesn't depend on \mathbf{z}

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

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

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

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

KL-divergence

$$\log \frac{Q}{P} = - \log \frac{P}{Q}$$

- Measures distance between distributions

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

if Q is "far" from $P \Rightarrow KL(Q||P)$ ^{very positive} large

if Q is "close" to $P \Rightarrow KL(Q||P)$ low

KL always ≥ 0

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

E-step also doesn't decrease potential function 2

- Fixing θ to $\theta^{(t)}$:

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

doesn't
depend on
Q

want to maximize
 \Rightarrow set $\text{KL} = 0$

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

E-step also doesn't decrease potential function 3

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

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

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

$\forall z$ write a table
(if z is discrete)

$\forall j$

z	$Q(z \mathbf{x}_j)$
true	0.8
false	0.2

\Rightarrow M step
weight
 $(\mathbf{x}_j, z=t), 0.8$
 $(\mathbf{x}_j, z=f), 0.2$

- 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 $\ell(\theta : \mathcal{D})$):

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

learn from weighted data

invariant

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

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

doesn't depend on Q

setting to ϕ

- **“Realigns” F with likelihood:**

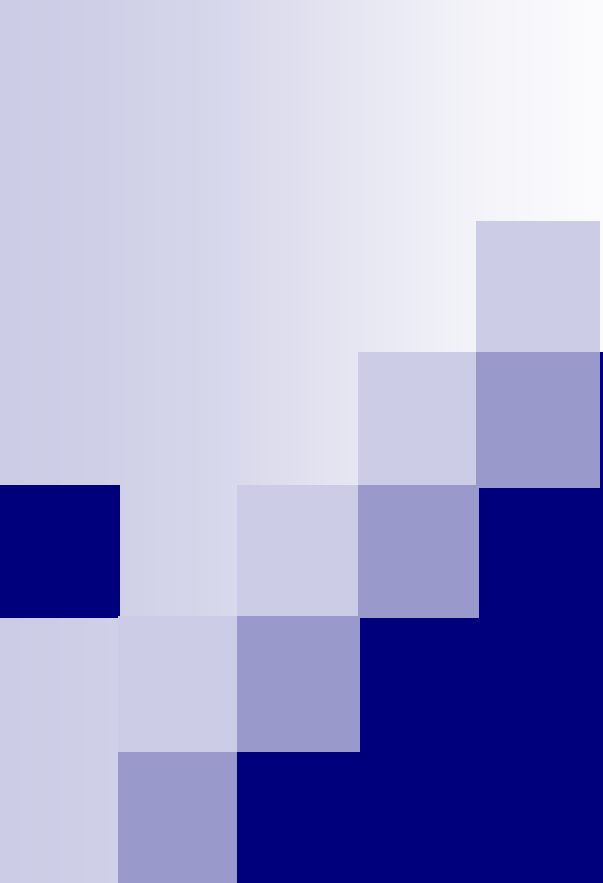
$$F(\theta^{(t)}, Q^{(t+1)}) = \ell(\theta^{(t)} : \mathcal{D})$$

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 *(random restarts)*
- 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>



EM for HMMs a.k.a. The Baum-Welch Algorithm

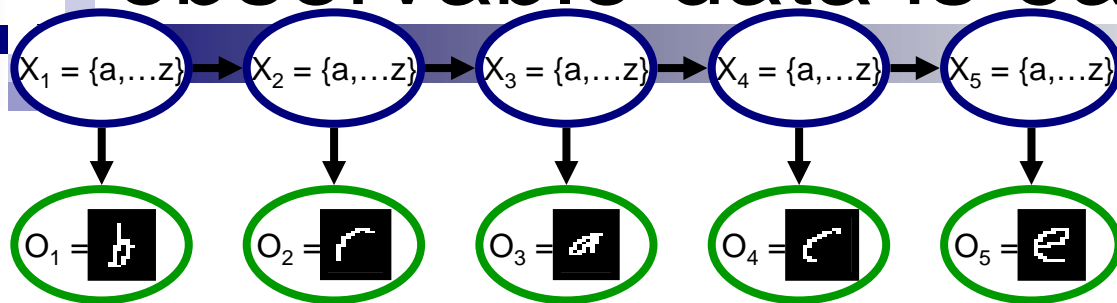
Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

April 10th, 2006

Learning HMMs from fully observable data is easy



Learn 3 distributions:

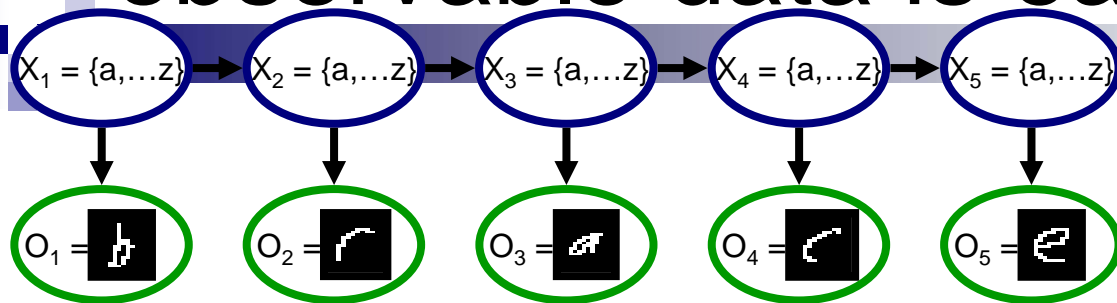
$$P(X_1^a) = \frac{\text{count}(\# \text{ first letter was } a)}{N = \text{dataset size}}$$

$$P(O_i^{\text{pixel } l \text{ is white}} | X_i^a) = \frac{\text{count}(\text{pixel } l \text{ was white, } X_i = a)}{\text{count}(X_i = a)}$$

$$P(X_i^a | X_{i-1}^b) = \frac{\text{count}(a \text{ appears after } b)}{\text{count}(\# \text{ of } b\text{'s that are not at the end of the word})}$$

select training data where letter was a

Learning HMMs from fully observable data is easy



Learn 3 distributions:

$$P(X_1^a) = \frac{\text{count}(\# \text{ first letter was } a)}{N = \text{dataset size}}$$

$$P(O_i^{\text{pixel } l \text{ is white}} | X_i^a) = \frac{\text{count}(\text{pixel } l \text{ was white, } X_i = a)}{\text{any } a \text{ in}}$$

select training data where letter was a

$$P(X_i^a | X_{i-1}^b)$$

What if **O** is observed, but **X** is hidden

Log likelihood for HMMs with hidden \mathbf{X}

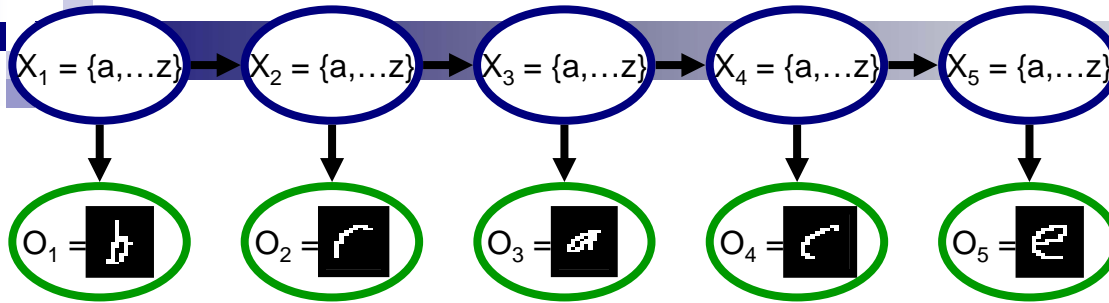
- Marginal likelihood – \mathbf{O} is observed, \mathbf{X} is missing
 - For simplicity of notation, we'll consider training data consists of only one sequence:

$$\begin{aligned}\ell(\theta : \mathcal{D}) &= \log P(\mathbf{o} | \theta) \\ &= \log \sum_{\mathbf{x}} P(\mathbf{x}, \mathbf{o} | \theta)\end{aligned}$$

- If there were m sequences:

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

E-step

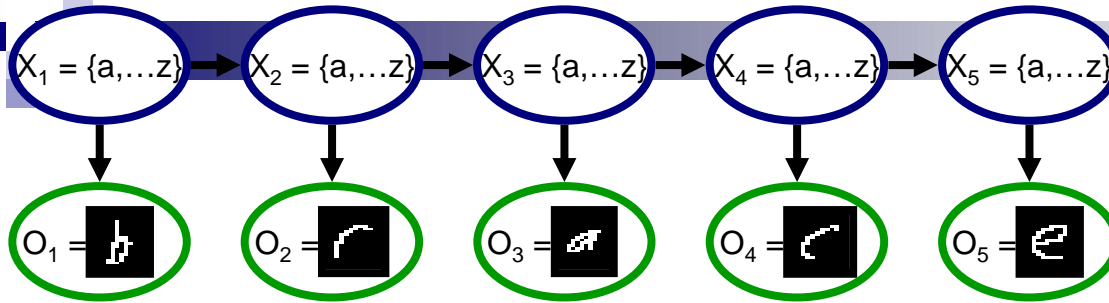


- E-step computes probability of hidden vars \mathbf{x} given \mathbf{o}

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

- Will correspond to inference
 - use forward-backward algorithm!

The M-step



■ Maximization step:

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

■ Use expected counts instead of counts:

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

Starting state probability $P(X_1)$

- Using expected counts

- $P(X_1=a) = \theta_{X_1=a}$

$$\theta_{X_1=a} = \frac{\sum_{j=1}^m Q(X_1 = a \mid \mathbf{o}^{(j)})}{m}$$

Transition probability $P(X_{t+1}|X_t)$

- Using expected counts

- $P(X_{t+1}=a|X_t=b) = \theta_{X_{t+1}=a|X_t=b}$

$$\theta_{X_{t+1}=a|X_t=b} = \frac{\sum_{j=1}^m \sum_{t=1}^{n-1} Q(X_{t+1} = a, X_t = b | \mathbf{o}^{(j)})}{\sum_{j=1}^m \sum_{t=1}^{n-1} \sum_{i=1}^k Q(X_{t+1} = i, X_t = b | \mathbf{o}^{(j)})}$$

Observation probability $P(O_t|X_t)$

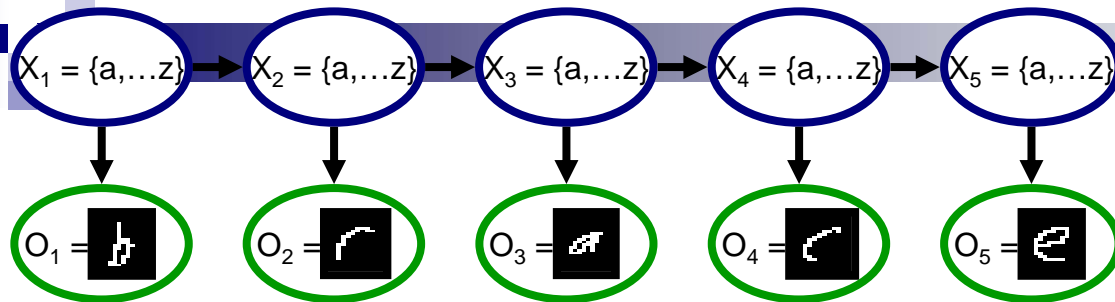
- Using expected counts

- $P(O_t=a|X_t=b) = \theta_{O_t=a|X_t=b}$

$$\theta_{O_t=a|X_t=b} = \frac{\sum_{j=1}^m \sum_{t=1}^n \delta(o_t^{(j)} = a) Q(X_t = b | o^{(j)})}{\sum_{j=1}^m \sum_{t=1}^n Q(X_t = b | o^{(j)})}$$

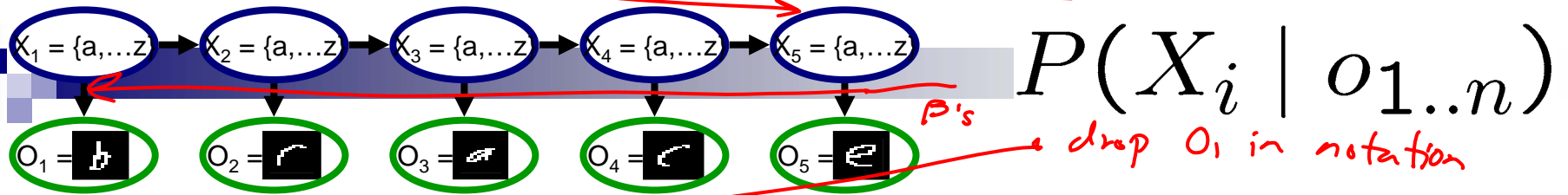
E-step revisited

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



- E-step computes probability of hidden vars \mathbf{x} given \mathbf{o}
- Must compute:
 - $Q(x_t = a | \mathbf{o})$ – marginal probability of each position
 - $Q(x_{t+1} = a, x_t = b | \mathbf{o})$ – joint distribution between pairs of positions

FS The forwards-backwards algorithm



■ Initialization: $\alpha_1(X_1) = P(X_1)P(o_1 | X_1)$

■ For $i = 2$ to n

□ Generate a forwards factor by eliminating X_{i-1}

sum out previous var prob obs

$$\alpha_i(X_i) = \sum_{x_{i-1}} P(o_i | X_i) P(X_i | X_{i-1} = x_{i-1}) \alpha_{i-1}(x_{i-1})$$

transition prob

■ Initialization: $\beta_n(X_n) = 1$

■ For $i = n-1$ to 1

□ Generate a backwards factor by eliminating X_{i+1}

$\forall x_i$

$$\beta_i(X_i) = \sum_{x_{i+1}} P(o_{i+1} | x_{i+1}) P(x_{i+1} | X_i) \beta_{i+1}(x_{i+1})$$

x_i

■ $\forall i$, probability is: $P(X_i | o_{1..n}) = \alpha_i(X_i) \beta_i(X_i)$

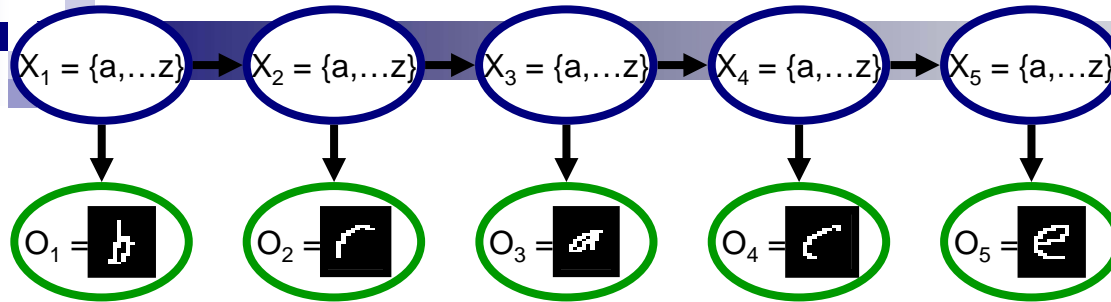
$\alpha_n(X_n)$
normalized
 $= P(X_n | O_{1:n})$

$\beta_1(X_1) \alpha_1(X_1)$
normalized
 $= P(X_1 | O_{1:n})$

$\alpha_5(a)$
 $\alpha_5(b)$
 \vdots
 $\alpha_5(z)$

E-step revisited

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



- E-step computes probability of hidden vars \mathbf{x} given \mathbf{o}
- Must compute:
 - $Q(x_t = a | \mathbf{o})$ – marginal probability of each position
 - Just forwards-backwards!
 - $Q(x_{t+1} = a, x_t = b | \mathbf{o})$ – joint distribution between pairs of positions
 - Homework! 😊