# Document Retrieval
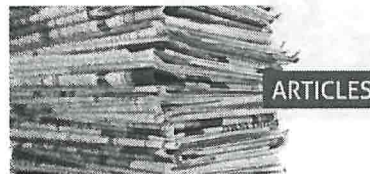
- **Goal:** Retrieve documents of interest
- **Challenges:**
  - Tons of articles out there
  - How should we measure similarity?



ARTICLES

29

---

# Task 1: Find Similar Documents

- **So far...**
  - **Input:** Query article  X
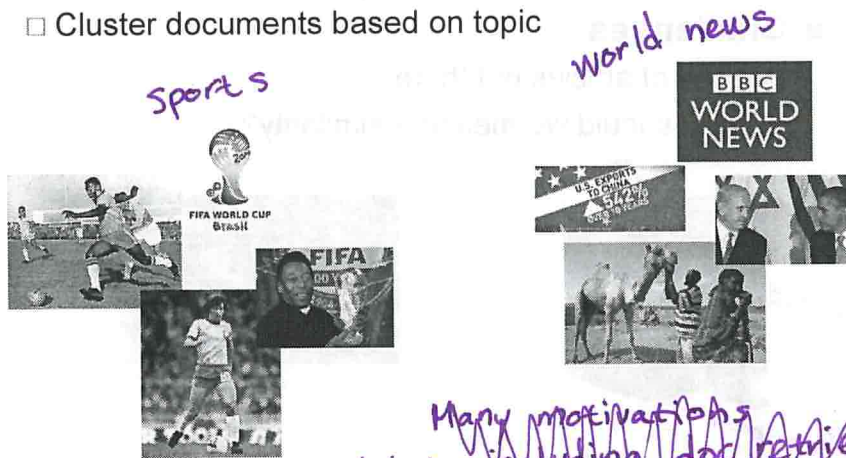  - **Output:** Set of k similar articles



FIFA WORLD CUP
Brasil

FIFA

30

1

# Task 2: Cluster Documents

- **Now:**
  - ☐ Cluster documents based on topic

*Sports*

*World news*

Many motivations including doc retrieval

More global description of corpus

---

# Document Representation

- Bag of words model

document *d*

previously:

$$x = \begin{bmatrix} \\ \\ \end{bmatrix}$$

vector fcn of word counts (e.g. tf-idf)

performed operations on this vector

now:

$$x = \{w_1, \ldots, w_N\}$$

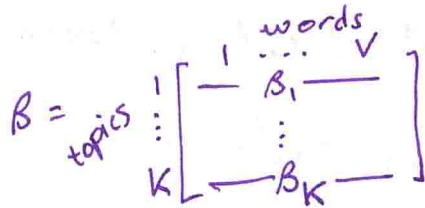unordered set of $N$ words $w_i \in V$ (vocab) in doc

# A Generative Model

*probabilistic model for simulating obs.*

- Documents: $x^1, \ldots, x^D$   with   $x^d = \{w_1^d, \ldots, w_{N_d}^d\}$
- Associated topics: $z^1, \ldots, z^D$   with   $z^d \in \{1, \ldots, K\}$
- Parameters: $\theta = \{\pi, \beta\}$

*total # of topics / clusters*

$$\pi = [\pi_1, \ldots, \pi_K] \quad \text{topic probabilities}$$

with

$$Pr(z^d = k) = \pi_k$$

$$\beta = \begin{array}{c} \text{topics} \\ \end{array} \begin{array}{c} 1 \\ \vdots \\ K \end{array} \left[ \begin{array}{c} \overset{\text{words}}{\overset{1 \cdots V}{- \beta_1 -}} \\ \vdots \\ - \beta_K - \end{array} \right]$$
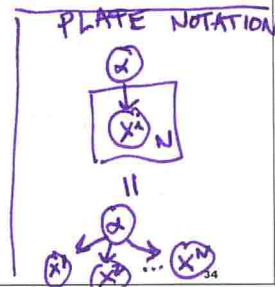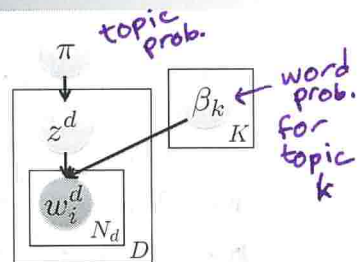
*word probabilities for each topic*

---

# A Generative Model

- Documents: $x^1, \ldots, x^D$
- Associated topics: $z^1, \ldots, z^D$
- Parameters: $\theta = \{\pi, \beta\}$
- Generative model:

*"drawn from"*

$$z^d \sim \pi$$

$$w_i^d \mid z^d \sim \beta_{z^d} \qquad i = 1, \ldots, N_d$$

*"given"*

Given topic $z^d = k$ for doc $d$,
draw each word ind. from $\beta_k$

$\pi$   *topic prob.*

$z^d$

$\beta_k$   *word prob. for topic k*

$K$

$w_i^d$

$N_d$

$D$

PLATE NOTATION

# Form of Likelihood

- Conditioned on topic...

$$p(x^d \mid z^d, \beta) = \prod_{i=1}^{N_d} p(w_i^d \mid z^d, \beta) = \prod_{i=1}^{N_d} \beta_{z_i^d w_i^d}$$

*unobserved/latent*

- Marginalizing latent topic assignment:

$$p(x^d \mid \beta, \pi) = \sum_{k=1}^{K} p(x^d, z^d = k \mid \beta, \pi)$$

*Convex comb. of $p(x^d \mid z^d, \beta)$*

$$= \sum_{k=1}^{K} p(x^d \mid z^d = k, \beta) \, p(z^d = k \mid \pi)$$

$$= \sum_{k=1}^{K} \pi_k \, p(x^d \mid z^d = k, \beta)$$

©Emily Fox 2013

35

---

# Gaussian Mixture Model

- Most commonly used mixture model
- Observations: $x^1, \ldots, x^N$

  with $x^i \in \mathbb{R}^d$

- Parameters:

  $$\pi = [\pi_1, \ldots, \pi_K]$$

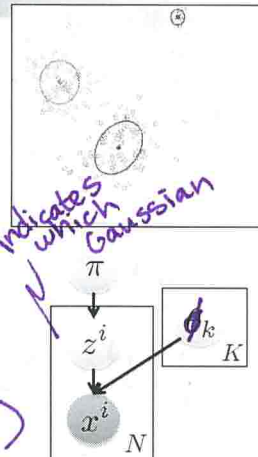  $$\theta = \{\theta_k\} = \{\mu_k, \Sigma_k\}$$

  *← params for cluster $k$*

- Likelihood:

  $$p(x^i \mid \theta) = \sum_{k=1}^{K} \pi_k \, N(x^i; \mu_k, \Sigma_k)$$

*indicates which Gaussian*

$\pi$

$z^i$    $\phi_k$   $K$

$x^i$   $N$

- Ex. $z^i$ = country of origin, $x^i$ = height of $i^{th}$ person
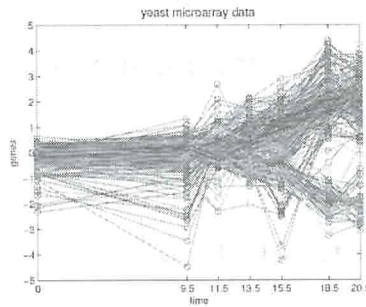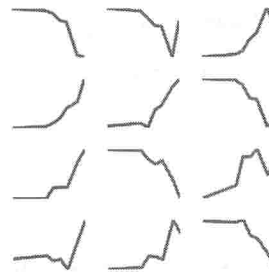  - $k^{th}$ mixture component = distribution of heights in country $k$

©Emily Fox 2013

36

2

# Another Example

(Taken from Kevin Murphy's ML textbook)
- Data: gene expression levels $\cdots / \backslash \, , / \rightarrow \begin{bmatrix} \\ \end{bmatrix} \in \mathbb{R}^7$
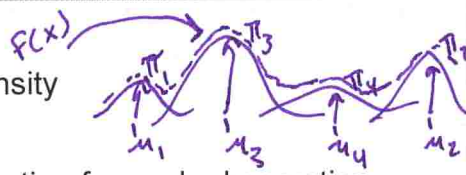- Goal: cluster genes with similar expression trajectories

cluster means

37

---

# 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:

$f(x)$   $\pi_3$   $\pi_2$   $\pi_1$   $\pi_4$

$\mu_1 \quad \mu_3 \quad \mu_4 \quad \mu_2$

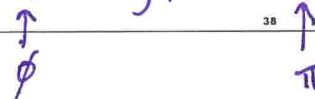"responsibility of point $i$ for cluster $k$"

$$p(z^i = k \mid x^i, \theta) = \frac{p(x^i \mid z^i = k, \phi)\, p(z^i = k \mid \pi)}{p(x^i \mid \theta)}$$

  - Hard clustering:

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

$$= \arg\max_k \log p(x^i \mid z^i = k, \theta) + \log p(z^i = k \mid \theta)$$

$\phi \qquad\qquad \pi$

38

---

3

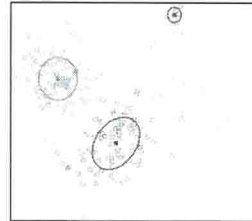# Issues

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

- Log likelihood is not convex in the parameters
  - No closed form gradient updates
  - Problem is simpler for "complete data likelihood"

  *IF we knew $z^i$*

- More on this next time...

---

# What you need to know

- Mixture model formulation
  - Generative model
  - Likelihood

# ML Estimate of Mixture Model Params

- Log likelihood

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

$$p(x|\theta) = \prod_i p(x_i | \theta)$$

- Want ML estimate

$$\hat{\theta}^{ML} = \arg\max_\theta L_x(\theta)$$

- Assume exponential family $p(x, z \mid \theta) = \dfrac{1}{Z(\theta)} e^{\theta' \phi(x,z)}$

$$L_x(\theta) = \sum_i \log \left( \sum_{z_i} e^{\theta^T \phi(z_i, x_i)} \right) - N \log Z(\theta)$$

- Neither convex nor concave and local optima

---

# If "complete" data were observed…

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

$$\pi_{z_i}$$

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

$$= \sum_k \sum_{i:z_i=k} \log p(x_i | z_i, \phi_k) + \sum_{j=1}^{K-1} N_j \log \pi_j + N_k \log\left(1 - \sum_{j=1}^{K-1} \pi_j\right)$$

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

$$[\{z_i : z_i = j\}] \qquad \sum \pi_j = 1$$

$$\hat{\phi}_k = \arg\max_{\phi_k} \sum_{i:z_i=k} \log p(x_i | z_i, \phi_k) \qquad \hat{\pi}_k = \frac{N_k}{N}$$

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

$$\hat{\mu}_k = \frac{1}{N_k} \sum_{i:z_i=k} x_i$$

$$\hat{\Sigma}_k = \frac{1}{N_k} \sum_{i:z_i=k} x_i x_i^T - \hat{\mu}_k \hat{\mu}_k^T \qquad \hat{\pi}_k = \frac{N_k}{N}$$

# 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

*[handwritten: estimate $z_i \leftarrow \theta$, max, "responsibility" of cluster $k$ for point $i$]*

- Example: MoG (derivation soon... + HW)
  1. Infer "responsibilities"

$$r_{ik} = p(z^i = k \mid x^i, \hat{\theta}^{(t-1)}) = \frac{\pi_k^{(t-1)} p(x_i \mid \phi_k^{(t-1)})}{\sum_j \pi_j^{(t-1)} p(x_i \mid \phi_j^{(t-1)})}$$

  2. Optimize parameters

max w.r.t. $\pi_k$ : $\quad \pi_k^{(t)} = \frac{1}{N} \sum r_{ik} = \frac{r_k}{N} \leftarrow$ *soft counts!*

max w.r.t. $\phi_k$ :

*[handwritten: $\mu_k^{(t)} = \frac{\sum r_{ik} x_i}{r_k} \leftarrow$ weighted mean$\qquad \Sigma_k^{(t)} = \frac{1}{r_k} \sum r_{ik} x_i x_i^T - \mu_k^{(t)} \mu_k^{(t)T}$]*

15

---

# Gaussian Mixture Example: Start

*[handwritten: Start with initial estimate of $\pi^{(0)}, \phi^{(0)}$]*

*[handwritten: → leads to initial "responsibilities"]*



*[handwritten labels near points: p=0.333, p=0.333, .333]*

16

# After first iteration

maximize
likelihood given
soft assignments

$\longrightarrow$ use new
$\pi^{(1)}, \phi^{(1)}$
to compute
new $r_{ix}$

p=0.06
p=0.272
=0.221

17

# After 2nd iteration

Iterate

p=0.3
p=0.306
=0.320

18

3

# After 3rd iteration

p=0.3
p=0.307

19



# After 4th iteration

p=0.331
p=0.286

20

4

# After 5th iteration

p=0.322

p=0.285

# After 6th iteration

p=0.315

p=0.287

# After 20th iteration

*Looks pretty good!*

23

# Some Bio Assay data

24

# GMM clustering of the assay data

25

# Resulting Density Estimator

*recall that GMMs can be used for density estimation*

26

7

# Expectation Maximization (EM) – Setup

- More broadly applicable than just to mixture models considered so far

- Model: $x$    observable – *"incomplete" data*    ← *what we actually have*

       $y$    not (fully) observable – *"complete" data*    ← *what we wish we had*

       $\theta$    parameters

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

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

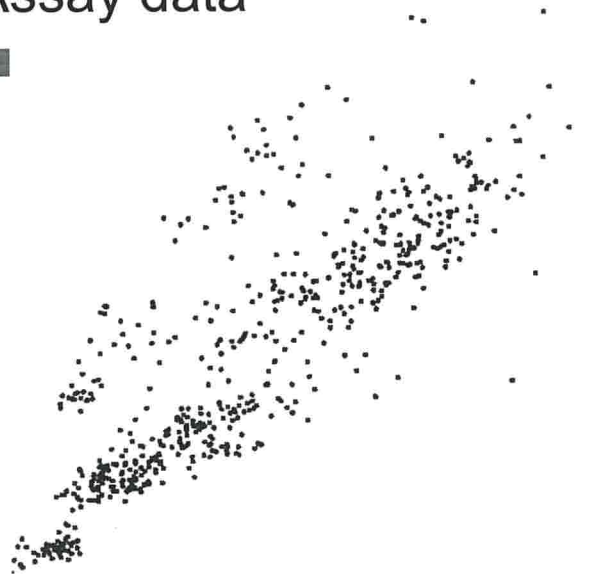- Special case:    *non-invertible determistic fun*

$$x = g(y)$$

*e.g.* $y = \begin{bmatrix} z \\ x \end{bmatrix}$ ← *class labels*    *in standard*
              ← *observations*    *mixture model*

27

---

# Expectation Maximization (EM) – Derivation

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

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

       ↑
       $x = g(y)$

$$\Rightarrow \log p(x \mid \theta) = \log p(y \mid \theta) - \log p(y \mid x, \theta)$$

       $\underbrace{\qquad}_{L_x(\theta)}$

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

$$L_x(\theta) = \underbrace{E[\log p(y \mid \theta) \mid x, \hat{\theta}]}_{U(\theta, \hat{\theta})} + \underbrace{E[-\log p(y \mid x, \theta) \mid x, \hat{\theta}]}_{V(\theta, \hat{\theta})}$$

28

8

# Expectation Maximization (EM) – Derivation

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

$$L_x(\theta) - L_x(\hat{\theta}) = \left[U(\theta,\hat{\theta}) - U(\hat{\theta},\hat{\theta})\right] + \left[V(\theta,\hat{\theta}) - V(\hat{\theta},\hat{\theta})\right]$$

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

  Proof: Use Jensen's Ineq $E[f(x)] \leq f[E[x]]$
  
  for any concave $f(\cdot)$

  Here:

  $E_p[\log q] - E_p[\log p] = E_p\left[\log \frac{q}{p}\right]$

  $\leq \log E_p\left[\frac{q}{p}\right] = \log \left(\int_x p(x) \frac{q(x)}{p(x)} dx = \log 1 = 0\right.$

29

---

# Expectation Maximization (EM) – Derivation

$$L_x(\theta) - L_x(\hat{\theta}) = [U(\theta,\hat{\theta}) - U(\hat{\theta},\hat{\theta})] - [V(\theta,\hat{\theta}) - V(\hat{\theta},\hat{\theta})]$$

$\geq 0$

- Step 4
  - Determine conditions under which log likelihood at $\theta$ exceeds that at $\hat{\theta}$

Using Gibbs inequality: $V(\theta,\hat{\theta}) = E[-\log p(y|x,\theta)|x,\hat{\theta}]$

$\geq E[-\log p(y|x,\hat{\theta})|x,\hat{\theta}]$

$= V(\hat{\theta},\hat{\theta})$ $\forall \theta$

If $U(\theta,\hat{\theta}) \geq U(\hat{\theta},\hat{\theta})$

Then

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

choosing a $\theta$ s.t. this is true means we're moving in the right direction (or not wrong!)

30

9

# Motivates EM Algorithm

- Initial guess: $\hat{\theta}^{(0)}$
- Estimate at iteration $t$: $\hat{\theta}^{(t)}$

- **E-Step**

  Compute $U(\theta, \hat{\theta}^{(t)}) = E[\log p(y|\theta) | x, \hat{\theta}^{(t)}]$

- **M-Step**

  Compute $\hat{\theta}^{(t+1)} = \arg\max_{\theta} U(\theta, \hat{\theta}^{(t)})$

From before, $U(\hat{\theta}^{(t+1)}, \hat{\theta}^{(t)}) \geq U(\hat{\theta}^{(t)}, \hat{\theta}^{(t+1)})$

$\Rightarrow L_x(\hat{\theta}^{(t+1)}) \geq L_x(\hat{\theta}^{(t)})$

©Emily Fox 2013    31

---

# 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.

$$p(x^i, z^i \mid \theta) = \pi_{z^i} p(x^i \mid \phi_{z^i}) = \prod_{k=1}^{K} (\pi_k p(x^i \mid \phi_k))^{I(z^i = k)}$$

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

$$U(\theta, \hat{\theta}^{(t)}) = \sum_i \sum_k E_{q_t}[I(z^i = k)] \log[\pi_k p(x^i \mid \phi_k)]$$

E-Step compute these

$$= \sum_i \sum_k p(z_i = k \mid x_i, \hat{\theta}^{(t)}) \log[\pi_k p(x^i \mid \phi_k)]$$

M-step maximize

$$= \sum_i \sum_k r_{ik} \log \pi_k + \sum_i \sum_k r_{ik} \log p(x_i \mid \phi_k)$$
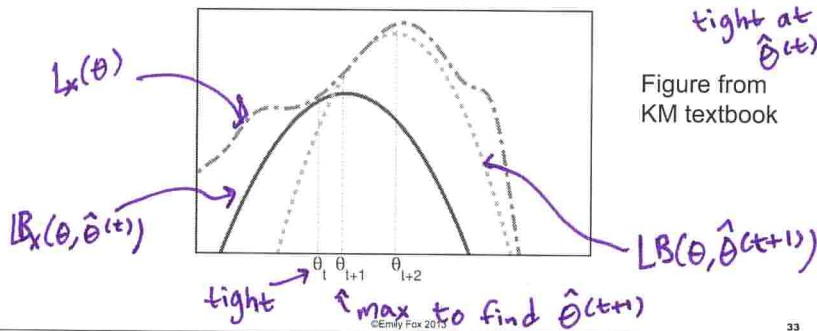
©Emily Fox 2013    32

1

# Coordinate Ascent Behavior

- Bound log likelihood:

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

*tight at $\hat{\theta}^{(t)}$*

$L_x(\theta)$

$LB_x(\theta, \hat{\theta}^{(t)})$

Figure from KM textbook

$\theta_t \quad \theta_{t+1} \quad \theta_{t+2}$

$LB(\theta, \hat{\theta}^{(t+1)})$

*tight*   ↑ *max to find $\hat{\theta}^{(t+1)}$*

©Emily Fox 2013

33

---

# Comments on EM

- Since Gibbs inequality is satisfied with equality only if *p=q*, any step that changes $\theta$ should strictly **increase likelihood**

  *assuming identifiability (i.e. $\nexists \theta \neq \theta'$ s.t. $p(x|\theta) = p(x|\theta')$)*

- In practice, can replace the **M-Step** with increasing *U* instead of maximizing it (**Generalized EM**)

  *exact max can be hard to compute*

- 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

  *like in mix model with $y = \{z, x\}$ "cluster assign."*

- If you want to choose any *y*, not necessarily *x=g(y)*, replace $p(y \mid \theta)$ in *U* with $p(y, x \mid \theta)$

©Emily Fox 2013

34

2

# 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

    *... many choices!*

- Can be quite important to convergence rates in practice

  *and quality of local optima found*

35

# MAP Estimation

- Bayesian approach:
  - ☐ Place **prior** $p(\theta)$ on parameters
  - ☐ Infer **posterior** $p(\theta \mid x)$ ≈ $\dfrac{p(x\mid\theta)\,p(\theta)}{p(x)}$

- Many, many, many motivations and implications
  - ☐ For the sake of this call, simplest motivation is to think of this as akin to regularization

$$\hat{\theta}^{MAP} = \arg\max_{\theta} \log p(\theta \mid x) = \arg\max_{\theta} \overbrace{\log p(x\mid\theta)}^{ML\ term} + \log p(\theta) \quad reg.$$

  - ☐ Saw importance of regularization in logistic regression (ML estimate can overfit data and lead to poor generalization)

36

3

# EM Algorithm – MAP Case

- Re-derive EM algorithm for $p(\theta \mid x)$

  *Prev* $E[\log p(y \mid \theta) \mid x, \hat{\theta}]$
  *Now* $E[\log p(y \mid \theta) \mid x, \hat{\theta}] + \log p(\theta)$

- Add $\log p(\theta)$ to $U(\theta, \hat{\theta}^{(t)})$
  - What must be computed in E-Step remains unchanged because this term does not depend on *y*.
  - M-Step becomes:

  $$\hat{\theta}^{(t+1)} = \arg \max_{\theta} U(\theta, \hat{\theta}^{(t)}) + \log p(\theta)$$

  *affects max w.r.t.* $\theta$

# MAP EM Example – MoG

$p(\theta \mid y)$ in same family as $p(\theta)$

- For mixture of Gaussians, conjugate priors are:

  $$\pi \sim \text{Dir}(\alpha_1, \ldots, \alpha_K) \quad \{\mu_k, \Sigma_k\} \sim \text{NIW}(m_0, \kappa_0, \nu_0, S_0)$$

- Results in following M-Step:

  $\hat{\mu}_k$ from before
  mean of pseudo-obs.
  pseudocounts of obs in cluster k

  $$\hat{\mu}_k = \frac{r_k \bar{x}_k + \kappa_0 m_0}{r_k + \kappa_0} \qquad \hat{\pi}_k = \frac{r_k + \alpha_k - 1}{N + \sum_k \alpha_k - K}$$

  $\hat{\Sigma}_k$ from before

  $$\hat{\Sigma}_k = \frac{S_0 + r_k S_k + \frac{\kappa_0 r_k}{\kappa_0 + r_k}(\bar{x}_k - m_0)(\bar{x}_k - m_0)'}{\nu_0 + r_k + d + 2}$$

  *dimension*

# What you need to know

- Mixture model formulation
  - Generative model
  - Likelihood
- Expectation Maximization (EM) Algorithm
  - Derivation
  - Concept of non-decreasing log likelihood
  - Application to standard mixture models

39

# Course Announcements

- Homework 2 will be posted on Thursday
  - Due 2 weeks later (Feb 14)

- Project proposals:
  - Initial ideas now posted
  - Deadline extended to Tues, Feb 5
  - 1 page, 1-2 people

- Recitation on Thursday (Linda)

- Office hours as normal

40