

# MIXTURE MODELS

ASSUMES OBSERVED VARIABLES ARE CORRELATED B/C ARISE FROM HIDDEN COMMON CAUSE. LATENT VARIABLE MODELS.

→ HARDER TO FIT

→ FEWER PARAMS

→ CAN LEAD TO BOTTLENECKS → COMPRESSED REPRESENTATION OF DATA. UNSUPERVISED LEARNING, THINK AUTOENCODERS

## MIXTURE MODELS

SIMPLEST FORM OF LVM. E.G.  $\{1 \dots K\}$  DISCRETE LATENT STATES. DISCRETE PRIOR  $P(z_i) = \text{CAT}(\pi)$  • LIKELIHOOD:  $P(x_i | z_i = u) = P_u(x_i)$

→ WE MIX TOGETHER B. DIST.  $P(x_i | \theta) = \sum \pi_u P_u(x_i | \theta)$ ; CONVEX COMBINATION OF  $P_u$ ,  $K$ th BASE DISTRIBUTION

$$0 \leq \pi_u < 1, \sum \pi_u = 1, \pi_u \text{ MIXING WEIGHTS}$$

DIFFERENT TYPES OF MIXTURE MODELS DEPEND ON TYPE OF  $P(z_i)$  AND  $P(x_i | z_i)$

## MIXTURE OF GAUSSIANS

• PRIOR = DISCRETE

• LIKELIHOOD = MVN

$$P(x_u | \theta) = \sum \pi_u \cdot N(x_i | \mu_u, \Sigma_u)$$

## MIXTURE OF MULTINOMIALS

• PRIOR = DISCRETE

• LIKELIHOOD = DISCRETE

$$P(x_u | \theta) = \prod_j \text{BER}(x_{ij} | \mu_{ju}) = \prod_j \mu_{ju}^{x_{ij}} (1 - \mu_{ju})^{1-x_{ij}} \quad \mu_{ju} = \text{PROB BIT } j \text{ IS ON IN CLUSTER } u$$

$$\Sigma[x] = \sum \pi_u \mu_u \quad \text{COV}[x] = \sum \pi_u [\Sigma_u + \mu_u \mu_u^T] - E[x] E[x]^T, \quad \Sigma_u = \text{DIAG}(\mu_{ju}(1 - \mu_{ju}))$$

## USING MIXTURE MODELS FOR CLUSTERING

MM USED MAINLY FOR TWO THINGS: BAYES BOX DENSITY MODELING  $P(x_i)$ , DATA COMPRESSION, OUTLIER DETECTION, MODEL CURS CONDITIONAL DENSITY CLASSIFICATION

CLUSTERING: FIRST FIT MODEL, COMPUTE  $P(z_i \neq u | x_i, \theta) \rightarrow$  POSTERIOR OF POINT  $i$  BELONGING TO CLUSTER  $u$

## RESPONSIBILITY

$$R_{ik} = P(z_i = k | x_i, \theta) = \frac{P(z_i = k | \theta) P(x_i | z_i = k, \theta)}{\sum_{u=1}^K P(z_i = u | \theta) P(x_i | z_i = u, \theta)} \rightarrow \text{SOFT CLUSTERING. LIKE GEN CLASSIFIERS BUT HERE WE NEVER OBSERVE } z_i$$

UNCERTAINTY  $1 - \max_u R_{ik}$ , IF SMALL  $\rightarrow z_i^* = \text{ARGMAX}_u R_{ik} = \text{ARGMAX}_u \cdot \log P(x_i | z_i = u, \theta) + \log P(z_i = u, \theta)$  HARD CLUSTERING

• WE CAN REPRESENT CLUSTERS USING PROTOTYPES/CENTROIDS  $\rightarrow \mu_u$  FOR THE  $K_s$  (MAPs)

## MIXTURE OF EXPERTS

DISCRIMINATING MODELS FOR CLASSIFICATION/REGRESSION. MIXTURE COMPONENT DEPEND ON INPUT VALUE, EACH IS AN EXPERT IN A PART OF THE INPUT SPACE

$$P(y_i | x_i, z_i = u, \theta) = N(y_i | w_u^T x_i, \sigma_u^2), \quad P(z_i | x_i, \theta) = \text{CAT}(z_i | S(V^T x_i)), \quad P(z_i = u | x_i, \theta) \text{ GATING FUNCTION}$$

$$\text{COMPLETE PREDICTION MODEL } P(y_i | x_i, \theta) = \sum P(z_i = u | x_i, \theta) P(y_i | x_i, z_i = u, \theta)$$

• ANY MODEL CAN BE EXPERT

• NN GATING, NN EXPERTS  $\rightarrow$  MIXTURE DENSITY NETWORK

• EACH EXPERT IS MOE ITSELF  $\rightarrow$  HIERARCHICAL MOE

• MORE USEFUL IN INVERSE PROBLEMS  $\rightarrow$  IE INVERSE KINEMATICS WITH A MANY-TO-ONE MAPPING

## PARAMETER ESTIMATION

HOW TO LEARN THE  $\theta$ s. • IF  $z_i$  WERE OBSERVED  $\rightarrow$  D-SEPARATION  $\rightarrow$  POSTERIOR FACTORIZES  $\rightarrow$  EASY BUT  $z_i$  HIDDEN!

### PROBLEM: PARAMETER UNIDENTIFIABILITY

- IF  $z_i$  WERE OBSERVED  $P(\theta|D) = \text{Dir}(\pi|D) = \prod_u \text{NIW}(\mu_u, \Sigma_u|D) \rightarrow$  FIND GLOBAL MAP/MLE
- $z_i$  HIDDEN  $\rightarrow$  WE GET DIFFERENT UNIMODAL LIKELIHOOD FOR EACH WAY OF FILLING IN  $z_i$ . WHEN WE MARGINALIZE OVER  $z_i$  WE GET A MULTIMODAL POSTERIOR  $P(\theta|D)$ 
  - $\rightarrow$  DIFFERENT LABELING OF CLUSTERS, NO UNIQUE MLE/MAP,  $K!$  POSSIBLE LABELINGS
  - $\rightarrow$  FINDING OPTIMAL MLE FOR A GMM IS NP-HARD!
- ISSUES FOR BAYESIAN INFERENCE: IT'S POINTLESS TO AVG TOGETHER (FOR APPROX. POSTERIOR MEAN) SAMPLES FROM A MULTIMODAL POSTERIOR.
  - $\rightarrow$  STILL OK TO AVG SAMPLES FROM POSTERIOR PREDICTIVE BECAUSE LIKELIHOOD IS INVARIANT TO MODE.
- HOW TO FIX
  - OPTIMAL: MCMC
  - QUICK FIX: COMPUTE SINGLE LOCAL MODE/MAP APPROX. POSTERIOR UNCERTAINTY ABOUT PARAMS LL POST. UNCERTAINTY ABOUT HIDDEN  $\rightarrow P(z_i|x_i, \hat{\theta})$ .

### PROBLEM: NONCONVEX MAP ESTIMATE

$$\log P(D|\theta) = \sum_i \log \left[ \sum_j P(x_i, z_i | \theta) \right] \text{ CAN'T PUSH LOG INSIDE SUM, BUT SINCE STUFF IS IN EXPONENTIAL FAMILY WE CAN DO}$$

$$\ell_c(\theta) = \sum_i \log P(x_i, z_i | \theta) = \theta^T \left( \sum \phi(x_i, z_i) \right) - N \log Z(\theta) \quad \text{COMPLETE DATA LOG LIKELIHOOD}$$

$$\ell(\theta) = \sum_i \log \left[ \sum_j e^{\theta^T \phi(z_i, x_i)} \right] - N \log Z(\theta) \quad \text{OBSERVED DATA LL (MISSING DATA)}$$

- THE TWO TERMS ARE CONVEX BUT NO GUARANTEES ON THE DIFFERENCE.
- SOLVE WITH HEURISTIC OPTIMIZATION ALGOS: SIMULATED ANNEALING, GENETIC ALGORITHMS, MULTIPLE RANDOM RESTARTS, EXPONENTIAL CONVEX METHOD BASED ON  $\ell_1$  CONVEX PENALTIES, UNSUPERVISED SPARSE KERNEL, LOGISTIC REGRESSION

# EXPECTATION - MAXIMIZATION ALGORITHM

WHEN MISSING DATA/LATENT VARIABLES ML/MAP IS HARD TO COMPUTE. GENERALLY WE CAN USE GRADIENT BASED OPTIMIZER AND FIND LOCAL MINIMUM OF  $NLL = -\frac{1}{N} \log p(D|\theta)$

WHEN THERE ARE CONSTRAINTS TO BE ENFORCED IT'S SIMPLY (NOT NEC. LY FASTER) TO USE EXPECTATION - MAXIMIZATION

• ITERATIVE, CLOSED-FORM UPDATES AT EACH STEP

• **E STEP** → INFERS MISSING DATA GIVEN PARAMS |  $\text{MAX LL OVER LATENT PARAMS}$

• **M STEP** → OPTIMIZES THE PARAMS GIVEN THE 'FILLED IN DATA' |  $\text{MAX OVER LL PARAMS}$  EM IS SPECIAL CASE OF BOUND OPTIMIZATION / MINIMIZE - MAXIMIZE ALGOS

## GENERAL

$$l(\theta) = \sum \log \left[ \sum p(x_i | z_i, \theta) \right] \quad \text{OBSERVED LL} \quad \text{CANT' PUSH LOG INSIDE SUM} \quad \left| \quad l_c(\theta) = \sum \log p(x_i, z_i | \theta) \quad \text{COMPLETE DATA L.L. NO CAN COMPUTE!} \right. \\ \left. Q(\theta, \theta^{t-1}) = E[l_c(\theta) | D, \theta^{t-1}] \quad \text{EXPECTED COMPLETE DATA LL} \right.$$

• **E STEP** COMPUTES TERMS INSIDE Q ON WHICH MLE DEPENDS ON → EXPECTED SUFFICIENT STATISTICS •  $Q(\theta | \theta^t) = E_{z|x} [\log L(\theta, x, z)]$  → AUXILIARY FCN

• **M STEP** OPTIMIZES Q WRT TO  $\theta \quad \theta^t = \underset{\theta}{\text{ARGMAX}} Q(\theta, \theta^{t-1}) \quad \theta^t = \underset{\theta}{\text{ARGMAX}} Q(\theta, \theta^{t-1}) + \log p(\theta)$  MAP MLE MAP

• EM MONOTONICALLY INCREASES LL OF OBSERVED DATA (OR STAYS SAME). IF OBJ GOES DOWN → SOMETHING IS WRONG!

## EM FOR GMM

• **AUXILIARY FCN:**  $Q(\theta, \theta^{t-1}) = E \left[ \sum \log p(x_i, z_i | \theta) \right] = \sum_u \sum_n R_{un} \log \pi_u + \sum_u \sum_n R_{un} \log p(x_i | \theta)$

$R_{un}$  = RESPONSIBILITY OF CLUSTER  $u$  FOR DATAPoint  $i$

• **E STEP:**  $R_{un} = \frac{\pi_u p(x_i | \theta_u^{t-1})}{\sum_u \pi_u p(x_i | \theta_u^{t-1})}$

• **M STEP:**  $\pi_u = \frac{R_u}{N} \rightarrow$  WEIGHTED NO OF POINTS  $\in u$ .  $l(\mu_u, \Sigma_u) = \sum_u \sum_n R_{un} \log p(x_i | \theta)$  SOME MATH LATER

$$\mu_u = \frac{\sum_i R_{un} x_i}{R_u}$$

$$\Sigma_u = \frac{\sum_i R_{un} x_i x_i^T}{R_u} - \mu_u \mu_u^T$$

• **FINAL:**  $\theta^t = (\pi_u, \mu_u, \Sigma_u)$  FOR  $u=1:N$ , REPEAT & REPEAT

• QUALITATIVELY MEAN OF CLUSTER  $u$  IS WEIGHTED AVG OF ALL POINTS  $\in u$ . COV IS PROPORTIONAL TO EMPIRICAL SCATTER MATRIX.

## K-MEANS

CAN BE SEEN AS VARIANT OF EM ON GMM WHERE  $\Sigma_u = \sigma^2 I$   $\pi_u = 1/K$  ARE FIXED SO ONLY  $\mu_u$  ARE TO BE ESTIMATED

• POSTERIOR APPROXIMATED TO  $p(z_i = u | x_i, \theta) \approx 1(K=u)$  → HARD EM EQUAL SPHERICAL COVARIANCE MATRIX →  $z_i = \underset{u}{\text{ARGMIN}} \|x_i - \mu_u\|_2^2$

•  $\mu_u = \frac{1}{N_u} \sum x_i$  •  $O(NKD)$  TIME • UNTIL CONVERGENCE

• INITIALIZATION: RANDOM OR FARTHEST POINT CLUSTERING / K-MEANS++ GUARANTEES DISTANCE  $\leq O(\log K)$

• CHOOSE CENTROIDS WITH PROBABILITY PROPORTIONAL TO DISTANCE OF ALREADY PICKED ONES

## VECTOR QUANTIZATION

CAN BE SEEN AS GREEDY ALGO FOR OPTIMIZING LOSS RELATED TO DATA COMPRESSION

• ENCODE VECTORS  $x_i$  (REAL VALUES) WITH A DISCRETE SYMBOL  $z_i = \{1..M\}$  INDEX INTO A CODEBOOK OF  $M$  PROTOTYPES  $\mu_k$ . EACH VECTOR ENCODED w/ NEAREST PROTOTYPE

$$\text{ENCODE}(x) = \underset{k}{\text{ARGMIN}} \|x_i - \mu_k\|_2^2$$

$$\text{RECONSTRUCTION ERROR: } d(\mu, z | K, X) = \frac{1}{N} \sum \|x_i - \text{DECODE}(\text{ENCODE}(x_i))\|^2$$

• TAKES LESS SPACE  $O(NDC) \rightarrow O(\log K)$



## MAP ESTIMATION

BECAUSE MLES MAY OVERFIT. COLLAPSING VARIANCE PROBLEM WHEN ONLY 1 POINT TO CLUSTER.

- Q BECOMES EXPECTED COMPLETE DATA LL + LOG PRIOR. E STEP STAYS THE SAME. M STEP CHANGES (PP. 356-359 MURPHY)
- MUCH MORE RESILIENT IN HIGH-DIMENSIONALITY SCENARIOS TO PROBLEMS DUE TO SINGULAR MATRICES
- HAS HYPERPARAMS

## EM FOR MIXTURE OF EXPERTS

$$Q(\theta, \theta^{old}) = \sum_i \sum_k R_{ik} \log [\pi_{ik} \cdot N(y_i | w_k^T x_i, \sigma_k^2)] \quad \pi_{ik} = S(V^T x_i)_k \quad R_{ik} = \pi_{ik}^{old} N(y_i | x_i^T w_k, (\sigma_k^{old})^2)$$

- E SAME AS GMM BUT  $\pi_{ik} \rightarrow \pi_{ik}$
- M MAXIMIZES  $Q(\theta, \theta^{old})$  WRT  $w_k, \sigma_k^2, V$ 
  - $-w_k = (x^T R_{ik} x)^{-1} x^T R_{ik} y$  IS ANALOGOUS TO WEIGHTED LEAST SQUARES
  - $-\sigma_k^2 = \frac{\sum_i R_{ik} (y_i - w_k^T x_i)^2}{\sum_i R_{ik}}$
  - $-l(V) = \sum_i \sum_k R_{ik} \log \pi_{ik}$  IS LOG-LIKELIHOOD FOR MULTINOMIAL LOGISTIC BUT WITH SOFT LABELS  $\rightarrow$  FITTING A LOG. REG. MODEL

## EM FOR DGM W/ HIDDEN VARS

FOR ARBITRARY GENERIC DGM.  $E[\log P(D|\theta)] = \sum_{t \in \mathcal{T}} \sum_{c \in \mathcal{C}} \sum_k \bar{N}_{tcu} \log \theta_{tcu}$   $\bar{N}_{tcu}$  = EMPIRICAL COUNTS =  $\sum_i P(x_{it}=k, x_{1i}, p_{1i}(1)=c | D)$

- FAMILY MARGINAL:  $P(x_{it}, x_1, p_{1i}(1) | D, \theta)$
- E PRODUCES  $\bar{N}_{tcu}$  EXPECTED SUFF STATS
- M  $\hat{\theta}_{tcu} = \frac{\bar{N}_{tcu}}{\sum_k \bar{N}_{tcu}}$

## EM FOR STUDENT DISTRIBUTION

- GAUSSIAN MM ARE SENSITIVE TO OUTLIERS; STUDENT'S T IS MORE ROBUST BUT NO CLOSED FORM MLE! SO WE MUST OPTIMIZE ITERATIVELY  $\rightarrow$  EM
- INTRODUCE 'ARTIFICIAL' HIDDEN VAR, WRITE STUDENT'S AS INFINITE MIXTURE OF GAUSSIANS W/ DIFF  $\Sigma$

$$\hookrightarrow \text{GAUSSIAN SCALE MIXTURE } p(x_i | \mu, \Sigma, \nu) = \int N(x_i | \mu, \Sigma/z_i) \cdot \text{GA}(z_i | \frac{\nu}{2}, \frac{\nu}{2}) dz_i \quad \text{TREAT } z_i \text{ AS MISSING DATA}$$

$$l_c(\theta) = L_N(\mu, \Sigma) + L_G(\nu) \quad \text{[PP. 359-360 MURPHY FOR MATH]}$$

$\left\{ \begin{array}{l} \text{V KNOWN} \rightarrow \text{IGNORE } L_G, \text{ ONLY DO } E[z_i] \text{ WRT OLD PARAMS} \\ \text{V UNKNOWN} \rightarrow \text{VERY COMPLEX EXPRESSION REQUIRING GRADIENT-BASED OPTIMIZATION (NO CLOSED-FORM UPDATE)} \end{array} \right.$

GENERALIZED EM

- MIXTURE OF STUDENTS IS A THING, YOI!

## EM FOR PROBIT

$$P(y_i=1|z_i) = 1(z_i > 0), \quad z_i \sim N(w^T x_i, 1) \text{ IS UNBIAS.} \quad l(z_i, w | \nu_0) = \log p(y_i|z_i) + \log N(z_i | w, 1) + \log N(w | 0, \nu_0)$$

$$-E \text{ POSTERIOR } P(z_i | y_i, x_i, w) = \begin{cases} N(z_i | w^T x_i, 1) & y_i=1 \\ N(z_i | w^T x_i, 1) & y_i=0 \end{cases} \quad \text{TRUNCATED GAUSSIAN W ONLY LIN. DEP ON } z_i$$

$$E[z_i | w, x_i] = \begin{cases} \mu_1 + \frac{\phi(\mu_1)}{\Phi(\mu_1)} & y_i=1 \\ \mu_1 - \frac{\phi(\mu_1)}{\Phi(\mu_1)} & y_i=0 \end{cases} \quad \mu_1 = w^T x_i$$

$$-M \text{ ESTIMATES } w \text{ WITH RIDGE REGRESSION } \mu = E[z] \text{ IS PREDICTION } \hat{w} = (V_0^{-1} + X^T X)^{-1} X^T \mu$$

- SLOWER THAN DIRECT GRADIENT B/C POSTERIOR ENTROPY IS HIGH. USE STRONGER REGULARIZER TO CONSTRAIN  $z$  VALUES  $\rightarrow$  SPEEDUP CONVERGENCE

# E-M THEORY

'EXPECTED COMPLETE DATA LL IS LOWER BOUND'

- $Q(z_i)$  ARBITRARY DISTRIBUTION OVER HIDDEN VARIABLES
- OBSERVED DATA LL LOWER BOUND  $Q(\theta, q) = \sum_i E_{q_i} [\log p(x_i, z_i | \theta)] + H(q_i)$   
 $\rightarrow$  PICK  $Q$  W/ TIGHTER LOWER BOUND ENTROPY
- BUT ALSO  $-KL(Q(z_i) || P(z_i | x_i, \theta)) + \log p(x_i | \theta)$
- $Q(\theta, q^t) = \sum E_{q^t} [\log p(x_i, z_i | \theta)] + H(q^t)$
- M-STEP:  $\theta^{t+1} = \text{ARGMAX} \sum E_{q^t} [\log p(x_i, z_i | \theta)]$
- E-STEP: IT MONOTONICALLY INCREASES OBS DATA LL UNTIL LOCAL MAX
- K-L = 0 ~~PROVIDE~~ FOR SOME REASON  $\rightarrow L(\theta^t, q_i) = \log p(x_i | \theta)$
- $Q(\theta^t, \theta^t) = \sum \log p(x_i | \theta^t) = \ell(\theta^t)$  \*
- LOWER BOUND TOUCHES FCN AT EVERY STEP. MAXIMIZING LB 'PUSHES UP' FCN  $\rightarrow$  M STEP GUARANTEES TO INCREASE DATA LIKELIHOOD

## EM MONOTONICALLY INCREASES OBSERVED DATA LL

- $\ell(\theta^{t+1}) \geq Q(\theta^{t+1}, \theta^t) \geq Q(\theta^t, \theta^t) = \ell(\theta^t)$  ;  $\text{II}$  IS LB ON  $\text{I}$ ;  $\text{II} = \max_{\theta} Q(\theta, \theta^t) \geq \text{III}$ ; \*
- I II III IV
- IS SAME AS VARIATIONAL APPROACH, BUT  $Q$  HERE IS TRACTABLE

## ONLINE EM

FOR STREAMING DATASETS

- BATCH EM: WE COMPUTE SECTOR OF SUFFICIENT STATISTICS, EXPECTATIONS, AND THEIR SUMS S Q M
- INCREMENTAL EM: WE KEEP TRACK OF SUFF. STATS AND THEIR SUMS AFTER EACH DATAPOINT.  $S_i$  ALL COMPUTED INITIALLY AND  $M$  INITIALIZED AS THE SUM THEN WE UPDATE  $S_{NEW}$
- STEPWISE EM: WHENEVER  $S_i$  IS UPDATE, MOVE  $M$  TOWARDS IT
- SPEED = STEPWISE > INCREMENTAL >> BATCH
- ACCURACY = STEPWISE => BATCH > INCREMENTAL

## OTHER EM VARIANTS

- ANNEALED EM: DETERMINISTIC ANNEALING TO INCREASE CHANCE OF GLOBAL MAXIMUM. TEMPERATURE ACTS ON POSTERIOR.
- VARIATIONAL EM: WHEN CANNOT EXACT INFERENCE FOR E STEP. APPROXIMATION STILL ENSURES LOWER BOUND TO LIKELIHOOD.
- MONTECARLO EM: \* \* \* \* \* SAMPLE FROM POSTERIOR  $z_i \sim P(z_i | x_i, \theta^t)$  AND THEN COMPUTE SUFFICIENT STATS  
MCMC FOR SAMPLING OR FOR FULL INFERENCE
- GENERALIZED EM: E STEP OK BUT M STEP NOT OK. PARTIAL M-STEP
- ECM(E): EXPECTATION - CONDITIONAL MAXIMIZATION. IF  $M$  PARAMS ARE DEPENDENT  $\rightarrow$  OPTIMIZE THEM SEQUENTIALLY
- OVER-RELAXED EM:  $\theta^{t+1} = \theta^t + \eta (M(\theta^t) - \theta^t)$ . CAN LEAD TO FASTER CONVERGENCE, <sup>WE LOTS OF</sup> MISSING DATA

## MODEL SELECTION FOR EM

HOW TO PICK NUMBER OF LATENT VARIABLES / NUMBER OF CLUSTERS

### — PROBABILISTIC MODELS

OPTIMAL  $K^* = \text{ARGMAX}_K P(D|K)$ : LARGEST MARGINAL LIKELIHOOD

- TOUGH FOR LVM  $\rightarrow$  BIC APPROXIMATION, X-VALIDATED LIKELIHOOD,
- LARGE MODEL SPACE  $\rightarrow$  BAYES PAZOZ, STOCHASTIC SAMPLING IE MCMC

### — NON-PROBABILISTIC MODELS

- NO LIKELIHOOD, IE NUMBER OF  $K$  IN  $K$ -MEANS  $\rightarrow$  RELY ON RECONSTRUCTION ERROR  $E(D, K) = \frac{1}{|D|} \sum \|x_i - \hat{x}_i\|^2$
- ERROR ON TEST SET DECREASES ALWAYS BECAUSE MORE MODELS  $\rightarrow$  LARGER CHANCE OF CLOSE PROTOTYPE
- TRY TO IDENTIFY  $WEE / MIN$  ON TRAINING DATA VS  $K$  PLOT, BECAUSE  $L_{K^*}$  DECREASES GREATLY,  $\therefore K^*$  NOT SO MUCH

## FITTING MODELS w/ MISSING DATA

WE WANT JOINT DENSITY MODEL BY MLE BUT WE HAVE MISSING DATA  $O_{ij} = 1$  IF THERE, 0 IF MISSING

$$X_V = \{x_{ij} : O_{ij} = 1\} \quad X_H = \{x_{ij} : O_{ij} = 0\} \quad \xrightarrow{\text{WE WANT}} \quad \theta = \text{ARGMAX}_{\theta} P(X_V | \theta, O)$$

MAR ASSUMPTIONS:  $P(X_V | \theta, O) = \prod_{i=1}^N P(x_{iV} | \theta)$

$$\log P(X_V | \theta) = \sum \log p(x_{iV} | \theta)$$

$$\log p(x_{iV} | \theta) = \sum p(x_{iV}, x_{iH} | \theta)$$

$\downarrow \rightarrow$  USUAL PROBLEM WITH NO CAN PUSH LOG INSIDE SUM  $\rightarrow$  EM

(PP 373-374 MURPHY)

### — FOR A MVN

- MLE FOR FULLY OBSERVED ROWS • E-STEP:  $Q(\theta | \theta^{t-1})$  • M STEP:  $\mu^t = \frac{1}{N} \sum [x_i^t]$ ,  $\Sigma^t = \frac{1}{N} \sum E[x_i x_i^T] - \mu^t (\mu^t)^T$
- NOT EQUIVALENT TO JUST AVERAGE VARS WITH EXPECTATIONS AND DO MLE.  $\rightarrow$  IGNORES POSTERIOR VARIANCE
- WE COMPUTE EXPECTATIONS OF SUFF. STATS AND USE THOSE FOR MLE
- CAN ALSO USE FOR MAP WITH EES INTO MAP EQUATIONS