## Artificial Intelligence

A Course About Foundations



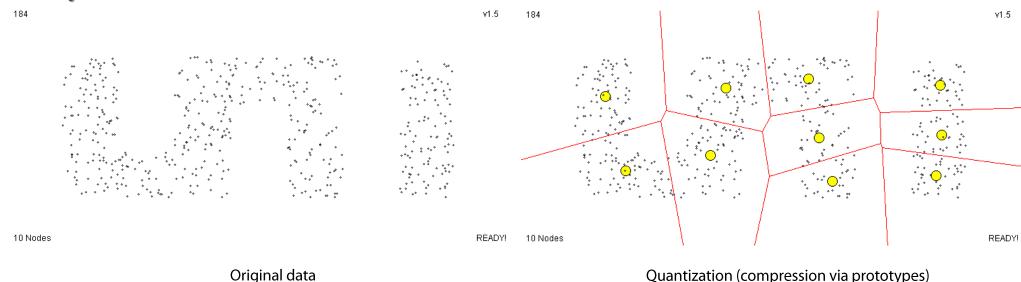
## Probabilistic Reasoning: Unsupervised Learning

Marco Piastra

Artificial Intelligence 2024–2025 Unsupervised Learning [1]

An aside: The K-means algorithm (alternate optimization)

#### Vector quantization



The basic idea is to replace each real-valued vector  $\mathbf{x} \in \mathbb{R}^d$  with a value  $\mathbf{w}_j \in \mathbb{R}^d$  which belongs to a finite <u>codebook</u> of k <u>prototypes</u>  $\theta := \{\mathbf{w}_1, \dots, \mathbf{w}_k\}$ 

Each data vector is encoded by using the index of the most similar prototype, where similarity is measured in terms, for instance, of Euclidean distance:

$$w(\boldsymbol{x}) := \operatorname{argmin}_{\boldsymbol{w}_j} \|\boldsymbol{x} - \boldsymbol{w}_j\|$$

For instance, part of mpeg4 and QuickTime (Apple) video compression algorithms work in this way and so does the Ogg Vorbis audio compression algorithm

Given a set  $D:=\{m{x}_1,\ldots,m{x}_N\}$  of observations (i.e. vectors in  $\mathbb{R}^d$ )

Clustering problem: given k, find a set of k prototypes  $\theta := \{w_1, \dots, w_k\}$  and an <u>assignment function</u>  $w : D \to \theta$  such that the objective (loss) function:

$$J(D, \theta) := \frac{1}{2} \sum_{i=1}^{N} \| \boldsymbol{x}_i - w(\boldsymbol{x}_i) \|^2$$

is minimized.

Artificial Intelligence 2024-2025

#### *k-means* algorithm:

- 1) Position the k prototypes at random
- 2) Assign each observation to its closest prototype

$$w(\boldsymbol{x}_i) := \operatorname{argmin}_{\boldsymbol{w}_i} \|\boldsymbol{x}_i - \boldsymbol{w}_j\|$$

3) Position each prototype at the *centroid* of the observations assigned to it

$$m{w}_j = rac{1}{|D(m{w}_j)|} \sum_{D(m{w}_j)} m{x}_i \qquad ext{where} \ \ D(m{w}_j) := \{m{x}_i \in D \mid w(m{x}_i) = m{w}_j\}$$

4) Unless no prototype was moved in step 3), go back to step 2)

This algorithm converges to a <u>local</u> minimum of  $J(D,\theta)$ 

Why does the algorithm work: *alternate optimization (also 'coordinate descent')* 

Step 2): Assign observations while keeping the k prototype fixed

$$w(m{x}_i) := \mathrm{argmin}_{m{w}_j} \|m{x}_i - m{w}_j\|$$
 which minimizes each of the terms in  $J(D, heta) := rac{1}{2} \sum_{i=1}^N \|m{x}_i - w(m{x}_i)\|^2$ 

Step 3): Reposition the k prototypes while keeping the assignments fixed

$$J(D, \theta) := \frac{1}{2} \sum_{i=1}^{N} \|\boldsymbol{x}_i - \boldsymbol{w}(\boldsymbol{x}_i)\|^2 = \frac{1}{2} \sum_{j} \sum_{D(\boldsymbol{w}_j)} (\boldsymbol{x}_i - \boldsymbol{w}_j)^2$$

$$\frac{\partial}{\partial \boldsymbol{w}_j} J(D, \theta) = \frac{\partial}{\partial \boldsymbol{w}_j} \frac{1}{2} \sum_{D(\boldsymbol{w}_j)} (\boldsymbol{x}_i - \boldsymbol{w}_j)^2 = \frac{\partial}{\partial \boldsymbol{w}_j} \frac{1}{2} \sum_{D(\boldsymbol{w}_j)} (\boldsymbol{x}_i - \boldsymbol{w}_j)^T (\mathbf{x}_i - \boldsymbol{w}_j)$$

$$= \frac{\partial}{\partial \boldsymbol{w}_j} \frac{1}{2} \sum_{D(\boldsymbol{w}_j)} (\boldsymbol{x}_i^2 + \boldsymbol{w}_j^2 - 2 \boldsymbol{x}_i^T \boldsymbol{w}_j) = \sum_{D(\boldsymbol{w}_j)} (\boldsymbol{w}_j - \boldsymbol{x}_i)$$

then, by imposing  $\frac{\partial}{\partial \boldsymbol{w}_i} J(D, \theta) = 0$  we obtain

$$oldsymbol{w}_j = rac{1}{|D(oldsymbol{w}_j)|} \sum_{D(oldsymbol{w}_j)} oldsymbol{x}_i$$

#### Discussion of the *k-means algorithm*

- a) At each step of the algorithm  $J(D,\theta)$  could not increase: it could only decrease or stay equal
- b) The algorithm is a variant of a *gradient descent*, in which at each step the *gradient descent* is performed on <u>one subset of variables</u> only
- c) It must reach a *fixed point*, where both gradients vanish
- d) But the only guarantee is that the algorithm reaches a local minimum (unless it gets stuck in a saddle point or a plateau)

Artificial Intelligence 2024–2025 Unsupervised Learning [7]

Given a set  $D:=\{oldsymbol{x}_1,\ldots,oldsymbol{x}_N\}$  of observations (i.e. vectors in  $\mathbb{R}^d$ ) and a set  $\theta := \{ \boldsymbol{w}_1, \dots, \boldsymbol{w}_k \}$  of k prototypes (i.e. vectors in  $\mathbb{R}^d$ )

Position each prototype at the *centroid* of the observations in its Voronoi cell

#### Voronoi cell:

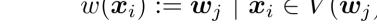
$$V(oldsymbol{w}_j) := \{oldsymbol{x} \in \mathbb{R}^d \mid \|oldsymbol{x} - oldsymbol{w}_j\| \leq \|oldsymbol{x} - oldsymbol{w}_l\| \ , orall \ l 
eq j \}$$

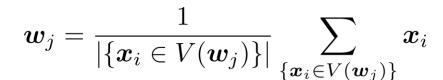
**Voronoi tesselation**: the complex of all Voronoi cells of  $\theta$ 

#### **Algorithm** (rewritten):

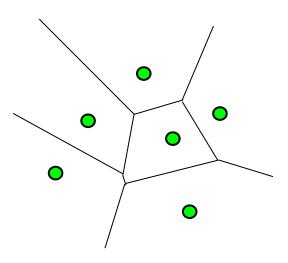
- 1) Position the k prototypes at random
- Assign each observation to its Voronoi cell

$$w(\boldsymbol{x}_i) := \boldsymbol{w}_j \mid \boldsymbol{x}_i \in V(\boldsymbol{w}_j)$$





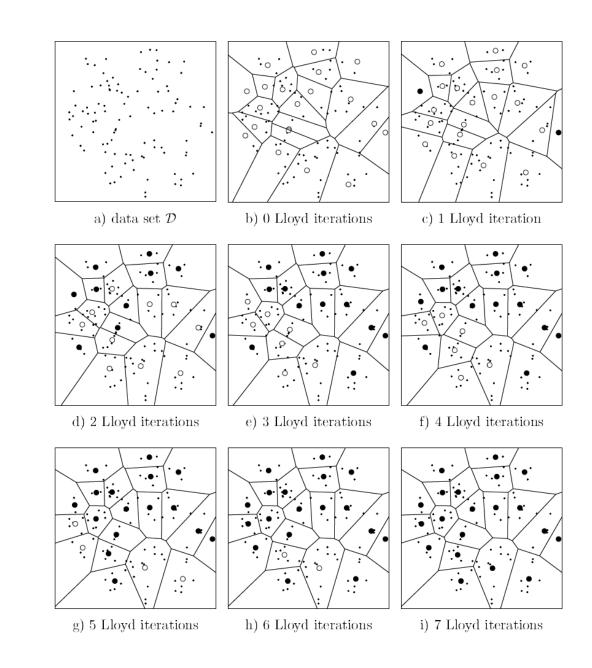
Unless no prototype was moved in step 3), go back to step 2)



#### k-means

*An example run of the algorithm* 

The landmarks (empty circles) become black when they cease to move



Artificial Intelligence 2024–2025 Unsupervised Learning [9]

## The Expectation–Maximization (EM) algorithm

### Expected value of a random variable

(also expectation)

**Basic definition** 

$$\mathbb{E}_X[X] := \sum_{x \in \mathcal{X}} x \ P(X = x)$$

A linear operator

$$\mathbb{E}[X+Y] = \mathbb{E}[X] + \mathbb{E}[Y]$$
$$\mathbb{E}[cX] = c\mathbb{E}[X]$$

More concise notation

$$\mathbb{E}[X] := \sum_{x \in \mathcal{X}} x \ P(x)$$

Continuous case

$$\mathbb{E}[X] := \int_{x \in \mathcal{X}} x \ p(x) dx$$

Conditional expectation

$$\mathbb{E}_X[X|Y=y] = \mathbb{E}[X|Y=y] := \sum_{x \in \mathcal{X}} x \ P(X=x|Y=y)$$

Iterated expectation (see Wikipedia)

$$\mathbb{E}_X[X] = \mathbb{E}_Y[\mathbb{E}_X[X|Y]]$$

#### Joint Expected Value

The **expected value** of a function f of a <u>set</u> of random variables  $\{X_i\}$  is

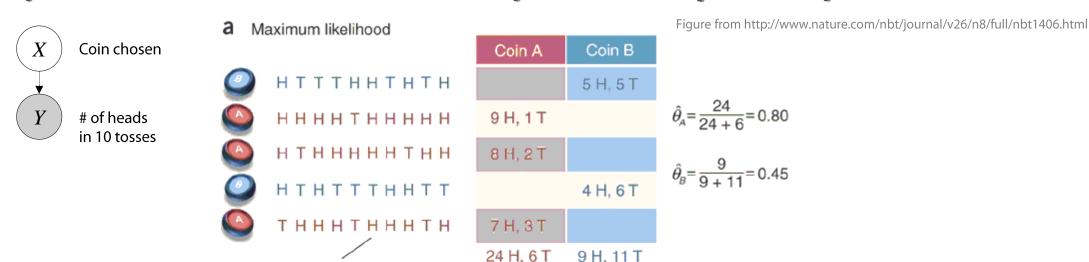
$$\mathbb{E}[f(\{X_i\})] := \sum_{\{X_i\}} f(\{X_i\}) \ P(\{X_i\})$$

the sum is over all possible combinations of values of the random variables

(Unless specified otherwise, the  $\mathbb{E}$  operator acts over *all* the random variables enclosed)

The extension to the continuous case is obvious

Artificial Intelligence 2024–2025 Unsupervised Learning [12]



An experiment with two coins

At each step, one coin is selected at random (with equal probability) and then tossed ten times

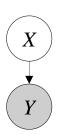
Random variables: Y number of heads, X selected coin (i.e A or B)

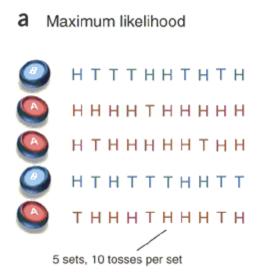
Parameters to be learnt:  $\theta = \{\theta_A, \; \theta_B\}$  probabilities of landing on head of A and B

When the results are fully observable, by MLE:

5 sets, 10 tosses per set

$$\theta_A^* = \frac{N_{Y=1,X=A}}{N_{X=A}}$$
  $\theta_B^* = \frac{N_{Y=1,X=B}}{N_{X=B}}$ 





Coin A	Coin B
	5 H, 5 T
9 H, 1 T	
8 H, 2 T	
	4 H, 6 T
7 H, 3 T	
24 H, 6 T	9 H, 11 T

Figure from http://www.nature.com/nbt/journal/v26/n8/full/nbt1406.html

 $\hat{\theta}_{A} = \frac{24}{24+6} = 0.80$ 

 $\hat{\theta}_{B} = \frac{9}{9 + 11} = 0.45$ 

An experiment with two coins

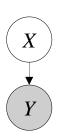
At each step, one coin is selected at random (with equal probability) and then tossed ten times

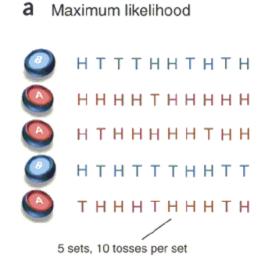
Random variables: Y number of heads, X selected coin (i.e A or B)

Parameters to be learnt:  $\theta = \{\theta_A, \theta_B\}$  probabilities of landing on head of A and B

■ What if *X* is hidden (= latent, = unobserved)?

The results of each sequence of coin tosses are known, but not the coin selected





Coin A	Coin B
	5 H, 5 T
9 H, 1 T	
8 H, 2 T	
	4 H, 6 T
7 H, 3 T	
24 H, 6 T	9 H, 11 T

Figure from http://www.nature.com/nbt/journal/v26/n8/full/nbt1406.html

$$\hat{\theta}_{A} = \frac{24}{24+6} = 0.80$$

$$\hat{\theta}_B = \frac{9}{9+11} = 0.45$$

• What if X is hidden (= latent, = unobserved)?

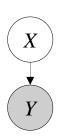
Likelihood

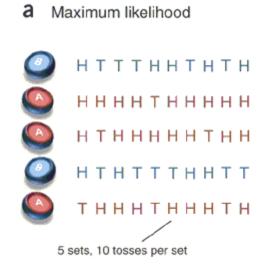
$$P(D \mid \theta) = P(\{Y^{(i)}\} \mid \theta) = \sum_{\{X^{(i)}\}} P(\{(Y^{(i)}, X^{(i))})\} \mid \theta)$$

MLE

$$\theta^* := \operatorname{argmax}_{\theta} \sum_{\{X^{(i)}\}} P(\{(Y^{(i)}, X^{(i)})\} \mid \theta)$$

\* This optimization problem is intractable, in general





Coin A	Coin B
	5 H, 5 T
9 H, 1 T	
8 H, 2 T	
	4 H, 6 T
7 H, 3 T	
24 H, 6 T	9 H, 11 T

Figure from http://www.nature.com/nbt/journal/v26/n8/full/nbt1406.html

- $\hat{\theta}_A = \frac{24}{24+6} = 0.80$
- $\hat{\theta}_B = \frac{9}{9+11} = 0.45$

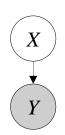
- What if X is hidden (= latent, = unobserved)?
  Intuitive idea: use expected values for unobserved variables
  - 1. Define an initial (random) guess  $\hat{\theta}^{(0)}$
  - 2. Compute  $Q_i(X^{(i)}) := P(X^{(i)} | Y^{(i)}; \hat{\theta}^{(t)})$

E-step

3. Maximize

$$\hat{\theta}^{(t+1)} = \operatorname{argmax}_{\theta} \ \sum \mathbb{E}_{Q_i(X^{(i)})}[Y^{(i)} \mid X^{(i)}; \hat{\theta}^{(t)}] \quad \textit{M-step}$$

4. Unless some convergence criterion has been met, go to step 2.



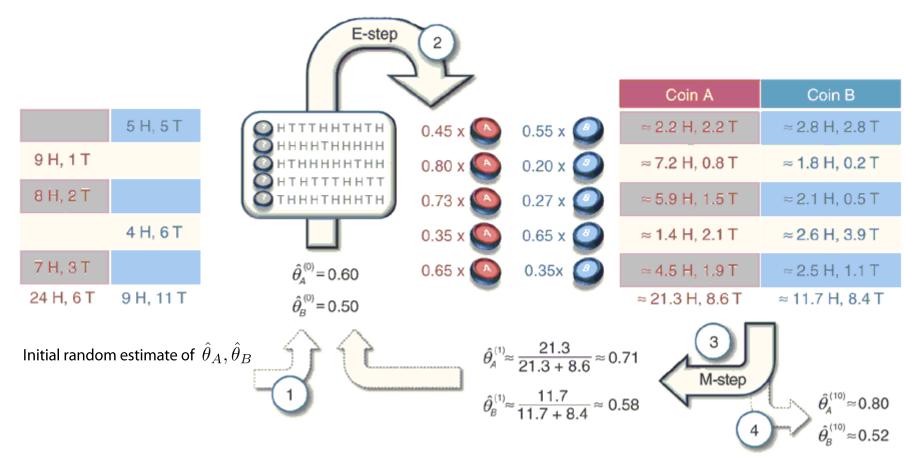
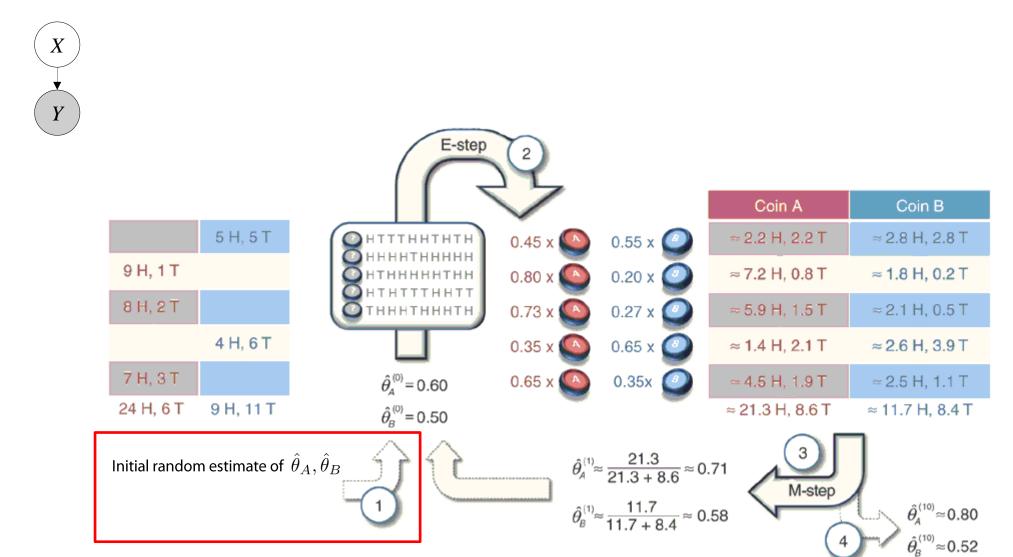
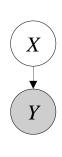


Figure from http://www.nature.com/nbt/journal/v26/n8/full/nbt1406.html

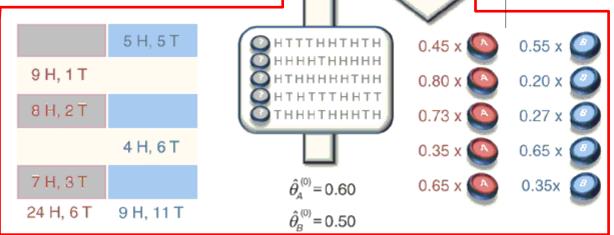


E-step



$$Q_i(X^{(i)}) := P(X^{(i)} \mid Y^{(i)}; \hat{\theta}^{(t)})$$

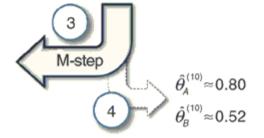
Compute the probability distribution of hidden observations probability of having used coin X for sequence igiven the result Y and current parameter estimate (see next)



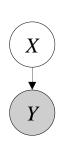
Coin A	Coin B
≈ 2.2 H, 2.2 T	≈ 2.8 H, 2.8 T
≈ 7.2 H, 0.8 T	≈ 1.8 H, 0.2 T
≈ 5.9 H, 1.5 T	≈ 2.1 H, 0.5 T
≈ 1.4 H, 2.1 T	≈ 2.6 H, 3.9 T
≈ 4.5 H, 1.9 T	≈ 2.5 H, 1.1 T
≈ 21.3 H, 8.6 T	≈ 11.7 H, 8.4 T

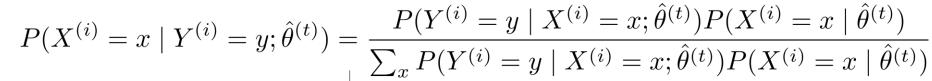
Initial random estimate of  $\hat{\theta}_A$ ,  $\hat{\theta}_B$ 

 $\hat{\theta}_{A}^{(1)} \approx \frac{21.3}{21.3 + 8.6} \approx 0.71$   $\hat{\theta}_{B}^{(1)} \approx \frac{11.7}{11.7 + 8.4} \approx 0.58$ 



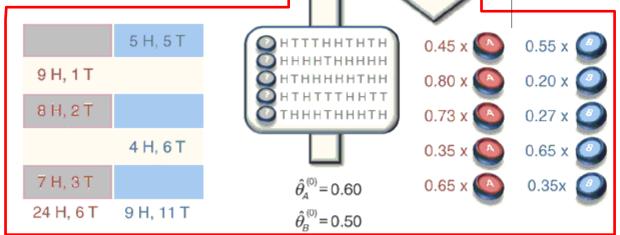
E-step





Compute the probability distribution of hidden observations

where:  $P(Y^{(i)} = y \mid X^{(i)} = x, \theta) = {10 \choose y} \theta_x^y (1 - \theta_x)^{10 - y}$  $P(X^{(i)} = x \mid \theta) = 0.5$ 

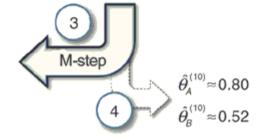


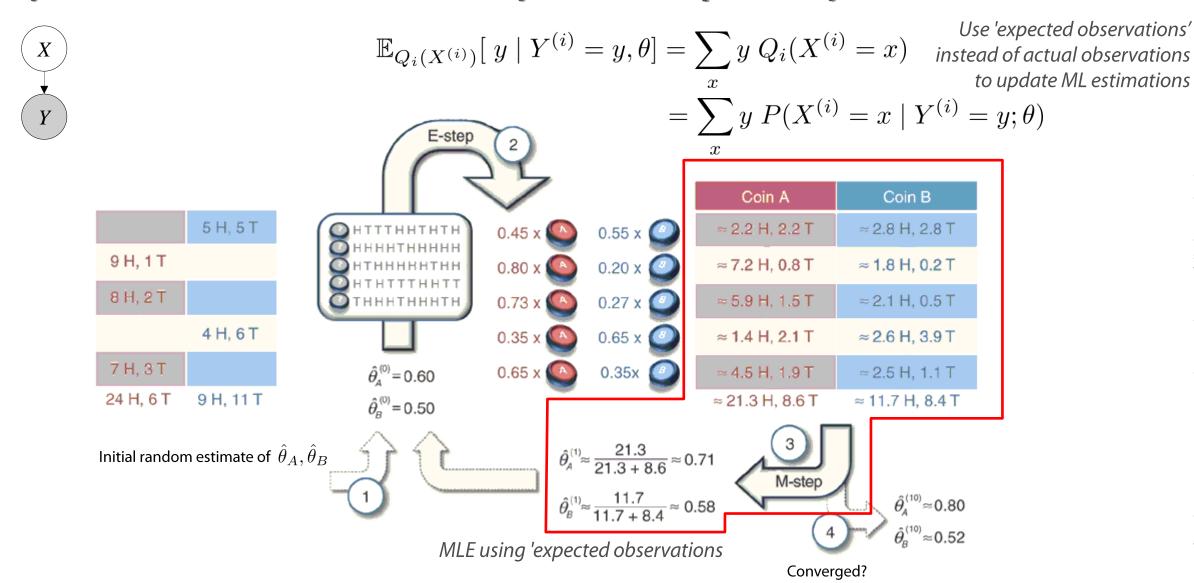
Coin A	Coin B
≈ 2.2 H, 2.2 T	≈ 2.8 H, 2.8 T
≈ 7.2 H, 0.8 T	≈ 1.8 H, 0.2 T
≈ 5.9 H, 1.5 T	≈ 2.1 H, 0.5 T
≈ 1.4 H, 2.1 T	≈ 2.6 H, 3.9 T
≈ 4.5 H, 1.9 T	≈ 2.5 H, 1.1 T
≈ 21.3 H, 8.6 T	≈ 11.7 H, 8.4 T

Initial random estimate of  $\ \hat{\theta}_A, \hat{\theta}_B$ 

$$\hat{\theta}_{A}^{(1)} \approx \frac{21.3}{21.3 + 8.6} \approx 0.71$$

$$\hat{\theta}_{B}^{(1)} \approx \frac{11.7}{11.7 + 8.4} \approx 0.58$$





-igure from http://www.nature.com/nbt/journal/v26/n8/full/nbt1406.html

# Convergence of the EM algorithm\* (in the discrete case)

#### An aside: Jensen's inequality

A relationship between probability and geometry

When f is convex function

$$f(\mathbb{E}[\{X_i\}]) \leq \mathbb{E}[f(\{X_i\})]$$

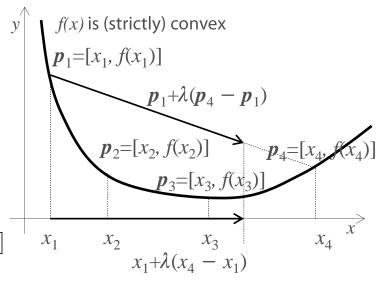
f is **convex** when for any two points  $p_i$  and  $p_j$  the segment  $(p_i - p_j)$  is not below f

That is, when

$$\lambda f(x_i) + (1 - \lambda)f(x_j) \ge f(\lambda x_i + (1 - \lambda)x_j), \, \forall \lambda \in [0, 1]$$

Furthermore, f is **strictly convex** when

$$\lambda f(x_i) + (1 - \lambda)f(x_j) \ge f(\lambda x_i + (1 - \lambda)x_j), \, \forall \lambda \in (0, 1)$$



#### Corollary:

when f is *strictly* convex, if and only if all the variables in  $\{X_i\}$  are <u>constant</u> it is true that

$$f(\mathbb{E}[\{X_i\}]) \leq \mathbb{E}[f(\{X_i\})]$$

Dual results also hold for *concave* functions

#### An aside: Jensen's inequality

A relationship between probability and geometry

When f is convex function

$$f(\mathbb{E}[\{X_i\}]) \leq \mathbb{E}[f(\{X_i\})]$$

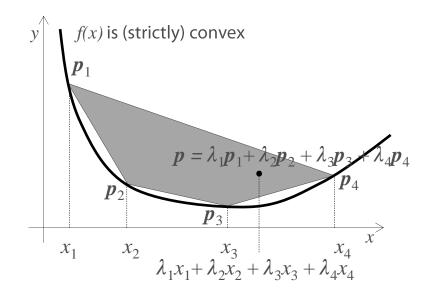
To see this, consider

$$\boldsymbol{p} = \lambda_1 \boldsymbol{p}_1 + \lambda_2 \boldsymbol{p}_2 + \lambda_3 \boldsymbol{p}_3 + \lambda_4 \boldsymbol{p}_4$$

i.e. a *linear combination* of  $p_i$  points

This is an *affine* combination if and it is a *convex* combination if also

$$\sum \lambda_i = 1$$
$$\lambda_i \ge 0, \ \forall i$$



When the  $\lambda_i$  define a probability, then p is a convex combination of  $p_i$  points

Any convex combination of  $p_i$  points lies inside their **convex hull** (see figure) and therefore above f:

$$\sum_{i} \lambda_{i} f(x_{i}) \geq f(\sum_{i} \lambda_{i} x_{i})$$

Corollary: the only way to make the convex hull be <u>on</u> f is to shrink it to a single point (i.e. the Jensen's corollary)

#### Incomplete observations

Likelihood function with hidden random variables

Artificial Intelligence 2024–2025 Unsupervised Learning [25]

#### Expectation-Maximization (EM) Algorithm

*Alternate optimization (coordinate ascent)* 

Log-likelihood function:

$$\ell(\theta|D) \geq \sum_{m} \sum_{\{X_i\}} Q^{(m)}(\{X_i\}) \log \frac{P(D^{(m)}, \{X_i\}|\theta)}{Q^{(m)}(\{X_i\})}$$
This inequality becomes equality when this term is constant (see Jensen's corollary)

1) Keep  $\theta$  constant, define  $Q^{(m)}(\{Z_i\})$  so that the right side of the inequality is maximized

$$Q^{(m)}(\{X_i\}) := \frac{P(D^{(m)}, \{X_i\}|\theta)}{\sum\limits_{\{X_i\}} P(D^{(m)}, \{X_i\}|\theta)} = \frac{P(D^{(m)}, \{X_i\}|\theta)}{P(D^{(m)}|\theta)} = P(\{X_i\}|D^{(m)}, \theta) =: p_{\{X_i\}}^{(m)}$$
These numbers can be computed from the graphical model (i.e. as an inference step)

2) Then maximize the log-likelihood while keeping  $\,Q^{(m)}(\{Z_i\})\,$  constant

$$\theta^* = \operatorname{argmax}_{\theta} \sum_{m} \sum_{\{X_i\}} p_{\{X_i\}}^{(m)} \log \frac{P(D^{(m)}, \{X_i\} | \theta)}{p_{\{X_i\}}^{(m)}}$$

$$= \operatorname{argmax}_{\theta} \sum_{m} \left( \sum_{\{X_i\}} p_{\{X_i\}}^{(m)} \log P(D^{(m)}, \{X_i\} | \theta) - \sum_{\{X_i\}} p_{\{X_i\}}^{(m)} \log p_{\{X_i\}}^{(m)} \log p_{\{X_i\}}^{(m)} \right) \right)$$

$$= \operatorname{argmax}_{\theta} \sum_{m} \sum_{\{X_i\}} p_{\{X_i\}}^{(m)} \log P(D^{(m)}, \{X_i\} | \theta)$$

$$= \operatorname{argmax}_{\theta} \sum_{m} \sum_{\{X_i\}} p_{\{X_i\}}^{(m)} \log P(D^{(m)}, \{X_i\} | \theta)$$

#### Expectation-Maximization (EM) Algorithm

Alternate optimization (coordinate ascent)

Log-likelihood function and its estimator:

$$\ell(\theta|D) \geq \sum_{m} \sum_{\{X_i\}} Q^{(m)}(\{X_i\}) \log \frac{P(D^{(m)}, \{X_i\}|\theta)}{Q^{(m)}(\{X_i\})}$$

#### **Algorithm:**

- 1) Assign the  $\theta$  at random
- 2) (*E-step*) Compute the probabilities

$$p_{\{X_i\}}^{(m)} = Q^{(m)}(\{X_i\}) = P(\{X_i\}|D^{(m)},\theta)$$

3) (*M-step*) Compute a new estimate of  $\theta$ 

$$\theta^* = \operatorname{argmax}_{\theta} \sum_{m} \sum_{\{X_i\}} p_{\{X_i\}}^{(m)} \log P(D^{(m)}, \{X_i\} | \theta)$$

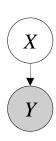
4) Go back to step 2) until some convergence criterion is met

The algorithm converges to a local maximum of the log-likelihood The effectiveness of algorithm depends on the form of  $P(\{X_i\}|D^{(m)},\theta)$  (see step 3) In particular, when this distribution is <u>exponential</u>... (e.g. Gaussian – see next slide)

Artificial Intelligence 2024–2025 Unsupervised Learning [27]

#### An aside: the EM algorithm in the continuous case\* (Mixture of Gaussians)

### EM Algorithm: mixture of Gaussians



#### **Model:**

The hidden variable X has k possible values, the observable variable Y is a point in  $\mathbb{R}^d$ 

$$P(X=k):=\phi_k$$
 Multivariate normal distribution 
$$P(Y=y|X=k)=\mathcal{N}(y;\mu_k,\Sigma_k):=(2\pi)^{-d/2}(\det\Sigma_k)^{-1/2}\exp\left(-\frac{1}{2}(y-\mu_k)^T\Sigma_k^{-1}(y-\mu_k)\right)$$
 i.e. the condition probabilities are normal distributions The observations are a set  $D=\{y^{(1)},\ldots,y^{(N)}\}$  of points in  $\mathbb{R}^d$ 

#### **Algorithm:**

- 1) For each value k, assign  $\phi_k$ ,  $\mu_k$  and  $\Sigma_k$  at random
- 2) (*E-step*) For all the  $y^{(m)}$  in D compute the probabilities  $p_k^{(m)} = P(X = k | y^{(m)}, \phi_k, \mu_k, \Sigma_k) = \phi_k \cdot \mathcal{N}(x^{(m)}; \mu_k, \Sigma_k)$
- 3) (*M-step*) Compute the new estimates for the parameters

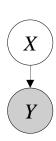
$$\phi_k = \frac{1}{n} \sum_{m} p_k^{(m)}$$

$$\mu_k = \frac{\sum_{m} p_k^{(m)} y^{(m)}}{\sum_{m} p_k^{(m)}}$$

$$\Sigma_k = \frac{\sum_{m} p_k^{(m)} (y - \mu_k) (y - \mu_k)^T}{\sum_{m} p_k^{(m)}}$$

4) Go back to step 2) until some convergence criterion is met

### EM Algorithm: mixture of Gaussians



#### **Model:**

The hidden variable X has k possible values, the variable Y is a point in  $\mathbb{R}^d$ 

$$P(X=k) := \phi_k$$

$$P(Y = y | X = k) = \mathcal{N}(y; \mu_k, \Sigma_k) := (2\pi)^{-d/2} (\det \Sigma_k)^{-1/2} \exp\left(-\frac{1}{2}(y - \mu_k)^T \Sigma_k^{-1} (y - \mu_k)\right)$$

i.e. the condition probabilities are <u>normal</u> distributions

The observations are a set  $D = \{y^{(1)}, \dots, y^{(N)}\}\$  of points in  $\mathbb{R}^d$ 

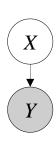
#### **Proof** (of the M-step):

$$\sum_{m} \sum_{k} p_{k}^{(m)} \log P(Y^{(m)}, X = k | \phi_{k}, \mu_{k}, \Sigma_{k})$$

$$= \sum_{m} \sum_{k} p_{k}^{(m)} \log P(Y^{(m)} | X = k, \mu_{k}, \Sigma_{k}) P(X = k | \phi_{k})$$

$$= \sum_{m} \sum_{k} p_{k}^{(m)} \left( \log \left( 2\pi^{-d/2} (\det \Sigma_{k})^{-1/2} \right) + \left( -\frac{1}{2} (y - \mu_{k})^{T} \Sigma_{k}^{-1} (y - \mu_{k}) \right) + \log \varphi_{k} \right)$$

### EM Algorithm: mixture of Gaussians



#### **Model:**

The hidden variable X has k possible values, the variable Y is a point in  $\mathbb{R}^d$ 

$$P(X=k) := \phi_k$$

$$P(Y = y | X = k) = \mathcal{N}(y; \mu_k, \Sigma_k) := (2\pi)^{-d/2} (\det \Sigma_k)^{-1/2} \exp\left(-\frac{1}{2}(y - \mu_k)^T \Sigma_k^{-1} (y - \mu_k)\right)$$

i.e. the condition probabilities are normal distributions

The observations are a set  $D = \{y^{(1)}, \dots, y^{(N)}\}\$  of points in  $\mathbb{R}^d$ 

#### **Proof** (of the M-step):

$$\begin{split} \frac{\partial}{\partial \mu_j} \sum_m \sum_k p_k^{(m)} \left( \log \left( (2\pi)^{-d/2} (\det \Sigma_k)^{-1/2} \right) + \left( -\frac{1}{2} (y^{(m)} - \mu_k)^T \Sigma_k^{-1} (y^{(m)} - \mu_k) \right) + \log \phi_k \right) \\ &= \frac{\partial}{\partial \mu_j} \sum_m \sum_k p_k^{(m)} \left( -\frac{1}{2} (y^{(m)} - \mu_k)^T \Sigma_k^{-1} (y^{(m)} - \mu_k) \right) \\ &= \frac{\partial}{\partial \mu_j} \sum_m \sum_k p_k^{(m)} \left( -\frac{1}{2} (y^{(m)}^T \Sigma_k^{-1} y^{(m)} + \mu_k^T \Sigma_k^{-1} \mu_k - 2 y^{(m)}^T \Sigma_k^{-1} \mu_k) \right) \\ &= \sum_m p_j^{(m)} \left( x^T \Sigma_j^{-1} - \mu_j^T \Sigma_j^{-1} \right) \\ &= \sum_m p_j^{(m)} \left( x^T \Sigma_j^{-1} - \mu_j^T \Sigma_j^{-1} \right) = 0 \qquad \Rightarrow \qquad \qquad \\ \text{In the web page for the derivations of other parameters} \dots \end{split}$$

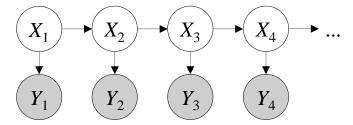
See the link in the web page for the derivations of other parameters ...

# The EM algorithm for learning with <u>missing data</u>

### Missing Values

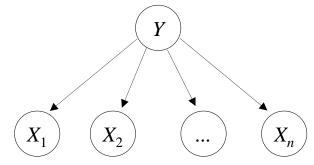
#### Hidden Variables

Some of the variables may be hidden, i.e., non observable 'by design' Example: 'Hidden Markov Model'



#### Incomplete Observations

Sometimes, however, observations may be missing 'by accident' and not 'by design' Example: 'Naïve Bayesian Classifier'



What if some classifications Y are missing, or a few features  $X_i$  are not available?

Artificial Intelligence 2024–2025 Unsupervised Learning [33]

### Missing Values: Observability Model

Observed and Unobserved Variables

Let's consider a graphical model with a set of random variables:

$$\boldsymbol{X} := \{X_1, \dots, X_n\}$$

In each actual observation (i.e., a data item)

$$\boldsymbol{X}^{(m)} := \{X_1^{(m)}, \dots, X_n^{(m)}\}$$

each value  $X_i^{(i)}$  may be either observed or unobserved (i.e., missing) determined by a <u>binary</u> random variable  $O_{X_i} \in \{0,1\}$ 

An **observability model** for a graphical model with random variables X is a set of binary random observability variables

$$O_X := \{O_{X_1}, \dots, O_{X_n}\}$$

with probability distribution

$$P_{missing}(\boldsymbol{X}, \boldsymbol{O}_{\boldsymbol{X}}) = P(\boldsymbol{X}) P_{missing}(\boldsymbol{O}_{\boldsymbol{X}} | \boldsymbol{X})$$

probability distribution with no missing values

### Missing Values: MCAR

Missing Completely at Random (MCAR)

MCAR assumption

$$\langle m{X} \perp m{O_X} 
angle$$

It entails that:

$$P_{missing}(\boldsymbol{X}, \boldsymbol{O}_{\boldsymbol{X}}) = P(\boldsymbol{X}) P_{missing}(\boldsymbol{O}_{\boldsymbol{X}})$$

This is tempting and it could ease all subsequent computations...

... but it is too strong, and hardly enforceable in many practical cases

Moral: we need a weaker assumption

### Missing Values: MAR

Missing at Random (MAR)

Consider a generic data item, possibly with missing values

$$oldsymbol{X}^{(m)} := oldsymbol{X}^{(m)}_{obs} \cup oldsymbol{X}^{(m)}_{hid}$$
observed missing

Missing values need NOT be for the same variables in each data item

MAR assumption, for each data item:

$$\langle oldsymbol{X}_{hid}^{(m)} \perp oldsymbol{O_{oldsymbol{X}}} \mid oldsymbol{X}_{obs}^{(m)} 
angle$$

Namely, the *values* of the missing variables are independent from their *observability* given the *values* of the observed variables

It is still a strong assumption, yet much more realistic...

Artificial Intelligence 2024–2025 Unsupervised Learning [36]

### Missing Values: MAR

Missing at Random (MAR)

MAR assumption, for each data item:

$$\langle oldsymbol{X}_{hid}^{(m)} \perp oldsymbol{O_{oldsymbol{X}}} \mid oldsymbol{X}_{obs}^{(m)} 
angle$$

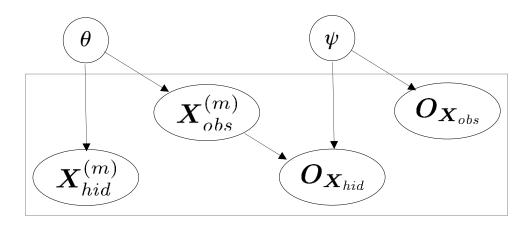
This entails:

$$\begin{split} P_{missing}(\boldsymbol{X}^{(m)}, \boldsymbol{O}_{\boldsymbol{X}}) &= P_{missing}(\boldsymbol{X}_{obs}^{(m)}, \boldsymbol{X}_{hid}^{(m)}, \boldsymbol{O}_{\boldsymbol{X}}) \\ &= P(\boldsymbol{X}_{obs}^{(m)}, \boldsymbol{X}_{hid}^{(m)}) \, P_{missing}(\boldsymbol{O}_{\boldsymbol{X}} | \boldsymbol{X}_{obs}^{(m)}, \boldsymbol{X}_{hid}^{(m)}) \\ &= P(\boldsymbol{X}_{obs}^{(m)}, \boldsymbol{X}_{hid}^{(m)}) \, P_{missing}(\boldsymbol{O}_{\boldsymbol{X}} | \boldsymbol{X}_{obs}^{(m)}) \\ P_{missing}(\boldsymbol{X}_{obs}^{(m)}, \boldsymbol{O}_{\boldsymbol{X}}) &= \sum_{\boldsymbol{X}_{hid}} P(\boldsymbol{X}_{obs}^{(m)}, \boldsymbol{X}_{hid}^{(m)}) \, P_{missing}(\boldsymbol{O}_{\boldsymbol{X}} | \boldsymbol{X}_{obs}^{(m)}) \end{split}$$

$$P_{missing}(m{X}_{obs}^{(m)},m{O_X}) = P(m{X}_{obs}^{(m)})\,P_{missing}(m{O_X}|m{X}_{obs}^{(m)})$$
 This is the relevant property

Artificial Intelligence 2024–2025 Unsupervised Learning [37]

#### Likelihood under MAR



$$\langle oldsymbol{X}_{hid}^{(m)} \perp oldsymbol{O_{oldsymbol{X}}} \mid oldsymbol{X}_{obs}^{(m)} 
angle$$

Variables and parameters in an observability model as a graphical model

#### Likelihood

$$L(\theta, \psi \mid D) = \prod_{m=1}^{N} P(D^{(m)} \mid \theta, \psi) \quad \text{where:} \quad D^{(m)} := \{\boldsymbol{X}_{obs}^{(m)}, \boldsymbol{X}_{hid}^{(m)}\}$$

$$l(\theta, \psi \mid D) = \sum_{m=1}^{N} \log P(D^{(m)} \mid \theta, \psi)$$

$$= \sum_{m=1}^{N} \log \left(P(\boldsymbol{X}_{obs}^{(m)}, \boldsymbol{X}_{hid}^{(m)} \mid \theta, \psi) P_{missing}(\boldsymbol{O}_{\boldsymbol{X}} | \boldsymbol{X}_{obs}^{(m)}, \boldsymbol{X}_{hid}^{(m)}, \theta, \psi)\right)$$

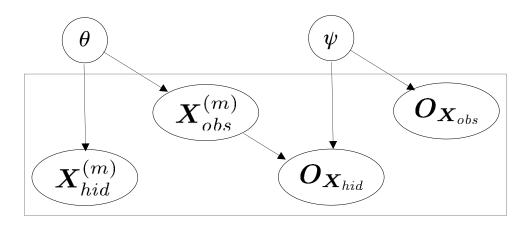
$$= \sum_{m=1}^{N} \log \left(P(\boldsymbol{X}_{obs}^{(m)}, \boldsymbol{X}_{hid}^{(m)} \mid \theta) P_{missing}(\boldsymbol{O}_{\boldsymbol{X}} | \boldsymbol{X}_{obs}^{(m)}, \psi)\right)$$

$$= \sum_{m=1}^{N} \log P(\boldsymbol{X}_{obs}^{(m)}, \boldsymbol{X}_{hid}^{(m)} \mid \theta) + \sum_{m=1}^{N} \log P_{missing}(\boldsymbol{O}_{\boldsymbol{X}} | \boldsymbol{X}_{obs}^{(m)}, \psi)$$

We are interested in optimizing  $\theta$  ...

yet we have only <u>observed</u> values

#### Likelihood under MAR



$$\langle oldsymbol{X}_{hid}^{(m)} \perp oldsymbol{O_{oldsymbol{X}}} \mid oldsymbol{X}_{obs}^{(m)} 
angle$$

Variables and parameters in an observability model as a graphical model

#### Likelihood (for observed values)

$$\begin{split} l(\theta \mid D) &:= \sum_{m=1}^{N} \log P(\boldsymbol{X}_{obs}^{(m)} \mid \theta) \\ &= \sum_{m=1}^{N} \log \sum_{\boldsymbol{X}_{hid}^{(m)}} P(\boldsymbol{X}_{obs}^{(m)}, \boldsymbol{X}_{hid}^{(m)} \mid \theta) \\ &= \sum_{m=1}^{N} \log \sum_{\boldsymbol{X}_{hid}^{(m)}} \left( P(\boldsymbol{X}_{obs}^{(m)} \mid \theta) P(\boldsymbol{X}_{hid}^{(m)} \mid \boldsymbol{X}_{obs}^{(m)}, \theta) \right) \end{split}$$

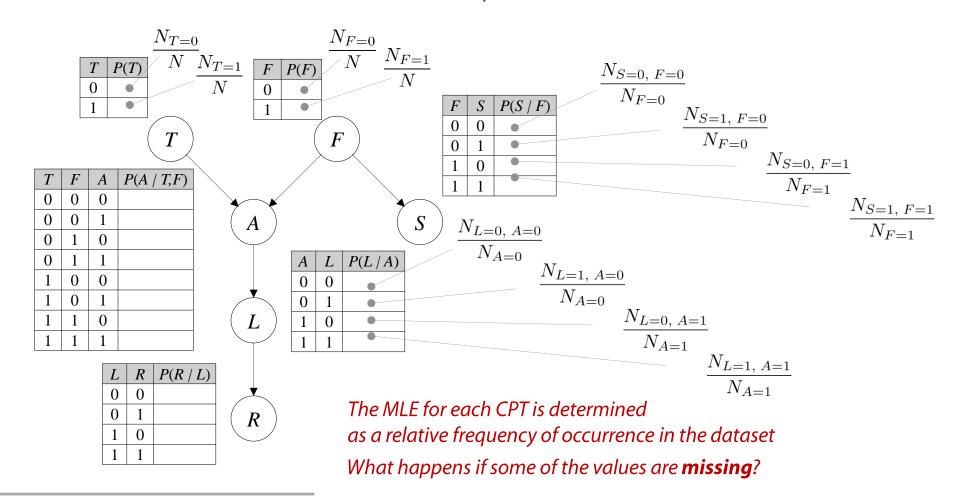
Looks promising: using probabilities instead of missing values ...

Artificial Intelligence 2024–2025 Unsupervised Learning [39]

**Model**: random variables plus the graph of dependencies

**Observations**: dataset of values, from <u>completely observed</u> outcomes

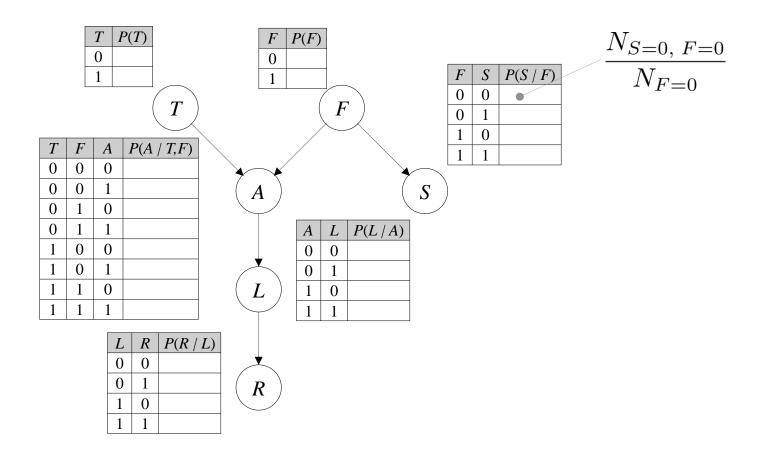
Parameters (to be determined): all conditional probabilities (i.e. all CPTs)



Artificial Intelligence 2024–2025 Unsupervised Learning [40]

#### Fundamental idea: using probabilities of observations

In the completely observed case: probabilities are estimated as frequencies of occurrence



Artificial Intelligence 2024-2025 Unsupervised Learning [41]

**Fundamental idea:** using probabilities of missing observations

Let's consider a dataset 
$$D = \left\{ oldsymbol{X}^{(i)} 
ight\}_{i=1}^{N}$$

Each data item  $oldsymbol{X}^{(i)}$  may contain some missing data Example:  $m{X}^{(i)} = (X_1^{(i)} = x_1, X_2^{(i)} = x_2, X_3^{(i)} = ?)$ 

Example: 
$$\boldsymbol{X}^{(i)} = (X_1^{(i)} = x_1, X_2^{(i)} = x_2, X_3^{(i)} = ?)$$

Define  $ilde{m{X}}^{(i)}$  as one possible **completion** of  $m{X}^{(i)}$  , namely one in which there are no missing data

Example: assuming that 
$$X_3 \in \{0,1\}$$
 ,  $\begin{array}{l} (X_1^{(i)} = x_1, X_2^{(i)} = x_2, X_3^{(i)} = 0) \\ (X_1^{(i)} = x_1, X_2^{(i)} = x_2, X_3^{(i)} = 1) \end{array}$  are the two possible completions of  $\boldsymbol{X}^{(i)}$ 

Note that there will be as many completions of  $m{X}^{(i)}$  as there are combinations of possible values for the missing data For any <u>complete</u> observation,  $ilde{X}^{(i)}=X^{(i)}$  i.e., there is only one possible completion that coincides with the data item itself

Likewise,  $m{X}_{obs}^{(i)}$  is the part of  $m{X}^{(i)}$  which contain the actual observations

Example: 
$$\boldsymbol{X}_{obs}^{(i)} = (X_1^{(i)} = x_1, X_2^{(i)} = x_2)$$

#### Fundamental idea: using probabilities of observations

In the completely observed case: probabilities are estimated as frequencies of occurrence More in general:

$$\frac{N_{X_i, \mathbf{Z}}}{N_{\mathbf{Z}}}$$
 where:  $\mathbf{Z} = parents(X_i)$ 

In the EM algorithm, use *estimated* occurrences:

EM algorithm, use *estimated* occurrences: 
$$\frac{\tilde{N}_{X_i, \mathbf{Z}}}{\tilde{N}_{\mathbf{Z}}} \qquad \text{where:} \qquad \tilde{N}_{\mathbf{X}} := \sum_{i=1}^{N} \sum_{\tilde{\boldsymbol{X}}^{(i)}} P(\tilde{\boldsymbol{X}}^{(i)} \mid \boldsymbol{X}_{obs}^{(i)}, \theta)$$

*In words, any incomplete observations 'splits up' and contributes with the probabilities of possible completions Note that, when all observations are complete:* 

$$\tilde{N}_{\boldsymbol{X}} = N_{\boldsymbol{X}}$$

Fundamental idea: using probabilities of observations

#### **Algorithm:**

- 1) Assign parameters  $\theta^{(0)}$  at random
- 2) Compute  $P(oldsymbol{X} \mid heta^{(t)})$

E-step

3) Update all parameters using estimated occurrences:

$$heta_{X_i \,|\, m{Z}}^{(t+1)} = rac{ ilde{N}_{X_i, m{Z}}}{ ilde{N}_{m{Z}}}$$
 where all estimations are made using  $P(m{X} \mid m{ heta}^{(t)})$  M-step

3) Go back to step 2) until some convergence criterion is met

Artificial Intelligence 2024-2025