### Lecture 18: 10 June, 2021

Madhavan Mukund

https://www.cmi.ac.in/~madhavan

Data Mining and Machine Learning April–July 2021

#### Mixture models

- Probabilistic process parameters ⊖
  - Tossing a coin with  $\Theta = \{Pr(H)\} = \{p\}$
- Perform an experiment
  - Toss the coin N times,  $H T H H \cdots T$
- Estimate parameters from observations
  - From h heads, estimate p = h/N
  - Maximum Likelihood Estimator (MLE)
- What if we have a mixture of two random processes
  - Two coins,  $c_1$  and  $c_2$ , with  $Pr(H) = p_1$  and  $p_2$ , respectively
  - Repeat N times: choose  $c_i$  with probability 1/2 and toss it
  - Outcome:  $N_1$  tosses of  $c_1$  interleaved with  $N_2$  tosses of  $c_2$ ,  $N_1 + N_2 = N$
  - $\blacksquare$  Can we estimate  $p_1$  and  $p_2$ ?

### Mixture models ...

- Two coins,  $c_1$  and  $c_2$ , with  $Pr(H) = p_1$  and  $p_2$ , respectively
- Sequence of *N* interleaved coin tosses *H T H H · · · · H H T*
- If the sequence is labelled, we can estimate  $p_1$ ,  $p_2$  separately
  - $\blacksquare$  H T T H H T
  - $p_1 = 8/12 = 2/3$ ,  $p_2 = 3/8$
- What the observation is unlabelled?
- Iterative algorithm to estimate the parameters
  - Make an initial guess for the parameters
  - Compute a (fractional) labelling of the outcomes
  - Re-estimate the parameters



- Iterative algorithm to estimate the parameters
  - Make an initial guess for the parameters
  - Compute a (fractional) labelling of the outcomes
  - Re-estimate the parameters

Madhavan Mukund Lecture 18: 10 June, 2021 DMML Apr-Jul 2021

- Iterative algorithm to estimate the parameters
  - Make an initial guess for the parameters
  - Compute a (fractional) labelling of the outcomes
  - Re-estimate the parameters
- $\blacksquare$  HTTHHTHTHHTHTHTHTHTHT
  - Initial guess:  $p_1 = 1/2$ ,  $p_2 = 1/4$

- Iterative algorithm to estimate the parameters
  - Make an initial guess for the parameters
  - Compute a (fractional) labelling of the outcomes
  - Re-estimate the parameters
- - Initial guess:  $p_1 = 1/2$ ,  $p_2 = 1/4$
  - $Pr(c_1 = T) = q_1 = 1/2, Pr(c_2 = T) = q_2 = 3/4.$

4/14

Madhavan Mukund Lecture 18: 10 June. 2021

- Iterative algorithm to estimate the parameters
  - Make an initial guess for the parameters
  - Compute a (fractional) labelling of the outcomes
- Re-estimate the parameters
  - НТТННТНТННТНТНТНТНТ
    - Initial guess:  $p_1 = 1/2$ ,  $p_2 = 1/4$
    - $Pr(c_1 = T) = q_1 = 1/2, Pr(c_2 = T) = q_2 = 3/4,$
    - For each H, likelihood it was  $c_i$ ,  $Pr(c_i \mid H)$ , is  $p_i/(p_1 + p_2)$

Madhavan Mukund Lecture 18: 10 June, 2021 DMML Apr-Jul 2021 4 / 14

- Iterative algorithm to estimate the parameters
  - Make an initial guess for the parameters
  - Compute a (fractional) labelling of the outcomes
  - Re-estimate the parameters
- - Initial guess:  $p_1 = 1/2$ ,  $p_2 = 1/4$
  - $Pr(c_1 = T) = q_1 = 1/2, Pr(c_2 = T) = q_2 = 3/4,$
  - For each H, likelihood it was  $c_i$ ,  $Pr(c_i \mid H)$ , is  $p_i/(p_1 + p_2)$
  - For each T, likelihood it was  $c_i$ ,  $Pr(c_i \mid T)$ , is  $q_i/(q_1+q_2)$

- Iterative algorithm to estimate the parameters
  - Make an initial guess for the parameters
  - Compute a (fractional) labelling of the outcomes
  - Re-estimate the parameters
- $\blacksquare$  H T T H H T H T H T H T H T H T H T
  - Initial guess:  $p_1 = 1/2$ ,  $p_2 = 1/4$
  - $Pr(c_1 = T) = q_1 = 1/2, Pr(c_2 = T) = q_2 = 3/4,$
  - For each H, likelihood it was  $c_i$ ,  $Pr(c_i \mid H)$ , is  $p_i/(p_1 + p_2)$
  - For each T, likelihood it was  $c_i$ ,  $Pr(c_i \mid T)$ , is  $q_i/(q_1 + q_2)$
  - Assign fractional count  $Pr(c_i \mid H)$  to each  $H: 2/3 \times c_1, 1/3 \times c_2$

4/14

Madhavan Mukund Lecture 18: 10 June, 2021 DMML Apr-Jul 2021

- Iterative algorithm to estimate the parameters
  - Make an initial guess for the parameters
  - Compute a (fractional) labelling of the outcomes
  - Re-estimate the parameters

initial guess: 
$$p_1 = 1/2$$
,  $p_2 = 1/4$ 

$$Pr(c_1 = T) = q_1 = 1/2, Pr(c_2 = T) = q_2 = 3/4,$$

- For each H, likelihood it was  $c_i$ ,  $Pr(c_i \mid H)$ , is  $p_i/(p_1 + p_2)$
- For each T, likelihood it was  $c_i$ ,  $Pr(c_i \mid T)$ , is  $q_i/(q_1+q_2)$
- Assign fractional count  $Pr(c_i \mid H)$  to each  $H: 2/3 \times c_1, 1/3 \times c_2$
- Likewise, assign fractional count  $Pr(c_i \mid T)$  to each  $T: 2/5 \times c_1, 3/5 \times c_2$

- $\blacksquare$  H T T H H T H T H T H T H T H T H T
- Initial guess:  $p_1 = 1/2$ ,  $p_2 = 1/4$
- Fractional counts: each H is  $2/3 \times c_1$ ,  $1/3 \times c_2$ , each T:  $2/5 \times c_1$ ,  $3/5 \times c_2$

Madhayan Mukund Lecture 18: 10 June. 2021 DMML Apr-Jul 2021 5 / 14

- *HTTHHTHTHTHTHTHTHT* 20 < 11 4
- Initial guess:  $p_1 = 1/2$ ,  $p_2 = 1/4$
- Fractional counts: each H is  $2/3 \times c_1$ ,  $1/3 \times c_2$ , each T:  $2/5 \times c_1$ ,  $3/5 \times c_2$
- Add up the fractional counts
  - $c_1$ :  $11 \cdot (2/3) = 22/3$  heads,  $9 \cdot (2/5) = 18/5$  tails
  - $c_2$ :  $11 \cdot (1/3) = 11/3$  heads,  $9 \cdot (3/5) = 27/5$  tails

5/14

Madhavan Mukund Lecture 18: 10 June, 2021 DMML Apr-Jul 2021

- $\blacksquare$  H T T H H T H T H T H T H T H T H T
- Initial guess:  $p_1 = 1/2$ ,  $p_2 = 1/4$
- Fractional counts: each H is  $2/3 \times c_1$ ,  $1/3 \times c_2$ , each T:  $2/5 \times c_1$ ,  $3/5 \times c_2$
- Add up the fractional counts
  - $c_1$ : 11 · (2/3) = 22/3 heads, 9 · (2/5) = 18/5 tails
  - $c_2$ :  $11 \cdot (1/3) = 11/3$  heads,  $9 \cdot (3/5) = 27/5$  tails

 $\frac{h}{N} = \frac{0.32}{0.3370.4}$ 

5/14

Re-estimate the parameters

- $\blacksquare$  HTTHHTHTHHTHTHTHTHTHT
- Initial guess:  $p_1 = 1/2$ ,  $p_2 = 1/4$
- Fractional counts: each H is  $2/3 \times c_1$ ,  $1/3 \times c_2$ , each T:  $2/5 \times c_1$ ,  $3/5 \times c_2$
- Add up the fractional counts
  - $c_1$ :  $11 \cdot (2/3) = 22/3$  heads,  $9 \cdot (2/5) = 18/5$  tails
  - $c_2$ :  $11 \cdot (1/3) = 11/3$  heads,  $9 \cdot (3/5) = 27/5$  tails
- Re-estimate the parameters

■ 
$$p_2 = \frac{11/3}{11/3 + 27/5} = 55/136 = 0.40$$
,  $q_2 = 1 - p_2 = 0.60$ 

■ Repeat until convergence



Lecture 18: 10 June, 2021

■ Mixture of probabilistic models  $(M_1, M_2, ..., M_k)$  with parameters  $\Theta = (\theta_1, \theta_2, ..., \theta_k)$ 



6/14

Madhavan Mukund Lecture 18: 10 June, 2021 DMML Apr-Jul 2021

- Mixture of probabilistic models  $(M_1, M_2, ..., M_k)$  with parameters  $\Theta = (\theta_1, \theta_2, ..., \theta_k)$
- Observation  $O = o_1 o_2 \dots o_N$

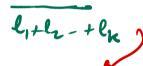


6/14

Madhavan Mukund Lecture 18: 10 June, 2021 DMML Apr-Jul 2021

- Mixture of probabilistic models  $(M_1, M_2, ..., M_k)$  with parameters  $\Theta = (\theta_1, \theta_2, \dots, \theta_k)$
- Observation  $O = o_1 o_2 \dots o_N$
- Expectation step
  - Compute likelihoods  $Pr(M_i|o_i)$  for each  $M_i$ ,  $o_i$

- Mixture of probabilistic models  $(M_1, M_2, ..., M_k)$  with parameters  $\Theta = (\theta_1, \theta_2, ..., \theta_k)$
- Observation  $O = o_1 o_2 \dots o_N$
- Expectation step
  - Compute likelihoods  $Pr(M_i|o_j)$  for each  $M_i$ ,  $o_j$
- li

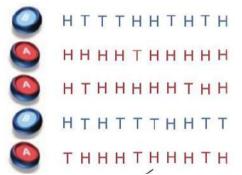


- Maximization step
  - Recompute MLE for each  $M_i$  using fraction of O assigned using likelihood

- Mixture of probabilistic models  $(M_1, M_2, ..., M_k)$  with parameters  $\Theta = (\theta_1, \theta_2, ..., \theta_k)$
- Observation  $O = o_1 o_2 \dots o_N$
- Expectation step
  - Compute likelihoods  $Pr(M_i|o_j)$  for each  $M_i$ ,  $o_j$
- Maximization step
  - **Recompute MLE** for each  $M_i$  using fraction of O assigned using likelihood
- Repeat until convergence
  - Why should it converge?
  - If the value converges, what have we computed?



Two biased coins, choose a coin and toss 10 times, repeat 5 times



Two biased coins, choose a coin and toss 10 times, repeat 5 times











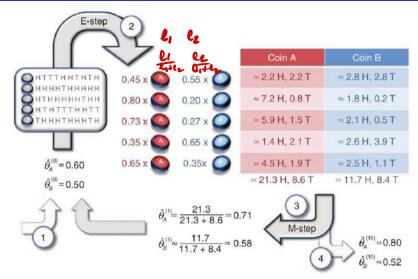
If we know the breakup, we can separately compute MLE for each coin

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

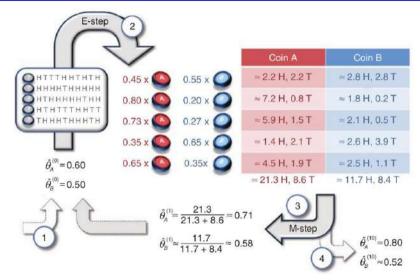
$$\hat{\theta}_A = \frac{24}{24 + 6} = 0.80$$

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

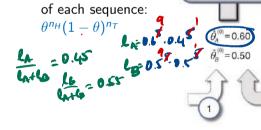
Expectation-Maximization

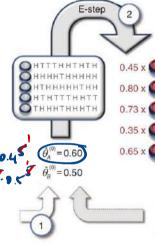


- Expectation-Maximization
- Initial estimates,  $\theta_A = 0.6$ ,  $\theta_B = 0.5$



- Expectation-Maximization
- Initial estimates.  $\theta_{A} = 0.6, \, \theta_{B} = 0.5$
- Compute likelihood of each sequence:

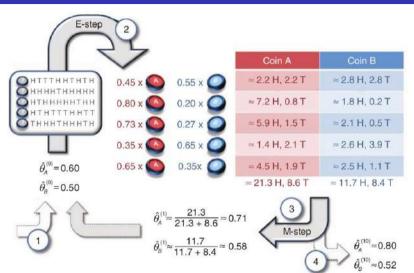




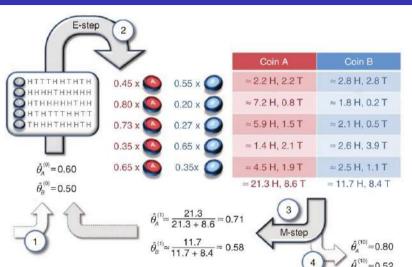
		Coin A	Coin B
( O	0.55 x	≈ 2.2 H, 2.2 T	≈ 2.8 H, 2.8 T
0	0.20 x	≈ 7.2 H, 0.8 T	≈ 1.8 H, 0.2 T
0	0.27 x	≈ 5.9 H, 1.5 T	≈ 2.1 H, 0.5 T
0	0.65 x	≈ 1.4 H, 2.1 T	= 2.6 H, 3.9 T
0	0.35x	≈ 4.5 H, 1.9 T	= 2.5 H, 1.1 T
_		≈ 21.3 H, 8.6 T	≈ 11.7 H, 8.4 T



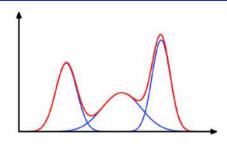
- Expectation-Maximization
- Initial estimates,  $\theta_A = 0.6$ ,  $\theta_B = 0.5$
- Compute likelihood of each sequence:  $\theta^{n_H}(1-\theta)^{n_T}$
- Assign each sequence proportionately



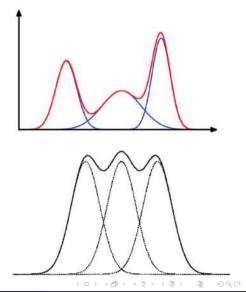
- Expectation-Maximization
- Initial estimates,  $\theta_A = 0.6$ ,  $\theta_B = 0.5$
- Compute likelihood of each sequence:  $\theta^{n_H}(1-\theta)^{n_T}$
- Assign each sequence proportionately
- Converge to  $\theta_A = 0.8$ ,  $\theta_B = 0.52$



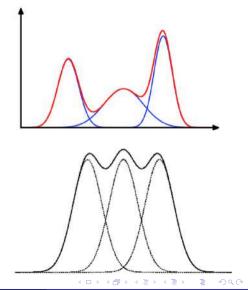
■ Sample uniformly from multiple Gaussians,  $\mathcal{N}(\mu_i, \sigma_i)$ 



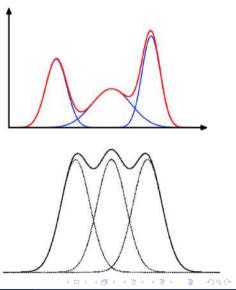
- Sample uniformly from multiple Gaussians,  $\mathcal{N}(\mu_i, \sigma_i)$
- For simplicity, assume all  $\sigma_i = \sigma$



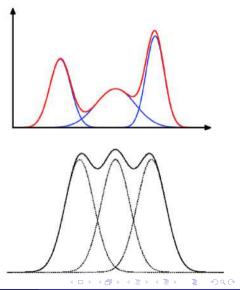
- Sample uniformly from multiple Gaussians,  $\mathcal{N}(\mu_i, \sigma_i)$
- For simplicity, assume all  $\sigma_i = \sigma$
- N sample points  $z_1, z_2, \ldots, z_N$



- Sample uniformly from multiple Gaussians,  $\mathcal{N}(\mu_i, \sigma_i)$
- For simplicity, assume all  $\sigma_i = \sigma$
- N sample points  $z_1, z_2, \ldots, z_N$
- lacksquare Make an initial guess for each  $\mu_j$



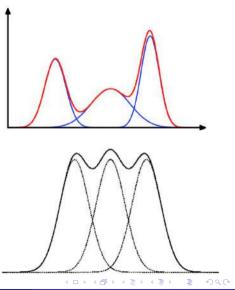
- Sample uniformly from multiple Gaussians,  $\mathcal{N}(\mu_i, \sigma_i)$
- For simplicity, assume all  $\sigma_i = \sigma$
- N sample points  $z_1, z_2, \ldots, z_N$
- lacksquare Make an initial guess for each  $\mu_j$
- $Pr(z_i \mid \mu_j) = exp(-\frac{1}{2\sigma^2}(z_i \mu_j)^2)$



- Sample uniformly from multiple Gaussians,  $\mathcal{N}(\mu_i, \sigma_i)$
- For simplicity, assume all  $\sigma_i = \sigma$
- *N* sample points  $z_1, z_2, ..., z_N$
- lacksquare Make an initial guess for each  $\mu_j$

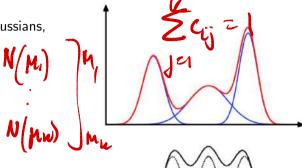
• 
$$Pr(z_i \mid \mu_j) = exp(-\frac{1}{2\sigma^2}(z_i - \mu_j)^2)$$

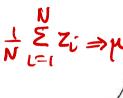
$$Pr(\mu_j \mid z_i) = c_{ij} = \frac{Pr(z_i \mid \mu_j)}{\sum_k Pr(z_i \mid \mu_k)}$$

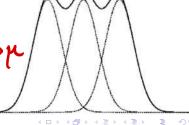


■ Sample uniformly from multiple Gaussians,  $\mathcal{N}(\mu_i, \sigma_i)$ 

- For simplicity, assume all  $\sigma_i = \sigma$
- N sample points  $z_1, z_2, \ldots, z_N$
- lacksquare Make an initial guess for each  $\mu_j$
- $Pr(z_i \mid \mu_j) = exp(-\frac{1}{2\sigma^2}(z_i \mu_j)^2)$
- $Pr(\mu_j \mid z_i) = c_{ij} = \frac{Pr(z_i \mid \mu_j)}{\sum_k Pr(z_i \mid \mu_k)}$
- MLE of  $\mu_j$  is sample mean,  $\sum_{j}$





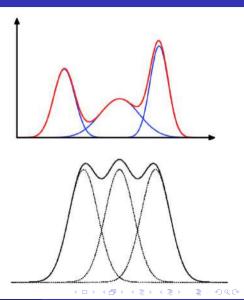


- Sample uniformly from multiple Gaussians,  $\mathcal{N}(\mu_i, \sigma_i)$
- For simplicity, assume all  $\sigma_i = \sigma$
- N sample points  $z_1, z_2, \ldots, z_N$
- lacksquare Make an initial guess for each  $\mu_j$

• 
$$Pr(z_i \mid \mu_j) = exp(-\frac{1}{2\sigma^2}(z_i - \mu_j)^2)$$

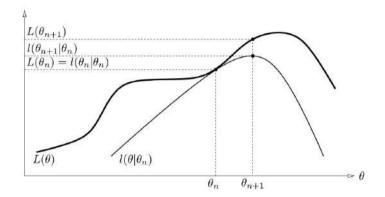
$$Pr(\mu_j \mid z_i) = c_{ij} = \frac{Pr(z_i \mid \mu_j)}{\sum_k Pr(z_i \mid \mu_k)}$$

- MLE of  $\mu_j$  is sample mean,  $\frac{\sum_i c_{ij} z_i}{\sum_i c_{ij}}$
- Update estimates for  $\mu_i$  and repeat



### Theoretical foundations of EM

■ Mixture of probabilistic models  $(M_1, M_2, ..., M_k)$  with parameters  $\Theta = (\theta_1, \theta_2, ..., \theta_k)$ 



10 / 14

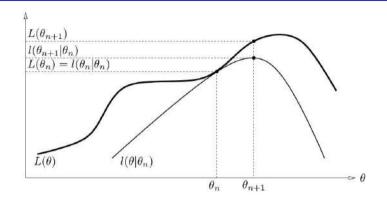
Madhavan Mukund Lecture 18: 10 June, 2021 DMML Apr–Jul 2021

### Theoretical foundations of EM

■ Mixture of probabilistic models  $(M_1, M_2, ..., M_k)$  with parameters  $\Theta = (\theta_1, \theta_2, ..., \theta_k)$ 

Observation

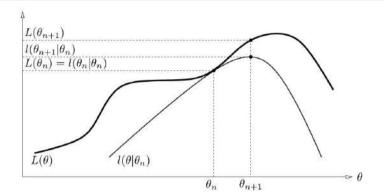
$$O = o_1 o_2 \dots o_N$$



- Mixture of probabilistic models  $(M_1, M_2, ..., M_k)$  with parameters  $\Theta = (\theta_1, \theta_2, ..., \theta_k)$
- Observation

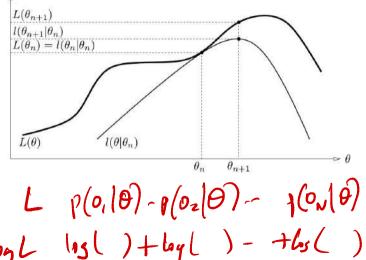
$$O = o_1 o_2 \dots o_N$$

■ EM builds a sequence of estimates  $\Theta_1, \Theta_2, \dots, \Theta_n$ 



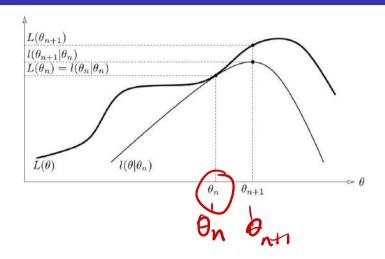
10 / 14

- Mixture of probabilistic models  $(M_1, M_2, \ldots, M_k)$ with parameters  $\Theta = (\theta_1, \theta_2, \dots, \theta_k)$
- Observation  $O = o_1 o_2 \dots o_N$
- EM builds a sequence of estimates  $\Theta_1, \Theta_2, \dots, \Theta_n$
- $L(\Theta_i)$  log-likelihood function,  $\ln Pr(O \mid \Theta_i)$



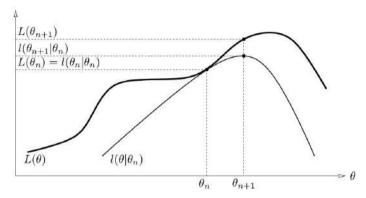
10 / 14

- Mixture of probabilistic models  $(M_1, M_2, \ldots, M_k)$ with parameters  $\Theta = (\theta_1, \theta_2, \dots, \theta_k)$
- Observation  $Q = Q_1 Q_2 \dots Q_N$
- EM builds a sequence of estimates  $\Theta_1, \Theta_2, \dots, \Theta_n$
- $L(\Theta_i)$  log-likelihood function,  $\ln Pr(O \mid \Theta_i)$
- Want to extend the sequence with  $\Theta_{n+1}$  such that  $L(\Theta_{n+1}) > L(\Theta_n)$



10 / 14

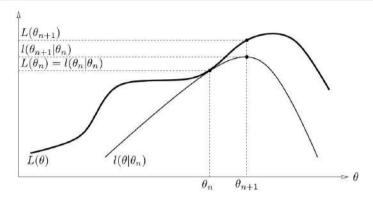
- Mixture of probabilistic models  $(M_1, M_2, ..., M_k)$ with parameters  $\Theta = (\theta_1, \theta_2, ..., \theta_k)$
- Observation  $O = o_1 o_2 \dots o_N$
- EM builds a sequence of estimates  $\Theta_1, \Theta_2, \dots, \Theta_n$
- $L(\Theta_j)$  log-likelihood function,  $\ln Pr(O \mid \Theta_j)$
- Want to extend the sequence with  $\Theta_{n+1}$  such that  $L(\Theta_{n+1}) > L(\Theta_n)$



■ EM performs a form of gradient descenct

Madhavan Mukund

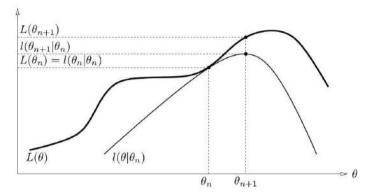
- Mixture of probabilistic models  $(M_1, M_2, ..., M_k)$  with parameters  $\Theta = (\theta_1, \theta_2, ..., \theta_k)$
- Observation  $O = o_1 o_2 \dots o_N$
- EM builds a sequence of estimates  $\Theta_1, \Theta_2, \dots, \Theta_n$
- $L(\Theta_j)$  log-likelihood function,  $\ln Pr(O \mid \Theta_j)$
- Want to extend the sequence with  $\Theta_{n+1}$  such that  $L(\Theta_{n+1}) > L(\Theta_n)$



- EM performs a form of gradient descenct
- If we update  $\Theta_n$  to  $\Theta'$  we get an new likelihood  $L(\Theta_n) + \Delta(\Theta', \Theta_n)$  which we call  $\ell(\Theta' \mid \Theta_n)$



- Mixture of probabilistic models  $(M_1, M_2, ..., M_k)$ with parameters  $\Theta = (\theta_1, \theta_2, ..., \theta_k)$
- Observation  $O = o_1 o_2 \dots o_N$
- EM builds a sequence of estimates  $\Theta_1, \Theta_2, \dots, \Theta_n$
- $L(\Theta_j)$  log-likelihood function,  $\ln Pr(O \mid \Theta_j)$
- Want to extend the sequence with  $\Theta_{n+1}$  such that  $L(\Theta_{n+1}) > L(\Theta_n)$



- EM performs a form of gradient descenct
- If we update  $\Theta_n$  to  $\Theta'$  we get an new likelihood  $L(\Theta_n) + \Delta(\Theta', \Theta_n)$  which we call  $\ell(\Theta' \mid \Theta_n)$
- Choose  $\Theta_{n+1}$  to maximize  $\ell(\Theta' \mid \Theta_n)$

Madhavan Mukund Lecture 18: 10 June. 2021 DMML Apr-Jul 2021 10 / 14

Supervised learning requires labelled training data

- Supervised learning requires labelled training data
- What if we don't have enough labelled data?

- Supervised learning requires labelled training data
- What if we don't have enough labelled data?
- For a probabilistic classifier we can apply EM

- Supervised learning requires labelled training data
- What if we don't have enough labelled data?
- For a probabilistic classifier we can apply EM
  - Use available training data to assign initial probabilities

- Supervised learning requires labelled training data
- What if we don't have enough labelled data?
- For a probabilistic classifier we can apply EM
  - Use available training data to assign initial probabilities
  - Label the rest of the data using this model fractional labels

11 / 14

- Supervised learning requires labelled training data
- What if we don't have enough labelled data?
- For a probabilistic classifier we can apply EM
  - Use available training data to assign initial probabilities
  - Label the rest of the data using this model fractional labels
  - Add up counts and re-estimate the parameters

11 / 14

■ Each document is a multiset or bag of words over a vocabulary

$$V = \{w_1, w_2, \dots, w_m\}$$

12 / 14

■ Each document is a multiset or bag of words over a vocabulary

$$V = \{w_1, w_2, \dots, w_m\}$$

■ Each topic c has probability Pr(c)

12 / 14

- Each document is a multiset or bag of words over a vocabulary  $V = \{w_1, w_2, \dots, w_m\}$
- Each topic c has probability Pr(c)
- Each word  $w_i \in V$  has conditional probability  $Pr(w_i \mid c_j)$ , for  $c_j \in C$ 
  - Note that  $\sum_{i=1}^{m} Pr(w_i \mid c_j) = 1$

Madhayan Mukund Lecture 18: 10 June. 2021 DMML Apr-Jul 2021 12 / 14

- Each document is a multiset or bag of words over a vocabulary  $V = \{w_1, w_2, \dots, w_m\}$
- Each topic c has probability Pr(c)
- Each word  $w_i \in V$  has conditional probability  $Pr(w_i \mid c_j)$ , for  $c_j \in C$ 
  - Note that  $\sum_{i=1}^{m} Pr(w_i \mid c_j) = 1$
- Assume document length is independent of the class

12 / 14

- Each document is a multiset or bag of words over a vocabulary  $V = \{w_1, w_2, \dots, w_m\}$
- Each topic c has probability Pr(c)
- Each word  $w_i \in V$  has conditional probability  $Pr(w_i \mid c_j)$ , for  $c_j \in C$ 
  - Note that  $\sum_{i=1}^{m} Pr(w_i \mid c_j) = 1$
- Assume document length is independent of the class
- Only a small subset of documents is labelled
  - Use this subset for initial estimate of  $P_r(c)$ ,  $P_r(w_i \mid c_j)$

12 / 14

■ Current model Pr(c),  $Pr(w_i | c_j)$ 

- Current model Pr(c),  $Pr(w_i | c_j)$
- Compute  $Pr(c_i \mid d)$  for each unlabelled document d
  - Normally we assign the maximum among these as the class for d
  - Here we keep fractional values

13 / 14

- Current model Pr(c),  $Pr(w_i \mid c_i)$
- Compute  $Pr(c_i \mid d)$  for each unlabelled document d
  - Normally we assign the maximum among these as the class for d
  - Here we keep fractional values
- Recompute  $Pr(c_j) = \frac{\sum_{d \in D} Pr(c_j \mid D)}{|D|}$ 
  - For labelled d,  $Pr(c_i \mid d) \in \{0, 1\}$

Pr(w/c)



■ For unlabelled d,  $Pr(c_i \mid d)$  is fractional value computed from current parameters

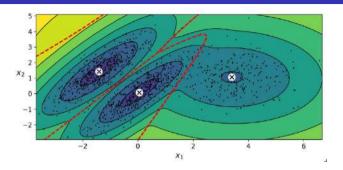
13 / 14

- Current model Pr(c),  $Pr(w_i | c_j)$
- Compute  $Pr(c_j \mid d)$  for each unlabelled document d
  - Normally we assign the maximum among these as the class for d
  - Here we keep fractional values
- Recompute  $Pr(c_j) = \frac{\sum_{d \in D} Pr(c_j \mid D)}{|D|}$ 
  - For labelled d,  $Pr(c_j \mid d) \in \{0, 1\}$
  - For unlabelled d,  $Pr(c_i \mid d)$  is fractional value computed from current parameters
- Recompute  $Pr(w_i \mid c_j)$  fraction of occurrences of  $w_i$  in documents labelled  $c_j$ 
  - $n_{id}$  occurrences of  $w_i$  in d
  - $Pr(w_i \mid c_j) = \frac{\sum_{d \in D} n_{id} Pr(c_j \mid d)}{\sum_{t=1}^{m} \sum_{d \in D} n_{td} Pr(c_j \mid d)}$

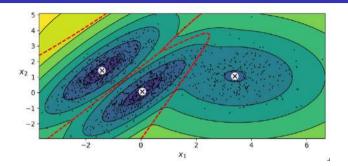


13 / 14

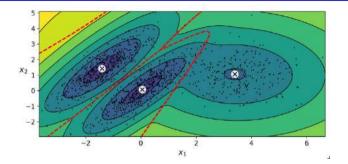
 Data points from a mixture of Gaussian distributions



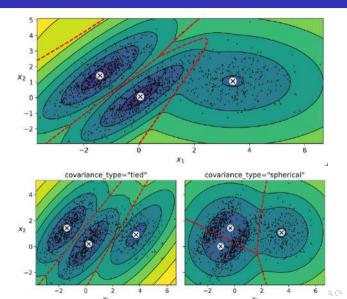
- Data points from a mixture of Gaussian distributions
- Use EM to estimate the parameters of each Gaussian distribution



- Data points from a mixture of Gaussian distributions
- Use EM to estimate the parameters of each Gaussian distribution
- Assign each point to "best"
  Gaussian



- Data points from a mixture of Gaussian distributions
- Use EM to estimate the parameters of each Gaussian distribution
- Assign each point to "best"
  Gaussian
- Can tweak the shape of the clusters by constraining the covariance matrix



- Data points from a mixture of Gaussian distributions
- Use EM to estimate the parameters of each Gaussian distribution
- Assign each point to "best"
  Gaussian
- Can tweak the shape of the clusters by constraining the covariance matrix
- Outliers are those that are outside  $k\sigma$  for all the Gaussians

