

Latent variable models

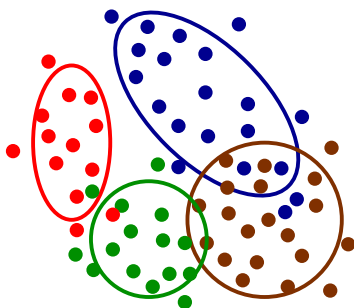
1 / 28

Mixture models

Unsupervised classification

Unsupervised classification

- **Input:** $x_1, x_2, \dots, x_n \in \mathbb{R}^d$, number of classes $K \in \mathbb{N}$.
- **Output:** function $f: \mathbb{R}^d \rightarrow \{1, 2, \dots, K\}$.
(Notation: $[K] := \{1, 2, \dots, K\}$.)
- **Typical semantics:** hidden subpopulation structure.



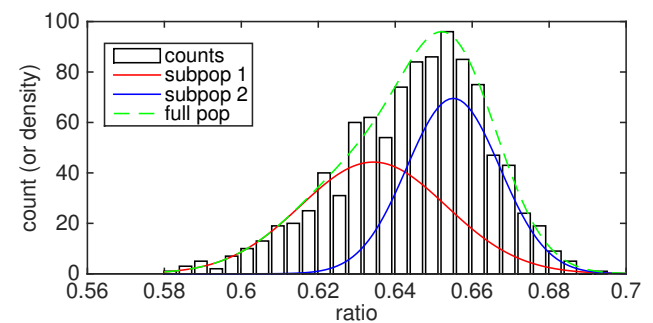
3 / 28

Example: Pearson's crabs (1894)

Data: ratio of forehead-width to body-length for 1000 crabs.



Maybe the sample is comprised of two different sub-species of crab?



4 / 28

Gaussian mixture model

Gaussian mixture model: statistical model $\mathcal{P} = \{P_{\theta} : \theta \in \Theta\}$ on $\mathcal{X} \times [K]$, where $(\mathbf{X}, Y) \sim P_{\theta}$ means that

$Y \sim (\pi_1, \dots, \pi_K)$ (discrete distribution over $[K]$; $P_{\theta}(Y = j) = \pi_j$)

$\mathbf{X} \mid Y = j \sim \mathcal{N}(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$ (Gaussian with mean $\boldsymbol{\mu}_j$ and covariance $\boldsymbol{\Sigma}_j$)

Parameter space Θ comprises all $\theta = (\pi_1, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \dots, \pi_K, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K)$ where $\boldsymbol{\mu}_j \in \mathbb{R}^d$, $\boldsymbol{\Sigma}_j \succeq \mathbf{0}$ (positive definite $d \times d$ matrix), $\pi_j \in [0, 1]$, and $\sum_{j=1}^K \pi_j = 1$.

Looks familiar?

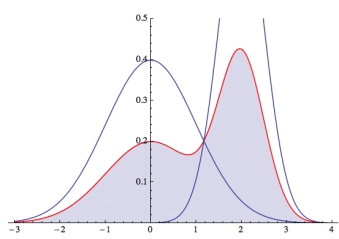
Even though this model is for data $(\mathbf{x}, y) \in \mathbb{R}^d \times [K]$, we declare only the \mathbf{x} part to be **observable**, and declare the y part to be **hidden** (or **latent**).

Models of this sort are called **mixture models**; this one in particular is called the **Gaussian mixture model**.

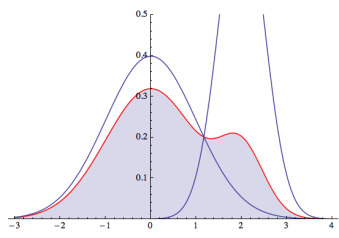
$$p_{\theta}(\mathbf{x}) = \sum_{j=1}^K \pi_j \cdot (2\pi)^{-d/2} \sqrt{\det(\boldsymbol{\Sigma}_j^{-1})} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_j)^{\top} \boldsymbol{\Sigma}_j^{-1}(\mathbf{x} - \boldsymbol{\mu}_j)\right)$$

Mixing weights π ; mixture components $\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1), \dots, \mathcal{N}(\boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K)$.

Gaussian mixtures in \mathbb{R}^1

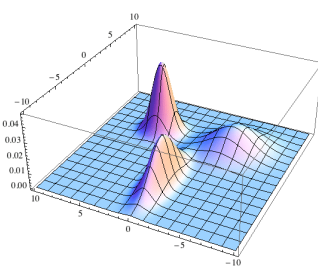


$$\frac{1}{2} \mathcal{N}(0, 1) + \frac{1}{2} \mathcal{N}(2, 1/4)$$

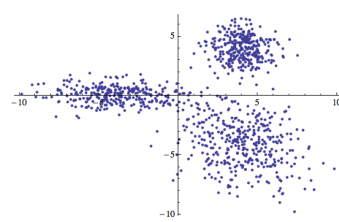


$$\frac{4}{5} \mathcal{N}(0, 1) + \frac{1}{5} \mathcal{N}(2, 1/4)$$

Gaussian mixtures in \mathbb{R}^2



Plot of the mixture density.



An iid sample of size 1000.

Soft assignments

Suppose you have the parameters $\theta = (\pi_1, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \dots, \pi_K, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K)$ of a Gaussian mixture distribution, and further that $(\mathbf{X}, Y) \sim P_{\theta}$.

Assignment variables $\Phi = (\Phi_1, \Phi_2, \dots, \Phi_K) \in \{0, 1\}^K$ (as in K -means):

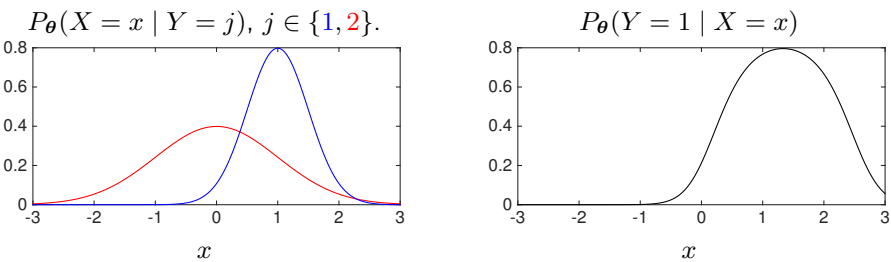
$$\Phi_j := \mathbb{1}\{Y = j\}.$$

Soft assignment of a data point $\mathbf{x} \in \mathbb{R}^d$ to component $j \in [K]$:

$$\begin{aligned} \mathbb{E}_{\theta}[\Phi_j \mid \mathbf{X} = \mathbf{x}] &= P_{\theta}(Y = j \mid \mathbf{X} = \mathbf{x}) \\ &= \frac{P_{\theta}(Y = j) \cdot p_{\theta}(\mathbf{x} \mid Y = j)}{p_{\theta}(\mathbf{x})} \\ &= \frac{\pi_j \cdot \sqrt{\det(\boldsymbol{\Sigma}_j^{-1})} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_j)^{\top} \boldsymbol{\Sigma}_j^{-1}(\mathbf{x} - \boldsymbol{\mu}_j)\right)}{\sum_{j'=1}^K \pi_{j'} \cdot \sqrt{\det(\boldsymbol{\Sigma}_{j'}^{-1})} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_{j'})^{\top} \boldsymbol{\Sigma}_{j'}^{-1}(\mathbf{x} - \boldsymbol{\mu}_{j'})\right)}. \end{aligned}$$

Soft clustering

Example: a Gaussian mixture distribution with $k = 2$ in \mathbb{R}^1 .



$$P_{\theta}(Y = 1 | X = x) = \frac{\pi_1 \cdot \frac{1}{\sigma_1} \exp\left(-\frac{(x - \mu_1)^2}{2\sigma_1^2}\right)}{\pi_1 \cdot \frac{1}{\sigma_1} \exp\left(-\frac{(x - \mu_1)^2}{2\sigma_1^2}\right) + \pi_2 \cdot \frac{1}{\sigma_2} \exp\left(-\frac{(x - \mu_2)^2}{2\sigma_2^2}\right)}.$$

Parameter estimation for Gaussian mixtures

Maximum likelihood estimation of $\theta = (\pi_1, \mu_1, \Sigma_1, \dots, \pi_K, \mu_K, \Sigma_K)$ given data x_1, x_2, \dots, x_n (regarded as an i.i.d. sample).

$$\begin{aligned} \hat{\theta} &:= \arg \max_{\theta \in \Theta} \sum_{i=1}^n \ln p_{\theta}(x_i) \\ &= \arg \max_{\theta \in \Theta} \sum_{i=1}^n \ln \left\{ \sum_{j=1}^K \pi_j \cdot \sqrt{\det(\Sigma_j^{-1})} \exp\left(-\frac{1}{2}(x - \mu_j)^{\top} \Sigma_j^{-1} (x - \mu_j)\right) \right\} \end{aligned}$$

Argh! The $\ln \left\{ \sum_{j=1}^K \dots \right\}$ does not simplify nicely!

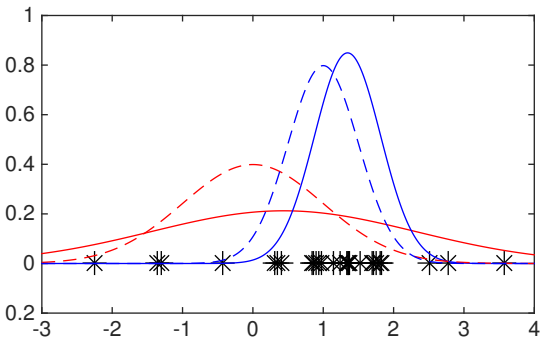
MLE for Gaussian mixtures: **not a convex optimization problem.**

$$\arg \min_{\theta \in \Theta} - \sum_{i=1}^n \ln \left\{ \sum_{j=1}^K \pi_j \cdot \sqrt{\det(\Sigma_j^{-1})} \exp\left(-\frac{1}{2}(x - \mu_j)^{\top} \Sigma_j^{-1} (x - \mu_j)\right) \right\}$$

Local optima not guaranteed to be global optima;
could be arbitrarily far from / worse than the MLE.

Local optimization

- ▶ For the purpose of modeling the density of X , a “good enough” local maximizer could be sufficient.
- ▶ If the data are actually generated by a Gaussian mixture distribution $P_{\theta_{\star}}$, then θ_{\star} may be close to some local maximizer of the likelihood.



Methods like gradient ascent would work, but there’s a much nicer local optimization method for this case: the E-M algorithm.

Expectation-Maximization for Gaussian mixtures

Motivating derivation

Suppose we had *softly labeled* data $\{(\mathbf{x}_i, \mathbf{w}_i)\}_{i=1}^n$ from $\mathbb{R}^d \times [0, 1]^K$.
(Each $\mathbf{w}_i = (w_{i,1}, w_{i,2}, \dots, w_{i,K})$ is a probability distribution on $[K]$.)

The “complete log-likelihood” of $\theta = (\pi_1, \mu_1, \Sigma_1, \dots, \pi_K, \mu_K, \Sigma_K)$ is

$$\sum_{i=1}^n \sum_{j=1}^K w_{i,j} \ln \left\{ \pi_j \cdot \sqrt{\det(\Sigma_j^{-1})} \exp \left(-\frac{1}{2} (\mathbf{x} - \mu_j)^\top \Sigma_j^{-1} (\mathbf{x} - \mu_j) \right) \right\} \\ = \sum_{i=1}^n \sum_{j=1}^K w_{i,j} \left(\ln \pi_j + \frac{1}{2} \ln \det(\Sigma_j^{-1}) - \frac{1}{2} (\mathbf{x} - \mu_j)^\top \Sigma_j^{-1} (\mathbf{x} - \mu_j) \right),$$

which can be easily maximized w.r.t. θ .

$$\hat{\pi}_j := \frac{1}{n} \sum_{i=1}^n w_{i,j}$$

$$\hat{\mu}_j := \frac{1}{n \hat{\pi}_j} \sum_{i=1}^n w_{i,j} \mathbf{x}_i$$

$$\hat{\Sigma}_j := \frac{1}{n \hat{\pi}_j} \sum_{i=1}^n w_{i,j} (\mathbf{x}_i - \hat{\mu}_j)(\mathbf{x}_i - \hat{\mu}_j)^\top.$$

13 / 28

Log-likelihood vs. complete log-likelihood

Three different functions of θ :

1. **Log-likelihood function**, given observed data:

$$\mathcal{L}(\theta) = \ln P_\theta(\text{observed data}).$$

2. **Complete log-likelihood function**, given observed **and unobserved** data:

$$\mathcal{L}_c(\theta) = \ln P_\theta(\text{observed data, unobserved data}).$$

- ★ Suppose we have some initial guess of parameters $\hat{\theta}$.
Treat **unobserved data** as random variables.

Conditional distribution of unobserved data given observed data $P_{\hat{\theta}}$:

$$P_{\hat{\theta}}(\text{unobserved data} \mid \text{observed data}).$$

(E.g., distribution of “soft labels” given the observed data.)

3. **Expected complete log-likelihood function** given observed data:

$$\mathbb{E}_{\hat{\theta}}[\mathcal{L}_c(\theta) \mid \text{observed data}] = \mathbb{E}_{\hat{\theta}}[\ln P_\theta(\text{observed \& unobserved data}) \mid \text{observed data}]$$

(Expectation $\mathbb{E}_{\hat{\theta}}[\dots]$ is with respect to $P_{\hat{\theta}}$ given observed data.)

14 / 28

Expected complete log-likelihood

Given initial guess of parameters $\hat{\theta}$,

$$w_{i,j} := \mathbb{E}_{\hat{\theta}}[\Phi_{i,j} \mid \mathbf{X}_i = \mathbf{x}_i \forall i \in [n]]$$

can be interpreted as predicted “soft labels”.

Using these $w_{i,j}$ to form expected complete log-likelihood function:

$$\mathbb{E}_{\hat{\theta}} \left[\sum_{i=1}^n \sum_{j=1}^K \Phi_{i,j} \left(\ln \pi_j + \frac{1}{2} \ln \det(\Sigma_j^{-1}) - \frac{1}{2} (\mathbf{x} - \mu_j)^\top \Sigma_j^{-1} (\mathbf{x} - \mu_j) \right) \mid \mathbf{X}_i = \mathbf{x}_i \forall i \in [n] \right] \\ = \sum_{i=1}^n \sum_{j=1}^K w_{i,j} \left(\ln \pi_j + \frac{1}{2} \ln \det(\Sigma_j^{-1}) - \frac{1}{2} (\mathbf{x} - \mu_j)^\top \Sigma_j^{-1} (\mathbf{x} - \mu_j) \right).$$

This function of θ is easy to maximize!

(But why do we care about the expected complete log-likelihood?)

15 / 28

Expectation-Maximization (E-M)

E-M algorithm for Gaussian mixtures

Initialize $\hat{\theta} = (\hat{\pi}_1, \hat{\mu}_1, \hat{\Sigma}_1, \dots, \hat{\pi}_K, \hat{\mu}_K, \hat{\Sigma}_K)$ somehow. Then repeat:

1. **E step:** expectation of “hidden variables” w.r.t. $P_{\hat{\theta}}$ conditioned on data.
For each $i \in \{1, 2, \dots, n\}$ and $j \in \{1, 2, \dots, K\}$,

$$w_{i,j} := \frac{\hat{\pi}_j \cdot \sqrt{\det(\hat{\Sigma}_j^{-1})} \exp \left(-\frac{1}{2} (\mathbf{x} - \hat{\mu}_j)^\top \hat{\Sigma}_j^{-1} (\mathbf{x} - \hat{\mu}_j) \right)}{\sum_{j'=1}^K \hat{\pi}_{j'} \cdot \sqrt{\det(\hat{\Sigma}_{j'}^{-1})} \exp \left(-\frac{1}{2} (\mathbf{x} - \hat{\mu}_{j'})^\top \hat{\Sigma}_{j'}^{-1} (\mathbf{x} - \hat{\mu}_{j'}) \right)}$$

2. **M step:** maximize “expected complete log-likelihood” w.r.t. parameters.
For each $j \in \{1, 2, \dots, K\}$,

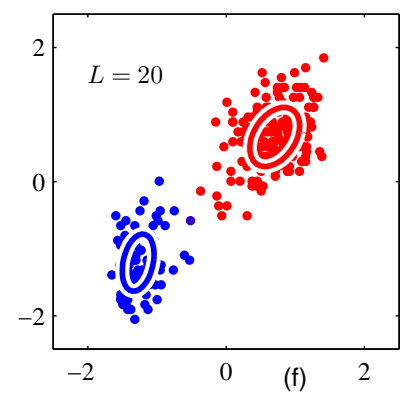
$$\hat{\pi}_j := \frac{1}{n} \sum_{i=1}^n w_{i,j}$$

$$\hat{\mu}_j := \frac{1}{n \hat{\pi}_j} \sum_{i=1}^n w_{i,j} \mathbf{x}_i$$

$$\hat{\Sigma}_j := \frac{1}{n \hat{\pi}_j} \sum_{i=1}^n w_{i,j} (\mathbf{x}_i - \hat{\mu}_j)(\mathbf{x}_i - \hat{\mu}_j)^\top.$$

16 / 28

Sample run of the E-M algorithm



After 20 rounds of E-M.

Using the E-M algorithm

E-M for Gaussian mixtures

- 1. **E step:** For each $i \in [n], j \in [K]$,

$$w_{i,j} \propto \hat{\pi}_j \cdot p_{\hat{\mu}_j, \hat{\Sigma}_j}(\mathbf{x}_i)$$

where $p_{\mu, \Sigma}$ is the $N(\mu, \Sigma)$ pdf.

- 2. **M step:** For each $j \in [K]$,

$$\hat{\pi}_j := \frac{1}{n} \sum_{i=1}^n w_{i,j}$$

$$\hat{\mu}_j := \frac{1}{n\hat{\pi}_j} \sum_{i=1}^n w_{i,j} \mathbf{x}_i$$

$$\hat{\Sigma}_j := \frac{1}{n\hat{\pi}_j} \sum_{i=1}^n w_{i,j} (\mathbf{x}_i - \hat{\mu}_j)(\mathbf{x}_i - \hat{\mu}_j)^\top.$$

Some details

- ▶ **No step sizes to tune!**
- ▶ **Initialization:** a bit of an art; both D^2 -sampling and Lloyd's algorithm are reasonable.
- ▶ **Starved clusters:** problems can occur if $\hat{\pi}_j$ becomes too small (e.g., $\hat{\Sigma}_j$ could be near singular).
Remove/replace such components.
- ▶ Log-likelihood of E-M iterates is **non-decreasing**; converges to a **stationary point**.
 \therefore Run E-M from many random initializations; pick the result with highest likelihood.

Derivation of E-M

E-M algorithm (Dempster, Laird, and Rubin, 1977)

E-M is a general algorithmic template for climbing log-likelihood objectives of models with **latent variables** (e.g., cluster assignments).

- ▶ What is the role of the expected complete log-likelihood?
- ▶ Why do parameters produced E-M iterations

$$\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(t)}$$

have increasing (or at least non-decreasing) log-likelihood

$$\mathcal{L}(\theta^{(1)}) \leq \mathcal{L}(\theta^{(2)}) \leq \dots \leq \mathcal{L}(\theta^{(t)})?$$

Statistical model $\mathcal{P} = \{P_{\theta} : \theta \in \Theta\}$, where each $P_{\theta} \in \mathcal{P}$ specifies distribution of observed variables \mathbf{X} and latent variables \mathbf{Y} .
(For simplicity, assume they are discrete random variables.)

Here, if we want to consider models for iid samples, then we shall take \mathcal{P} to be an appropriate “lifted” model (in the n -fold product form $\mathcal{P} = \mathcal{P}_0^n$ for a base model \mathcal{P}_0).

- So, if \mathcal{P} is a “lifted” model, then \mathbf{X} and \mathbf{Y} are observed and latent variables for *all* n data points in the sample.

Initialize parameters $\theta^{(1)} \in \Theta$ somehow.
For $t = 1, 2, \dots$:

- **E step:** Construct expected complete log-likelihood function

$$\theta \mapsto \mathbb{E}_{\theta^{(t)}}[\mathcal{L}_c(\theta) \mid \mathbf{X} = \mathbf{x}]$$

where expectation (of latent variables) is with respect to distribution $P_{\theta^{(t)}}$ given observed data $\mathbf{X} = \mathbf{x}$.

- **M step:** Choose $\theta^{(t+1)}$ to maximize that function,

$$\theta^{(t+1)} := \arg \max_{\theta \in \Theta} \mathbb{E}_{\theta^{(t)}}[\mathcal{L}_c(\theta) \mid \mathbf{X} = \mathbf{x}].$$

- Let \mathbf{x} be the (observed) data.
- Let $q^{(t)}$ denote distribution of \mathbf{Y} given $\mathbf{X} = \mathbf{x}$ when $(\mathbf{X}, \mathbf{Y}) \sim P_{\theta^{(t)}}$:

$$q^{(t)}(\mathbf{y}) := \frac{P_{\theta^{(t)}}(\mathbf{x}, \mathbf{y})}{\sum_{\mathbf{y}' \in \mathcal{Y}} P_{\theta^{(t)}}(\mathbf{x}, \mathbf{y}')}.$$

- Log-likelihood function \mathcal{L} :

$$\mathcal{L}(\theta) = \ln \left(\sum_{\mathbf{y} \in \mathcal{Y}} q^{(t)}(\mathbf{y}) \cdot \frac{P_{\theta}(\mathbf{x}, \mathbf{y})}{q^{(t)}(\mathbf{y})} \right) = \ln \left(\mathbb{E}_{\mathbf{Y} \sim q^{(t)}} \left[\frac{P_{\theta}(\mathbf{x}, \mathbf{Y})}{q^{(t)}(\mathbf{Y})} \right] \right).$$

- Define **log-likelihood lower-bound** function $\mathcal{L}_{\text{LB}}^{(t)}$ by

$$\mathcal{L}_{\text{LB}}^{(t)}(\theta) := \mathbb{E}_{\mathbf{Y} \sim q^{(t)}} \left[\ln \left(\frac{P_{\theta}(\mathbf{x}, \mathbf{Y})}{q^{(t)}(\mathbf{Y})} \right) \right].$$

(Here, $\mathbb{E}_{\mathbf{Y} \sim q^{(t)}}[\dots]$ means expectation conditioned on $\mathbf{X} = \mathbf{x}$ under $q^{(t)}$.)

Log-likelihood lower-bound function:

$$\begin{aligned} \mathcal{L}_{\text{LB}}^{(t)}(\theta) &:= \mathbb{E}_{\mathbf{Y} \sim q^{(t)}} \left[\ln \left(\frac{P_{\theta}(\mathbf{x}, \mathbf{Y})}{q^{(t)}(\mathbf{Y})} \right) \right] \\ &= \underbrace{\mathbb{E}_{\mathbf{Y} \sim q^{(t)}} [\ln P_{\theta}(\mathbf{x}, \mathbf{Y})]}_{\substack{\text{Expected complete log-likelihood} \\ \mathbb{E}_{\theta^{(t)}}[\mathcal{L}_c(\theta) \mid \text{observed data}]} } + (\text{stuff not depending on } \theta). \end{aligned}$$

Therefore, maximizing $\mathcal{L}_{\text{LB}}^{(t)}$ is the same as maximizing **expected complete log-likelihood**

$$\theta \mapsto \mathbb{E}_{\theta^{(t)}}[\mathcal{L}_c(\theta) \mid \text{observed data}].$$

Lower-bound property

Jensen's inequality: for any *concave* function g and random variable Z ,

$$g(\mathbb{E}[Z]) \geq \mathbb{E}[g(Z)].$$

Since natural logarithm is concave,

$$\mathcal{L}(\theta) = \ln \left(\mathbb{E}_{\mathbf{Y} \sim q^{(t)}} \left[\frac{P_{\theta}(\mathbf{x}, \mathbf{Y})}{q^{(t)}(\mathbf{Y})} \right] \right) \geq \mathbb{E}_{\mathbf{Y} \sim q^{(t)}} \left[\ln \left(\frac{P_{\theta}(\mathbf{x}, \mathbf{Y})}{q^{(t)}(\mathbf{Y})} \right) \right] = \mathcal{L}_{\text{LB}}^{(t)}(\theta).$$

Moreover, we have $\mathcal{L}(\theta^{(t)}) = \mathcal{L}_{\text{LB}}^{(t)}(\theta^{(t)})$:

$$\begin{aligned} \mathbb{E}_{\mathbf{Y} \sim q^{(t)}} \left[\ln \left(\frac{P_{\theta^{(t)}}(\mathbf{x}, \mathbf{Y})}{q^{(t)}(\mathbf{Y})} \right) \right] &= \mathbb{E}_{\mathbf{Y} \sim q^{(t)}} \left[\ln \left(\sum_{\mathbf{y} \in \mathcal{Y}} P_{\theta^{(t)}}(\mathbf{x}, \mathbf{y}) \right) \right] \\ &= \ln \left(\sum_{\mathbf{y} \in \mathcal{Y}} P_{\theta^{(t)}}(\mathbf{x}, \mathbf{y}) \right) = \mathcal{L}(\theta^{(t)}). \end{aligned}$$

E-M in terms of log-likelihood lower-bound

In t -th iteration of E-M:

- **E step:** Construct log-likelihood lower-bound function $\mathcal{L}_{\text{LB}}^{(t)}$ via distribution $q^{(t)}$,

$$q^{(t)}(\mathbf{y}) := \frac{P_{\theta^{(t)}}(\mathbf{x}, \mathbf{y})}{\sum_{\mathbf{y}' \in \mathcal{Y}} P_{\theta^{(t)}}(\mathbf{x}, \mathbf{y}')},$$

so that lower-bound is tight at $\theta^{(t)}$.

- **M step:** Choose $\theta^{(t+1)}$ to maximize $\mathcal{L}_{\text{LB}}^{(t)}$:

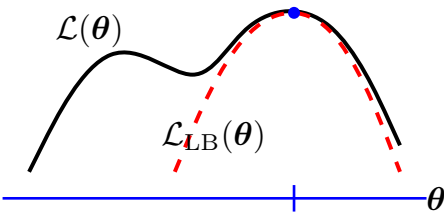
$$\theta^{(t+1)} := \arg \max_{\theta \in \Theta} \mathcal{L}_{\text{LB}}^{(t)}(\theta).$$

Theorem

In t -th iteration of E-M, we have

$$\mathcal{L}(\theta^{(t)}) = \mathcal{L}_{\text{LB}}^{(t)}(\theta^{(t)}) \leq \mathcal{L}_{\text{LB}}^{(t)}(\theta^{(t+1)}) \leq \mathcal{L}(\theta^{(t+1)}).$$

Constructing and maximizing \mathcal{L}_{LB}



M step: choose $\hat{\theta}$ to maximize \mathcal{L}_{LB} .

Key takeaways

1. **Mixture models:** similar to generative models for classification, except class labels are not observed.
 - Important example: **Gaussian mixture models**.
2. **E-M algorithm for Gaussian mixture models.**
3. **Recipe to derive E-M algorithm** for a general latent variable model:
 - **E step:** use conditional distribution of latent variables (given observed variables) under $P_{\theta^{(t)}}$ to form “expected complete log-likelihood” function.
 - **M step:** maximize the “expected complete log-likelihood” function.
 - (We’ll see more examples next time.)
4. **Key properties of E-M:**
 - “Expected complete log-likelihood” function is **lower-bound on log-likelihood function** that is tight at $\theta^{(t)}$.
 - Log-likelihoods of E-M iterates are **non-decreasing**.

More latent variable models

- ▶ Statistical model $\mathcal{P} = \{P_{\theta} : \theta \in \Theta\}$ with both *observed variables* \mathbf{X} and *latent variables* \mathbf{Z} .
- ▶ Would like to use MLE for parameter estimation, but this is often computationally difficult.
- ▶ Typically resort to using E-M to find parameters with high likelihood (though maybe not the highest possible).

(Actually, there are many options besides MLE ...)

1 / 22

2 / 22

Log-likelihood vs. complete log-likelihood

1. **Log-likelihood function**, given observed data \mathbf{x} :

$$\mathcal{L}(\theta) = \ln \left(\sum_{\mathbf{z} \in \mathcal{Z}} P_{\theta}(\mathbf{x}, \mathbf{z}) \right).$$

(If \mathbf{Z} is continuous r.v., replace sum with integral ...)

2. **Complete log-likelihood function**, given observed data \mathbf{x} and unobserved data \mathbf{z} :

$$\mathcal{L}_c(\theta) = \ln P_{\theta}(\mathbf{x}, \mathbf{z}).$$

3. **Expected complete log-likelihood function** given observed data \mathbf{x} , treating unobserved data as random \mathbf{Z} :

$$\mathbb{E}_{\hat{\theta}}[\mathcal{L}_c(\theta)] = \mathbb{E}_{\hat{\theta}}[\ln P_{\theta}(\mathbf{x}, \mathbf{Z}) \mid \mathbf{X} = \mathbf{x}].$$

(Expectation $\mathbb{E}_{\hat{\theta}}[\dots]$ is with respect to $P_{\hat{\theta}}$ given $\mathbf{X} = \mathbf{x}$.)

3 / 22

Expectation-Maximization

Initialize parameters $\theta^{(1)} \in \Theta$ somehow.

For $t = 1, 2, \dots$:

- ▶ **E step**: Construct expected complete log-likelihood function

$$\theta \mapsto \mathbb{E}_{\theta^{(t)}}[\mathcal{L}_c(\theta) \mid \mathbf{X} = \mathbf{x}]$$

where expectation (of latent variables) is with respect to distribution $P_{\theta^{(t)}}$ given observed data $\mathbf{X} = \mathbf{x}$.

- ▶ **M step**: Choose $\theta^{(t+1)}$ to maximize that function,

$$\theta^{(t+1)} := \arg \max_{\theta \in \Theta} \mathbb{E}_{\theta^{(t)}}[\mathcal{L}_c(\theta) \mid \mathbf{X} = \mathbf{x}].$$

4 / 22

- 1. Mixture models (last time)
- 2. Mechanical Turk model
- 3. Conditional mixture models
- 4. ...

Mechanical Turk model

Amazon Mechanical Turk

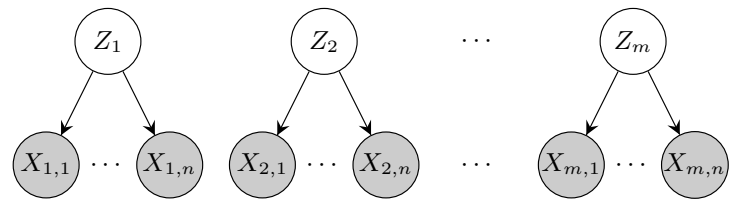
- Passonneau and Carpenter (TACL 2014):** asked Mechanical Turk workers (presumably humans) to label words by their “word sense”.
- **Items** (“human intelligence tasks”): English words with multiple possible meanings, as they appear in natural sentences.
 - **Labels:** the correct meanings in the given contexts.
 - Some workers are **adept** at this task (likely to give correct answer), some are **inept** (as good as random guessing), others are **malicious** (likely to give wrong answer).
 - For binary labels, can be difficult to distinguish **adept** from **malicious**, but can at least hope to ignore **inept** workers.
- If assume more **adept** workers than **malicious** workers, can use (weighted) majority vote over non-**inept** workers’ labels.
- Modeling worker accuracies → labeled data set with label “reliabilities”.

Mechanical Turk model (Dawid and Skene, 1979)

- **Observed:** predicted labels on m items from n workers $\{x_{i,j}\}_{i \in [m], j \in [n]}$.
- **Hidden:** correct labels $\{z_i\}_{i=1}^m$ for all m items.
- **Model:** $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$, where $\theta = (\pi_1, \pi_2, \dots, \pi_m, p_1, p_2, \dots, p_n)$.
 - Binary labels $\{0, 1\}$.
 - Data for items $\{(Z_i, X_{i,1}, X_{i,2}, \dots, X_{i,n})\}_{i=1}^m$ are independent.
 - Nature determines correct label for item i :
$$P_\theta(Z_i = 1) = 1 - P_\theta(Z_i = 0) = \pi_i.$$
 - Conditioned on Z_i , predicted labels $\{X_{i,j}\}_{j=1}^n$ from workers on item i are independent.
 - Worker j is correct with probability p_j :
$$P_\theta(X_{i,j} = z \mid Z_i = z) = p_j.$$
- **Notation:**
 - $\mathbf{Z} = (Z_1, \dots, Z_m)$ (all correct labels);
 - $\mathbf{X} = (X_{1,1}, \dots, X_{m,n})$ (all predicted labels);
 - $\mathbf{X}_i = (X_{i,1}, \dots, X_{i,n})$ (predicted labels for item i).

Aside: graphical models

(Directed) graphical model: use (directed) graph structure over random variables to represent conditional independence assumptions.



Semantics

- ▶ Observed variables shaded, latent variables not shaded.
- ▶ **Disjoint connected components are independent** of each other.
E.g., $(Z_1, \mathbf{X}_1) \perp\!\!\!\perp (Z_2, \mathbf{X}_2) \perp\!\!\!\perp \dots \perp\!\!\!\perp (Z_m, \mathbf{X}_m)$.
- ▶ If connected component is **directed tree** (i.e., every node has ≤ 1 parent), **conditioning on r.v.** is akin to **removing it from the graph**.
E.g., $X_{i,1} \perp\!\!\!\perp X_{i,2} \perp\!\!\!\perp \dots \perp\!\!\!\perp X_{i,n} \mid Z_i = z$.
- ▶ Other rules more involved (e.g., for vertices with multiple parents).

Complete log-likelihood

Complete log-likelihood of $\theta = (\pi_1, \pi_2, \dots, \pi_m, p_1, p_2, \dots, p_n)$ given unobserved data \mathbf{z} and observed data \mathbf{x} :

(Use *conditional independence* properties to *factor* the likelihood.)

$$\begin{aligned} \mathcal{L}_c(\theta) &= \ln \prod_{i=1}^m P_{\theta}(Z_i = z_i) \prod_{j=1}^n P_{\theta}(X_{i,j} = x_{i,j} \mid Z_i = z_i) \\ &= \ln \prod_{i=1}^m \pi_i^{z_i} (1 - \pi_i)^{1-z_i} \prod_{j=1}^n p_j^{(1-x_{i,j})(1-z_i) + x_{i,j}z_i} (1 - p_j)^{(1-x_{i,j})z_i + x_{i,j}(1-z_i)} \\ &= \sum_{i=1}^m [z_i \ln \pi_i + (1 - z_i) \ln(1 - \pi_i)] \\ &\quad + \sum_{i=1}^m \sum_{j=1}^n [(1 - x_{i,j})(1 - z_i) + x_{i,j}z_i] \ln p_j \\ &\quad + \sum_{i=1}^m \sum_{j=1}^n [(1 - x_{i,j})z_i + x_{i,j}(1 - z_i)] \ln(1 - p_j) \end{aligned}$$

Expected complete log-likelihood

Assume we have some initial parameters $\hat{\theta} = (\hat{\pi}_1, \hat{\pi}_2, \dots, \hat{\pi}_m, \hat{p}_1, \hat{p}_2, \dots, \hat{p}_n)$.
Expected complete log-likelihood of θ w.r.t. $P_{\hat{\theta}}$ conditioned on $\mathbf{X} = \mathbf{x}$:

$$\begin{aligned} \mathbb{E}_{\hat{\theta}}[\mathcal{L}_c(\theta) \mid \mathbf{X} = \mathbf{x}] &= \sum_{i=1}^m [w_i \ln \pi_i + (1 - w_i) \ln(1 - \pi_i)] \\ &\quad + \sum_{i=1}^m \sum_{j=1}^n [(1 - x_{i,j})(1 - w_i) + x_{i,j}w_i] \ln p_j \\ &\quad + \sum_{i=1}^m \sum_{j=1}^n [(1 - x_{i,j})w_i + x_{i,j}(1 - w_i)] \ln(1 - p_j) \end{aligned}$$

where

$$\begin{aligned} w_i &:= P_{\hat{\theta}}(Z_i = 1 \mid \mathbf{X}_i = \mathbf{x}_i) = \frac{P_{\hat{\theta}}(Z_i = 1, \mathbf{X}_i = \mathbf{x}_i)}{P_{\hat{\theta}}(Z_i = 1, \mathbf{X}_i = \mathbf{x}_i) + P_{\hat{\theta}}(Z_i = 0, \mathbf{X}_i = \mathbf{x}_i)} \\ &= \frac{\hat{\pi}_i \prod_{j=1}^n \hat{p}_j^{x_{i,j}} (1 - \hat{p}_j)^{1-x_{i,j}}}{\hat{\pi}_i \prod_{j=1}^n \hat{p}_j^{x_{i,j}} (1 - \hat{p}_j)^{1-x_{i,j}} + (1 - \hat{\pi}_i) \prod_{j=1}^n \hat{p}_j^{1-x_{i,j}} (1 - \hat{p}_j)^{x_{i,j}}} \end{aligned}$$

Interpreting w_i

$$\frac{w_i}{1 - w_i} = \frac{P_{\hat{\theta}}(Z_i = 1 \mid \mathbf{X}_i = \mathbf{x}_i)}{P_{\hat{\theta}}(Z_i = 0 \mid \mathbf{X}_i = \mathbf{x}_i)} = \frac{\hat{\pi}_i}{1 - \hat{\pi}_i} \prod_{j=1}^n \left(\frac{\hat{p}_j}{1 - \hat{p}_j} \right)^{2x_{i,j} - 1}$$

- ▶ Start with $\frac{\hat{\pi}_i}{1 - \hat{\pi}_i}$: “prior” odds ratio.
- ▶ **Perfectly inept workers** $\hat{p}_j = 1/2$: **contributes nothing**.
- ▶ **Adept workers** $\hat{p}_j > 1/2$:
 - ▶ If $x_{i,j} = 1$: **increase ratio**.
 - ▶ If $x_{i,j} = 0$: **decrease ratio**.
- ▶ **Malicious workers** $\hat{p}_j < 1/2$:
 - ▶ If $x_{i,j} = 1$: **decrease ratio**.
 - ▶ If $x_{i,j} = 0$: **increase ratio**.

Function $\theta \mapsto \mathbb{E}_{\hat{\theta}}[\mathcal{L}_c(\theta) \mid \mathbf{X} = \mathbf{x}]$ is concave w.r.t. θ ;

$$\begin{aligned} \frac{\partial}{\partial \pi_i} \mathbb{E}_{\hat{\theta}}[\mathcal{L}_c(\theta) \mid \mathbf{X} = \mathbf{x}] &= \frac{w_i}{\pi_i} - \frac{1 - w_i}{1 - \pi_i}, \\ \frac{\partial}{\partial p_j} \mathbb{E}_{\hat{\theta}}[\mathcal{L}_c(\theta) \mid \mathbf{X} = \mathbf{x}] &= \frac{\sum_{i=1}^m (1 - x_{i,j})(1 - w_i) + x_{i,j} w_i}{p_j} \\ &\quad - \frac{\sum_{i=1}^m (1 - x_{i,j}) w_i + x_{i,j} (1 - w_i)}{1 - p_j}. \end{aligned}$$

Partial derivatives are zero when

$$\begin{aligned} \pi_i &= w_i, \\ p_j &= \frac{1}{m} \sum_{i=1}^m \left\{ w_i x_{i,j} + (1 - w_i)(1 - x_{i,j}) \right\}. \end{aligned}$$

Maximizers of expected complete log-likelihood function:

$$\begin{aligned} \pi_i &= w_i, \\ p_j &= \frac{1}{m} \sum_{i=1}^m \left\{ w_i x_{i,j} + (1 - w_i)(1 - x_{i,j}) \right\}. \end{aligned}$$

- ▶ $\pi_i = w_i$: of course.
- ▶ Now pretend w_i are true “soft” labels of items.
- ▶ p_j : “soft” accuracy of worker j based on predicted labels across all items.

Initialize $\hat{\theta} = (\hat{\pi}_1, \hat{\pi}_2, \dots, \hat{\pi}_m, \hat{p}_1, \hat{p}_2, \dots, \hat{p}_n)$.

Repeat:

▶ **E step:**

$$w_i := \frac{\hat{\pi}_i \prod_{j=1}^n \hat{p}_j^{x_{i,j}} (1 - \hat{p}_j)^{1-x_{i,j}}}{\hat{\pi}_i \prod_{j=1}^n \hat{p}_j^{x_{i,j}} (1 - \hat{p}_j)^{1-x_{i,j}} + (1 - \hat{\pi}_i) \prod_{j=1}^n \hat{p}_j^{1-x_{i,j}} (1 - \hat{p}_j)^{x_{i,j}}}$$

for all $i = 1, 2, \dots, m$.

▶ **M step:**

$$\hat{\pi}_i := w_i \quad \text{for all } i = 1, 2, \dots, m,$$

$$\hat{p}_j := \frac{1}{m} \sum_{i=1}^m \left\{ w_i x_{i,j} + (1 - w_i)(1 - x_{i,j}) \right\} \quad \text{for all } j = 1, 2, \dots, n.$$

Can extend to handle multi-class labels, label-dependent worker abilities, etc.

Conditional mixture model

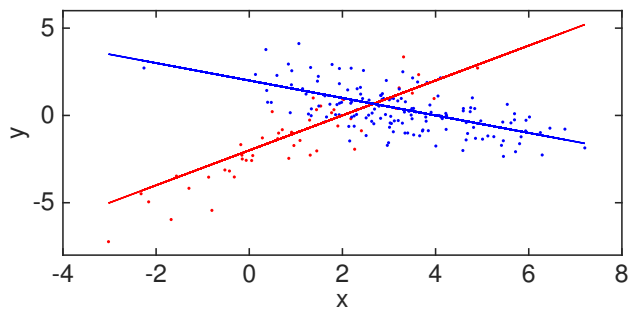
Conditional mixture model

Value of \mathbf{x} (probabilistically) determines which regression model y follows: e.g.,

$$y = \langle \beta_0, \mathbf{x} \rangle + \text{noise},$$

or

$$y = \langle \beta_1, \mathbf{x} \rangle + \text{noise}.$$



Mixture of two linear regressions

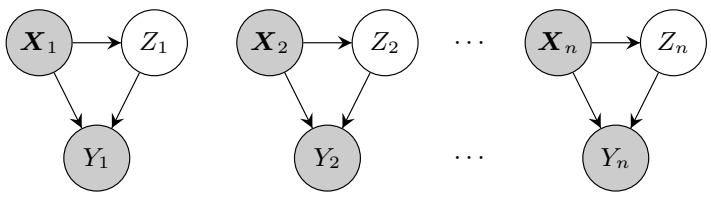
- **Observed:** labeled data $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ from $\mathbb{R}^d \times \mathbb{R}$.
- **Hidden:** hidden bits $\{z_i\}_{i=1}^n$ from $\{0, 1\}$.
- **Model:** $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$, where $\theta = (\alpha, \beta_0, \beta_1) \in \mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R}^d$.
 - $\{(\mathbf{X}_i, Y_i, Z_i)\}_{i=1}^n$ are iid.
 - Hidden bit for i -th data point depends on \mathbf{X}_i :

$$Z_i \mid \mathbf{X}_i = \mathbf{x}_i \sim \text{Bern}(\text{logistic}(\langle \alpha, \mathbf{x}_i \rangle))$$

where $\text{logistic}(t) = 1/(1 + e^{-t})$. (Note: $1 - \text{logistic}(t) = \text{logistic}(-t)$.)

- Hidden bit determines which regression coefficients to use:

$$Y_i \mid \mathbf{X}_i = \mathbf{x}_i, Z_i = z_i \sim N(\langle (1 - z_i)\beta_0 + z_i\beta_1, \mathbf{x}_i \rangle, 1).$$



Complete log-(conditional) likelihood

$$\begin{aligned} \mathcal{L}_c(\theta) &= \ln \prod_{i=1}^n P_\theta(Z_i = z_i \mid \mathbf{X}_i = \mathbf{x}_i) \cdot p_\theta(y_i \mid \mathbf{X}_i = \mathbf{x}_i, Z_i = z_i) \\ &= \sum_{i=1}^n \left(z_i \ln \text{logistic}(\langle \alpha, \mathbf{x}_i \rangle) + (1 - z_i) \ln \text{logistic}(-\langle \alpha, \mathbf{x}_i \rangle) \right. \\ &\quad \left. - \frac{1}{2} (y_i - \langle (1 - z_i)\beta_0 + z_i\beta_1, \mathbf{x}_i \rangle)^2 \right) + (\text{stuff not involving } \theta). \end{aligned}$$

Given $\hat{\theta} = (\hat{\alpha}, \hat{\beta}_0, \hat{\beta}_1)$,

$$\begin{aligned} w_i &:= \mathbb{E}_{\hat{\theta}}[Z_i \mid \mathbf{X}_i = \mathbf{x}_i, Y_i = y_i] = P_{\hat{\theta}}(Z_i = 1 \mid \mathbf{X}_i = \mathbf{x}_i, Y_i = y_i) \\ &= \frac{\text{logistic}(\langle \hat{\alpha}, \mathbf{x}_i \rangle) \cdot e^{-\frac{1}{2}(y_i - \langle \hat{\beta}_1, \mathbf{x}_i \rangle)^2}}{\text{logistic}(\langle \hat{\alpha}, \mathbf{x}_i \rangle) \cdot e^{-\frac{1}{2}(y_i - \langle \hat{\beta}_1, \mathbf{x}_i \rangle)^2} + \text{logistic}(-\langle \hat{\alpha}, \mathbf{x}_i \rangle) \cdot e^{-\frac{1}{2}(y_i - \langle \hat{\beta}_0, \mathbf{x}_i \rangle)^2}}. \end{aligned}$$

Maximizing $\mathbb{E}_{\hat{\theta}}[\mathcal{L}_c(\theta) \mid \mathbf{X} = \mathbf{x}, \mathbf{Y} = \mathbf{y}]$

$$\begin{aligned} \mathcal{L}_c(\theta) &= \sum_{i=1}^n \left(z_i \ln \text{logistic}(\langle \alpha, \mathbf{x}_i \rangle) + (1 - z_i) \ln \text{logistic}(-\langle \alpha, \mathbf{x}_i \rangle) \right. \\ &\quad \left. - \frac{1}{2} (y_i - \langle (1 - z_i)\beta_0 + z_i\beta_1, \mathbf{x}_i \rangle)^2 \right) + (\text{ignorable}). \end{aligned}$$

Replace z_i with r.v. Z_i , and take expectations:

$$\begin{aligned} &\mathbb{E}_{\hat{\theta}}[\mathcal{L}_c(\theta) \mid \mathbf{X} = \mathbf{x}, \mathbf{Y} = \mathbf{y}] \\ &= \sum_{i=1}^n \left(w_i \ln \text{logistic}(\langle \alpha, \mathbf{x}_i \rangle) + (1 - w_i) \ln \text{logistic}(-\langle \alpha, \mathbf{x}_i \rangle) \right. \\ &\quad \left. - \frac{1 - w_i}{2} (y_i - \langle \beta_0, \mathbf{x}_i \rangle)^2 - \frac{w_i}{2} (y_i - \langle \beta_1, \mathbf{x}_i \rangle)^2 \right) + (\text{ignorable}). \end{aligned}$$

- Maximizing w.r.t. α : **weighted logistic regression**.
- Maximizing w.r.t. β_0 : **weighted linear regression**.
- Maximizing w.r.t. β_1 : **weighted linear regression**.

Initialize $\hat{\theta} = (\hat{\alpha}, \hat{\beta}_0, \hat{\beta}_1)$.

Repeat:

► **E step:**

$$w_i := \frac{\text{logistic}(\langle \hat{\alpha}, \mathbf{x}_i \rangle) \cdot e^{-\frac{1}{2}(y_i - \langle \hat{\beta}_1, \mathbf{x}_i \rangle)^2}}{\text{logistic}(\langle \hat{\alpha}, \mathbf{x}_i \rangle) \cdot e^{-\frac{1}{2}(y_i - \langle \hat{\beta}_1, \mathbf{x}_i \rangle)^2} + \text{logistic}(-\langle \hat{\alpha}, \mathbf{x}_i \rangle) \cdot e^{-\frac{1}{2}(y_i - \langle \hat{\beta}_0, \mathbf{x}_i \rangle)^2}}$$

for all $i = 1, 2, \dots, n$.

► **M step:** Solve a weighted logistic regression problem, and two weighted least squares problems to get new $\hat{\theta} = (\hat{\alpha}, \hat{\beta}_0, \hat{\beta}_1)$.

1. More examples of E-M algorithm for latent variable models.
 - Need to mind conditional independence assumptions. (Graphical models can help with this.)
2. Sometimes E- and M-steps are not available in closed-form.