

# **Expectation Maximization**

**STATS 305C: Applied Statistics**

Scott Linderman

April 20, 2022

# Recall our Bayesian Mixture Model

1. Sample the proportions from a Dirichlet prior:

$$\underline{\pi} \sim \text{Dir}(\alpha) \quad (1)$$

*cluster mean*

2. Sample the parameters for each component:

$$\theta_k \stackrel{\text{iid}}{\sim} p(\theta \mid \phi, \nu) \quad \text{for } k = 1, \dots, K \quad (2)$$

3. Sample the assignment of each data point:

$$z_n \stackrel{\text{iid}}{\sim} \pi \quad \text{for } n = 1, \dots, N \quad (3)$$

4. Sample data points given their assignments:

$$\mathbf{x}_n \sim p(\mathbf{x} \mid \theta_{z_n}) \quad \text{for } n = 1, \dots, N \quad (4)$$

# Joint distribution

- This generative model corresponds to the following factorization of the joint distribution,

$$p(\pi, \{\theta_k\}_{k=1}^K, \{(z_n, \mathbf{x}_n)\}_{n=1}^N \mid \phi, \nu, \alpha) = p(\pi \mid \alpha) \prod_{k=1}^K p(\theta_k \mid \phi, \nu) \prod_{n=1}^N p(z_n \mid \pi) p(\mathbf{x}_n \mid z_n, \{\theta_k\}_{k=1}^K) \quad (5)$$

- Equivalently,

$$p(\pi, \{\theta_k\}_{k=1}^K, \{(z_n, \mathbf{x}_n)\}_{n=1}^N \mid \phi, \nu, \alpha) = p(\pi \mid \alpha) \prod_{k=1}^K p(\theta_k \mid \phi, \nu) \prod_{n=1}^N \prod_{k=1}^K [\Pr(z_n = k \mid \pi) p(\mathbf{x}_n \mid \theta_k)]^{\mathbb{I}[z_n=k]} \quad (6)$$

- Substituting in the assumed forms

$$p(\pi, \{\theta_k\}_{k=1}^K, \{(z_n, \mathbf{x}_n)\}_{n=1}^N \mid \phi, \nu, \alpha) = \text{Dir}(\pi \mid \alpha) \prod_{k=1}^K p(\theta_k \mid \phi, \nu) \prod_{n=1}^N \prod_{k=1}^K [\pi_k p(\mathbf{x}_n \mid \theta_k)]^{\mathbb{I}[z_n=k]} \quad (7)$$

# Exponential family mixture models

What about  $p(\mathbf{x} \mid \boldsymbol{\theta}_k)$  and  $p(\boldsymbol{\theta}_k \mid \boldsymbol{\phi}, \nu)$ ?

Let's assume an **exponential family** likelihood,

$$p(\mathbf{x} \mid \boldsymbol{\theta}_k) = h(\mathbf{x}_n) \exp \left\{ \langle t(\mathbf{x}_n), \boldsymbol{\theta}_k \rangle - A(\boldsymbol{\theta}_k) \right\}. \quad (8)$$

Then assume a **conjugate prior**,

$$p(\boldsymbol{\theta}_k \mid \boldsymbol{\phi}, \nu) \propto \exp \left\{ \langle \boldsymbol{\phi}, \boldsymbol{\theta}_k \rangle - \nu A(\boldsymbol{\theta}_k) \right\}. \quad (9)$$

The hyperparameters  $\boldsymbol{\phi}$  are **pseudo-observations** of the sufficient statistics (like statistics from fake data points) and  $\nu$  is a **pseudo-count** (like the number of fake data points).

Note that the product of prior and likelihood remains in the same family as the prior. That's why we call it conjugate.

## Example: Gaussian mixture model

Assume the conditional distribution of  $\mathbf{x}_n$  is a Gaussian with mean  $\boldsymbol{\theta}_k \in \mathbb{R}^D$  and identity covariance,

$$p(\mathbf{x}_n \mid \boldsymbol{\theta}_k) = \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\theta}_k, I) \quad (10)$$

$$= (2\pi)^{-D/2} \exp \left\{ -\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\theta}_k)^\top (\mathbf{x}_n - \boldsymbol{\theta}_k) \right\} \quad (11)$$

$$= (2\pi)^{-D/2} \exp \left\{ -\frac{1}{2} \mathbf{x}_n^\top \mathbf{x}_n + \mathbf{x}_n^\top \boldsymbol{\theta}_k - \frac{1}{2} \boldsymbol{\theta}_k^\top \boldsymbol{\theta}_k \right\}, \quad (12)$$

which is an exponential family distribution with base measure  $h(\mathbf{x}_n) = (2\pi)^{-D/2} e^{-\frac{1}{2} \mathbf{x}_n^\top \mathbf{x}_n}$ , sufficient statistics  $t(\mathbf{x}_n) = \mathbf{x}_n$ , and log normalizer  $A(\boldsymbol{\theta}_k) = \frac{1}{2} \boldsymbol{\theta}_k^\top \boldsymbol{\theta}_k$ .

The conjugate prior is a Gaussian prior on the mean,

$$p(\boldsymbol{\theta}_k \mid \boldsymbol{\phi}, \nu) = \mathcal{N}(\nu^{-1} \boldsymbol{\phi}, \nu^{-1} I) \propto \exp \left\{ \boldsymbol{\phi}^\top \boldsymbol{\theta}_k - \frac{\nu}{2} \boldsymbol{\theta}_k^\top \boldsymbol{\theta}_k \right\} = \exp \left\{ \boldsymbol{\phi}^\top \boldsymbol{\theta}_k - \nu A(\boldsymbol{\theta}_k) \right\}. \quad (13)$$

Note that  $\boldsymbol{\phi}$  sets the location and  $\nu$  sets the precision (i.e. inverse variance).

## EM in the Gaussian mixture model

K-Means made **hard assignments** of data points to clusters in each iteration. What if we used **soft assignments** instead?

Instead of assigning  $z_n^*$  to the closest cluster, we compute *responsibilities* for each cluster:

1. For each data point  $n$  and component  $k$ , set the *responsibility* to,

$$\omega_{nk} = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\theta}_k, I)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\theta}_j, I)}.$$

likelihood of  $\mathbf{x}_n$  belonging to cluster  $k$

(14)

2. For each component  $k$ , set the new mean to

$$\boldsymbol{\theta}_k^* = \frac{1}{N_k} \sum_{n=1}^K \omega_{nk} \mathbf{x}_n,$$
(15)

where  $N_k = \sum_{n=1}^N \omega_{nk}$ .

This is called the **expectation maximization (EM)** algorithm.

# What is EM doing?

Rather than maximizing the **joint probability**, EM is maximizing the **marginal probability**,

$$\log p(\mathbf{X}, \boldsymbol{\theta}) = \log p(\boldsymbol{\theta}) + \log \sum_{\mathbf{z}} p(\mathbf{X}, \mathbf{z} \mid \boldsymbol{\theta}) \quad (16)$$

*← assignment of data points*

$$= \log p(\boldsymbol{\theta}) + \log \prod_{n=1}^N \sum_{z_n} p(\mathbf{x}_n, z_n \mid \boldsymbol{\theta}) \quad (17)$$

*$z_n$  are conditionally independent*

$$= \log p(\boldsymbol{\theta}) + \sum_{n=1}^N \log \sum_{z_n} p(\mathbf{x}_n, z_n \mid \boldsymbol{\theta}) \quad (18)$$

For discrete mixtures (with small enough  $K$ ) we can evaluate the log marginal probability (with what complexity?).

We can usually evaluate its gradient too, so we could just do gradient ascent to find  $\boldsymbol{\theta}^*$ .

However, EM typically obtains faster convergence rates.

# What is EM doing? II

**Idea:** Obtain a lower bound on the marginal probability,

$$\log p(\mathbf{X}, \boldsymbol{\theta}) = \log p(\boldsymbol{\theta}) + \sum_{n=1}^N \log \sum_{z_n} p(\mathbf{x}_n, z_n \mid \boldsymbol{\theta}) \quad (19)$$

$$= \log p(\boldsymbol{\theta}) + \sum_{n=1}^N \log \sum_{z_n} q(z_n) \frac{p(\mathbf{x}_n, z_n \mid \boldsymbol{\theta})}{q(z_n)} \quad (20)$$

$$= \log p(\boldsymbol{\theta}) + \sum_{n=1}^N \log \mathbb{E}_{q(z_n)} \left[ \frac{p(\mathbf{x}_n, z_n \mid \boldsymbol{\theta})}{q(z_n)} \right] \quad (21)$$

↳ true when  $q(z_n) \neq 0$

where  $q(z_n)$  is any distribution on  $z_n \in \{1, \dots, K\}$  such that  $p(\mathbf{x}_n, z_n \mid \boldsymbol{\theta})$  is **absolutely continuous** w.r.t.  $q(z_n)$ .



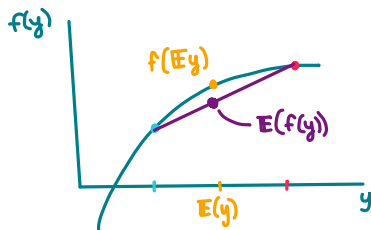
# Jensen's Inequality

Jensen's inequality states that,

$$f(\mathbb{E}_{p(y)}[y]) \geq \mathbb{E}_{p(y)}[f(y)] \quad (22)$$

if  $f$  is a **concave function**, with equality iff  $f$  is linear.

[Picture]



## What is EM doing? III

Applied to the log marginal probability, Jensen's inequality yields,

$$\log p(\mathbf{X}, \boldsymbol{\theta}) = \log p(\boldsymbol{\theta}) + \sum_{n=1}^N \log \mathbb{E}_{q_n(z_n)} \left[ \frac{p(\mathbf{x}_n, z_n \mid \boldsymbol{\theta})}{q_n(z_n)} \right] \quad (23)$$

$$\geq \log p(\boldsymbol{\theta}) + \sum_{n=1}^N \mathbb{E}_{q_n(z_n)} [\log p(\mathbf{x}_n, z_n \mid \boldsymbol{\theta}) - \log q_n(z_n)] \quad (24)$$

$$\triangleq \mathcal{L}[\boldsymbol{\theta}, \mathbf{q}] \quad (25)$$

ELBO

where  $\mathbf{q} = (q_1, \dots, q_N)$  is a tuple of densities.

This is called the **evidence lower bound**, or **ELBO** for short.

It is a function of  $\boldsymbol{\theta}$  and a **functional** of  $\mathbf{q}$ , since each  $q_n$  is a probability density function.   
 PDFs! not numbers.

function of functions

We can think of **EM** as coordinate ascent on the **ELBO**.

## M-step: Maximizing the ELBO wrt $\theta$ (Gaussian case)

Suppose we fix  $q$ . Since each  $z_n$  is a discrete latent variable,  $q_n$  must be a probability mass function. Let it be denoted by,

$$q_n(z_n) = [q_n(z_n = 1), \dots, q_n(z_n = K)]^\top = [\omega_{n1}, \dots, \omega_{nK}]^\top. \quad (26)$$

(These will be the **responsibilities** from before.)

Now, recall our basic model,  $\mathbf{x}_n \sim \mathcal{N}(\boldsymbol{\theta}_{z_n}, I)$ , and assume a prior  $\boldsymbol{\theta}_k \sim \mathcal{N}(\boldsymbol{\phi}, \nu^{-1}I)$ , Then,

$$\mathcal{L}[\boldsymbol{\theta}, q] = \log p(\boldsymbol{\theta}) + \sum_{n=1}^N \mathbb{E}_{q_n(z_n)} [\log p(\mathbf{x}_n, z_n | \boldsymbol{\theta})] + c \quad (27)$$

$$= \log p(\boldsymbol{\theta}) + \sum_{n=1}^N \sum_{k=1}^K \omega_{nk} \log p(\mathbf{x}_n, z_n = k | \boldsymbol{\theta}) + c \quad (28)$$

$$= \sum_{k=1}^K [\boldsymbol{\phi}^\top \boldsymbol{\theta}_k - \frac{\nu}{2} \boldsymbol{\theta}_k^\top \boldsymbol{\theta}_k] + \sum_{n=1}^N \sum_{k=1}^K \omega_{nk} [\mathbf{x}_n^\top \boldsymbol{\theta}_k - \frac{1}{2} \boldsymbol{\theta}_k^\top \boldsymbol{\theta}_k] + c \quad (29)$$

find  $\boldsymbol{\theta}_k^*$  maximizing  $\mathcal{L}(\boldsymbol{\theta}, q)$ .

## M-step: Maximizing the ELBO wrt $\theta$ (Gaussian case) II

Zooming in on just  $\theta_k$ ,

$$\mathcal{L}[\theta, q] = \phi_{N,k}^\top \theta_k - \frac{1}{2} \nu_{N,k} \theta_k^\top \theta_k \quad (30)$$

where

$$\phi_{N,k} = \phi + \sum_{n=1}^N \omega_{nk} \mathbf{x}_n \quad \nu_{N,k} = \nu + \sum_{n=1}^N \omega_{nk} \quad (31)$$

Taking derivatives and setting to zero yields,

$$\theta_k^* = \frac{\phi_{N,k}}{\nu_{N,k}} = \frac{\phi + \sum_{n=1}^N \omega_{nk} \mathbf{x}_n}{\nu + \sum_{n=1}^N \omega_{nk}}. \quad (32)$$

*weighted average of data points plus a bias term.*

In the improper uniform prior limit where  $\phi \rightarrow 0$  and  $\nu \rightarrow 0$ , we recover the EM updates shown on slide 6.

## E-step: Maximizing the ELBO wrt $q$ (Gaussian case)

As a function of  $q_n$ , for discrete Gaussian mixtures with identity covariance,

$$\mathcal{L}[\boldsymbol{\theta}, \boldsymbol{q}] = \mathbb{E}_{q_n(z_n)} [\log p(\mathbf{x}_n, z_n | \boldsymbol{\theta}) - \log q_n(z_n)] + c \quad (33)$$

$$\begin{aligned} & \text{keep terms concerned with } \omega_{nk}! \quad = \sum_{k=1}^K \omega_{nk} [\log \mathcal{N}(\mathbf{x}_n | \boldsymbol{\theta}_k, \boldsymbol{I}) + \log \underbrace{\pi_k}_{\text{prior}} - \log \omega_{nk}] + c \end{aligned} \quad (34)$$

where  $\boldsymbol{\pi} = [\pi_1, \dots, \pi_K]^\top$  is the vector of cluster probabilities.

We also have two constraints:  $\omega_{nk} \geq 0$  and  $\sum_k \omega_{nk} = 1$ . Let's ignore the non-negative constraint for now (it will automatically be satisfied anyway) and write the Lagrangian with the simplex constraint,

$$\mathcal{J}(\boldsymbol{\omega}_n, \lambda) = \sum_{k=1}^K \omega_{nk} [\log \mathcal{N}(\mathbf{x}_n | \boldsymbol{\theta}_k, \boldsymbol{I}) + \log \pi_k - \log \omega_{nk}] - \lambda \left( 1 - \sum_{k=1}^K \omega_{nk} \right) \quad (35)$$

## E-step: Maximizing the ELBO wrt $q$ (Gaussian case) II

Taking the partial derivative wrt  $\omega_{nk}$  and setting to zero yields,

$$\frac{\partial}{\partial \omega_{nk}} \mathcal{J}(\omega_n, \lambda) = \log \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\theta}_k, I) + \log \pi_k - \log \omega_{nk} - 1 + \lambda = 0 \quad (36)$$

$$\Rightarrow \log \omega_{nk}^* = \log \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\theta}_k, I) + \log \pi_k + \lambda - 1 \quad (37)$$

$$\Rightarrow \omega_{nk}^* \propto \pi_k \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\theta}_k, I) \quad (38)$$

Enforcing the simplex constraint yields,

$$\sum_{k=1}^K \omega_{nk} = 1$$

$$\omega_{nk}^* = \frac{\overset{\text{responsibility}}{\downarrow} \pi_k \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\theta}_k, I)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\theta}_j, I)}, \quad (39)$$

just like on slide 6.

Note that

$$\omega_{nk}^* \propto p(z_n = k) p(\mathbf{x}_n \mid z_n = k, \boldsymbol{\theta}) = \overset{\text{posterior distribution}}{p(z_n = k \mid \mathbf{x}_n, \boldsymbol{\theta})} \quad (40)$$

## The ELBO is tight after the E-step

Equivalently,  $q_n$  equals the posterior,  $p(z_n | \mathbf{x}_n, \boldsymbol{\theta})$ . At that point, the ELBO simplifies to,

$$\mathcal{L}[\boldsymbol{\theta}, q] = \log p(\boldsymbol{\theta}) + \sum_{n=1}^N \mathbb{E}_{q_n(z_n)} [\log p(\mathbf{x}_n, z_n | \boldsymbol{\theta}) - \log q_n(z_n)] \quad (41)$$

$$= \log p(\boldsymbol{\theta}) + \sum_{n=1}^N \mathbb{E}_{p(z_n | \mathbf{x}_n, \boldsymbol{\theta})} [\log p(\mathbf{x}_n, z_n | \boldsymbol{\theta}) - \log p(z_n | \mathbf{x}_n, \boldsymbol{\theta})] \quad (42)$$

$$= \log p(\boldsymbol{\theta}) + \sum_{n=1}^N \mathbb{E}_{p(z_n | \mathbf{x}_n, \boldsymbol{\theta})} [\log p(\mathbf{x}_n | \boldsymbol{\theta})] \quad (43)$$

(Baye's rule):  $\frac{p(\mathbf{x}_n, z_n | \boldsymbol{\theta})}{p(z_n | \mathbf{x}_n, \boldsymbol{\theta})} = p(\mathbf{x}_n | \boldsymbol{\theta})$

$$= \log p(\boldsymbol{\theta}) + \sum_{n=1}^N \log p(\mathbf{x}_n | \boldsymbol{\theta}) \quad (44)$$

$$= \log p(\mathbf{X}, \boldsymbol{\theta}) \leftarrow \text{marginal prob!} \quad (45)$$

In other words, **after the E step, the bound is tight!**

# EM as a minorize-maximize (MM) algorithm

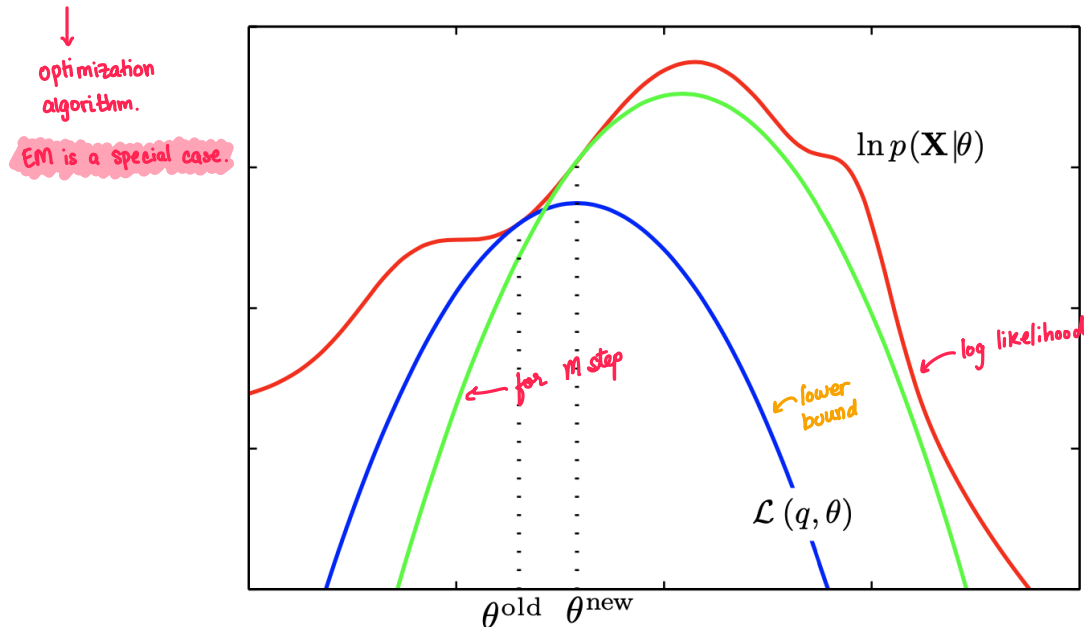


Figure: Bishop, Figure 9.14: EM alternates between constructing a lower bound (minorizing) and finding new parameters that maximize it.



## M-step: Maximizing the ELBO wrt $\theta$ (generic exp. fam.)

Now let's consider the general Bayesian mixture with exponential family likelihoods and conjugate priors. As a function of  $\theta$ ,

$$\mathcal{L}[\theta, q] = \log p(\theta) + \sum_{n=1}^N \mathbb{E}_{q_n(z_n)} [\log p(\mathbf{x}_n, z_n | \theta)] + c \quad (46)$$

$$= \log p(\theta) + \sum_{n=1}^N \sum_{k=1}^K \omega_{nk} \log p(\mathbf{x}_n, z_n = k | \theta) + c \quad (47)$$

$$= \sum_{k=1}^K [\phi^\top \theta_k - \nu A(\theta_k)] + \sum_{n=1}^N \sum_{k=1}^K \omega_{nk} [t(\mathbf{x}_n)^\top \theta_k - A(\theta_k)] + c \quad (48)$$

## M-step: Maximizing the ELBO wrt $\theta$ (generic exp. fam.) II

Zooming in on just  $\theta_k$ ,

$$\nabla_{\theta_k} \mathcal{L}(\theta, q) = \phi_{N,k} - v_{N,k} \cdot \nabla_{\theta} A(\theta) = 0 \Rightarrow 0$$

$$\mathcal{L}[\theta, q] = \phi_{N,k}^\top \theta_k - v_{N,k} A(\theta_k) \quad (49)$$

where

$$\phi_{N,k} = \phi + \sum_{n=1}^N \omega_{nk} \overbrace{t(\mathbf{x}_n)}^{\text{convert to natural parameters}} \quad v_{N,k} = v + \sum_{n=1}^N \omega_{nk} \quad (50)$$

Taking derivatives and setting to zero yields,

$$\theta_k^* = [\nabla A]^{-1} \left( \frac{\phi_{N,k}}{v_{N,k}} \right) \quad (51)$$

sufficient stats:  
everything you know  
about the data

## M-step: Maximizing the ELBO wrt $\theta$ (generic exp. fam.) III

What is the gradient of the log normalizer? We have,

$$\nabla A(\theta_k) = \nabla_{\theta_k} \log \int h(\mathbf{x}) \exp \{ \langle t(\mathbf{x}), \theta_k \rangle \} d\mathbf{x} \quad (52)$$

$$= \frac{\int h(\mathbf{x}) \exp \{ \langle t(\mathbf{x}), \theta_k \rangle \} t(\mathbf{x}) d\mathbf{x}}{\int h(\mathbf{x}) \exp \{ \langle t(\mathbf{x}), \theta_k \rangle \} d\mathbf{x}} \leftarrow \exp\{A(\theta_k)\} \quad (53)$$

$$= \int \underbrace{h(\mathbf{x}) \exp \{ \langle t(\mathbf{x}), \theta_k \rangle - A(\theta_k) \}}_{p(\mathbf{x}|\theta_k)} t(\mathbf{x}) d\mathbf{x} \quad (54)$$

$$= \mathbb{E}_{p(\mathbf{x}|\theta_k)}[t(\mathbf{x})] \quad p(\mathbf{x}|\theta_k) \quad (55)$$

Gradients of the log normalizer yield **expected sufficient statistics!**   
 want to find  $\theta_k$  that gives us  $\Rightarrow$  that's what  $A^{-1}$  is doing

## M-step: Maximizing the ELBO wrt $\theta$ (generic exp. fam.) IV

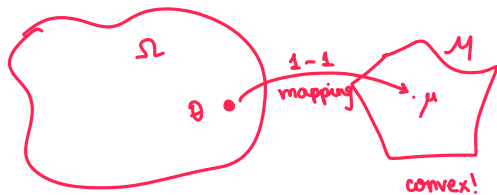
The gradient  $\nabla A$  is a map from the set of valid natural parameters  $\Omega$  (those for which the log normalizer is finite) to the set of realizable mean parameters  $\mathcal{M}$ , *marginal polytope*

$$\mathcal{M} = \{\mu \in \mathbb{R}^D : \exists p \text{ s.t. } \mathbb{E}_p[t(\mathbf{x})] = \mu\} \quad (56)$$

An exponential family is **minimal** if its sufficient statistics are linearly independent.

**Fact:** The gradient mapping  $\nabla A : \Omega \rightarrow \mathcal{M}$  is one-to-one (and hence invertible) if and only if the exponential family is minimal.

<Picture>



## M-step: Maximizing the ELBO wrt $\theta$ (generic exp. fam.) $\mathbf{V}$

Thus, the generic M-step in eq. 51 amounts to finding the natural parameters  $\theta_k^*$  that yield the expected sufficient statistics  $\phi_{N,k}/\nu_{N,k}$  by inverting the gradient mapping.

*Note: There is a longer and much more technical story about exponential families, maximum likelihood, convex analysis, and conjugate duals that you can read about in [Wainwright et al., 2008, Ch. 3] if you are interested.*

## E-step: Maximizing the ELBO wrt $q$ (generic exp. fam.)

In our first pass, we assumed  $q_n$  was a finite pmf. More generally,  $q_n$  will be a probability density function, and optimizing over functions usually requires the **calculus of variations**. (Ugh!)

However, note that we can write the ELBO in a slightly different form,

$$\mathcal{L}[\theta, q] = \log p(\theta) + \sum_{n=1}^N \mathbb{E}_{q_n(z_n)} [\log p(\mathbf{x}_n, z_n | \theta) - \log q_n(z_n)] \quad (57)$$

$$p(\mathbf{x}, z | \theta) = p(\mathbf{x} | \theta) \cdot p(z | \theta)$$

$$= \log p(\theta) + \sum_{n=1}^N \mathbb{E}_{q_n(z_n)} [\log p(z_n | \mathbf{x}_n, \theta) + \log p(\mathbf{x}_n | \theta) - \log q_n(z_n)] \quad (58)$$

constant wrt  $\mathbb{E}_{q_n(\cdot)}$

$$= \log p(\theta) + \sum_{n=1}^N [\log p(\mathbf{x}_n | \theta) - D_{\text{KL}}(q_n(z_n) \| p(z_n | \mathbf{x}_n, \theta))] \quad (59)$$

$$= \log p(\mathbf{X}, \theta) - \sum_{n=1}^N D_{\text{KL}}(q_n(z_n) \| p(z_n | \mathbf{x}_n, \theta)) \quad (60)$$

where  $D_{\text{KL}}(\cdot \| \cdot)$  denote the **Kullback-Leibler divergence**.

# Kullback-Leibler (KL) divergence

The KL divergence is defined as,

$$D_{\text{KL}}(q(z) \parallel p(z)) = \int q(z) \log \frac{q(z)}{p(z)} dz. \quad (61)$$

It gives a notion of how similar two distributions are, but it is **not a metric!** (It is not symmetric, e.g.)  
Still, it has some intuitive properties:

- ▶ It is non-negative,  $D_{\text{KL}}(q(z) \parallel p(z)) \geq 0$ .
- ▶ It equals zero iff the distributions are the same,  $D_{\text{KL}}(q(z) \parallel p(z)) = 0 \iff q(z) = p(z)$  almost everywhere.

↪ minimum.

doesn't satisfy triangle inequality.

## E-step: Maximizing the ELBO wrt $q$ (generic exp. fam.) II

Maximizing the ELBO wrt  $q_n$  amounts to minimizing the KL divergence to the posterior  $p(z_n | \mathbf{x}_n, \boldsymbol{\theta})$ ,

$$\mathcal{L}[\boldsymbol{\theta}, \mathbf{q}] = \log p(\boldsymbol{\theta}) + \sum_{n=1}^N [\log p(\mathbf{x}_n | \boldsymbol{\theta}) - D_{\text{KL}}(q_n(z_n) \| p(z_n | \mathbf{x}_n, \boldsymbol{\theta}))] \quad (62)$$

$$= -D_{\text{KL}}(q_n(z_n) \| p(z_n | \mathbf{x}_n, \boldsymbol{\theta})) + c \quad (63)$$

As we said, the KL is minimized when  $q_n(z_n) = p(z_n | \mathbf{x}_n, \boldsymbol{\theta})$ , so the optimal update is,

$$q_n^*(z_n) = p(z_n | \mathbf{x}_n, \boldsymbol{\theta}), \quad (64)$$

just like we found on slide 14.

*EM: special case of variational inference*



# References I

Martin J Wainwright, Michael I Jordan, et al. Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning*, 1(1-2):1-305, 2008.

$$\begin{aligned} \text{NIW}(\mu, \Sigma) &= \mathcal{N}(\mu | \mu_0, \kappa_0^{-1} \Sigma) \cdot \text{IW}(\Sigma | \nu_0, \Sigma_0) \\ &= |\kappa_0^{-1} \Sigma|^{1/2} \cdot \exp\left\{-\frac{\kappa_0}{2} (\mu - \mu_0)^T \Sigma^{-1} (\mu - \mu_0)\right\} \cdot |\Sigma|^{-\frac{\nu_0}{2}} \cdot \exp\left\{-\frac{1}{2} \text{tr}(\Sigma^{-1} \Sigma_0)\right\} \end{aligned}$$

$$\begin{aligned} \log \text{NIW}(\mu, \Sigma) &= -\frac{1}{2} \log |\Sigma| - \frac{\kappa_0}{2} \mu^T \Sigma^{-1} \mu + \kappa_0 \mu_0^T \Sigma^{-1} \mu - \frac{\kappa_0}{2} \mu_0^T \Sigma^{-1} \mu_0 - \frac{\nu_0}{2} \log |\Sigma| - \frac{1}{2} \text{tr}(\Sigma^{-1} \Sigma_0) \\ &= -\frac{(\nu_0+1)}{2} \log |\Sigma| - \frac{\kappa_0}{2} \mu^T \Sigma^{-1} \mu + \kappa_0 \mu_0^T \Sigma^{-1} \mu - \frac{1}{2} \text{tr}(\Sigma^{-1} \Sigma_0 + \kappa_0 \Sigma \mu_0 \mu_0^T) \\ &= \left\langle -\frac{(\nu_0+1)}{2}, \log |\Sigma| \right\rangle - \left\langle \frac{\kappa_0}{2}, \mu^T \Sigma^{-1} \mu \right\rangle + \left\langle \kappa_0 \mu_0, \Sigma^{-1} \mu \right\rangle + \underbrace{\left\langle -\frac{1}{2} (\Sigma_0 + \kappa_0 \mu_0 \mu_0^T), \Sigma^{-1} \right\rangle}_{\text{inner product of matrices!}} \end{aligned}$$

overcomplete representation!

→ line these terms in MVN.

→ EM algo w/ unknown means/covariances will combine these terms.