

LECTURE 10: THE EM ALGORITHM (CONTD)

STAT 545: INTRO. TO COMPUTATIONAL STATISTICS

Vinayak Rao

Purdue University

September 22, 2019

EXPONENTIAL FAMILY MODELS

Consider a space \mathbb{X} . E.g. \mathbb{R} , \mathbb{R}^d or \mathbb{N} .

We want to specify a probability distribution $p(x|\boldsymbol{\theta})$ on \mathbb{X} .

$\boldsymbol{\phi}(x) = [\phi_1(x), \dots, \phi_D(x)]$: (feature) vector of sufficient statistics

$\boldsymbol{\theta} = [\theta_1, \dots, \theta_D]$: vector of natural parameters

EXPONENTIAL FAMILY MODELS

Consider a space \mathbb{X} . E.g. \mathbb{R} , \mathbb{R}^d or \mathbb{N} .

We want to specify a probability distribution $p(x|\boldsymbol{\theta})$ on \mathbb{X} .

$\boldsymbol{\phi}(x) = [\phi_1(x), \dots, \phi_D(x)]$: (feature) vector of sufficient statistics

$\boldsymbol{\theta} = [\theta_1, \dots, \theta_D]$: vector of natural parameters

Exponential family distribution:

$$p(x|\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} h(x) \exp(\boldsymbol{\theta}^\top \boldsymbol{\phi}(x))$$

EXPONENTIAL FAMILY MODELS

Consider a space \mathbb{X} . E.g. \mathbb{R} , \mathbb{R}^d or \mathbb{N} .

We want to specify a probability distribution $p(x|\boldsymbol{\theta})$ on \mathbb{X} .

$\boldsymbol{\phi}(x) = [\phi_1(x), \dots, \phi_D(x)]$: (feature) vector of sufficient statistics

$\boldsymbol{\theta} = [\theta_1, \dots, \theta_D]$: vector of natural parameters

Exponential family distribution:

$$p(x|\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} h(x) \exp(\boldsymbol{\theta}^\top \boldsymbol{\phi}(x))$$

$h(x)$ is the base-measure or base distribution.

$Z(\boldsymbol{\theta}) = \int h(x) \exp(\boldsymbol{\theta}^\top \boldsymbol{\phi}(x)) dx$ is the normalization constant.

The normal distribution:

$$\begin{aligned} p(x|\mu, \sigma^2) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right) \\ &= \frac{\exp\left(-\frac{\mu^2}{2\sigma^2}\right)}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}x^2 + \frac{\mu}{\sigma^2}x\right) \end{aligned}$$

The normal distribution:

$$\begin{aligned} p(x|\mu, \sigma^2) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right) \\ &= \frac{\exp\left(-\frac{\mu^2}{2\sigma^2}\right)}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}x^2 + \frac{\mu}{\sigma^2}x\right) \end{aligned}$$

The Poisson distribution:

$$\begin{aligned} p(x|\lambda) &= \frac{\lambda^x \exp(-\lambda)}{x!} \\ &= \exp(-\lambda) \frac{1}{x!} \exp(\log(\lambda)x) \end{aligned}$$

MINIMAL EXPONENTIAL FAMILY

Sufficient statistics are linearly independent

Consider a K -component discrete distribution $\pi = (\pi_1, \dots, \pi_K)$

$$p(X) = \prod_{c=1}^K \pi_c^{\delta(X=c)} = \exp\left(\sum_{c=1}^K \delta(X=c) \log \pi_c\right)$$

MINIMAL EXPONENTIAL FAMILY

Sufficient statistics are linearly independent

Consider a K -component discrete distribution $\pi = (\pi_1, \dots, \pi_K)$

$$p(X) = \prod_{c=1}^K \pi_c^{\delta(X=c)} = \exp\left(\sum_{c=1}^K \delta(X=c) \log \pi_c\right)$$

Is it minimal?

$$\begin{aligned} p(X) &= \pi_K \exp\left(\sum_{c=1}^{K-1} \delta(X=c) \log \pi_c / \pi_K\right) \\ &= \frac{1}{Z} \exp\left(\sum_{c=1}^{K-1} \delta(X=c) \theta_c\right) \end{aligned}$$

Given N i.i.d. observations $X \equiv \{x_1, \dots, x_N\}$, the likelihood is

$$\begin{aligned}\mathcal{L}(X|\boldsymbol{\theta}) &= \prod_{i=1}^N \frac{1}{Z(\boldsymbol{\theta})} h(x_i) \exp(\boldsymbol{\theta}^\top \boldsymbol{\phi}(x_i)) \\ &= \left(\frac{1}{Z(\boldsymbol{\theta})}\right)^N \left(\prod_{i=1}^N h(x_i)\right) \exp(\boldsymbol{\theta}^\top \sum_{i=1}^N \boldsymbol{\phi}(x_i))\end{aligned}$$

MAXIMUM-LIKELIHOOD ESTIMATION

Given N i.i.d. observations $X \equiv \{x_1, \dots, x_N\}$, the likelihood is

$$\begin{aligned}\mathcal{L}(X|\boldsymbol{\theta}) &= \prod_{i=1}^N \frac{1}{Z(\boldsymbol{\theta})} h(x_i) \exp(\boldsymbol{\theta}^\top \boldsymbol{\phi}(x_i)) \\ &= \left(\frac{1}{Z(\boldsymbol{\theta})}\right)^N \left(\prod_{i=1}^N h(x_i)\right) \exp(\boldsymbol{\theta}^\top \sum_{i=1}^N \boldsymbol{\phi}(x_i))\end{aligned}$$

The log-likelihood $\ell(X|\boldsymbol{\theta}) = \log \mathcal{L}(X|\boldsymbol{\theta})$ is

$$\ell(X|\boldsymbol{\theta}) = \boldsymbol{\theta}^\top \left(\sum_{i=1}^N \boldsymbol{\phi}(x_i)\right) - N \log Z(\boldsymbol{\theta}) + \sum_{i=1}^N \log h(x_i)$$

MAXIMUM-LIKELIHOOD ESTIMATION

Given N i.i.d. observations $X \equiv \{x_1, \dots, x_N\}$, the likelihood is

$$\begin{aligned}\mathcal{L}(X|\boldsymbol{\theta}) &= \prod_{i=1}^N \frac{1}{Z(\boldsymbol{\theta})} h(x_i) \exp(\boldsymbol{\theta}^\top \boldsymbol{\phi}(x_i)) \\ &= \left(\frac{1}{Z(\boldsymbol{\theta})}\right)^N \left(\prod_{i=1}^N h(x_i)\right) \exp(\boldsymbol{\theta}^\top \sum_{i=1}^N \boldsymbol{\phi}(x_i))\end{aligned}$$

The log-likelihood $\ell(X|\boldsymbol{\theta}) = \log \mathcal{L}(X|\boldsymbol{\theta})$ is

$$\ell(X|\boldsymbol{\theta}) = \boldsymbol{\theta}^\top \left(\sum_{i=1}^N \boldsymbol{\phi}(x_i)\right) - N \log Z(\boldsymbol{\theta}) + \sum_{i=1}^N \log h(x_i)$$

To calculate a maximum likelihood estimate, we only need the sum of the suff. statistics.

$$\ell(X|\boldsymbol{\theta}) = \boldsymbol{\theta}^\top \left(\sum_{i=1}^N \boldsymbol{\phi}(x_i) \right) - N \log Z(\boldsymbol{\theta}) + \sum_{i=1}^N \log h(x_i)$$

$$\ell(X|\boldsymbol{\theta}) = \boldsymbol{\theta}^\top \left(\sum_{i=1}^N \boldsymbol{\phi}(x_i) \right) - N \log Z(\boldsymbol{\theta}) + \sum_{i=1}^N \log h(x_i)$$

At MLE of θ_d , the d th component of $\boldsymbol{\theta}$: $\frac{\partial \ell(X|\boldsymbol{\theta})}{\partial \theta_d} = 0$.

$$\ell(X|\boldsymbol{\theta}) = \boldsymbol{\theta}^\top \left(\sum_{i=1}^N \boldsymbol{\phi}(x_i) \right) - N \log Z(\boldsymbol{\theta}) + \sum_{i=1}^N \log h(x_i)$$

At MLE of θ_d , the d th component of $\boldsymbol{\theta}$: $\frac{\partial \ell(X|\boldsymbol{\theta})}{\partial \theta_d} = 0$.

$$\sum_{i=1}^N \phi_d(x_i) = N \frac{\partial \log Z(\boldsymbol{\theta})}{\partial \theta_d}$$

$$\ell(X|\boldsymbol{\theta}) = \boldsymbol{\theta}^\top \left(\sum_{i=1}^N \boldsymbol{\phi}(x_i) \right) - N \log Z(\boldsymbol{\theta}) + \sum_{i=1}^N \log h(x_i)$$

At MLE of θ_d , the d th component of $\boldsymbol{\theta}$: $\frac{\partial \ell(X|\boldsymbol{\theta})}{\partial \theta_d} = 0$.

$$\sum_{i=1}^N \phi_d(x_i) = N \frac{\partial \log Z(\boldsymbol{\theta})}{\partial \theta_d}$$

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N \phi_d(x_i) &= \frac{1}{Z(\boldsymbol{\theta})} \frac{\partial}{\partial \theta_d} Z(\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \frac{\partial}{\partial \theta_d} \left(\int h(x) \exp(\boldsymbol{\theta}^\top \boldsymbol{\phi}(x)) dx \right) \\ &= \frac{1}{Z(\boldsymbol{\theta})} \int h(x) \frac{\partial \exp(\boldsymbol{\theta}^\top \boldsymbol{\phi}(x))}{\partial \theta_d} dx \\ &= \int \frac{1}{Z(\boldsymbol{\theta})} h(x) \exp(\boldsymbol{\theta}^\top \boldsymbol{\phi}(x)) \phi_d(x) dx \end{aligned}$$

Match empirical and population averages of $\phi(x)$:

$$\frac{1}{N} \sum_{i=1}^N \phi_d(x_i) = \mathbb{E}_{\theta_{MLE}}[\phi_d(x)]$$

Match empirical and population averages of $\phi(x)$:

$$\frac{1}{N} \sum_{i=1}^N \phi_d(x_i) = \mathbb{E}_{\theta_{MLE}}[\phi_d(x)]$$

RHS: 'moment parameters' of the exponential distribution.

Thus: θ_{MLE} are natural parameters corresponding to empirical moment parameters ('moment matching').

MLE FOR EXPONENTIAL FAMILIES

Match empirical and population averages of $\phi(x)$:

$$\frac{1}{N} \sum_{i=1}^N \phi_d(x_i) = \mathbb{E}_{\theta_{MLE}}[\phi_d(x)]$$

RHS: 'moment parameters' of the exponential distribution.

Thus: θ_{MLE} are natural parameters corresponding to empirical moment parameters ('moment matching').

Is this a maximum?

- is second derivative (Hessian) negative (negative definite)?

We can show $\frac{\partial^2}{\partial \theta_i \partial \theta_j} \log Z(X|\theta) = \text{Cov}(\phi_i, \phi_j)$, and the Hessian of $\ell(X|\theta)$ is $-N$ times the feature covariance matrix

EXAMPLE

The 1-d Gaussian: $\phi = [x \quad x^2]$

Moment parameters are mean and mean squared

Easy to find corresponding natural parameters

Quite often, it is not the case.

However, we will restrict ourselves to cases where it is.

EXAMPLE

The 1-d Gaussian: $\phi = [x \quad x^2]$

Moment parameters are mean and mean squared

Easy to find corresponding natural parameters

Quite often, it is not the case.

However, we will restrict ourselves to cases where it is.

Can you do it for the Poisson?

$$p(x|\lambda) = \frac{1}{x!} \lambda^x \exp(-\lambda)$$

Let samples from the exponential family have two parts: $[x \ y]$.

Feature vector $\phi([x \ y]) := \phi(x, y)$.

$$P(x, y | \theta) = \frac{h(x, y)}{Z(\theta)} \exp(\theta^\top \phi(x, y))$$

Let samples from the exponential family have two parts: $[x \ y]$.

Feature vector $\phi([x \ y]) := \phi(x, y)$.

$$P(x, y | \theta) = \frac{h(x, y)}{Z(\theta)} \exp(\theta^\top \phi(x, y))$$

We observe only x . What is the conditional over y ?

$$P(y|x, \theta) = \frac{P(x, y | \theta)}{P(x | \theta)} = \frac{h(x, y)}{P(x | \theta) Z(\theta)} \exp(\theta^\top \phi(x, y))$$

MISSING DATA IN EXPONENTIAL FAMILY DISTRIBUTIONS

Let samples from the exponential family have two parts: $[x \ y]$.

Feature vector $\phi([x \ y]) := \phi(x, y)$.

$$P(x, y | \theta) = \frac{h(x, y)}{Z(\theta)} \exp(\theta^\top \phi(x, y))$$

We observe only x . What is the conditional over y ?

$$P(y | x, \theta) = \frac{P(x, y | \theta)}{P(x | \theta)} = \frac{h(x, y)}{P(x | \theta) Z(\theta)} \exp(\theta^\top \phi(x, y))$$

An exponential family distrib. over y (remember x is fixed) with:

- feature vector $\phi_x(y) = \phi(x, y)$
- base distribution $h_x(y) = h(x, y)$

MISSING DATA IN EXPONENTIAL FAMILY DISTRIBUTIONS

Let samples from the exponential family have two parts: $[x \ y]$.

Feature vector $\phi([x \ y]) := \phi(x, y)$.

$$P(x, y | \theta) = \frac{h(x, y)}{Z(\theta)} \exp(\theta^\top \phi(x, y))$$

We observe only x . What is the conditional over y ?

$$P(y | x, \theta) = \frac{P(x, y | \theta)}{P(x | \theta)} = \frac{h(x, y)}{P(x | \theta) Z(\theta)} \exp(\theta^\top \phi(x, y))$$

An exponential family distrib. over y (remember x is fixed) with:

- feature vector $\phi_x(y) = \phi(x, y)$
- base distribution $h_x(y) = h(x, y)$

Not necessarily easy to work with, but will restrict to this case.
 y_i with different x_i belong to different exp. fam. distrbs.

What about the marginal likelihood $P(x|\boldsymbol{\theta}) = \int P(x, y|\boldsymbol{\theta})dy$?

What about the marginal likelihood $P(x|\boldsymbol{\theta}) = \int P(x, y|\boldsymbol{\theta})dy$?

- Not an exponential family distribution!
- Calculating derivatives is messy
- No nice closed form expression for MLE

What about the marginal likelihood $P(x|\boldsymbol{\theta}) = \int P(x, y|\boldsymbol{\theta})dy$?

- Not an exponential family distribution!
- Calculating derivatives is messy
- No nice closed form expression for MLE

One approach: the EM algorithm.

An algorithm for MLE in exp. fam. distribs. with missing data

What about the marginal likelihood $P(x|\boldsymbol{\theta}) = \int P(x, y|\boldsymbol{\theta})dy$?

- Not an exponential family distribution!
- Calculating derivatives is messy
- No nice closed form expression for MLE

One approach: the EM algorithm.

An algorithm for MLE in exp. fam. distrib. with missing data

Problem: Given observations i.i.d. $X = \{x_1, \dots, x_N\}$ from $P(x|\boldsymbol{\theta})$ where $P(X, Y|\boldsymbol{\theta})$ is exponential family, maximize w.r.t. $\boldsymbol{\theta}$

$$\ell(X|\boldsymbol{\theta}) = \log P(X|\boldsymbol{\theta}) = \sum_{i=1}^N \log P(x_i|\boldsymbol{\theta})$$

MLE IN LATENT VARIABLE MODELS

$$\ell(X|\boldsymbol{\theta}) = \sum_{i=1}^N \log P(x_i|\boldsymbol{\theta}) = \sum_{i=1}^N \log \int P(x_i, y_i|\boldsymbol{\theta}) dy_i$$

MLE IN LATENT VARIABLE MODELS

$$\begin{aligned}\ell(X|\boldsymbol{\theta}) &= \sum_{i=1}^N \log P(x_i|\boldsymbol{\theta}) = \sum_{i=1}^N \log \int P(x_i, y_i|\boldsymbol{\theta}) dy_i \\ &= \sum_{i=1}^N \log \int q_i(y_i) \frac{P(x_i, y_i|\boldsymbol{\theta})}{q_i(y_i)} dy_i \quad (\text{for arbitrary densities } q_i(y_i))\end{aligned}$$

MLE IN LATENT VARIABLE MODELS

$$\begin{aligned}\ell(X|\boldsymbol{\theta}) &= \sum_{i=1}^N \log P(x_i|\boldsymbol{\theta}) = \sum_{i=1}^N \log \int P(x_i, y_i|\boldsymbol{\theta}) dy_i \\ &= \sum_{i=1}^N \log \int q_i(y_i) \frac{P(x_i, y_i|\boldsymbol{\theta})}{q_i(y_i)} dy_i \quad (\text{for arbitrary densities } q_i(y_i)) \\ &\geq \sum_{i=1}^N \int q_i(y_i) \log \frac{P(x_i, y_i|\boldsymbol{\theta})}{q_i(y_i)} dy_i \quad (\text{Jensen's inequality})\end{aligned}$$

MLE IN LATENT VARIABLE MODELS

$$\begin{aligned}\ell(X|\boldsymbol{\theta}) &= \sum_{i=1}^N \log P(x_i|\boldsymbol{\theta}) = \sum_{i=1}^N \log \int P(x_i, y_i|\boldsymbol{\theta}) dy_i \\ &= \sum_{i=1}^N \log \int q_i(y_i) \frac{P(x_i, y_i|\boldsymbol{\theta})}{q_i(y_i)} dy_i \quad (\text{for arbitrary densities } q_i(y_i)) \\ &\geq \sum_{i=1}^N \int q_i(y_i) \log \frac{P(x_i, y_i|\boldsymbol{\theta})}{q_i(y_i)} dy_i \quad (\text{Jensen's inequality}) \\ &= \sum_{i=1}^N \int q_i(y_i) \log P(x_i, y_i|\boldsymbol{\theta}) dy_i - \sum_{i=1}^N \int q_i(y_i) \log q_i(y_i) dy_i\end{aligned}$$

MLE IN LATENT VARIABLE MODELS

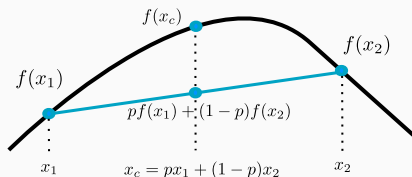
$$\begin{aligned}\ell(X|\boldsymbol{\theta}) &= \sum_{i=1}^N \log P(x_i|\boldsymbol{\theta}) = \sum_{i=1}^N \log \int P(x_i, y_i|\boldsymbol{\theta}) dy_i \\ &= \sum_{i=1}^N \log \int q_i(y_i) \frac{P(x_i, y_i|\boldsymbol{\theta})}{q_i(y_i)} dy_i \quad (\text{for arbitrary densities } q_i(y_i)) \\ &\geq \sum_{i=1}^N \int q_i(y_i) \log \frac{P(x_i, y_i|\boldsymbol{\theta})}{q_i(y_i)} dy_i \quad (\text{Jensen's inequality}) \\ &= \sum_{i=1}^N \int q_i(y_i) \log P(x_i, y_i|\boldsymbol{\theta}) dy_i - \sum_{i=1}^N \int q_i(y_i) \log q_i(y_i) dy_i \\ &= \sum_{i=1}^N \int q_i(y_i) \log P(x_i|\boldsymbol{\theta}) dy_i + \sum_{i=1}^N \int q_i(y_i) \log \frac{P(y_i|x_i, \boldsymbol{\theta})}{q_i(y_i)} dy_i\end{aligned}$$

MLE IN LATENT VARIABLE MODELS

$$\begin{aligned}\ell(X|\boldsymbol{\theta}) &= \sum_{i=1}^N \log P(x_i|\boldsymbol{\theta}) = \sum_{i=1}^N \log \int P(x_i, y_i|\boldsymbol{\theta}) dy_i \\ &= \sum_{i=1}^N \log \int q_i(y_i) \frac{P(x_i, y_i|\boldsymbol{\theta})}{q_i(y_i)} dy_i \quad (\text{for arbitrary densities } q_i(y_i)) \\ &\geq \sum_{i=1}^N \int q_i(y_i) \log \frac{P(x_i, y_i|\boldsymbol{\theta})}{q_i(y_i)} dy_i \quad (\text{Jensen's inequality}) \\ &= \sum_{i=1}^N \int q_i(y_i) \log P(x_i, y_i|\boldsymbol{\theta}) dy_i - \sum_{i=1}^N \int q_i(y_i) \log q_i(y_i) dy_i \\ &= \sum_{i=1}^N \int q_i(y_i) \log P(x_i|\boldsymbol{\theta}) dy_i + \sum_{i=1}^N \int q_i(y_i) \log \frac{P(y_i|x_i, \boldsymbol{\theta})}{q_i(y_i)} dy_i \\ &= \ell(X|\boldsymbol{\theta}) - \sum_{i=1}^N \text{KL}(q_i(y_i) \| P(y_i|X, \boldsymbol{\theta}))\end{aligned}$$

JENSEN'S INEQUALITY

Let $f(x)$ be a concave real-valued function defined on \mathbb{X} .



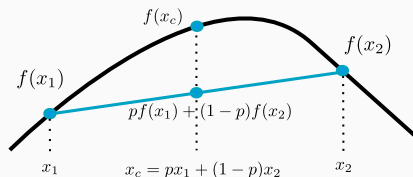
Concave: Non-positive 2nd-derivative (non-increasing deriv.)

A chord always lies below the function.

E.g. logarithm (defined on \mathbb{R}^+).

JENSEN'S INEQUALITY

Let $f(x)$ be a concave real-valued function defined on \mathbb{X} .



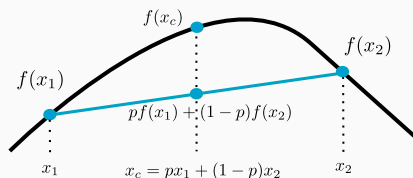
Concave: Non-positive 2nd-derivative (non-increasing deriv.)
A chord always lies below the function.

E.g. logarithm (defined on \mathbb{R}^+).

Jensen: for any prob. vector $p = (p_1, \dots, p_K)$ and any set of points (x_1, \dots, x_K) , $f(\sum_{i=1}^K p_i x_i) \geq \sum_{i=1}^K p_i f(x_i)$

JENSEN'S INEQUALITY

Let $f(x)$ be a concave real-valued function defined on \mathbb{X} .



Concave: Non-positive 2nd-derivative (non-increasing deriv.)
A chord always lies below the function.

E.g. logarithm (defined on \mathbb{R}^+).

Jensen: for any prob. vector $p = (p_1, \dots, p_K)$ and any set of points (x_1, \dots, x_K) , $f(\sum_{i=1}^K p_i x_i) \geq \sum_{i=1}^K p_i f(x_i)$

In fact, for a prob. density $p(x)$, $f(\int_{\mathbb{X}} xp(x)dx) \geq \int_{\mathbb{X}} f(x)p(x)dx$

Defining $Q(Y) = \prod_{i=1}^N q_i(y_i)$,

$$\begin{aligned}\ell(X|\boldsymbol{\theta}) &\geq \sum_{i=1}^N \int q_i(y_i) \log P(x_i, y_i|\boldsymbol{\theta}) dy_i + \sum_{i=1}^N H(q_i) \\ &= \sum_{i=1}^N \mathbb{E}_{q_i}[\log P(x_i, y_i|\boldsymbol{\theta})] + \sum_{i=1}^N H(q_i) \\ &= \ell(X|\boldsymbol{\theta}) - \sum_{i=1}^N \text{KL}(q_i(y_i) \| P(y_i|X, \boldsymbol{\theta})) \\ &:= \mathcal{F}_X(\boldsymbol{\theta}, Q(\cdot))\end{aligned}$$

Defining $Q(Y) = \prod_{i=1}^N q_i(y_i)$,

$$\begin{aligned}
 \ell(X|\boldsymbol{\theta}) &\geq \sum_{i=1}^N \int q_i(y_i) \log P(x_i, y_i|\boldsymbol{\theta}) dy_i + \sum_{i=1}^N H(q_i) \\
 &= \sum_{i=1}^N \mathbb{E}_{q_i}[\log P(x_i, y_i|\boldsymbol{\theta})] + \sum_{i=1}^N H(q_i) \\
 &= \ell(X|\boldsymbol{\theta}) - \sum_{i=1}^N \text{KL}(q_i(y_i) \| P(y_i|X, \boldsymbol{\theta})) \\
 &:= \mathcal{F}_X(\boldsymbol{\theta}, Q(\cdot))
 \end{aligned}$$

$\mathcal{F}_X(\boldsymbol{\theta}, Q(\cdot))$ is a lower bound to the log-likelihood $\ell(X|\boldsymbol{\theta})$.

Sometimes called ‘variational free energy’ and is function of $\boldsymbol{\theta}$ and the ‘variational distribution’ $Q(Y)$ (X is fixed).

OPTIMIZING THE VARIATIONAL LOWER BOUND

Our original goal was to maximize the log-likelihood:

$$\theta_{MLE} = \operatorname{argmax} \ell(X|\theta)$$

EM algorithm: maximize the lower-bound instead

$$(\theta^*, Q^*) = \operatorname{argmax} \mathcal{F}_X(\theta, Q(\cdot))$$

Hopefully easier, since all summations are outside logarithms.

Strategy: Coordinate ascent.

Alternately maximize w.r.t Q and θ

First find best lower-bound given the current θ_s .

Optimize this lower-bound to find θ_{s+1} .

Maximizing $\mathcal{F}_X(\boldsymbol{\theta}, Q)$ with $\boldsymbol{\theta}$ fixed:

- Recall $\mathcal{F}_X(\boldsymbol{\theta}, Q) = \ell(X|\boldsymbol{\theta}) - \sum_{i=1}^N \text{KL}(q_i(y_i) \| P(y_i|x_i, \boldsymbol{\theta}))$

Maximizing $\mathcal{F}_X(\boldsymbol{\theta}, Q)$ with $\boldsymbol{\theta}$ fixed:

- Recall $\mathcal{F}_X(\boldsymbol{\theta}, Q) = \ell(X|\boldsymbol{\theta}) - \sum_{i=1}^N \text{KL}(q_i(y_i) \| P(y_i|x_i, \boldsymbol{\theta}))$
- Solution: set $q_i(y_i) = P(y_i|x_i, \boldsymbol{\theta})$ for $i = 1, \dots, N$
- Recall: $P(\cdot|x_i, \boldsymbol{\theta})$ is an exponential family distribution with natural parameters $\boldsymbol{\theta}$ and feature vector $\boldsymbol{\phi}(x_i, \cdot)$

Maximizing $\mathcal{F}_X(\boldsymbol{\theta}, Q)$ with Q fixed:

$$\mathcal{F}_X(\boldsymbol{\theta}, Q) = \sum_{i=1}^N \mathbb{E}_{q_i}[\log P(x_i, y_i | \boldsymbol{\theta})] + \sum_{i=1}^N H(q_i)$$

The entropy terms $H(q_i)$ don't depend on $\boldsymbol{\theta}$. Ignore.

Maximizing $\mathcal{F}_X(\boldsymbol{\theta}, Q)$ with Q fixed:

$$\mathcal{F}_X(\boldsymbol{\theta}, Q) = \sum_{i=1}^N \mathbb{E}_{q_i}[\log P(x_i, y_i | \boldsymbol{\theta})] + \sum_{i=1}^N H(q_i)$$

The entropy terms $H(q_i)$ don't depend on $\boldsymbol{\theta}$. Ignore.

$$\log P(x_i, y_i | \boldsymbol{\theta}) = \boldsymbol{\theta}^\top \boldsymbol{\phi}(x_i, y_i) + \log h(x_i) - \log Z(\boldsymbol{\theta})$$

Maximizing $\mathcal{F}_X(\boldsymbol{\theta}, Q)$ with Q fixed:

$$\mathcal{F}_X(\boldsymbol{\theta}, Q) = \sum_{i=1}^N \mathbb{E}_{q_i}[\log P(x_i, y_i | \boldsymbol{\theta})] + \sum_{i=1}^N H(q_i)$$

The entropy terms $H(q_i)$ don't depend on $\boldsymbol{\theta}$. Ignore.

$$\log P(x_i, y_i | \boldsymbol{\theta}) = \boldsymbol{\theta}^\top \boldsymbol{\phi}(x_i, y_i) + \log h(x_i) - \log Z(\boldsymbol{\theta})$$

$$\mathcal{F}_X(\boldsymbol{\theta}, Q) = \sum_{i=1}^N \boldsymbol{\theta}^\top \mathbb{E}_{q_i}[\boldsymbol{\phi}(x_i, y_i)] - N \log Z(\boldsymbol{\theta}) + \text{const}$$

THE EXPECTATION-MAXIMIZATION ALGORITHM

$$\mathcal{F}_X(\boldsymbol{\theta}, Q) = \boldsymbol{\theta}^\top \sum_{i=1}^N \mathbb{E}_{q_i}[\boldsymbol{\phi}(x_i, y_i)] - N \log Z(\boldsymbol{\theta}) + \text{const}$$

THE EXPECTATION-MAXIMIZATION ALGORITHM

$$\mathcal{F}_X(\boldsymbol{\theta}, Q) = \boldsymbol{\theta}^\top \sum_{i=1}^N \mathbb{E}_{q_i}[\boldsymbol{\phi}(x_i, y_i)] - N \log Z(\boldsymbol{\theta}) + \text{const}$$

To maximize \mathcal{F}_X w.r.t. θ_d , solve $\frac{\partial}{\partial \theta_d} \mathcal{F}_X(\boldsymbol{\theta}, Q) = 0$:

- Solution: set θ^* to match moments (compare w. fully observed case)

$$\mathbb{E}_{\theta^*}[\phi_d(x, y)] = \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{q_i}[\phi_d(x_i, y_i)]$$

THE EXPECTATION-MAXIMIZATION ALGORITHM

$$\mathcal{F}_X(\boldsymbol{\theta}, Q) = \boldsymbol{\theta}^\top \sum_{i=1}^N \mathbb{E}_{q_i}[\boldsymbol{\phi}(x_i, y_i)] - N \log Z(\boldsymbol{\theta}) + \text{const}$$

To maximize \mathcal{F}_X w.r.t. θ_d , solve $\frac{\partial}{\partial \theta_d} \mathcal{F}_X(\boldsymbol{\theta}, Q) = 0$:

- Solution: set θ^* to match moments (compare w. fully observed case)

$$\mathbb{E}_{\theta^*}[\phi_d(x, y)] = \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{q_i}[\phi_d(x_i, y_i)]$$

$q_i(y_i) = P(y_i|X, \boldsymbol{\theta}^{old})$, an exponential family distribution whose moment parameters can be calculated (by assumption).

STEP S OF THE EM ALGORITHM

Current parameters: $\theta^s, Q^s(Y) = \prod_{i=1}^N q_i^s(y_i)$

STEP S OF THE EM ALGORITHM

Current parameters: $\theta^s, Q^s(Y) = \prod_{i=1}^N q_i^s(y_i)$

E-step:

For $i = 1, \dots, N$:

- Set $q_i^{s+1}(y_i) = P(y_i|X, \theta^s)$.

Exp. fam. distrib. with suff. stats $\phi(x_i, \cdot)$, natural params θ^s

- Calculate $\mathbb{E}_{q_i^{s+1}}[\phi(x_i, y_i)]$ (Expectation)

STEP S OF THE EM ALGORITHM

Current parameters: $\theta^s, Q^s(Y) = \prod_{i=1}^N q_i^s(y_i)$

E-step:

For $i = 1, \dots, N$:

- Set $q_i^{s+1}(y_i) = P(y_i | X, \theta^s)$.

Exp. fam. distrib. with suff. stats $\phi(x_i, \cdot)$, natural params θ^s

- Calculate $\mathbb{E}_{q_i^{s+1}}[\phi(x_i, y_i)]$ (Expectation)

M-step (Maximization):

- Set θ^{s+1} so that $\mathbb{E}_{\theta^{s+1}}[\phi(x, y)] = \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{q_i^s}[\phi(x_i, y_i)]$.

EM ALGORITHM NEVER DECREASES LOG-LIKELIHOOD

Recall $\mathcal{F}_X(\boldsymbol{\theta}^{s-1}, Q^{s-1}) = \ell(X|\boldsymbol{\theta}^{s-1}) - \text{KL}(Q(Y)||P(Y|X, \boldsymbol{\theta}^{s-1}))$.

After the E-step, $Q^s(Y) = P(Y|X, \boldsymbol{\theta}^{s-1})$

$$\ell(X|\boldsymbol{\theta}^{s-1}) = \mathcal{F}_X(\boldsymbol{\theta}^{s-1}, Q^s)$$

EM ALGORITHM NEVER DECREASES LOG-LIKELIHOOD

Recall $\mathcal{F}_X(\boldsymbol{\theta}^{s-1}, Q^{s-1}) = \ell(X|\boldsymbol{\theta}^{s-1}) - \text{KL}(Q(Y)||P(Y|X, \boldsymbol{\theta}^{s-1}))$.

After the E-step, $Q^s(Y) = P(Y|X, \boldsymbol{\theta}^{s-1})$

$$\begin{aligned}\ell(X|\boldsymbol{\theta}^{s-1}) &= \mathcal{F}_X(\boldsymbol{\theta}^{s-1}, Q^s) \\ &\leq \mathcal{F}_X(\boldsymbol{\theta}^s, Q^s)\end{aligned}$$

EM ALGORITHM NEVER DECREASES LOG-LIKELIHOOD

Recall $\mathcal{F}_X(\boldsymbol{\theta}^{s-1}, Q^{s-1}) = \ell(X|\boldsymbol{\theta}^{s-1}) - \text{KL}(Q(Y)||P(Y|X, \boldsymbol{\theta}^{s-1}))$.

After the E-step, $Q^s(Y) = P(Y|X, \boldsymbol{\theta}^{s-1})$

$$\begin{aligned}\ell(X|\boldsymbol{\theta}^{s-1}) &= \mathcal{F}_X(\boldsymbol{\theta}^{s-1}, Q^s) \\ &\leq \mathcal{F}_X(\boldsymbol{\theta}^s, Q^s) \\ &\leq \mathcal{F}_X(\boldsymbol{\theta}^s, Q^{s+1})\end{aligned}$$

EM ALGORITHM NEVER DECREASES LOG-LIKELIHOOD

Recall $\mathcal{F}_X(\boldsymbol{\theta}^{s-1}, Q^{s-1}) = \ell(X|\boldsymbol{\theta}^{s-1}) - \text{KL}(Q(Y)||P(Y|X, \boldsymbol{\theta}^{s-1}))$.

After the E-step, $Q^s(Y) = P(Y|X, \boldsymbol{\theta}^{s-1})$

$$\begin{aligned}\ell(X|\boldsymbol{\theta}^{s-1}) &= \mathcal{F}_X(\boldsymbol{\theta}^{s-1}, Q^s) \\ &\leq \mathcal{F}_X(\boldsymbol{\theta}^s, Q^s) \\ &\leq \mathcal{F}_X(\boldsymbol{\theta}^s, Q^{s+1}) \\ &= \ell(X|\boldsymbol{\theta}^s)\end{aligned}$$

EM ALGORITHM NEVER DECREASES LOG-LIKELIHOOD

Recall $\mathcal{F}_X(\boldsymbol{\theta}^{s-1}, Q^{s-1}) = \ell(X|\boldsymbol{\theta}^{s-1}) - \text{KL}(Q(Y)||P(Y|X, \boldsymbol{\theta}^{s-1}))$.

After the E-step, $Q^s(Y) = P(Y|X, \boldsymbol{\theta}^{s-1})$

$$\begin{aligned}\ell(X|\boldsymbol{\theta}^{s-1}) &= \mathcal{F}_X(\boldsymbol{\theta}^{s-1}, Q^s) \\ &\leq \mathcal{F}_X(\boldsymbol{\theta}^s, Q^s) \\ &\leq \mathcal{F}_X(\boldsymbol{\theta}^s, Q^{s+1}) \\ &= \ell(X|\boldsymbol{\theta}^s)\end{aligned}$$

Can also show that local maxima of \mathcal{F}_X are local maxima of ℓ .

EM ALGORITHM NEVER DECREASES LOG-LIKELIHOOD

Recall $\mathcal{F}_X(\boldsymbol{\theta}^{s-1}, Q^{s-1}) = \ell(X|\boldsymbol{\theta}^{s-1}) - \text{KL}(Q(Y)||P(Y|X, \boldsymbol{\theta}^{s-1}))$.

After the E-step, $Q^s(Y) = P(Y|X, \boldsymbol{\theta}^{s-1})$

$$\begin{aligned}\ell(X|\boldsymbol{\theta}^{s-1}) &= \mathcal{F}_X(\boldsymbol{\theta}^{s-1}, Q^s) \\ &\leq \mathcal{F}_X(\boldsymbol{\theta}^s, Q^s) \\ &\leq \mathcal{F}_X(\boldsymbol{\theta}^s, Q^{s+1}) \\ &= \ell(X|\boldsymbol{\theta}^s)\end{aligned}$$

Can also show that local maxima of \mathcal{F}_X are local maxima of ℓ .

Variants of E and M-steps.

EM ALGORITHM NEVER DECREASES LOG-LIKELIHOOD

Recall $\mathcal{F}_X(\boldsymbol{\theta}^{s-1}, Q^{s-1}) = \ell(X|\boldsymbol{\theta}^{s-1}) - \text{KL}(Q(Y)||P(Y|X, \boldsymbol{\theta}^{s-1}))$.

After the E-step, $Q^s(Y) = P(Y|X, \boldsymbol{\theta}^{s-1})$

$$\begin{aligned}\ell(X|\boldsymbol{\theta}^{s-1}) &= \mathcal{F}_X(\boldsymbol{\theta}^{s-1}, Q^s) \\ &\leq \mathcal{F}_X(\boldsymbol{\theta}^s, Q^s) \\ &\leq \mathcal{F}_X(\boldsymbol{\theta}^s, Q^{s+1}) \\ &= \ell(X|\boldsymbol{\theta}^s)\end{aligned}$$

Can also show that local maxima of \mathcal{F}_X are local maxima of ℓ .

Variants of E and M-steps.

Partial M-step: Update $\boldsymbol{\theta}$ to increase (rather than maximize) \mathcal{F}_X

EM ALGORITHM NEVER DECREASES LOG-LIKELIHOOD

Recall $\mathcal{F}_X(\boldsymbol{\theta}^{s-1}, Q^{s-1}) = \ell(X|\boldsymbol{\theta}^{s-1}) - \text{KL}(Q(Y)||P(Y|X, \boldsymbol{\theta}^{s-1}))$.

After the E-step, $Q^s(Y) = P(Y|X, \boldsymbol{\theta}^{s-1})$

$$\begin{aligned}\ell(X|\boldsymbol{\theta}^{s-1}) &= \mathcal{F}_X(\boldsymbol{\theta}^{s-1}, Q^s) \\ &\leq \mathcal{F}_X(\boldsymbol{\theta}^s, Q^s) \\ &\leq \mathcal{F}_X(\boldsymbol{\theta}^s, Q^{s+1}) \\ &= \ell(X|\boldsymbol{\theta}^s)\end{aligned}$$

Can also show that local maxima of \mathcal{F}_X are local maxima of ℓ .

Variants of E and M-steps.

Partial M-step: Update $\boldsymbol{\theta}$ to increase (rather than maximize) \mathcal{F}_X

Partial E-step: Update Q to decrease $\text{KL}(Q(Y)||P(Y|X, \boldsymbol{\theta}))$ (rather reduce to 0). E.g. update just one or a few q_i .

A TOY EXAMPLE:

A mixture of two Gaussians, $\mathcal{N}(x|m, 1)$ and $\mathcal{N}(x|5 - m, 2)$.
First has probability 0.6, the second 0.4.

A TOY EXAMPLE:

A mixture of two Gaussians, $\mathcal{N}(x|m, 1)$ and $\mathcal{N}(x|5 - m, 2)$.
First has probability 0.6, the second 0.4.

We observe a single data point $x = 1$.

What is the ML estimate of m ?

A TOY EXAMPLE:

A mixture of two Gaussians, $\mathcal{N}(x|m, 1)$ and $\mathcal{N}(x|5 - m, 2)$.
First has probability 0.6, the second 0.4.

We observe a single data point $x = 1$.

What is the ML estimate of m ?

What is the hidden variable?

A TOY EXAMPLE:

A mixture of two Gaussians, $\mathcal{N}(x|m, 1)$ and $\mathcal{N}(x|5 - m, 2)$.
First has probability 0.6, the second 0.4.

We observe a single data point $x = 1$.

What is the ML estimate of m ?

What is the hidden variable?

Is the overall model exponential family?

A TOY EXAMPLE:

A mixture of two Gaussians, $\mathcal{N}(x|m, 1)$ and $\mathcal{N}(x|5 - m, 2)$.
First has probability 0.6, the second 0.4.

We observe a single data point $x = 1$.

What is the ML estimate of m ?

What is the hidden variable?

Is the overall model exponential family?

What is the posterior distribution over the hidden variable?

A TOY EXAMPLE:

A mixture of two Gaussians, $\mathcal{N}(x|m, 1)$ and $\mathcal{N}(x|5 - m, 2)$.
First has probability 0.6, the second 0.4.

We observe a single data point $x = 1$.

What is the ML estimate of m ?

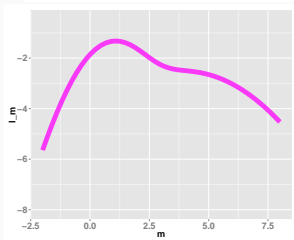
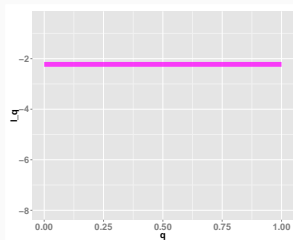
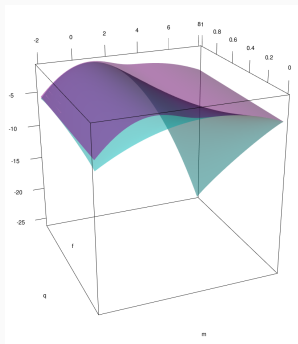
What is the hidden variable?

Is the overall model exponential family?

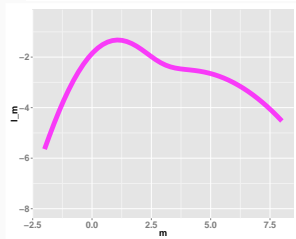
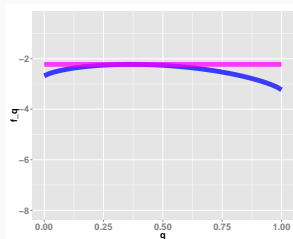
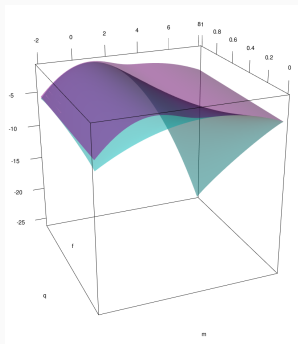
What is the posterior distribution over the hidden variable?

If we knew the hidden variable, what is the MLE?

THE EM ALGORITHM

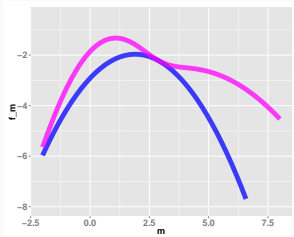
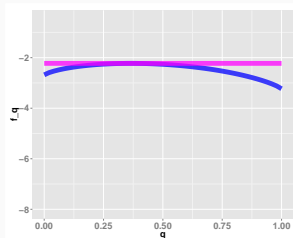
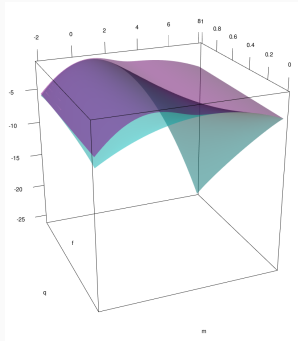


THE EM ALGORITHM



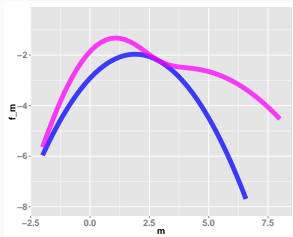
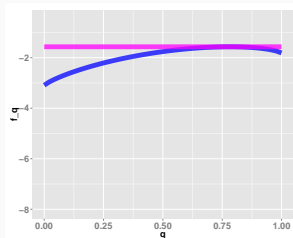
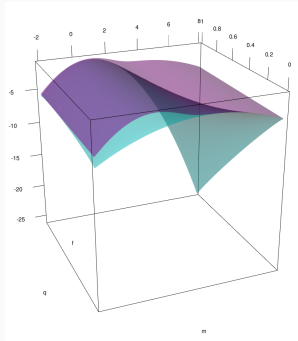
Initialize $m = 2.9$.

THE EM ALGORITHM



Initialize $m = 2.9$.
Set $q = 0.37$.

THE EM ALGORITHM

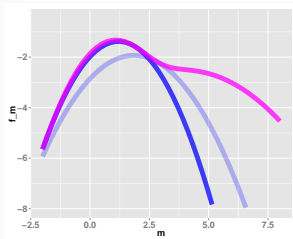
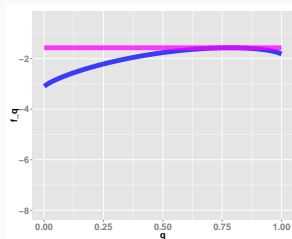
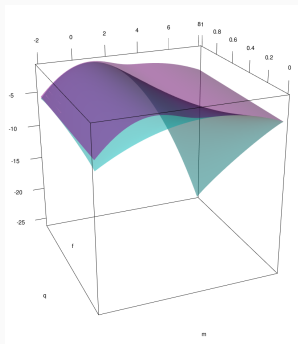


Initialize $m = 2.9$.

Set $q = 0.37$.

Set $m = 1.88$.

THE EM ALGORITHM



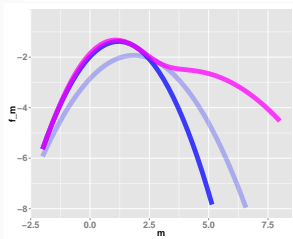
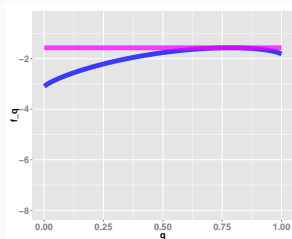
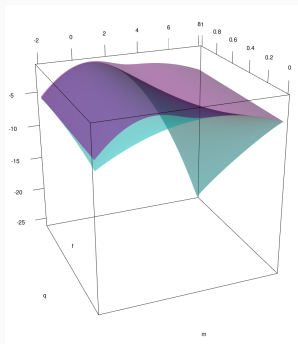
Initialize $m = 2.9$.

Set $q = 0.37$.

Set $m = 1.88$.

Set $q = 0.775$.

THE EM ALGORITHM



Initialize $m = 2.9$.

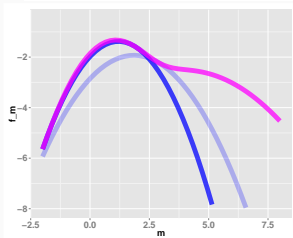
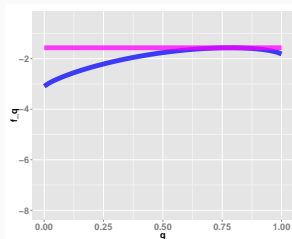
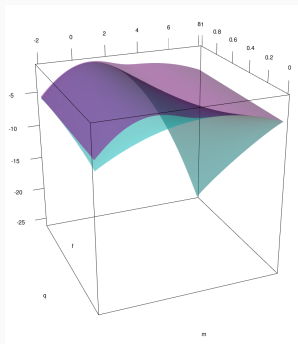
Set $q = 0.37$.

Set $m = 1.88$.

Set $q = 0.775$.

Set $m = 1.265$.

THE EM ALGORITHM



Initialize $m = 2.9$.

Set $q = 0.37$.

Set $m = 1.88$.

Set $q = 0.775$.

Set $m = 1.265$.

Repeat till convergence