

Exponential Families

①

the maximum entropy problem

X is a random variable with some distribution $p(x)$.

assume X_{i1}, \dots, X_{in} iid p , $i=1, 2, \dots, n$;

given a function $\phi_\alpha: \mathcal{X} \rightarrow \mathbb{R}$ ($X \in \mathcal{X}$)

we can compute the empirical expectation

$$\hat{\mu}_\alpha = \frac{1}{n} \sum_{i=1}^n \phi_\alpha(X_{i1}).$$

Now assume $\alpha \in I$ where I is some index set, and based on the $|I|$ -dim. empirical expectation vector

$$\hat{\mu} = (\hat{\mu}_\alpha, \alpha \in I),$$

we want to solve the following maximum entropy problem

$$p^* = \arg \max_{p \in \mathcal{P}} H(p) \quad \text{s.t.} \quad E_p[\phi_\alpha(X)] = \hat{\mu}_\alpha, \forall \alpha \in I$$

where \mathcal{P} is the set of all probability distributions over X .

Solution

We define Lagrange multipliers $\{\theta_\alpha\}$ and modify the objective:

$$p^*, \theta^* = \arg \max_{p \in \mathcal{P}, \theta \in \mathbb{R}^d} H(p) + \sum_\alpha \theta_\alpha \{ E_p[\phi_\alpha(X)] - \hat{\mu}_\alpha \} + \lambda (\sum p - 1)$$

where we denote $\theta = (\theta_\alpha, \alpha \in I)$ and the objective as $L(p, \theta, \lambda)$

Then we want $\nabla_p L(p, \theta, \lambda) = 0$

$$\Leftrightarrow p \propto \exp[\sum_\alpha \theta_\alpha \phi_\alpha(X)].$$

λ is to make p a valid distribution, and by solving θ we add in the constraints specified by the empirical expectations $\hat{\mu}$.

Formal Definition

②

Given a random variable $X = (X_1, X_2, \dots, X_m) \in \mathcal{X}^m$,
Let $\phi = (\phi_\alpha, \alpha \in I)$ be a collection of functions $\phi_\alpha: \mathcal{X}^m \rightarrow \mathbb{R}$,
and $\theta = (\theta_\alpha, \alpha \in I)$ be an associated vector. We define the
exponential family associated with ϕ as the set of the
following parameterized density functions

$$p_\theta(x_1, x_2, \dots, x_m) = \exp \{ \langle \theta, \phi(x) \rangle - A(\theta) \},$$

$$\theta \in \mathcal{L} := \{ \theta \in \mathbb{R}^d \mid A(\theta) = \log \int_{\mathcal{X}^m} \exp \{ \langle \theta, \phi(x) \rangle \} \nu dx < +\infty \}.$$

Short-hands: $\phi(x) = (\phi_\alpha(x), \alpha \in I) \in \mathbb{R}^d$, $p_\theta(x) = p_\theta(x_1, \dots, x_m)$,
 $x = (x_1, \dots, x_m)$

names of quantities:

ϕ : (log) potential functions, sufficient statistics

θ : canonical / exponential / natural parameters

$\mu = E_\theta[\phi(x)]$: mean parameters / moments

Other notions:

regular families: \mathcal{L} is an open set

minimal: there is a unique parameter θ associated with each distribution.

overcomplete: $\exists \theta \in \mathbb{R}^d$ s.t. $\langle \theta, \phi(x) \rangle$ is a constant.

Properties of $A(\theta)$

Proposition 3.1 If $A(\theta)$ is associated with some regular family, then:

(a) $\frac{\partial A}{\partial \theta_\alpha} = E_\theta[\phi_\alpha(x)]$ (we short-hand $E_{p_\theta(x)}[\cdot]$ as $E_\theta[\cdot]$)

$\frac{\partial^2 A}{\partial \theta_\alpha \partial \theta_\beta} = \text{cov}(\phi_\alpha(x), \phi_\beta(x))$

(b) A is a convex function of θ on \mathcal{L} , and strictly so if the representation is minimal.

Proof. (a) is straight forward by calculus.

(b): from (a) we have the full Hessian $\nabla^2 A$ the covariance matrix, so $\nabla^2 A \succeq 0$; also if the family is minimal, then $\text{Var}_\theta[\langle a, \phi(x) \rangle] = a^T \nabla^2 A(\theta) a > 0, \forall a \in \mathbb{R}^d, \theta \in \mathcal{L} \Rightarrow \nabla^2 A \succ 0. \quad \square$

Marginal Polytope

Recall the mean parameters' definition

$$\mu = (\mu_\alpha, \alpha \in \mathcal{I}), \mu_\alpha = E_p[\phi_\alpha(X)] \quad (\text{for any valid distribution } p)$$

We define the marginal polytope as those μ 's realizable:

$$\mathcal{M} = \{ \mu \in \mathbb{R}^d \mid \exists p, \text{ s.t. } E_p[\phi(X)] = \mu \}$$

properties

- (a) $\mathcal{M} = \text{conv} \{ \phi(X), X \in \mathcal{X}^m \}$ (convex hull) (proof by definition)
- (b) $\exists \{ (a_j, b_j) \in \mathbb{R}^d \times \mathbb{R} \mid j \in \mathcal{J} \}$ with $|\mathcal{J}| < +\infty$, s.t.
 $\mathcal{M} = \{ \mu \in \mathbb{R}^d \mid \langle a_j, \mu \rangle \geq b_j, \forall j \in \mathcal{J} \}$

example: consider an Ising model

$$p(x_1, x_2) \propto \exp \{ b_1 x_1 + b_2 x_2 + W x_1 x_2 \}, \quad x_1, x_2 \in \{0, 1\}$$

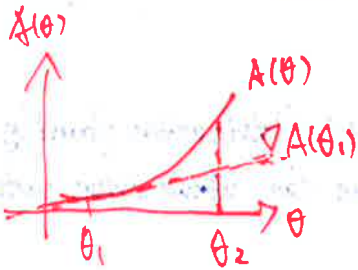
then (a) $\mu = \{ E_p[x_1], E_p[x_2], E_p[x_1 x_2] \}$
 $= \{ p(x_1=1), p(x_2=1), p(x_1=1, x_2=1) \}$
 $= p(x_1=0, x_2=0) \cdot (0, 0, 0) + p(x_1=1, x_2=0) \cdot (1, 0, 0)$
 $+ p(x_1=0, x_2=1) \cdot (0, 1, 0) + p(x_1=1, x_2=1) \cdot (1, 1, 1)$

- (b) also we need to constraint μ by
 $0 \leq p(x_1=1, x_2=1) \leq p(x_i=1), i=1, 2$
 and $1 + p(x_1=1, x_2=1) - p(x_1=1) - p(x_2=1) \geq 0$

(c) Proposition 3.2 from proposition 3.1 we have the gradient $\nabla A = \bar{E}_\theta[\phi(X)]$ defines a mapping $\nabla A: \mathcal{L} \rightarrow \mathcal{M}$; this mapping is one-to-one iff. the family is minimal.

Proof. 1) If the family is not minimal, then $\exists \gamma \in \mathbb{R}^d \setminus \{0\}$, s.t. $\langle \gamma, \phi(x) \rangle$ is a constant. We can choose some t s.t. for a given $\theta_1, \theta_2 = \theta_1 + t\gamma \in \mathcal{L}$. Then we have p_{θ_1} and p_{θ_2} induce the same distribution, and so $\nabla A(\theta_1) = \nabla A(\theta_2)$.

2) Conversely if the family is minimal, then from proposition 3.1 A is strictly convex, so we have $\forall \theta_1 \neq \theta_2$, $A(\theta_2) > A(\theta_1) + \langle \nabla A(\theta_1), \theta_2 - \theta_1 \rangle$, and by symmetry, $\langle \nabla A(\theta_1), \theta_1 - \theta_2 \rangle > \langle \nabla A(\theta_2), \theta_1 - \theta_2 \rangle > 0 \Rightarrow \nabla A(\theta_1) \neq \nabla A(\theta_2) \quad \square$



d) Theorem 3.3 If the family is minimal, then the gradient map ∇A is onto the interior M° , i.e. $\forall \mu \in M^\circ$, $\exists \theta = \theta(\mu) \in \Lambda$ s.t. $E_\theta[\phi(X)] = \mu$. ④

Conjugate Duality

From Theorem 3.3 we know that ∇A is invertible on M° , let's consider the form of this inverse mapping. Before that some more notations are needed:

Conjugate dual of $A(\theta)$:

$$A^*(\mu) = \sup_{\theta \in \Lambda} \{ \langle \mu, \theta \rangle - A(\theta) \}, \quad \mu \in \mathbb{R}^d$$

connections to the maximum likelihood problem on exponential families: consider

$D = \{X_{(1)}, X_{(2)}, \dots, X_{(N)}\}$, which induces

$\mu = \sum_{i=1}^N \phi(X_{(i)})$, then the ~~max~~ ML problem is

$$\max_{\theta \in \Lambda} \sum_{i=1}^N \log p_\theta(X_{(i)}) = N \sup_{\theta \in \Lambda} \{ \langle \mu, \theta \rangle - A(\theta) \}.$$

Solving the ML problem returns optimum $A^*(\mu)$, and the optimizer θ satisfies

$$E_\theta[\phi(X)] = \nabla A(\theta) = \mu$$

Writing $\theta = \theta(\mu)$ for the corresponded optimizer as Thm. 3.3 suggested,

Theorem 3.4 (a) $\forall \mu \in M^\circ$ and the corresponded $\theta(\mu)$

$$A^*(\mu) = -H(p_{\theta(\mu)})$$

$$(b) A(\theta) = \sup_{\mu \in M} \{ \langle \theta, \mu \rangle - A^*(\mu) \}$$

(c) the supremum in (b) is attained uniquely at $\mu = E_\theta[\phi(X)] \in M^\circ$

Remark ① We can solve the ML problem with any distribution family, but the solution in exponential family gives the ~~the~~ maximum entropy.
 ② When the family is minimal & regular,
 $(\nabla A)^{-1} = \nabla A^*$

Mean Field Methods

(5)

Consider an exponential family with $\phi = (\phi_\alpha, \alpha \in I)$ on graph $G = (V, E)$. In ^{Inference} ~~learning~~ we are interested in getting ~~$\theta = \theta(\mu)$~~ μ where θ is ~~calculated from data~~ given

Difficulties : ① from Thm. 3.4 we have $\theta = \nabla A^*(\mu)$, $\mu = \nabla A(\theta)$ but we generally don't know A^* and ∇A^* (also A)
② $A(\theta) = \sup_{\mu \in \mathcal{M}} \{ \langle \theta, \mu \rangle - A^*(\mu) \}$, but we have little idea on the nature of \mathcal{M} .

Mean field methods approximates the distribution with tractable distributions.

let F a subgraph of G ; $I(F) \subseteq I$ the subset of sufficient statistics evaluated on F , then we can define an exponential family with ϕ on graph F , where the parameters belong to

$$\mathcal{L}(F) = \{ \theta \in \mathcal{L} \mid \theta_\alpha = 0, \forall \alpha \in I \setminus I(F) \};$$

and the new "marginal" polytope defined on F :

$$\mathcal{M}_F(G; \phi) = \{ \mu \in \mathbb{R}^d \mid \mu = E_\theta[\phi(x)], \forall \theta \in \mathcal{L}(F) \}.$$

$$\star \mathcal{M}_F(G; \phi) \subseteq \mathcal{M}^\circ(G; \phi) \text{ since } \begin{cases} \nabla A(\mathcal{L}) = \mathcal{M}^\circ(G; \phi) \\ \nabla A(\mathcal{L}(F)) = \mathcal{M}_F(G; \phi) \end{cases}$$

Proposition 5.1 (Mean field lower bound)

$A(\theta) \geq \langle \theta, \mu \rangle - A^*(\mu)$ for $\forall \mu \in \mathcal{M}^\circ$; equality holds iff. $\mu = E_\theta[\phi(x)]$.

Proof.
$$\begin{aligned} A(\theta) &= \log \int x^\mu \exp \{ \langle \theta, \phi(x) \rangle \} dx \\ &\geq \int x^\mu q(x) [\langle \theta, \phi(x) \rangle - \log q(x)] dx \quad (\text{Jensen's}) \\ &= \langle \theta, \mu \rangle + H(q) \quad (\mu = E_q[\phi(x)]) \\ &= \langle \theta, \mu \rangle - A^*(\mu) \quad (\text{Thm. 3.4}) \quad \square \end{aligned}$$

It says the log-partition function is obtained by optimizing the RHS objective. However as we don't know $A^*(\mu)$ where $\mu \in \mathcal{M}^\circ$ (also the structure of \mathcal{M}°), we restrict the searching space to $\mathcal{M}_F(G; \phi)$ (short-handed as $\mathcal{M}_F(G)$). Denote $A_F^* = A^*|_{\mathcal{M}_F(G)}$, we get the best lower bound within $\mathcal{M}_F(G)$:

$$\max_{\mu \in \mathcal{M}_F(G)} \{ \langle \mu, \theta \rangle - A_F^*(\mu) \}$$

KL-divergence

(6)

We write $D(\theta_1 \| \theta_2) = D(P_{\theta_1} \| P_{\theta_2})$

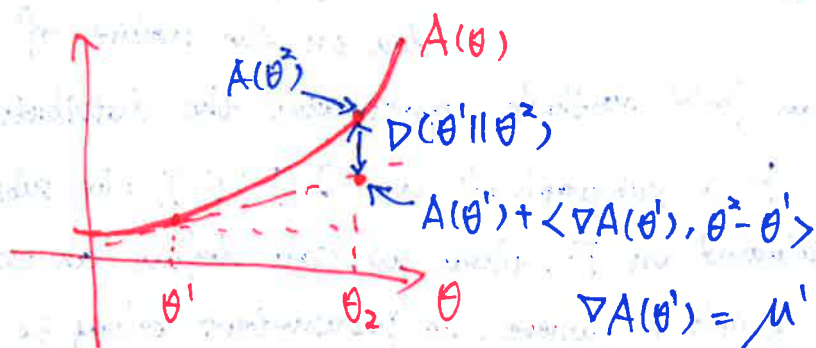
then little calculation reveals that: ($\mu^i = E_{P_{\theta^i}}[\phi(x)]$, $i=1,2$)

$$D(\theta^1 \| \theta^2) = A(\theta^2) - A(\theta^1) - \langle \mu^1, \theta^2 - \theta^1 \rangle \quad (\text{primal form})$$

$$= A(\theta^2) + A^*(\mu^1) - \langle \mu^1, \theta^2 \rangle \quad (\text{mixed form})$$

$$= A^*(\mu^1) - A^*(\mu^2) - \langle \theta^2, \mu^1 - \mu^2 \rangle \quad (\text{dual form})$$

primal form:



mixed form: as $D(\theta_1 \| \theta_2) \geq 0$, this also verifies that

$$A(\theta) \geq \langle \theta, \mu \rangle - A^*(\mu), \quad \forall \mu \in \mathcal{M}^0$$

dual form: change A to A^* in the figure of primal form.

back to mean field

$$\max_{\mu} \langle \mu, \theta \rangle - A_{F_0}^*(\mu) \Leftrightarrow \min_{\theta} D(\mu \| \theta) = A(\theta) + A_{F_0}^*(\mu) - \langle \mu, \theta \rangle$$

all subject to $\mu \in \mathcal{M}_F(G)$

example: Naive mean field for Ising model

short-hand $\mu_s = E[X_s] = P(X_s=1)$, $\mu_{st} = E[X_s X_t] = P(X_s=1, X_t=1)$

then $\mathcal{M}_{F_0}(G) = \{ \mu \in \mathbb{R}^{M+|E|} \mid 0 \leq \mu_s \leq 1, \forall s \in V, \text{ and } \mu_{st} = \mu_s \mu_t, \forall (s,t) \in E \}$

$(F_0 = (V, \phi))$

and now $-A_{F_0}^*(\mu) = H(P_{\theta}(\mu)) = -\sum_{s \in V} H_s(\mu_s)$.

the optimization problem reduces to

$$\max_{\mu} \left\{ \sum_{s \in V} \theta_s \mu_s + \sum_{(s,t) \in E} \theta_{st} \mu_s \mu_t + \sum_{s \in V} H_s(\mu_s) \right\}$$

solution by coordinate descent:

$$\mu_s \leftarrow \sigma \left(\theta_s + \sum_{t \in N(s)} \theta_{st} \mu_t \right)$$

neighbour set

Non-convexity of mean field

7

from definition of the marginal polytope

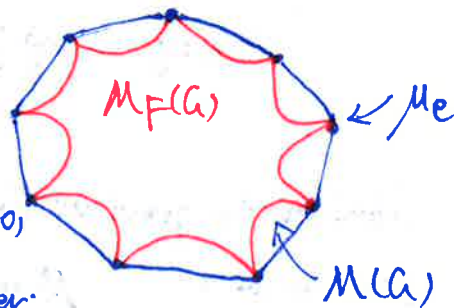
$$\mathcal{M}(\mathcal{G}) = \text{conv} \{ \phi(e), e \in \mathcal{X}^m \}$$

where e is the basis vector $(0, \dots, 0, 1, 0, \dots, 0)$

we know that μ_e is the mean parameter that the corresponded p_θ put all mass on a single element X_m .

Such a point belongs to $\overline{\mathcal{M}_F(\mathcal{G})}$ (easy to prove)

so if $\mathcal{M}_F(\mathcal{G}) \not\subseteq \mathcal{M}(\mathcal{G})$, then $\mathcal{M}_F(\mathcal{G})$ should be non-convex.



Variational Methods

maximum likelihood estimation

Given observation $D = \{X_{(1)}, \dots, X_{(N)}\}$, we can compute the empirical expectation

$$\hat{\mu} = \hat{E}[\phi(X)] = \frac{1}{N} \sum_{i=1}^N \phi(X_{(i)})$$

and the MLE problem reduces to

$$L(\theta; D) = \langle \theta, \hat{\mu} \rangle - A(\theta).$$

solution: $\hat{\theta}$ satisfies $E_{\hat{\theta}}[\phi(X)] = \hat{\mu}$

we know that $\lim_{N \rightarrow \infty} \hat{\mu} = \mu$ (Law of Large Numbers)

and MLE is an unbiased estimator:

$$\lim_{N \rightarrow \infty} \hat{\theta} = \theta \text{ where } \theta \text{ satisfies } E_{\theta}[\phi(X)] = \mu.$$

exact EM

define a joint exponential family distribution for (X, Y) :

$$p_{\theta}(x, y) = \exp \{ \langle \theta, \phi(x, y) \rangle - A(\theta) \}$$

where Y is observed and X is some hidden/latent variable:

given an observation $Y=y$, we get the conditional distribution

$$p_{\theta}(x|y) = \exp \{ \langle \theta, \phi(x, y) \rangle - A_y(\theta) \}$$

$$A_y(\theta) = \log \int_{\mathcal{X}^m} \exp \{ \langle \theta, \phi(x, y) \rangle \} dx$$

MLE problem on models with hidden variables:

$$L(\theta; y) = \log P_\theta(y) = \log \int_{\mathcal{X}} \exp\{\langle \theta, \phi(x, y) \rangle - A(\theta)\} dx$$

$$= A_y(\theta) - A(\theta)$$

Similar to previous definitions of mean parameter, marginal polytope, etc.

We define: $\mathcal{M}_y = \{\mu \in \mathbb{R}^d \mid \mu = E_p[\phi(X, y)] \text{ for some } p\}$

$$\mu_y = E_{p_\theta}[\phi(X, y)]$$

$$\Rightarrow \text{(dual)} \quad A_y(\theta) = \sup_{\mu_y \in \mathcal{M}_y} \{\langle \theta, \mu_y \rangle - A_y^*(\mu_y)\}$$

$$A_y^*(\mu) = \sup_{\theta \in \text{dom}(A_y)} \{\langle \mu_y, \theta \rangle - A_y(\theta)\}$$

\Rightarrow extension of proposition 5.1:

$$A_y(\theta) \geq \langle \mu_y, \theta \rangle - A_y^*(\mu_y) \quad , \quad \forall \mu_y \in \mathcal{M}_y$$

so we put a lower bound to the MLE problem as

$$L(\theta; y) \geq \langle \mu_y, \theta \rangle - A_y^*(\mu_y) - A(\theta) \triangleq \tilde{L}(\theta; y)$$

and update μ_y, θ with EM:

E-step: $\mu_y^{(t+1)} = \arg \max_{\mu_y \in \mathcal{M}_y} L(\mu_y, \theta^{(t)})$

$(\max_{\mu_y \in \mathcal{M}_y} \{\langle \mu_y, \theta^{(t)} \rangle - A_y^*(\mu_y)\})$ solution: $\mu_y^{(t+1)} = E_{\theta^{(t)}}[\phi(X, y)]$ on the conditional

which makes the bound tight: $L(\theta^{(t+1)}; y) = L(\theta^{(t)}; y)$

M-step: $\theta^{(t+1)} = \arg \max_{\theta \in \mathcal{A}} L(\mu_y^{(t+1)}, \theta)$

$(\max_{\theta \in \mathcal{A}} \{\langle \mu_y^{(t+1)}, \theta \rangle - A(\theta)\})$ solution: $\theta^{(t+1)} = \theta(\mu_y^{(t+1)})$, i.e. $\mu_y^{(t+1)} = E_{\theta^{(t+1)}}[\phi(X, Y)]$

increase L but also make $L \geq L$ again

important: expectation on the joint distribution

Variational EM: what if we cannot compute μ_y ?

Recall the lower bound $L(\theta; y)$, where we search optimizers in \mathcal{M}_y in E-step. Now \mathcal{M}_y is generally unknown, but we can use the subset $\mathcal{M}_{F, y}(G) \subsetneq \mathcal{M}_y$ to approximately solve E-step.

since the problem is solved in a subset, we generally have $L(\mu_y^{(t+1)}, \theta^{(t)}) \leq L(\theta^{(t)}; y)$ where $\mu_y^{(t+1)} = \arg \max_{\mu_y \in \mathcal{M}_{F, y}(G)} L(\mu_y, \theta^{(t)})$

Variational Bayes

9

In previous notes we assume θ an unknown parameter to estimate, now we introduce Bayesian inference and consider it as a random variable.

Model: let Y be the observed random variable and X be some hidden variable. Assume the joint distribution / likelihood is $p(x, y | \theta) = \exp\{\langle \eta(\theta), \phi(x, y) \rangle - A(\eta(\theta))\}$ ($\eta: \mathbb{R}^d \rightarrow \mathbb{R}^d$)

and the prior distribution is

$$P_{\xi, \lambda}(\theta) = \exp\{\langle \xi, \eta(\theta) \rangle - \lambda A(\eta(\theta)) - B(\xi, \lambda)\}$$

sufficient statistics log-partition function

Given a parametrization of the prior (ξ^*, λ^*) , Bayesian inference compute the marginal likelihood

$$\log P_{\xi^*, \lambda^*}(y) = \log \int P_{\xi^*, \lambda^*}(\theta) p(y | \theta) d\theta$$

$$= \log \int P_{\xi^*, \lambda^*}(\theta) p(y | \theta) \frac{P_{\xi, \lambda}(\theta)}{P_{\xi, \lambda}(\theta)} d\theta$$

$$\geq E_{\xi, \lambda}[\log p(y | \theta)] + E_{\xi, \lambda}\left[\log \frac{P_{\xi^*, \lambda^*}(\theta)}{P_{\xi, \lambda}(\theta)}\right]$$

$$= E_{\xi, \lambda}[A_y(\eta(\theta)) - A(\eta(\theta))] + E_{\xi, \lambda}\left[\log \frac{P_{\xi^*, \lambda^*}(\theta)}{P_{\xi, \lambda}(\theta)}\right]$$

(proposition 5.1)

$$\geq E_{\xi, \lambda}[\langle \mu(\theta), \eta(\theta) \rangle - A_y^*(\mu(\theta)) - A(\eta(\theta))] + \downarrow$$

Now we short-hand: $\bar{\eta} = E_{\xi, \lambda}[\eta(\theta)]$, $\bar{A} = E_{\xi, \lambda}[A(\eta(\theta))]$, $\mu = \mu(\theta)$ then the lower bound reduces to:

$$[\langle \mu, \bar{\eta} \rangle - A_y^*(\mu) - \bar{A}] + \langle \bar{\eta}, \xi^* - \xi \rangle + \langle -\bar{A}, \lambda^* - \lambda \rangle - B(\xi^*, \lambda) + B(\xi, \lambda),$$

also since $B^*(\bar{\eta}, \bar{A}) = \langle \bar{\eta}, \xi \rangle + \langle -\bar{A}, \lambda \rangle - B(\xi, \lambda)$, we have the lower bound changes to

$$\langle \mu + \xi^*, \bar{\eta} \rangle - A_y^*(\mu) + \langle \lambda^* + 1, -\bar{A} \rangle - B^*(\bar{\eta}, \bar{A}) \triangleq L(\mu, \bar{\eta}, \bar{A})$$

and now we can do EM on this objective.

$$\text{VB-E-step: } \mu^{(t+1)} = \arg \max_{\mu \in \mathcal{M}_y} L(\mu, \bar{\eta}^{(t)}, \bar{A}^{(t)}) = \arg \max_{\mu \in \mathcal{M}_y} \langle \mu, \bar{\eta}^{(t)} \rangle - A_y^*(\mu)$$

$$\text{solution: } \mu^{(t+1)} = E_{\bar{\eta}^{(t)}}[\phi(X, y)]$$

(10)

VB-M-step:
$$\begin{aligned}
 (\bar{\eta}^{(t+1)}, \bar{A}^{(t+1)}) &= \underset{(\bar{\eta}, \bar{A})}{\operatorname{argmax}} L(\mu^{(t+1)}, \bar{\eta}, \bar{A}) \\
 &= \underset{(\bar{\eta}, \bar{A})}{\operatorname{argmax}} \{ \langle \mu^{(t+1)} + \xi^*, \bar{\eta} \rangle - (1 + \lambda^*) \bar{A} - B^*(\bar{\eta}, \bar{A}) \}
 \end{aligned}$$

\Rightarrow then $(\bar{\eta}^{(t+1)}, \bar{A}^{(t+1)}) = E_{\xi^{(t+1)}, \lambda^{(t+1)}} [L(\bar{\eta}, \bar{A})]$

where we can get $(\xi^{(t+1)}, \lambda^{(t+1)}) = (\mu^{(t+1)} + \xi^*, \lambda^* + 1)$

In practice:

$P_{\xi, \lambda}(\theta)$ is interpreted as an approximation of $p(\theta | y)$, and we often assume family $\{P_{\xi, \lambda}(\theta)\}$ is minimal, hence the VB-M-step is valid (i.e. one-to-one correspondence between (ξ, λ) and $(\bar{\eta}, \bar{A})$).

However the solution for $\lambda \equiv \lambda^* + 1$, which means if we choose $\lambda^* = 0$ (as we often do), the posterior approximation is also intractable.

Also in computations of μ , we need a tractable $\bar{\eta} = E_{\xi, \lambda}[\eta(\theta)]$, i.e. we may want $\lambda^* + 1 = 0$, i.e. set $\lambda^* = -1$. But this is impossible: as we extend VB to N data points, update will change to $\lambda^{(t+1)} \leftarrow \lambda^* + N$, i.e. λ^* set as $-N \Rightarrow$ the prior depends on the observations, which is counter intuitive!