

# Lecture 2: Exponential Families and Variational Representation

Tianpei Xie

Aug. 25th., 2022

## Contents

<b>1</b>	<b>Background knowledge</b>	<b>2</b>
<b>2</b>	<b>Exponential family via maximum entropy</b>	<b>4</b>
2.1	Maximum entropy estimation . . . . .	4
2.2	Properties of log-partition function . . . . .	5
2.3	Conjugate Duality: Maximum Likelihood and Maximum Entropy . . . . .	6
2.4	Challenges in high dimensional setting . . . . .	8
2.5	Primal-dual formulation of KL divergence . . . . .	8

# 1 Background knowledge

- The commonly used function representation for distributions are the exponential family. The joint distribution  $p(\mathbf{x})$  follows the canonical form ***exponential famlity*** of distribution

$$\begin{aligned} p(x_1, \dots, x_m) &= p(\mathbf{x}; \boldsymbol{\eta}) = \exp(\langle \boldsymbol{\eta}, \boldsymbol{\phi}(\mathbf{x}) \rangle - A(\boldsymbol{\eta})) h(\mathbf{x}) \nu(d\mathbf{x}) \\ &= \exp\left(\sum_{\alpha} \eta_{\alpha} \phi_{\alpha}(\mathbf{x}) - A(\boldsymbol{\eta})\right) \end{aligned} \quad (1)$$

where  $\phi$  is a feature map and  $\boldsymbol{\phi}(\mathbf{x})$  defines a set of ***sufficient statistics*** (or ***potential functions***). The normalization factor is defined as

$$A(\boldsymbol{\eta}) := \log \int \exp(\langle \boldsymbol{\eta}, \boldsymbol{\phi}(\mathbf{x}) \rangle) h(\mathbf{x}) \nu(d\mathbf{x}) = \log Z(\boldsymbol{\eta})$$

$A(\boldsymbol{\eta})$  is also referred as ***log-partition function*** or *cumulant function*. The parameters  $\boldsymbol{\eta} = (\eta_{\alpha})$  are called ***natural parameters*** or *canonical parameters*.  $\boldsymbol{\eta} \in \{\boldsymbol{\eta} \in \mathbb{R}^d : A(\boldsymbol{\eta}) < \infty\}$ , which is called *natural parameter space*. Note that  $A(\boldsymbol{\eta})$  is a convex function.

In exponential family, due to property of exponent, we can formulate the (unnormalized) local functions as a exponential family too

$$\phi_C(\mathbf{x}_C; \boldsymbol{\eta}_C) = \exp\left(\sum_{k_C} \eta_{k_C} \phi_{k_C}(\mathbf{x}_C)\right)$$

Commonly known distribution and there natural parameterization:

- Bernoulli distribution  $B(x; p)$ :  $\nu =$  Counting measure,  $\eta = \log(p/(1-p))$ ,  $\phi(x) = x$

$$\begin{aligned} \langle \boldsymbol{\eta}, \boldsymbol{\phi}(x) \rangle &= \log\left(\frac{p}{1-p}\right) x \\ A(\eta) &= -\log(1-p) = \log(1 + \exp(\eta)) \end{aligned}$$

- Gaussian distribution  $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ :  $\nu =$  Lebesgue measure  $\mathbb{R}^d$ ,  $h(\mathbf{x}) = \frac{1}{(2\pi)^d}$ ,

$$\boldsymbol{\eta} = \left(\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}, -\frac{1}{2}\text{vec}(\boldsymbol{\Sigma}^{-1})\right) := \left(\boldsymbol{\theta}, -\frac{1}{2}\text{vec}(\boldsymbol{\Theta})\right) \quad (2)$$

$$\boldsymbol{\phi}(\mathbf{x}) = (\mathbf{x}, \text{vec}(\mathbf{x}\mathbf{x}^T)) \quad (3)$$

$$\begin{aligned} \langle \boldsymbol{\eta}, \boldsymbol{\phi}(x) \rangle &= \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} - \frac{1}{2} \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} = \mathbf{x}^T \boldsymbol{\theta} - \frac{1}{2} \mathbf{x}^T \boldsymbol{\Theta} \mathbf{x} \\ A(\boldsymbol{\eta}) &= \frac{1}{2} (\boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} + \log \det |\boldsymbol{\Sigma}|) = \frac{1}{2} (\boldsymbol{\theta}^T \boldsymbol{\Theta}^{-1} \boldsymbol{\theta} - \log \det |\boldsymbol{\Theta}|) \end{aligned} \quad (4)$$

- Poisson distribution  $\text{Poisson}(\lambda)$ :  $\nu =$  Counting measure  $h(x) = 1/(x!)$ ,  $\eta = \log(\lambda)$ ,  $\phi(x) = x$

$$\begin{aligned} \langle \boldsymbol{\eta}, \boldsymbol{\phi}(x) \rangle &= \log(\lambda) x \\ A(\eta) &= \lambda = \exp(\eta) \end{aligned}$$

- Gamma distribution  $\Gamma(\alpha, \lambda)$ :  $\nu = \text{Lebesgue measure } (0, \infty)$ ,  $\boldsymbol{\eta} = (-\lambda, \alpha - 1)$  and  $\boldsymbol{\phi}(x) = (x, \log(x))$

$$\begin{aligned}\langle \boldsymbol{\eta}, \boldsymbol{\phi}(x) \rangle &= -\lambda x + (\alpha - 1) \log(x) \\ A(\boldsymbol{\eta}) &= \log(\Gamma(\alpha)) - \alpha \log(\lambda) = \log(\Gamma(\eta_2 + 1)) - (\eta_2 + 1) \log(-\eta_1)\end{aligned}$$

- We can *re-parameterize* the exponential family by choosing  $\boldsymbol{\theta}$  as parameter when  $\boldsymbol{\eta} := \boldsymbol{\eta}(\boldsymbol{\theta})$ . In (1), if  $\boldsymbol{\eta} := \boldsymbol{\eta}(\boldsymbol{\theta}) = \boldsymbol{\theta}$ , we call it canonical form.

$$\begin{aligned}p(x_1, \dots, x_m) &= p(\mathbf{x}; \boldsymbol{\theta}) = \exp(\langle \boldsymbol{\eta}(\boldsymbol{\theta}), \boldsymbol{\phi}(\mathbf{x}) \rangle - A(\boldsymbol{\eta}(\boldsymbol{\theta}))) \\ &= \exp\left(\sum_{\alpha} \eta_{\alpha}(\boldsymbol{\theta}) \phi_{\alpha}(\mathbf{x}) - A(\boldsymbol{\eta}(\boldsymbol{\theta}))\right)\end{aligned}\tag{5}$$

From [Wainwright et al., 2008] we can see that the form in (1) and (5) are both *conjugate* to each other based on convex analysis.

- A special form of exponential family is the *generalized linear models (GLMs)*, when  $p(x_s|x_C)$  follows exponential family,  $\boldsymbol{\phi}(\mathbf{x}) = \mathbf{x}$ ,

$$\boldsymbol{\theta} = \mathbb{E}[\boldsymbol{\phi}(\mathbf{x})] = g^{-1}(\langle \boldsymbol{\eta}, \mathbf{x} \rangle)\tag{6}$$

where  $g$  is called the *link function*,  $\langle \boldsymbol{\eta}, \mathbf{x} \rangle$  is referred as linear predictor or system components.

- **Minimal:** It is typical to define an exponential family with a vector of sufficient statistics  $\boldsymbol{\phi}(\mathbf{x})$  for which there **does not exist** a nonzero vector  $\mathbf{a} \in \mathbb{R}^d$  such that the linear combination

$$\sum_{\alpha \in \mathcal{I}} a_{\alpha} \phi_{\alpha}(x) = \text{const.} \quad (\nu\text{-almost everywhere})$$

This condition gives rise to a so-called *minimal representation*, in which there is a unique parameter vector  $\boldsymbol{\mu}$  associated with each distribution.

## 2 Exponential family via maximum entropy

### 2.1 Maximum entropy estimation

The exponential family (1) is the unique solution to the following maximum entropy estimation problem:

$$\min_{q \in \Delta} \text{KL}(q \parallel p_0) \quad (7)$$

$$\text{s.t. } \mathbb{E}_q[\phi_\alpha(X)] = \mu_\alpha \quad \forall \alpha \in \mathcal{I} \quad (8)$$

where  $\text{KL}(q \parallel p_0) = \int \log(\frac{q}{p_0}) q dx = \mathbb{E}_q \left[ \log \frac{q}{p_0} \right]$  is the relative entropy or the Kullback-Leibler divergence of  $q$  w.r.t.  $p_0$ . To see this, we have the Lagrangian function

$$\mathcal{L}(q, \{\eta_\alpha\}) := \text{KL}(q \parallel p_0) - \sum_{\alpha} \eta_\alpha [\mathbb{E}_q[\phi_\alpha(X)] - \mu_\alpha]$$

$$\frac{\partial \mathcal{L}}{\partial q} = \log \left( \frac{q}{p_0} \right) + 1 - \sum_{\alpha} \eta_\alpha \phi_\alpha(x) = 0$$

The equation gives the exponential family in canonical form

$$q(x) = \exp \left( \sum_{\alpha} \eta_\alpha \phi_\alpha(x) - A(\boldsymbol{\eta}) \right)$$

Also note that  $\text{KL}(q \parallel p_0)$  is **convex** w.r.t.  $q$ , therefore the optimal solution is unique.

The canonical parameter  $\{\eta_\alpha\}$  forms a **canonical parameter space**

$$\Omega = \left\{ \boldsymbol{\eta} \in \mathbb{R}^d : A(\boldsymbol{\eta}) < \infty \right\} \quad (9)$$

The mean constraint (**moment matching** conditions) (8) defines a set of **mean parameters**  $\{\mu_\alpha\}_{\alpha \in \mathcal{I}}$  one for each of the  $|\mathcal{I}| = d$  sufficient statistics  $\phi_\alpha$ , with respect to an arbitrary density  $q$ . An interesting object is the set of all such vectors  $\boldsymbol{\mu} \in \mathbb{R}^d$  traced out as the underlying density  $q$  is varied. More formally, we define the **mean parameter space**

$$\mathcal{M} := \left\{ \boldsymbol{\mu} \in \mathbb{R}^d : \exists q \text{ s.t. } \mathbb{E}_q[\phi_\alpha(X)] = \mu_\alpha \quad \forall \alpha \in \mathcal{I} \right\} \quad (10)$$

We see that  $\mathcal{M}$  is the *feasible region* of the maximum entropy optimization (7). We see that  $\mathcal{M}$  is a *convex hull* spanned by sufficient statistics  $\{\phi_\alpha\}_{\alpha \in \mathcal{I}}$

$$\begin{aligned} \mathcal{M} &= \left\{ \boldsymbol{\mu} \in \mathbb{R}^d : \sum_{x \in \mathcal{X}^m} q(x) \phi_\alpha(x) = \mu_\alpha \text{ for some } \mathbf{q} \in \Delta_{|\mathcal{X}|}, \forall \alpha \in \mathcal{I} \right\} \\ &= \text{conv} \{ \phi_\alpha(x), x \in \mathcal{X}, \alpha \in \mathcal{I} \} \end{aligned}$$

It is thus a **convex polytope**.

Note that if the sufficient statistics are chosen as indicator functions of variables, the expectation constraint (8) becomes marginal distribution constraint.

$$\phi_{s;j}(x_s) = \mathbb{1} \{x_s = j\}$$

Moreover, for each edge  $(s, t)$  and pair of values  $(j, k) \in \mathcal{X} \times \mathcal{X}$ , define the sufficient statistics

$$\phi_{st;jk}(x_s, x_t) = \mathbb{1}\{x_s = j \wedge x_t = k\}$$

Thus the mean constraints becomes

$$\begin{aligned} \mathbb{E}_q[\mathbb{1}\{x_s = j \wedge x_t = k\}] &= \mu_{st;jk} = \mathbb{Q}(X_s = j \wedge X_t = k), \quad \forall s, t, j, k \\ \mathbb{E}_q[\mathbb{1}\{x_s = j\}] &= \mu_{s;j} = \mathbb{Q}(X_s = j), \quad \forall s, j \end{aligned}$$

Thus  $\mathcal{M}$  defined in (10) is also referred as the **marginal polytope** associated with the graph  $\mathcal{G}$ .

## 2.2 Properties of log-partition function

For the log-partition function (or cumulant function)  $A(\boldsymbol{\eta})$  we have the following theorem:

**Theorem 2.1** *The log-partition function  $A(\boldsymbol{\eta})$  is defined as*

$$A(\boldsymbol{\eta}) := \log \int \exp(\langle \boldsymbol{\eta}, \boldsymbol{\phi}(\mathbf{x}) \rangle) h(\mathbf{x}) \nu(d\mathbf{x}) \quad (11)$$

*associated with any regular exponential family has the following properties:*

- It has derivatives of **all orders** on its domain  $\Omega$ . The first two derivatives yield the **cumulants** of the random vector  $\boldsymbol{\phi}(X)$  as follows:

$$\frac{\partial A}{\partial \eta_\alpha} = \mathbb{E}_{\boldsymbol{\eta}}[\phi_\alpha(X)] := \int_{\mathcal{X}^m} \phi_\alpha(\mathbf{x}) q(\mathbf{x}; \boldsymbol{\eta}) d\mathbf{x} \quad (12)$$

$$\frac{\partial^2 A}{\partial \eta_\alpha \partial \eta_\beta} = \mathbb{E}_{\boldsymbol{\eta}}[\phi_\alpha(X) \phi_\beta(X)] - \mathbb{E}_{\boldsymbol{\eta}}[\phi_\alpha(X)] \mathbb{E}_{\boldsymbol{\eta}}[\phi_\beta(X)] \quad (13)$$

- Moreover,  $A$  is a **convex** function of  $\boldsymbol{\eta}$  on its domain  $\Omega$ , and strictly so if the representation is **minimal**.

From (12), we see that the **gradient of log-partition function** define a mapping  $\nabla A : \Omega \rightarrow \mathcal{M}$  from canonical parameters  $\boldsymbol{\eta}$  to the mean parameters  $\boldsymbol{\mu}$ . This mapping is called **forward mapping** [Wainwright et al., 2008]. The forward mapping  $\nabla A$  is **one-to-one mapping** when the exponential family is minimal.

**Proposition 2.2** *The gradient mapping  $\nabla A : \Omega \rightarrow \mathcal{M}$  is one-to-one if and only if the exponential representation is **minimal**.*

In particular, we say that the pair  $(\boldsymbol{\eta}, \boldsymbol{\mu})$  are **dual coupled** if  $\nabla A(\boldsymbol{\eta}) = \boldsymbol{\mu}$ .

We now consider the image  $\nabla A(\Omega)$  of the domain of valid canonical parameters  $\Omega$  under the gradient mapping  $\nabla A$ . We have the following theorem:

**Theorem 2.3** *In a minimal exponential family, the gradient map  $\nabla A$  is onto the **interior** of  $\mathcal{M}$ , denoted by  $\mathcal{M}^\circ$ . Consequently, for each  $\boldsymbol{\mu} \in \mathcal{M}^\circ$ , there exists some  $\boldsymbol{\eta} = \boldsymbol{\eta}(\boldsymbol{\mu}) \in \Omega$  such that  $\mathbb{E}_{\boldsymbol{\eta}}[\boldsymbol{\phi}(X)] = \boldsymbol{\mu}$*

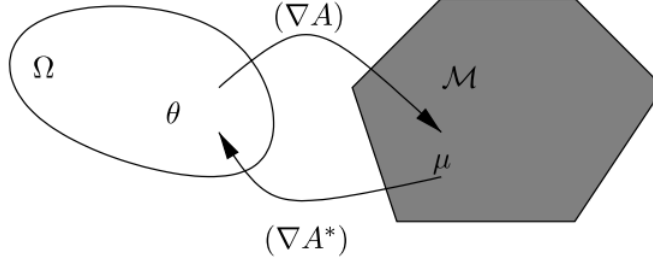


Fig. 3.8 Idealized illustration of the relation between the set  $\Omega$  of valid canonical parameters, and the set  $\mathcal{M}$  of valid mean parameters. The gradient mappings  $\nabla A$  and  $\nabla A^*$  associated with the conjugate dual pair  $(A, A^*)$  provide a bijective mapping between  $\Omega$  and the interior  $\mathcal{M}^\circ$ .

**Figure 1: The forward and backward mapping between the canonical parameter region to the marginal polytope.**

This implies that besides the *canonical parameterization* via  $\eta$ , the exponential family also has an equivalent parameterization: the *mean parameterization*:

$$p(\mathbf{x}; \boldsymbol{\mu}) := p(\mathbf{x}; \boldsymbol{\eta}(\boldsymbol{\mu})) = \exp(\langle \boldsymbol{\eta}(\boldsymbol{\mu}), \boldsymbol{\phi}(\mathbf{x}) \rangle - A(\boldsymbol{\eta}(\boldsymbol{\mu})))$$

For instance  $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$  is a mean parameterization of the normal distribution.

This fact is remarkable: *it means that (disregarding boundary points) all mean parameters  $\mathcal{M}$  that are realizable by some distribution can be realized by a member of the exponential family.* From this point of view, the maximum entropy problem (7) is just to **project** the prior distribution  $p_0$  into the space of exponential families  $\mathcal{M}$ . The moment matching conditions (8) are identical to those defining the maximum likelihood problem.

## 2.3 Conjugate Duality: Maximum Likelihood and Maximum Entropy

The convex *conjugate dual* of log-partition function  $A$  is defined as

$$A^*(\boldsymbol{\mu}) := \sup_{\boldsymbol{\eta} \in \Omega} \{ \langle \boldsymbol{\mu}, \boldsymbol{\eta} \rangle - A(\boldsymbol{\eta}) \} \quad (14)$$

Here  $\boldsymbol{\mu} \in \mathbb{R}^d$  is a fixed vector of so-called **dual variables** of the same dimension as  $\boldsymbol{\eta}$ .

The conjugate dual function (14) is the **negative entropy**. When  $\boldsymbol{\mu} \in \mathcal{M}^\circ$ , then

$$A^*(\boldsymbol{\mu}) = \int \log q(\mathbf{x}; \boldsymbol{\eta}(\boldsymbol{\mu})) q(\mathbf{x}; \boldsymbol{\eta}(\boldsymbol{\mu})) \nu(d\mathbf{x}) := -H(q_{\boldsymbol{\eta}(\boldsymbol{\mu})}),$$

where the  $q$  is the exponential family (1) with canonical parameter  $\boldsymbol{\eta}(\boldsymbol{\mu})$  and the moment matching conditions are met

$$\mathbb{E}_{\boldsymbol{\eta}(\boldsymbol{\mu})} [\boldsymbol{\phi}(X)] = \boldsymbol{\mu} \quad (15)$$

This fact is essential in the use of **variational methods**: it guarantees that any optimization problem involving the dual function can be reduced to an optimization problem over  $\mathcal{M}$ .

**Theorem 2.4** [Wainwright et al., 2008]

1. For any  $\boldsymbol{\mu} \in \mathcal{M}^\circ$ , denote by  $\boldsymbol{\eta}(\boldsymbol{\mu})$  the unique canonical parameter satisfying the dual matching condition (19). The conjugate dual function  $A^*$  takes the form

$$A^*(\boldsymbol{\mu}) = \begin{cases} -H(q_{\boldsymbol{\eta}(\boldsymbol{\mu})}) & \text{if } \boldsymbol{\mu} \in \mathcal{M}^\circ \\ +\infty & \text{if } \boldsymbol{\mu} \notin \overline{\mathcal{M}} \end{cases} \quad (16)$$

For any boundary point  $\boldsymbol{\mu} \in \partial\mathcal{M} := \overline{\mathcal{M}} \setminus \mathcal{M}^\circ$  we have

$$A^*(\boldsymbol{\mu}) = \lim_{n \rightarrow \infty} A^*(\boldsymbol{\mu}_n) \quad (17)$$

taken over any sequence  $\{\boldsymbol{\mu}_n\} \subset \mathcal{M}^\circ$  converging to  $\boldsymbol{\mu}$ .

2. In terms of this dual, the log-partition function has the **variational representation**

$$A(\boldsymbol{\eta}) = \sup_{\boldsymbol{\mu} \in \mathcal{M}} \{\langle \boldsymbol{\eta}, \boldsymbol{\mu} \rangle - A^*(\boldsymbol{\mu})\} \quad (18)$$

3. For all  $\boldsymbol{\eta} \in \Omega$ , the supremum in Equation (18) is attained uniquely at the vector  $\boldsymbol{\mu} \in \mathcal{M}^\circ$  specified by the moment matching conditions

$$\mathbb{E}_{\boldsymbol{\eta}}[\phi(X)] = \int_{\mathcal{X}^m} \phi(\mathbf{x}) q(\mathbf{x}; \boldsymbol{\eta}) \nu(d\mathbf{x}) = \boldsymbol{\mu} \quad (19)$$

The above theorem establishes the duality between log-partition function  $A$  and negative entropy  $A^*$ . Note that  $A^*$  is a function of mean parameter  $\boldsymbol{\mu}$  not like usual entropy as function of distribution. Moreover, the value  $-A^*(\boldsymbol{\mu})$  corresponds to the **optimal value** of maximum entropy optimization (7). Thus (18) formulate the maximum entropy estimation problem in (7). Third, above theorem also clarifies the precise nature of the **bijection** between the sets  $\Omega$  and  $\mathcal{M}^\circ$ , which holds for any minimal exponential family. In particular, the gradient mapping  $\nabla A$  maps  $\Omega$  in a one-to-one manner onto  $\mathcal{M}^\circ$ , whereas the **inverse mapping** from  $\mathcal{M}^\circ$  to  $\Omega$  is given by the gradient  $\nabla A^*$  of the dual function. The mapping  $\nabla A^* : \mathcal{M}^\circ \rightarrow \Omega$  is called **backward mapping**. See Figure 1 for an idealized illustration of this **bijection** correspondence based on the gradient mappings  $(\nabla A, \nabla A^*)$ .

With conjugate dual, we see that the **maximum likelihood estimation** problem is essentially the **dual problem** of the maximum entropy estimation (7).

$$\begin{aligned} & \max_{\boldsymbol{\eta}} \frac{1}{N} \sum_{n=1}^N \log q_{\boldsymbol{\eta}}(X_n) \\ \Rightarrow & \max_{\boldsymbol{\eta}} \langle \bar{\boldsymbol{\mu}}, \boldsymbol{\eta} \rangle - A(\boldsymbol{\eta}) \end{aligned} \quad (20)$$

where  $\bar{\boldsymbol{\mu}} = \hat{\mathbb{E}}[\phi(X)] = \frac{1}{N} \sum_{n=1}^N \phi(X_n)$  fits the moment matching conditions. (20) is the right-hand side of the conjugate dual of log-partition function  $A^*$  in (14). Thus we have one statistical interpretation of this variational problem (14):  $A^*$  is the **optimal value of the rescaled log likelihood** (20).

Also see that the gradient of log-likelihood function

$$\nabla_{\boldsymbol{\eta}} \frac{1}{N} \sum_{n=1}^N \log q_{\boldsymbol{\eta}}(X_n) = \nabla_{\boldsymbol{\eta}} (\langle \bar{\boldsymbol{\mu}}, \boldsymbol{\eta} \rangle - A(\boldsymbol{\eta}))$$

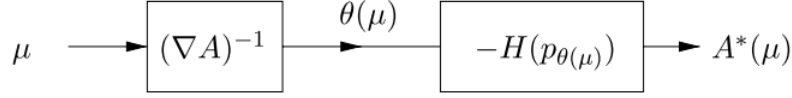


Fig. 3.9 A block diagram decomposition of  $A^*$  as the composition of two functions. Any mean parameter  $\mu \in \mathcal{M}^\circ$  is first mapped back to a canonical parameter  $\theta(\mu)$  in the inverse image  $(\nabla A)^{-1}(\mu)$ . The value of  $A^*(\mu)$  corresponds to the negative entropy  $-H(p_{\theta(\mu)})$  of the associated exponential family density  $p_{\theta(\mu)}$ .

Figure 2: The computation of  $A^*$ .

$$\begin{aligned}
&= \hat{\mathbb{E}}[\phi(X)] - \mathbb{E}_\eta[\phi(X)] \\
&= \bar{\mu} - \mu = \text{sample mean} - \text{model mean}
\end{aligned} \tag{21}$$

## 2.4 Challenges in high dimensional setting

For general multivariate exponential families, there are two **primary challenges** associated with the **variational representation**:

1. In many cases, the constraint set  $\mathcal{M}$  of realizable mean parameters is extremely difficult to characterize in an **explicit** manner. Note that even in discrete cases, the number of constraints defining  $\mathcal{M}$  grows exponentially with respect to dimension of sample space  $\mathcal{X}^m$ .
2. The negative entropy function  $A^*$  is defined indirectly in a variational manner so that it too typically **lacks an explicit form**.

To understand the complexity inherent in evaluating the dual value  $A^*(\mu)$ , note that Theorem 2.4 provides only an implicit characterization of  $A^*$  as the **composition** of mappings: first, the **inverse mapping**  $(\nabla A)^{-1} : \mathcal{M}^\circ \rightarrow \Omega$ , in which  $\mu$  maps to  $\eta(\mu)$ , corresponding to the *exponential family member* with **mean parameters**  $\mu$ ; and second, the mapping from  $\eta(\mu)$  to the negative entropy  $-H(q_{\eta(\mu)})$  of the associated exponential family density. This decomposition of the value  $A^*(\mu)$  is illustrated in Figure 2. Computing the inverse mapping  $(\nabla A)^{-1}$  as well as entropy  $-H$  are both challenging in high dimensional setting. These difficulties motivate the use of **approximations** to  $\mathcal{M}$  and  $A^*$ .

## 2.5 Primal-dual formulation of KL divergence

The conjugate duality between  $A$  and  $A^*$ , as characterized in Theorem (2.4), leads to several alternative forms of the KL divergence for *exponential family members*.

$$\text{KL}(q \parallel p) := \int_{\mathcal{X}^m} q(x) \log \left[ \frac{q(x)}{p(x)} \right] \nu(dx)$$

Consider two canonical parameter vectors  $\eta_1, \eta_2 \in \Omega$ , we formulate the KL divergence between two distributions in exponential family

$$\begin{aligned}
\text{KL}(p_{\eta_1} \parallel p_{\eta_2}) &\equiv \text{KL}(\eta_1 \parallel \eta_2) \\
&:= \mathbb{E}_{\eta_1}[A(\eta_2) - A(\eta_1) - \langle \phi(X), \eta_2 - \eta_1 \rangle]
\end{aligned}$$



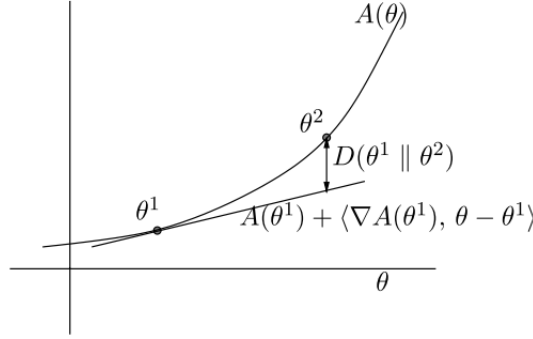


Fig. 5.1 The hyperplane  $A(\theta^1) + \langle \nabla A(\theta^1), \theta - \theta^1 \rangle$  supports the epigraph of  $A$  at  $\theta^1$ . The Kullback–Leibler divergence  $D(\theta^1 \parallel \theta^2)$  is equal to the difference between  $A(\theta^2)$  and this tangent approximation.

Figure 3: The geometrical intepretation using primal  $A$  and dual  $A^*$

$$\begin{aligned} &= A(\eta_2) - A(\eta_1) - \langle \mu_1, \eta_2 - \eta_1 \rangle \\ &\equiv A(\eta_2) - A(\eta_1) - \langle \nabla A(\eta_1), \eta_2 - \eta_1 \rangle \end{aligned} \quad (22)$$

where  $\mu_1 = \mathbb{E}_{\eta_1} [\phi(X)] = \nabla A(\eta_1)$  is the mean parameters for  $p_{\eta_1}$ . This is the **primal form of the KL divergence** at  $p_{\eta_1}$ . As illustrated in Figure 3, this form of the KL divergence can be interpreted as the difference between  $A(\eta_2)$  and the *hyperplane tangent* to  $A$  at  $\eta_1$  with normal  $\nabla A(\eta_1) = \mu_1$ .

A second form of the KL divergence can be obtained by using the strong duality condition (18) for dually coupled parameters.

$$\begin{aligned} \text{KL}(\eta_1 \parallel \eta_2) &\equiv \text{KL}(\mu_1 \parallel \eta_2) = A(\eta_2) - A(\eta_1) - \langle \mu_1, \eta_2 - \eta_1 \rangle \\ &= A(\eta_2) - (\langle \mu_1, \eta_1 \rangle - A^*(\mu_1)) - \langle \mu_1, \eta_2 - \eta_1 \rangle \\ &= A(\eta_2) + A^*(\mu_1) - \langle \mu_1, \eta_2 \rangle \end{aligned} \quad (23)$$

This is the **primal-dual mixed form of the KL divergence**.

Finally, we have the **dual form of KL divergence**:

$$\begin{aligned} \text{KL}(\mu_1 \parallel \mu_2) &= A(\eta_2) + A^*(\mu_1) - \langle \mu_1, \eta_2 \rangle \\ &= \langle \mu_2, \eta_2 \rangle - A^*(\mu_2) + A^*(\mu_1) - \langle \mu_1, \eta_2 \rangle \\ &= A^*(\mu_1) - A^*(\mu_2) - \langle \eta_2, \mu_1 - \mu_2 \rangle \\ &\equiv A^*(\mu_1) - A^*(\mu_2) - \langle \nabla A^*(\mu_2), \mu_1 - \mu_2 \rangle \end{aligned} \quad (24)$$

The dual form is related to the *Bregman divergence*, which induce the **projection operation**:

**Definition** Let  $F : \mathcal{X} \rightarrow \mathbb{R}$  be a *continuously-differentiable, strictly convex* function defined on a convex set  $\mathcal{X}$ . The **Bregman divergence** associated with  $F$  for points  $p, q \in \mathcal{X}$  is the difference between the value of  $F$  at point  $p$  and the value of the *first-order Taylor expansion* of  $F$  around point  $q$  evaluated at point  $p$ :

$$\mathbb{D}^F(p \parallel q) = F(p) - F(q) - \langle \nabla F(q), p - q \rangle \quad (25)$$

we see that dual form  $\text{KL}(\mu_1 \parallel \mu_2) = \mathbb{D}^{A^*}(\mu_1 \parallel \mu_2)$ , where  $F = A^*$  is the negative entropy.

## References

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.