



Exponential Family

Once we choose global structure and local structure of the network, we define a family of all distributions that can be attained by different parameters for this specific choice of CPDs

Let X be a set of variables. An exponential family P over X is specified by 4 components

A sufficient statistics function $T(\cdot)$ from assignments to X to \mathbb{R}^k

Sufficient statistics function

A parameter space that is a convex set $\Theta \subset \mathbb{R}^m$ of legal parameters

parameter space

A natural parameter function $t(\cdot)$: $\Theta \rightarrow \mathbb{R}^k$

legal parameter

An auxiliary measure A over X

each parameter $\theta \in \Theta$ specifies a distribution P_θ in this family as

$$P_\theta(x) = \frac{1}{Z(\theta)} A(x) \exp[\langle t(\theta), T(x) \rangle]$$

$$Z(\theta) = \sum_x A(x) \exp[\langle t(\theta), T(x) \rangle] \quad \text{the partition function, must be finite}$$

The parametric family P is defined as

$$P = \{P_\theta \mid \theta \in \Theta\}$$

The exponential family is a concise representation of a class of probability distribution that share a similar function form. A member of this family is specified by a parameter vector $\theta \in \Theta$

The sufficient statistic function $T(\cdot)$ summarizes the aspects of an instance that are relevant for assigning it a probability

The function $t(\cdot)$ maps the parameters to space of the sufficient statistics

A assigns additional preferences among instances that do not depend on the parameter (A is constant in most cases)

e.g. Bernoulli:

$$T(x) = \begin{bmatrix} 1 & \{x=1\} \\ 0 & \{x=0\} \end{bmatrix} \quad t(\theta) = \begin{bmatrix} \ln(\theta) \\ \ln(1-\theta) \end{bmatrix}$$

$$\exp[\langle T(x), t(\theta) \rangle] = \exp \left[\begin{bmatrix} 1 & \{x=1\} \\ 0 & \{x=0\} \end{bmatrix} \begin{bmatrix} \ln(\theta) \\ \ln(1-\theta) \end{bmatrix} \right] = \begin{cases} \theta & x=1 \\ 1-\theta & x=0 \end{cases}$$

e.g. Gaussian

$$P(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{(x-\mu)^2}{2\sigma^2} \right]$$

$$T(x) = \begin{bmatrix} x \\ x^2 \end{bmatrix} \quad t(\mu, \sigma^2) = \begin{bmatrix} \mu \\ \frac{\mu^2 + 1}{\sigma^2} \end{bmatrix} \quad Z(\mu, \sigma^2) = \sqrt{2\pi}\sigma \exp \left(\frac{\mu^2}{2\sigma^2} \right)$$

• Linear Exponential Family

(parameters have the same dimension as Jacob's representation $T(x)$)

θ is called the natural parameters for the given sufficient statistic function

$$p_\theta(x) = \frac{1}{Z(\theta)} \exp\{\langle \theta, T(x) \rangle\}$$

require that each $\theta \in \Theta$ gives a legal distribution

define the set of allowable natural parameters, the parameter space, to be

$$\Theta = \left\{ \theta \in \mathbb{R}^k : \int \exp(\langle \theta, T(x) \rangle) dx < \infty \right\}$$

(every choice of θ leads to a valid distribution if x is discrete)

An exponential family over the natural parameter parameter space, and for which the natural parameter space is open and convex, is called a linear exponential family
(we need to define only the function $T()$ to specify such a family)

g. Bernoulli:

in the previous example, $t(\theta) = \begin{bmatrix} \ln \theta \\ \ln(1-\theta) \end{bmatrix}$ is not a convex set

$$\text{let } T(x) = 1\{x=1\} \quad t(\theta) = \ln \frac{\theta}{1-\theta}$$

$$\exp(\langle t(\theta), T(x) \rangle) = \exp \left[\ln \frac{\theta}{1-\theta} \cdot 1\{x=1\} \right]$$

$$\exp[\langle t(\theta), T(1) \rangle] = \frac{\theta}{1-\theta} \quad Z(\theta) = \frac{\theta}{1-\theta} + 1 = \frac{1}{1-\theta}$$

$$p_\theta(x) = (\theta)^x \exp\left[\ln \frac{\theta}{1-\theta} \cdot 1\{x=1\}\right]$$

$$p_\theta(1) = \theta$$

$$p_\theta(0) = 1-\theta$$

g. Multinomial

$$t(\theta) = \langle \ln \theta_1, \ln \theta_2, \dots, \ln \theta_K \rangle$$

$$T(x) = \langle 1\{x=1\}, 1\{x=2\}, \dots, 1\{x=k\} \rangle$$

$$p_\theta(x) = \exp(\langle t(\theta), T(x) \rangle)$$

$t(\theta)$ is not a convex set

$$t(\theta) = \left\langle \ln \frac{\theta_1}{\theta_K}, \ln \frac{\theta_2}{\theta_K}, \dots, \ln \frac{\theta_{K-1}}{\theta_K} \right\rangle$$

$$T(x) = (1\{x=1\}, 1\{x=2\}, \dots, 1\{x=k\})$$

$$Z(\theta) = \sum \exp(\langle t(\theta), T(x) \rangle) = \frac{\theta_1}{\theta_K} + \frac{\theta_2}{\theta_K} + \dots + \frac{\theta_{K-1}}{\theta_K} + 1 = \frac{1}{1-\theta_K} + 1 = \frac{1}{1-\theta_K}$$

$$p_\theta(x=i) = \begin{cases} \frac{\theta_i}{\theta_K} \cdot \theta_K = \theta_i \\ \theta_K \cdot 1 = \theta_K \end{cases}$$

• Factored Exponential Family

the log-linear model is

$$p(x_1, x_2, \dots, x_n) \propto \exp \left\{ \sum_i \theta_i f_i(x_i) \right\}$$

it's a linear exponential family $\theta = (\theta_1, \dots, \theta_n)$ $T(x) = \langle f(x), -, f_1(x), \dots, f_n(x) \rangle$

By choosing appropriate features, we can derive a log-linear model to represent a given discrete Markov network structure
discrete Markov networks are linear exponential families

• Factored Exponential Family: Product Distribution

An (unnormalized) exponential factor families \mathcal{E} is defined by A, T, t, A and θ

A factor in this family is

$$\phi_i(x) = A(x) \exp \left\{ \langle t(\theta_i), T(x) \rangle \right\}$$

Let $\mathcal{E}_1, \dots, \mathcal{E}_k$ be exponential families. the family composition of $\mathcal{E}_1, \dots, \mathcal{E}_k$ is the family $\mathcal{E}_1 \times \dots \times \mathcal{E}_k$, parametrized by $\theta = \theta_1 \oplus \theta_2 \oplus \dots \oplus \theta_k \in D_1 \times D_2 \times \dots \times D_k$

$$p_\theta(x) \propto \prod_i \phi_i(x) = \left(\prod_i A_i(x) \right) \exp \left\{ \sum_i \langle t_i(\theta_i), T_i(x) \rangle \right\}$$

where ϕ_i is a factor in family \mathcal{E}_i

The composition of exponential factors is an exponential family

$$T(x) = T_1(x) \circ T_2(x) \circ \dots \circ T_k(x) \quad t(\theta) = t_1(\theta_1) \circ t_2(\theta_2) \circ \dots \circ t_k(\theta_k)$$

• Factored Exponential Family: Bayesian Network

If we have a set of CPDs from an exponential family, then their product is also in the exponential family

$$P_{\theta_i}(x_i | x_{-i}) \propto A_i(x_i) \cdot \exp \left\{ \langle t_i(\theta_i), T_i(x_i | x_{-i}) \rangle \right\}$$

$$P_\theta(x_1, \dots, x_n) \propto \prod_i P_{\theta_i}(x_i | x_{-i}) \cdot \exp \left\{ \sum_i \langle t_i(\theta_i), T_i(x_i | x_{-i}) \rangle \right\} \\ = \prod_i A_i(x_i) \cdot \exp \left\{ \langle t(\theta), T(x) \rangle \right\}$$

A Bayesian network with exponential CPDs defines an exponential family

by table-CPD $P(x|u)$

$$T_{x|u}(x|u) = \langle \{x=x, u=u\} : x \in \text{val}(x), u \in \text{val}(u) \rangle$$

$$t_{x|u}(\theta) = \langle \ln P(x|u) : x \in \text{val}(x), u \in \text{val}(u) \rangle$$

$$P(\theta|u) = \exp \left\{ \langle t_{x|u}(\theta), T_{x|u}(\theta) \rangle \right\}$$

can use similar representation to capture any table CPD / tree CPD

g. Linear Gaussian

$$X = \beta_0 + \beta_1 U_1 + \cdots + \beta_k U_k + \epsilon \quad \epsilon \sim N(0, \sigma^2)$$

$$p(x|u) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2\sigma^2} (x - (\beta_0 + \beta_1 u_1 + \cdots + \beta_k u_k))^2\right\}$$

$$\Gamma_{X|U}(x) = \langle 1, x, u_1, \dots, u_k, xu_1, \dots, xu_k, u_1^2, u_2^2, \dots, u_k^2 \rangle$$

A Bayesian network that is the product of exponential CPDs defines an exponential family

We cannot construct a Bayesian network from a product of unnormalized exponential factors

g. A CPD $P(B|A)$ with both A and B binary

$A \sim \text{Bernoulli}$:

$$T(A) = \langle \sum A = a \rangle \quad t(\theta_A) = \ln \frac{P_A}{P_{A^c}}$$

$$\exp\{\langle t(\theta_A), T(A) \rangle\} = \frac{P_A}{P_{A^c}}$$

$$\exp\{\langle t(\theta_A), T(a^*) \rangle\} = 1$$

$$Z_A(\theta_A) = 1 + \frac{P_A}{P_{A^c}} = \frac{1}{P_{A^c}}$$

$$P_{A^c}(a^*) = P_A \quad P_{A^c}(a^*) = 1 - P_A$$

$$\left. \begin{aligned} & \begin{array}{c} A \\ \downarrow \\ B \end{array} \quad \begin{array}{c} \theta_A \\ \theta_{A^c} \\ \theta_B \\ \theta_{B^c} \end{array} \\ & T_A(\theta_A) = \frac{P_A}{P_{A^c}} \\ & T_A(A) = \langle \sum A = a \rangle \end{aligned} \right\}$$

$$\begin{array}{c} A \\ \downarrow \\ B \end{array} \quad \begin{array}{c} \theta_A \\ \theta_{A^c} \\ \theta_B \\ \theta_{B^c} \end{array}$$

θ_A	B	θ_{A^c}
a^*	b^*	$B_{a^*b^*}$
a^*	b^*	$\theta_{a^*b^*}$
a^*	b^*	$B_{a^*b^*}$
a^*	b^*	$\theta_{a^*b^*}$

$B|A \sim \text{Bernoulli}$:

$$T(B|A) = \langle \sum B = b^*, A = a^* \rangle, \langle \sum B = b^*, A = a^* \rangle \rangle$$

$$T(B|A) = \langle T(B|a^*), T(B|a^*) \rangle$$

$$t(\theta) = \langle \ln \frac{P_{B|a^*}}{P_{B|a^c}}, \ln \frac{P_{B|a^*}}{P_{B|a^c}} \rangle$$

$$t(\theta) = \langle t(\theta|a^*), t(\theta|a^*) \rangle$$

$$\exp\{\langle T(B|a^*), t(\theta) \rangle\} = 1$$

$$\exp\{\langle T(B|a^*), t(\theta) \rangle\} = \frac{P_{B|a^*}}{P_{B|a^c}}$$

$$Z_{B|A}(B) = 1 + \frac{P_{B|a^*}}{P_{B|a^c}} = \frac{1}{1 - P_{B|a^c}}$$

$$P(B|a^*) = P_{B|a^*} \cdot \frac{P_{B|a^*}}{P_{B|a^c}} = P_{B|a^*}$$

$$P(B|a^*) = P_{B|a^*}|a^*$$

$$t_{AB}(\theta) = \left[\ln \frac{P_A}{P_{A^c}}, \ln \frac{P_{B|a^*}}{P_{B|a^c}}, \ln \frac{P_{B|a^*}}{P_{B|a^c}} \right]$$

$$T_{AB}(A, B) = \langle \sum A = a^*, \sum B = b^*, A = a^* \rangle, \langle \sum B = b^*, A = a^* \rangle \rangle$$

$$\exp\{\langle T_{AB}(A), T_{AB}(a^*, b^*) \rangle\} = 1$$

$$\exp\{\langle T_{AB}(A), T_{AB}(a^*, b^*) \rangle\} = \frac{P_{AB}}{Z_{AB}}$$

$$\exp\{\langle T_{AB}(A), T_{AB}(a^*, b^*) \rangle\} = \frac{P_{AB}}{Z_{AB}}$$

$$\exp\{\langle T_{AB}(A), T_{AB}(a^*, b^*) \rangle\} = \frac{P_{AB}}{Z_{AB}}$$

$$t_{B|A}(\theta_{B|A}) = \left[\ln \frac{P_{B|a^*}}{P_{B|a^c}}, \ln \frac{P_{B|a^*}}{P_{B|a^c}} \right]$$

$$T_{B|A}(B|A) = \langle \sum B = b^*, A = a^* \rangle, \langle \sum B = b^*, A = a^* \rangle \rangle$$

$\exists \theta$ is different for 2 CPDs

$$P(B|A=a^*) \quad P(B|A=a^*)$$

$$P(a^*, b^*) = \frac{1}{Z_A(B)} \cdot \frac{1}{Z_{B|A}(B)} \exp\{\langle \theta_B, b^* \rangle\} = P_{B|a^*} \cdot P_{a^*}$$

$$P(a^*, b^*) = P_{a^*} \cdot \frac{1}{Z_{B|A}(B)} \cdot \frac{1}{Z_A(B)} \exp\{\langle \theta_B, b^* \rangle\} = P_{B|a^*} \cdot P_{a^*} \cdot P_{a^*}$$

$$P(a^*, b^*) = \frac{1}{Z_A(B)} \cdot \frac{1}{Z_{B|A}(B)} \exp\{\langle \theta_B, b^* \rangle\} = P_{a^*} P_{B|a^*} \cdot \frac{1}{Z_{B|A}(B)} = P_{B|a^*} \cdot P_{a^*}$$

$$P(a^*, b^*) = P_{a^*} \cdot \frac{1}{Z_{B|A}(B)} \cdot \frac{1}{Z_A(B)} \exp\{\langle \theta_B, b^* \rangle\} = P_{B|a^*} \cdot P_{a^*}$$

To have an exponential representation of a Bayesian network, need to ensure that each CPD is locally normalized
 Can have an extra 1 in $T(x)$ and $\log \frac{1}{Z_{B|A}}$ in $t(\theta)$

$\text{g} \quad A \sim \text{Bernoulli}(p_A)$

$$\left. \begin{aligned} T(A) &= 1\{\delta = a\} & t(\theta_{01}) &= \ln \frac{\theta_{01}}{F_{01}} \\ \exp \{ \langle T(a), t(\theta_{01}) \rangle \} &= \frac{\theta_{01}}{F_{01}} \\ \exp \{ \langle T(a), t(\theta_{01}) \rangle \} &= 1 \\ Z_\theta(\theta_{01}) &= 1 + \frac{\theta_{01}}{F_{01}} = \frac{1}{F_{01}} \end{aligned} \right\} \Rightarrow \begin{aligned} T(A) &= \langle 1\{\delta = a\}, 1 \rangle \\ t(\theta_{01}) &= \langle \ln \frac{\theta_{01}}{F_{01}}, \ln(1-\theta_{01}) \rangle \\ \exp \{ \langle T(a), t(\theta_{01}) \rangle \} &= \frac{\theta_{01}}{F_{01}} \cdot (1-\theta_{01}) = \theta_{01} \\ \exp \{ \langle T(a), t(\theta_{01}) \rangle \} &= 1 \cdot (1-\theta_{01}) = 1-\theta_{01} \end{aligned}$$

$B/A \sim \text{Bernoulli};$

$$\left. \begin{aligned} T(A|B) &= \langle 1\{B=b\}, A=a \rangle, 1\{B=b\}, A=a \rangle \\ t(\theta_{01|0}, \theta_{11|0}) &= \langle \ln \frac{\theta_{01|0}}{F_{01|0}}, \ln \frac{\theta_{11|0}}{F_{11|0}} \rangle \\ \left\{ \begin{aligned} \exp \{ \langle t(\theta), T(b^*, a^*) \rangle \} &= 1 \\ \exp \{ \langle t(\theta), T(b^*, a^*) \rangle \} &= \frac{\theta_{b^*|a^*}}{F_{b^*|a^*}} \\ Z_{\theta|01|0}(B) &= 1 + \frac{\theta_{01|0}}{F_{01|0}} = \frac{1}{1-\theta_{01|0}} \end{aligned} \right. \end{aligned} \right\} \Rightarrow \begin{aligned} T_{01|0}(A|B) &= \langle 1\{B=b\}, 1\{A=a\}, 1\{B=b\}, A=a \rangle \\ t_{01|0}(\theta_{01|0}, \theta_{11|0}) &= \langle \ln \frac{\theta_{01|0}}{F_{01|0}}, \ln(1-\theta_{01|0}), \ln \frac{\theta_{11|0}}{F_{11|0}}, \ln(1-\theta_{11|0}) \rangle \\ \exp \{ \langle t_{01|0}(\cdot), T_{01|0}(a^*, b^*) \rangle \} &= 1 - \theta_{01|0} = \theta_{11|0} \\ \exp \{ \langle t_{01|0}(\cdot), T_{01|0}(a^*, b^*) \rangle \} &= \frac{\theta_{b^*|a^*}}{F_{b^*|a^*}} \cdot (1-\theta_{01|0}) = \theta_{b^*|a^*} \\ \exp \{ \langle t_{01|0}(\cdot), T_{01|0}(a^*, b^*) \rangle \} &> 1 - \theta_{01|0} = \theta_{b^*|a^*} \\ \exp \{ \langle t_{01|0}(\cdot), T_{01|0}(a^*, b^*) \rangle \} &= \frac{\theta_{01|0}}{F_{01|0}} \cdot (1-\theta_{01|0}) = \theta_{01|0} \end{aligned}$$

$$T_{AB}(A|B) = \langle 1\{\delta = a\}, 1\{B=b\}, A=a \rangle, 1\{A=a\}, 1\{B=b\}, A=a \rangle$$

$$t_{AB}(\theta_{01}, \theta_{11|0}, \theta_{01|0}) = \langle \ln \frac{\theta_{01}}{F_{01}}, \ln(1-\theta_{01}), \ln \frac{\theta_{11|0}}{F_{11|0}}, \ln(1-\theta_{11|0}), \ln \frac{\theta_{01|0}}{F_{01|0}}, \ln(1-\theta_{01|0}) \rangle$$

$$P_{AB}(A|B) = \exp \{ \langle T_{AB}(A|B), t_{AB}(\theta_{01}, \theta_{11|0}, \theta_{01|0}) \rangle \}$$

Though a Bayesian network with suitable CPDs defines an exponential family, this family is not generally a linear one.

Any Bayesian network that contain immorality does not induce a linear exponential family