

# exponential families

## ● Definition (Exponential family for measures in different spaces):

Let

- $(\Omega, \mathcal{F})$  be a measurable space
- $(\mathbb{R}^n, \mathcal{B}_{\mathbb{R}^n}, \mu)$  be a measure space with  $\mathcal{B}_{\mathbb{R}^n}$  being the Borel  $\sigma$ -algebra of  $\mathbb{R}^n$  and  $\mu$  being a  $\sigma$ -finite measure (usually,  $\mu$  is the Lebesgue or counting measures)
- $\Theta \subseteq \mathbb{R}^s$
- $\mathcal{P}_\Theta = \{P_\theta : \theta \in \Theta\}$  be a family probability measures in  $\Omega$  indexed by  $\Theta$ , ie, the map  $\Theta \rightarrow \mathcal{P}_\Theta, \theta \mapsto P_\theta$  is bijective.

Suppose, for each  $\theta \in \Theta, \exists$  a random vector  $\bar{X} : \Omega \rightarrow \mathbb{R}^n$  such that  $P_\theta = \bar{X}^* \mu$  and  $\bar{X}_* P_\theta \ll \mu$ .

We call  $\mathcal{P}_\Theta$  an Exponential family if exist

- $\eta : \Theta \rightarrow \mathbb{R}^s$  (of class  $C^1(\Theta)$  when  $\Theta$  is continuous)
- $\zeta : \Theta \rightarrow \mathbb{R}$
- $T : \mathbb{R}^n \rightarrow \mathbb{R}^s$  measurable
- $h : \mathbb{R}^n \rightarrow \mathbb{R}$  measurable such that  $h(x) > 0 \forall x \in \mathbb{R}^n$ .

such that the density function may be written as

$$p_x(x; \theta) = e^{\langle \eta(\theta), T(x) \rangle - \zeta(\theta)} h(x).$$

Remarks: • When  $\Theta$  is closed and continuous,  $\eta, \zeta \in C^1(\text{Int}(\Theta))$  and this condition concerns the applicability of the ML method

- $\eta$  is called natural function
- $T$  is called natural sufficient statistics
- $\zeta$  is the log-cumulant function, responsible by the normalization of  $P_\theta$  over  $\Omega$

- $h$  is called base function, which induces a new measure  $h\mu$  (usually denoted by  $h d\mu$ ) and which connects  $P_\theta$  and  $\mu$  by the exponential term in the calculation of the probability of an event in  $\mathcal{F}$ , as represented below,

$$P_\theta(\vec{X}(A)) = \int_{\vec{X}(A)} e^{\langle \eta(\theta), T(x) \rangle - \psi(\theta)} h(x) d\mu(x), \quad \forall A \in \mathcal{F}$$

where the integral is reduced to a sum in the discrete case.

- This definition specifies the probability measures in  $\Omega$  such that the Radon-Nikodym derivative of them and  $\mu$  are given by the exponential expression.
- Usually, Exponential families are defined in  $\mathbb{R}^n$ . Then, as both  $\mu$  and  $P_\theta$  are in  $\mathbb{R}^n$ , the comparison  $P_\theta \ll \mu$  makes sense. In our case,  $\Omega$  isn't necessarily equal to  $\mathbb{R}^n$  and, consequently, we need a measurable function  $(X, \vec{X} \text{ or } [X])$  connecting them.

### Definition (Exponential families): Let

- $(\Omega, \mathcal{F}, \mu)$  measure space such that  $\mu$  is  $\sigma$ -finite and  $\Omega \subset \mathbb{R}^n$
- $h: \Omega \rightarrow \mathbb{R}$  be a measurable function,  $h(x) > 0$  almost everywhere
- $T: \Omega \rightarrow \mathbb{R}^d$  be a measurable function
- $\eta: \square \rightarrow \mathbb{R}^d$  be a measurable function ( $\in C^1(\square)$  when  $\square$  is continuous), where
  - $\square = \left\{ \theta \in \mathbb{R}^d : \psi(\theta) := \log \int_{\Omega} e^{\langle T(x), \eta(\theta) \rangle} h(x) d\mu(x) < \infty \right\}$

The family of probability densities  $\{\rho_\eta := \rho(\cdot; \theta) : \theta \in \square\}$  which elements are defined as

$$\rho(\cdot; \theta) : \Omega \longrightarrow \mathbb{R}$$

$$x \longmapsto \rho(x; \theta) = \int_{\mathbb{R}^n} e^{\langle T(x), \eta(\theta) \rangle - \psi(\theta)} h(x) d\mu(x)$$

is called a  $d$ -parameter exponential family.  $\Xi$  is called a parameter space,  $\eta$  is called log-cumulant function,  $h$  is called base function (sometimes, it presented as  $e^{K(x)}$ , where  $K$  is called carrier measure) and  $T$  is a sufficient statistics for  $\eta$ .  
 When  $\eta = \text{id}_{\Xi}$ ,  $\{\rho_{\theta} : \theta \in \Xi\}$  is called a  $d$ -parameter exponential family in the canonical form and  $\Xi$  is called the natural parameter space for that family

↳ When the support of the density depends on the parameter, we call that family "irregular". They aren't exponential families

● Why is  $T$  sufficient statistics?

● By the Fisher-Neyman factorization criteria, it's easy to see it

● A family of probability measures  $\mathcal{P} = \{P\}$  in  $\Omega \subseteq \mathbb{R}^n$  is called a  $d$ -parametric exponential family of probability measures if

- $P \ll \mu \quad \forall P \in \mathcal{P}$
- $\exists h: \Omega \rightarrow \mathbb{R}$  a measurable function,  $h(x) > 0$  almost everywhere
- $\exists T: \Omega \rightarrow \mathbb{R}^d$  be a measurable function for some  $d \in \mathbb{N}$
- $\exists \eta: \Xi \subset \mathbb{R}^d \rightarrow \mathbb{R}$  ( $\in C^1(\Xi)$ ) measurable function such that
  - $\Xi = \left\{ \theta \in \mathbb{R}^d : \eta(\theta) := \log \int_{\Omega} e^{\langle T(x), \eta(\theta) \rangle} h(x) d\mu(x) < \infty \right\}$
- $\exists$  bijection  $\Xi \rightarrow \mathcal{P}$ ,  $\theta \mapsto P_{\theta} \in \mathcal{P}$

Such that the Radon-Nykodim derivatives of  $\mu$  and  $P_{\eta}$ , for  $\theta \in \Xi$ , assume the exponential form

$$\rho(x; \theta) = e^{\langle T(x), \eta(\theta) \rangle - \eta(\theta)} h(x) \text{ almost everywhere}$$

● Examples:

● Binomial distribution:  $\theta = p$ ,  $p \in (0, 1)$ ,  $\Omega = \mathbb{N}$ , for  $n \in \mathbb{N}$

$$\rho(x; p) = \binom{n}{x} p^x (1-p)^{n-x} = \binom{n}{x} e^{x \ln\left(\frac{p}{1-p}\right) + n \ln(1-p)}$$

$\eta(p) = \ln\left(\frac{p}{1-p}\right)$	$T(x) = x$	$\eta(p) = -n \log(1-p)$	$h(x) = \binom{n}{x}$	$\mu = \mu_c$
---	------------	--------------------------	-----------------------	---------------

- Normal distribution:  $\theta = (\mu, \sigma^2)$ ,  $\mu \in \mathbb{R}$ ,  $\sigma^2 > 0$ ,  $\Omega = \mathbb{R}$

$$p(x; \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} = \frac{1}{\sqrt{2\pi}} \exp \left[ \frac{\mu}{\sigma^2} x - \frac{1}{2\sigma^2} x^2 - \left( \frac{\mu^2}{2\sigma^2} + \ln \sigma \right) \right]$$

$\eta(\theta) = \left( \frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2} \right)$	$T(x) = (x, x^2)$	$\zeta(\theta) = \frac{\mu^2}{2\sigma^2} + \ln \sigma$	$h(x) = \frac{1}{\sqrt{2\pi}}$	$\mu = \mu_{\text{Leb}}$
--	-------------------	--	--------------------------------	--------------------------

- Gamma:  $\theta = (\alpha, \beta)$ ,  $\alpha, \beta > 0$ ,  $\Omega = (0, +\infty)$

$$p(x; \theta) = \frac{x^{\alpha-1} e^{-x/\beta}}{\Gamma(\alpha) \beta^\alpha} = e^{-x/\beta} \exp \left[ \ln \left( \frac{x^{\alpha-1}}{\Gamma(\alpha) \beta^\alpha} \right) \right]$$

$$= \exp \left[ (\alpha-1) \ln(x) - \ln(\Gamma(\alpha)) - \alpha \ln \beta - x/\beta \right]$$

$$= \exp \left[ (\alpha-1) \ln(x) - x/\beta - (\ln(\Gamma(\alpha)) + \alpha \ln \beta) \right]$$

$\eta(\theta) = (\alpha-1, 1/\beta)$	$T(x) = (\ln(x), -x)$	$\zeta(\theta) = \ln(\Gamma(\alpha)) + \alpha \ln \beta$	$h(x) = 1$	$\mu = \mu_{\text{Leb}}$
--------------------------------------	-----------------------	--	------------	--------------------------

↳ Distribution maximizes the entropy with the constraints  
 $\mathbb{E}[X] = \alpha \beta$ ,  $\mathbb{E}[\ln X] = \psi(\alpha) + \ln(\beta)$ ,  $\psi$  is called digamma function

- Erlang:  $\theta = (k, \lambda)$ ,  $k \in \mathbb{N}$ ,  $\lambda \in (0, +\infty)$ ,  $\Omega = [0, +\infty)$

$$p(x; \theta) = \frac{\lambda^k x^{k-1} e^{-\lambda x}}{(k-1)!} = \exp \left[ k \ln \lambda + (k-1) \ln x - \lambda x - \ln((k-1)!) \right]$$

$\eta(\theta) = (k-1, \lambda)$	$T = (\ln x, -x)$	$\zeta(\theta) = \ln((k-1)!) - k \ln \lambda$	$h(x) = 1$	$\mu = \mu_{\text{Leb}}$
---------------------------------	-------------------	---	------------	--------------------------

↳ It's the gamma for  $k = \alpha \in \mathbb{N}$ ,  $\lambda = 1/\beta$



- Chi-squared:  $\theta = k$ ,  $k \in \mathbb{N}$ ,  $\Omega = (0, +\infty)$  ↖ degrees of freedom

$$p(x; \theta) = \frac{x^{\frac{k}{2}-1} e^{-\frac{x}{2}}}{2^{k/2} \Gamma(k/2)} = \exp \left[ \frac{k}{2} \ln(x) - \frac{k}{2} \ln(2) + \ln \Gamma(k/2) \right] e^{-\ln(x) - x/2}$$

$\eta = k/2$	$T = \ln x$	$\varphi(k) = \frac{k}{2} \ln(2) - \ln \Gamma(k/2)$	$h(x) = e^{-\ln(x) - x/2}$	$\mu = \mu_{\text{leb}}$
--------------	-------------	---	----------------------------	--------------------------

↳ It's a gamma distribution for  $\alpha = k/2$ ,  $\beta = 2$

- Beta:  $\theta = (\alpha, \beta)$ ,  $\alpha > 0$ ,  $\beta > 0$ ,  $\Omega = (0, 1)$

$$p(x; \theta) = \frac{x^{\alpha-1} (1-x)^{\beta-1}}{B(\alpha, \beta)} = \exp \left[ (\alpha-1) \ln x + (\beta-1) \ln(1-x) - \ln B(\alpha, \beta) \right]$$

$$= \exp \left[ \alpha \ln x - \ln x + \beta \ln(1-x) - \ln(1-x) - \ln B(\alpha, \beta) \right]$$

$B(\alpha, \beta) = \frac{\Gamma(\alpha) \Gamma(\beta)}{\Gamma(\alpha+\beta)}$

$\eta(\theta) = (\alpha, \beta)$	$T(x) = (\ln x, \ln(1-x))$	$\varphi(\theta) = \ln B(\alpha, \beta)$	$h(x) = x(1-x)$	$\mu = \mu_{\text{leb}}$
----------------------------------	----------------------------	--	-----------------	--------------------------

- Bernoulli:  $\theta = p$ ,  $p \in [0, 1]$ ,  $\Omega = \{0, 1\}$

$$p(x; \theta) = p^x (1-p)^{1-x} = \exp(x \ln(p) + \ln(1-p) - x \ln(1-p))$$

$\eta(\theta) = \ln(p/1-p)$	$T(x) = x$	$\varphi(p) = -\ln(1-p)$	$h(x) = 1$	$\mu = \mu_c$
-----------------------------	------------	--------------------------	------------	---------------

↳ Binomial for  $n=1$

- Poisson:  $\theta = \lambda$ ,  $\lambda \in (0, +\infty)$ ,  $\Omega = \mathbb{Z}_{\geq 0}$

$$p(x; \lambda) = \frac{\lambda^x e^{-\lambda}}{x!} = \exp[x \ln \lambda - \lambda - \ln(x!)]$$

$\eta(\theta) = \ln(\lambda)$	$T(x) = x$	$\zeta(\theta) = \lambda$	$h(x) = 1/x!$	$\mu = \mu_0$
-------------------------------	------------	---------------------------	---------------	---------------

- Exponential:  $\theta = \lambda$ ,  $\lambda > 0$ ,  $\Omega = [0, +\infty)$

$$p(x; \theta) = \lambda e^{-\lambda x} = e^{\ln \lambda - \lambda x}$$

$\eta(\theta) = -\lambda$	$T(x) = x$	$\zeta(\theta) = -\ln \lambda$	$h(x) = 1$	$\mu = \mu_{\text{lab}}$
---------------------------	------------	--------------------------------	------------	--------------------------

↳ Given a mean  $\frac{1}{\lambda}$ , it's the continuous distribution of maximum entropy

- Dirichlet:  $\theta = \vec{\alpha}$ , given  $K \in \mathbb{N} \setminus \{1\}$ ,  $\vec{\alpha} = (\alpha_1, \dots, \alpha_K) \in \mathbb{R}^K$ ,  $\alpha_i > 0 \forall 1 \leq i \leq K$   
 $\Omega = \text{Simp}(K-1) := \{x \in (0, 1]^K : \sum_{i=1}^K x_i = 1\}$

$$\begin{aligned}
 p(x; \theta) &= \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{i=1}^K x_i^{\alpha_i - 1} = \frac{1}{B(\vec{\alpha})} \prod_{i=1}^K x_i^{\alpha_i - 1} \\
 &= \exp\left(-\ln B(\vec{\alpha}) + \sum_{i=1}^K (\alpha_i - 1) \ln x_i\right) \\
 &= \exp\left(-\ln B(\vec{\alpha}) + \sum_{i=1}^K \alpha_i \ln x_i - \sum_{i=1}^K \ln x_i\right)
 \end{aligned}$$

$\eta(\theta) = \vec{\alpha}$	$T(x) = (\ln x_1, \dots, \ln x_K)$	$\zeta(\theta) = \ln B(\vec{\alpha})$	$h(x) = \prod_{i=1}^K \frac{1}{x_i}$	$\mu = \mu_{\text{lab}}$
-------------------------------	------------------------------------	---------------------------------------	--------------------------------------	--------------------------

↳ multivariate beta distribution

- Categorical:  $\theta = \vec{p}$ , given  $K \in \mathbb{N}$ ,  $\vec{p} = (p_1, \dots, p_K) \in \mathbb{R}^K$ ,  $p_i \geq 0 \forall 1 \leq i \leq K$ ,  $\sum_{i=1}^K p_i = 1$   
 $\Omega = \{1, \dots, K\}$

$$p(x; \theta) = \prod_{i=1}^K p_i^{[x=i]} = \exp \left[ \sum_{i=1}^K [x=i] \ln(p_i) \right]$$

Inversion bracket  
 $[x=i] = \begin{cases} 1, & \text{if it is true} \\ 0, & \text{otherwise} \end{cases}$

$\eta(\theta) = (\ln(p_1), \dots, \ln(p_K))$	$T(x) = ([x=1], \dots, [x=K])$	$\zeta(\theta) = 0$	$h(x) = 1$	$\mu = \mu_0$
--	--------------------------------	---------------------	------------	---------------

- Geometrical:  $\theta = p$ ,  $p \in (0, 1)$ ,  $\Omega = \mathbb{N}$   
success probability of each trial  $K \in \mathbb{N}$   
usually, it's defined for  $p \in (0, 1]$ . In that case, the geometric distributions are not an exponential family. Then, we remove the degenerated case of sure success of all the trials

$$p(x; \theta) = (1-p)^{x-1} p = \exp [x \ln(1-p) - \ln(1-p) + \ln p]$$

$\eta(\theta) = \ln(1-p)$	$T(x) = x$	$\zeta(\theta) = \ln\left(\frac{1-p}{p}\right)$	$h(x) = 1$	$\mu = \mu_0$
---------------------------	------------	---	------------	---------------

↳ Given a mean  $1/p$ , it's the continuous distribution of maximum entropy

- Von Mises:  $\theta = (\mu, \kappa)$ ,  $\mu \in \mathbb{R}$ ,  $\kappa > 0$ ,  $\Omega = [t, t+2\pi)$  for some  $t \in \mathbb{R}$  ( $\Omega \cong \mathbb{S}^1$ )

$$p(x; \theta) = \frac{e^{\kappa \cos(x-\mu)}}{2\pi I_0(\kappa)} = \frac{e^{\kappa (\cos x \cdot \cos \mu - \sin x \cdot \sin \mu)}}{2\pi I_0(\kappa)}$$

$\eta(\theta) = (\kappa \cos \mu, -\kappa \sin \mu)$	$T(x) = (\cos x, \sin x)$	$\zeta(\theta) = \ln(2\pi I_0(\kappa))$	$h(x) = 1$	$\mu = \mu_{\text{leb}}$
--	---------------------------	---	------------	--------------------------

- Inverse Gaussian:  $\theta = (\nu, \lambda)$ ,  $\nu, \lambda > 0$ ,  $\Omega = (0, +\infty)$   
mean  
shape

$$p(x; \theta) = \sqrt{\frac{\lambda}{2\pi x^3}} \exp \left[ -\frac{\lambda(x-\nu)^2}{2\nu^2 x} \right]$$

$$= \exp \left[ \frac{1}{2} \ln \lambda - \frac{1}{2} \ln(2\pi x^3) - \frac{\lambda x}{2\nu^2} - \frac{\lambda}{2x} + \frac{\lambda}{\nu} \right]$$

$\eta(\theta) = \left(\frac{\lambda}{2\nu^2}, \frac{\lambda}{2}\right)$	$T(x) = \left(-x, -\frac{1}{x}\right)$	$\psi(\theta) = -\frac{1}{2} \ln \lambda - \frac{\lambda}{\nu}$	$h(x) = \frac{1}{2} \ln(2\pi x^3)$	$\mu = \mu_{\text{Leb}}$
---	--	---	------------------------------------	--------------------------

↳ it name arises from the inverse role of these distributions in the Brownian movement

- Wishart:  $\theta = ([V], p)$ , given  $n \in \mathbb{N}$ ,  $[V] \in S_n^{++}(\mathbb{R}) \subset \mathbb{R}^{n \times n}$ ,  $p \geq n$ ,  $\Omega = \mathcal{C}_n$  the cone of symmetric positive semidefinite matrices in  $\mathbb{R}^{n \times n}$

$$p([x]; \theta) = \frac{\det([x])^{\frac{p-n-1}{2}} e^{-\frac{1}{2} \text{tr}([V]^{-1}[x])}}{2^{np/2} \det([V])^{p/2} \Gamma_n\left(\frac{p}{2}\right)}$$

$$= \exp \left[ \left( \frac{p-n-1}{2} \right) \ln(\det[x]) - \frac{np}{2} \ln 2 - \frac{p}{2} \ln(\det[V]) - \ln \left[ \Gamma_n\left(\frac{p}{2}\right) \right] - \frac{1}{2} \text{tr}([V]^{-1}[x]) \right]$$

$$= \exp \left[ \frac{p}{2} \ln(\det[x]) - \frac{n}{2} \ln(\det[x]) - \frac{\ln(\det[x])}{2} - \frac{p}{2} (n \ln 2 + \ln(\det[V])) - \ln \left[ \Gamma_n\left(\frac{p}{2}\right) \right] - \frac{1}{2} \text{tr}([V]^{-1}[x]) \right]$$

$\eta(\theta) = \left(p, -\frac{1}{2}[V]\right)$	$T([x]) = \left(\ln(\det[x]), [x]\right)$	$\psi(\theta) = \frac{p}{2} (n \ln 2 + \ln(\det[V])) + \ln \left[ \Gamma_n\left(\frac{p}{2}\right) \right]$
$h([x]) = \det[x]^{-\frac{(n+1)}{2}}$		$\mu = \mu_{\text{Leb induced}}$

$$\hookrightarrow \langle (a, [A]), (b, [B]) \rangle = ab + \text{tr}([A]^T [B])$$

- Inverse Wishart:  $\theta = ([\Psi], p)$ , given  $n \in \mathbb{N}$ ,  $[\Psi] \in S_n^{++}(\mathbb{R}) \subset \mathbb{R}^{n \times n}$ ,  $p \geq n$ ,  $\Omega = \mathcal{C}_n$  cone of symmetric positive semidefinite matrices in  $\mathbb{R}^{n \times n}$

$$p([x]; \theta) = \frac{\det([\Psi])^{p/2} \det([x])^{-\frac{(p+n+1)}{2}} e^{-\frac{1}{2} \text{tr}([\Psi][x]^{-1})}}{2^{\frac{pn}{2}} \Gamma_n\left(\frac{p}{2}\right)}$$

$$= \exp \left[ \frac{p}{2} \ln[\det([\Psi])] - \frac{p}{2} \ln(\det[x]) - \frac{(n+1)}{2} \ln(\det[x]) - \frac{pn}{2} \ln 2 - \ln \left[ \Gamma_n\left(\frac{p}{2}\right) \right] - \frac{1}{2} \text{tr}([\Psi][x]^{-1}) \right]$$

$\eta(\theta) = \left(-\frac{\mathbb{P}}{2}, \frac{1}{2}[\Psi]\right)$	$T([x]) = (\ln(\det[x]), [x])$	$h([x]) = \det[x]^{-\frac{(n+1)}{2}}$
$\varphi(\theta) = \frac{\mathbb{P}}{2} \left( n \ln 2 - \ln[\det([\Psi])] \right) + \ln\left[\Gamma_n\left(\frac{\mathbb{P}}{2}\right)\right]$		$\mu = \mu_{\text{Leb induced}}$

● **Isn't exponential**: Let's see an example of a family of probability distributions which have an "exponential format", but isn't an exponential family

● **Wrapped normal**:  $\theta = (\mu, \sigma)$ ,  $\mu \in \mathbb{R}$ ,  $\sigma > 0$ ,  $\Omega = [t, t + 2\pi) \equiv \mathbb{S}^1$  for  $t \in \mathbb{R}$

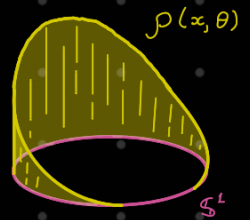
$$\rho(x; \theta) = \frac{1}{2\pi} \mathcal{V}\left(\frac{\theta - \mu}{2\pi}, \frac{i\sigma^2}{2\pi}\right) = \frac{1}{\sigma\sqrt{2\pi}} \sum_{k=-\infty}^{\infty} \exp\left[-\frac{(x - \mu + 2\pi k)^2}{2\sigma^2}\right]$$

↳ Jacobi theta function

$$\mathcal{V}(\alpha, \beta) = \sum_{k=-\infty}^{+\infty} \left[ (e^{i\pi\alpha})^{2k} \right] (e^{i\pi\beta})^{k^2}$$

↳ normal distribution wrapping  $\mathbb{S}^1$

To see that it doesn't admit the format of the exponential family, just observe the infinite sum over  $k$ . As the product  $\langle \cdot, \cdot \rangle$  is always finite, that family of density are not exponential



## Differential identities

● **Theorem**: Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be a measurable function and let  $\Xi_f$  be the set of values for  $\theta \in \mathbb{R}^d$  where

$$\int |f(x)| e^{\langle \theta, T(x) \rangle} h(x) d\mu(x) < \infty$$

for  $\mu$  measure in  $\mathbb{R}^n$  and  $h$  and  $T$  following the suppositions in the definition of exponential families. Then, the function

$$g(\theta) = \int f(x) \exp^{\langle \theta, T(x) \rangle} h(x) d\mu(x) \in C^\infty(\text{Int } \Xi_f)$$

and the derivatives can be computed by differentiation under the integral sign

**Interpretation:** Notice that we are not working necessarily with exponential families, because we don't have any condition about normalization. Otherwise, we can connect the log-cumulant term of normalization with the natural sufficient when the exponential family admits canonical form

For that, consider  $f(x) = 1 \forall x \in \Omega \subset \mathbb{R}^n$  and, by definition,  $\Xi_f = \Xi$ . Then, consider the function

$$g(\theta) = e^{\eta(\theta)} = \int e^{\langle \theta, T(x) \rangle} h(x) d\mu(x)$$

Differentiating under the integral sign if  $\theta \in \text{Int}(\Xi)$ , we obtain

$$e^{\eta(\theta)} \frac{\partial \eta(\theta)}{\partial \theta_i} = \int T_i(x) e^{\langle \theta, T(x) \rangle} h(x) d\mu(x)$$

dividing by  $e^{\eta(\theta)} \neq 0 \forall \theta$

$$\Rightarrow \frac{\partial \eta(\theta)}{\partial \theta_i} = \int T_i(x) e^{\langle \theta, T(x) \rangle - \eta(\theta)} h(x) d\mu(x)$$

We see the exponential family associated to  $\eta, h, T$  and  $\mu$  appearing.

$$\frac{\partial \eta(\theta)}{\partial \theta_i} = \int T_i(x) \rho(x; \theta) d\mu(x)$$

Consequently, if a random variable  $X: \mathbb{R}^n \rightarrow \mathbb{R}$  has density  $\rho_\theta$  with respect to  $\mu$ , then

$$\therefore \frac{\partial \eta(\theta)}{\partial \theta_i} = E_\theta[T_i(X)] \quad \forall \theta \in \text{Int } \Xi$$



Sometimes, we use the Dominated Convergence Theorem. For that, there are some useful bounds

$$\bullet |e^t - 1| \leq |t| e^{|t|} \quad \forall t \in \mathbb{R}$$

$$\bullet |t| < e^{|t|} \quad \forall t \in \mathbb{R}$$

## Moments and Cumulants

● Consider the random vector  $\vec{T} = (T_1, \dots, T_d) : \mathbb{R}^n \rightarrow \mathbb{R}^d$ . The moment-generating function, in this case, is

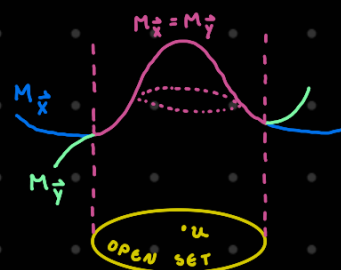
$$M_{\vec{T}}(u) = \mathbb{E}[e^{\langle u, \vec{T} \rangle}], \quad u \in \mathbb{R}^d$$

If  $M_{\vec{T}}(u)$  exists, the cumulant generating function is

$$K_{\vec{T}}(u) = \ln M_{\vec{T}}$$

**Lemma:** If the moment generating functions  $M_{\vec{X}}$  and  $M_{\vec{Y}}$  for  $\vec{X}, \vec{Y}$  random vectors are finite and coincide in some nonempty open set, then  $p_{\vec{X}} = p_{\vec{Y}}$

→ The moments, locally, determine the probability density globally



Now, let  $\vec{X} : \mathbb{R}^n \rightarrow \mathbb{R}^d$  be a random vector with exponential probability density in the canonical form. Then, considering the random vector  $T(\vec{X})$ , we have

$$\begin{aligned} \mathbb{E}_{\theta} [e^{\langle u, T(x) \rangle}] &= \int_{\Omega} e^{\langle u, T(x) \rangle} p_{\theta}(x) d\mu(x) \\ &= \int_{\Omega} e^{\langle u, T(x) \rangle} e^{\langle \theta, T(x) \rangle - \psi(\theta)} h(x) d\mu(x) \\ &= \int_{\Omega} e^{\langle u + \theta, T(x) \rangle - \psi(\theta)} h(x) d\mu(x) \end{aligned}$$



$$= \int_{\Omega} e^{\langle u, T(x) \rangle} p_{\theta}(x) d\mu(x)$$

$$= \mathbb{E}_{\theta} [e^{\langle u, T(\tilde{X}) \rangle}]$$

It express the fact of the distribution of  $T(\tilde{X})$  be induced by the distribution of  $\tilde{X}$ .

Then, in fact, the cumulant-generating function of a random vector which follows a exponential density (assuming that it admits canonical form) is a function of the log-cumulant function.

$$K_T(u) = \psi(u+\theta) - \psi(\theta), \quad u+\theta \in \text{Int}(\mathbb{Q}_d)$$

Then, a  $K$ -th cumulant for  $T$  is

$$\chi_K(T(x)) = \frac{\partial^{k_1}}{\partial u^{k_1}} \dots \frac{\partial^{k_d}}{\partial u^{k_d}} (\psi(u+\theta) - \psi(\theta))$$

$$= \frac{\partial^{k_1}}{\partial u^{k_1}} \dots \frac{\partial^{k_d}}{\partial u^{k_d}} \psi(u+\theta)$$

$$= \frac{\partial^{k_1}}{\partial \theta^{k_1}} \dots \frac{\partial^{k_d}}{\partial \theta^{k_d}} \psi(\theta), \quad K = \sum_{i=1}^d k_i$$

) changing the variables

Then,

$$\chi_K(T(\tilde{X})) = \frac{\partial^{k_1}}{\partial \theta^{k_1}} \dots \frac{\partial^{k_d}}{\partial \theta^{k_d}} \psi(\theta), \quad K = \sum_{i=1}^d k_i$$

$K$ -th cumulant for random vectors are not unique, so sometimes people write  $\chi_{k_1, \dots, k_d}$  to make explicit the which term we are taking in the Maclaurin series for multiple variables

# Curved exponential families

**Definition (Full rank exponential family):** An exponential family of densities  $p_\theta(x) = e^{\langle \eta(\theta), T(x) \rangle - \psi(\theta)} h(x)$ ,  $\theta \in \Xi$ , is said to be of full rank if

- $\text{Int } \eta(\Xi) \neq \emptyset$
  - $\nexists v \in \mathbb{R}^d \setminus \{0\}, \nexists c \in \mathbb{R}$  such that  $\langle v, T(x) \rangle = c$  a.e. with respect to  $\mu$
- ↳ The coordinates of the statistic  $T$  are "linear independent", i.e., no-one is a linear combination of the other plus a cte
- ↳ Each component of  $T(x)$  contributes with additional information about  $\theta$
- If  $\text{Int } \eta(\Xi) = \emptyset$ ,  $\eta$  can not vary freely in one open set
- ↳ problems of inference (can not apply ML method, for example)

**Theorem:** In an exponential family  $\{p_\theta = e^{\langle \eta(\theta), T(x) \rangle - \psi(\theta)} h(x) : \theta \in \Xi\}$  of full rank,  $T$  is always a complete statistic

**Definition (Curved exponential families):** Let  $\mathcal{P} = \{p_\theta : \theta \in \Xi\}$  be a full rank  $d$ -parameter canonical exponential family with complete statistic  $T$ . Consider the subfamily  $\mathcal{P}_0 \subsetneq \mathcal{P}$  of  $\mathcal{P}$  not necessarily in the canonical form, but parametrized by  $\tilde{\Xi} \subset \mathbb{R}^s$ , with  $\eta : \tilde{\Xi} \rightarrow \Xi$  being the parameter change function and  $s < d$ . Then, we can write

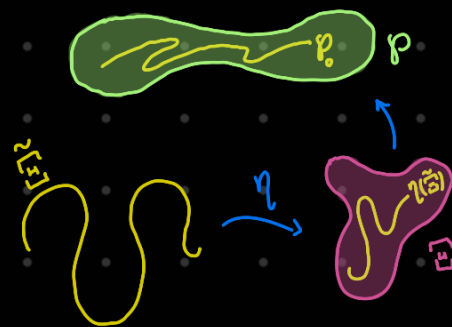
$$\mathcal{P}_0 = \{p_{\eta(\xi)} : p_{\eta(\xi)} \in \mathcal{P}, \xi \in \tilde{\Xi}\}.$$

Usually, that is the common definition for curved exponential family, but we will put more one constraint which justifies the name.

If there is one  $(d-s) \times s$  matrix  $M$  and no vector  $\alpha \in \mathbb{R}^{d-s} \setminus \{0\}$  such that

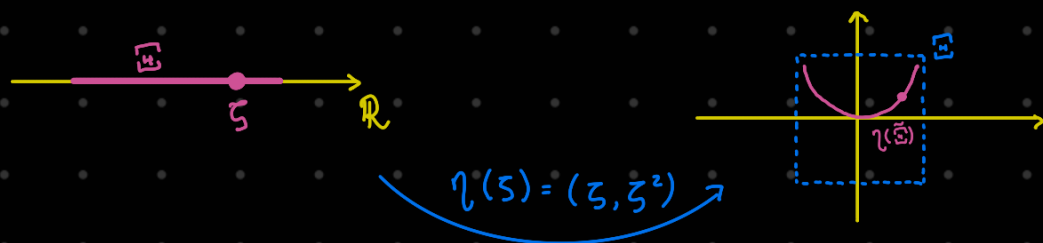
$$\eta(\tilde{\Xi}) = \{\eta(\theta) : M\eta(\theta) - \alpha = 0 \forall \theta \in \Xi\},$$

$\mathcal{P}_0$  is called a  $s$ -curved exponential family.



Interpretation:  $\eta(\tilde{\Xi})$  can't be described as a set of solutions of a linear system of equations. Consequently,  $\eta(\tilde{\Xi})$  is not a affine subset of  $\Xi$ . The consequence of the projection of  $T$  over  $\eta(\Theta)$  is that  $T$  may not be complete anymore, but it's still minimal sufficient. There is no guarantee of any function of  $T$  with cte expectation is constant.

Example of :  $\tilde{\Xi} = (-a, a), a > 1$        $\Xi = (a^2 + \epsilon, a^2 + \epsilon) \times (a^2 + \epsilon, a^2 + \epsilon), \epsilon > 0$   
 non-linearity



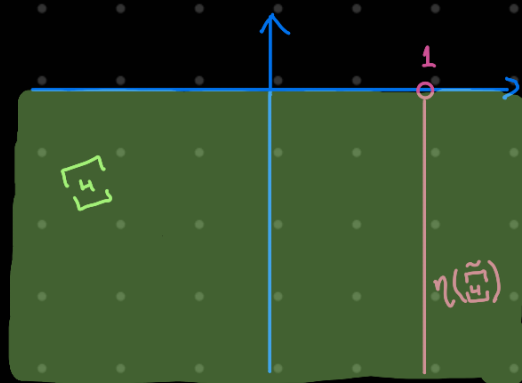
Example of linear exponential family: Consider  $\{\mathcal{N}(\mu, \sigma^2)\}$ , which the canonical parameters are  $\eta = \left(\frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2}\right)$  and it admits the complete sufficient statistic  $T(X) = (X, X^2)$ .

Consider the subfamily  $\{\mathcal{N}(\theta, \theta)\}, \theta > 0$ . Then,

$$\bullet \text{ if } \mu = \theta \text{ and } \sigma^2 = \theta \Rightarrow \begin{cases} \frac{\mu}{\sigma^2} = \frac{\theta}{\theta} = 1 \\ -\frac{1}{2\sigma^2} = -\frac{1}{2\theta} \end{cases}$$

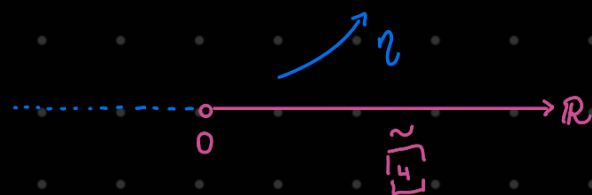
Consequently,

$$\eta : \tilde{\Xi} = \{1\} \times (0, +\infty) \rightarrow \Xi = \mathbb{R} \times (-\infty, 0) \\ (1, \theta) \mapsto \left(1, -\frac{1}{2\theta^2}\right)$$



Note it's linear. For all  $\theta \in \tilde{\Xi}$ ,  $\eta(\tilde{\Xi})$  is subset of

$$(1 \ 0) \begin{pmatrix} \tilde{\eta}_1 \\ \tilde{\eta}_2 \end{pmatrix} = 1$$

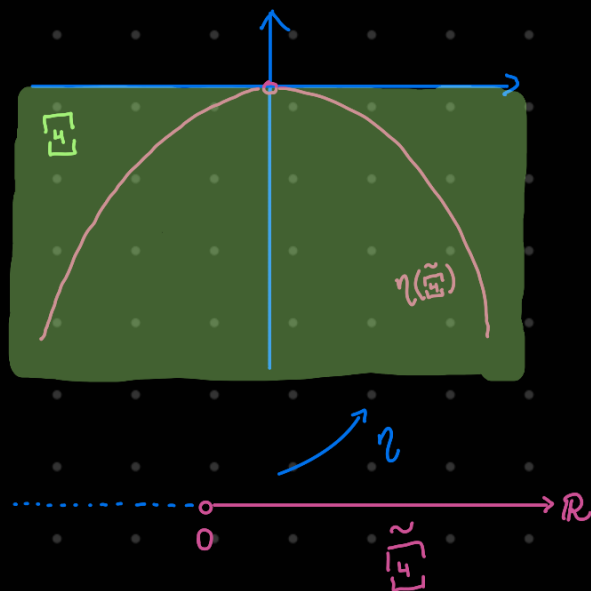


Example of 1-curved exponential family: Consider now the subfamily  $\{N(\theta, \theta^2)\}$ ,  $\theta > 0$ . Then, in this case,

$$\bullet \text{ if } \mu = \theta \text{ and } \sigma^2 = \theta \Rightarrow \begin{cases} \frac{\mu}{\sigma^2} = \frac{\theta}{\theta^2} = \frac{1}{\theta} \\ -\frac{1}{2\sigma^2} = -\frac{1}{2\theta^2} \end{cases}$$

Consequently,

$$\eta: \tilde{\mathcal{U}} = (0, +\infty) \longrightarrow \mathcal{U} = \mathbb{R} \times (-\infty, 0) \\ \theta \longmapsto \left( \frac{1}{\theta}, -\frac{1}{2\theta^2} \right)$$



## Convexity

**Theorem:** Given an exponential family  $\{p_\theta(x) = e^{\langle \theta, T(x) \rangle - \psi(\theta)} h(x)\}$  in the canonical form. The natural parameter space  $\mathcal{U}$  of this family is convex and  $\psi$  is a convex function on  $\mathcal{U}$ , i.e.,

$$\psi(\alpha\theta_1 + (1-\alpha)\theta_2) \leq \alpha\psi(\theta_1) + (1-\alpha)\psi(\theta_2), \quad \forall \alpha \in (0, 1), \forall \theta_1, \theta_2 \in \mathcal{U}, \theta_1 \neq \theta_2$$

**Proposition:** Let  $\{p_\theta(x) = e^{\langle \theta, T(x) \rangle - \psi(\theta)} h(x)\}$  be an exponential family such that  $\text{Var}_\theta [T(X)] > 0 \quad \forall \theta \in \text{Int}(\mathcal{U})$ , where  $\mathcal{U}$  is the natural space of parameters. Then,  $\mathbb{E}_\theta [T(X)]$  is a monotonic function on  $\mathcal{U}$ , i.e.,  $\text{grad}_\theta \mathbb{E}_\theta [T(X)] > 0$ .



## Maximum entropy under linear constraints

**Lemma:** Let  $(\Omega, \mathcal{F})$  be a measurable space and  $(V, \|\cdot\|_V)$  be a normed vector space. Consider  $S(\Omega)$  the vector space of all finite measures of  $\Omega$  and endowed with the total variation norm  $\|\cdot\|_S$  (it will be clear in the proof). Let  $\varphi: S(\Omega) \rightarrow V$  be a continuous linear operator. Then, exists a  $\mathcal{F}$ -measurable function  $T: \Omega \rightarrow V$  such that, for all  $\rho \in S(\Omega)$ ,

$$\varphi(\rho) = \int_{\Omega} T(\omega) d\rho(\omega)$$

**Proof:** Define the function

$$\begin{aligned} m: \mathcal{F} &\longrightarrow V \\ A &\longmapsto \varphi(\delta_A), \end{aligned}$$

where  $\delta_A$  is the characteristic function of  $A \in \mathcal{F}$ . By linearity of  $\varphi$ ,

$$m(A \cup B) = \varphi(\delta_{A \cup B}) = \varphi(\delta_A + \delta_B) = \varphi(\delta_A) + \varphi(\delta_B) = m(A) + m(B) \quad \forall A, B \in \mathcal{F} \text{ such that } A \cap B = \emptyset$$

Generalizing it, we have

$$m\left(\bigcup_{i \in \mathbb{N}} A_i\right) = \sum_{i \in \mathbb{N}} m(A_i) \quad \forall \{A_i\}_{i \in \mathbb{N}} \subset \mathcal{F} \text{ such that } A_j \cap A_k = \emptyset \quad \forall j, k \in \mathbb{N},$$

and we know the right side converges since  $\mathcal{F}$  is closed by union and  $m$  is finite. Then, in other words,  $m$  is countable additive.

Now, for each  $A \in \mathcal{F}$ , consider all the countable decomposition  $\{A_i: A = \bigcup_i A_i, A_i \in \mathcal{F}\}$  of  $A$  by measurable sets. For each one, define

$$S(A; \{A_i\}) = \sum_i \|m(A_i)\|_V = \sum_i \|\varphi(\delta_{A_i})\|_V$$

It induces the measure in  $\Omega$  above, called total variation of  $m$ :

$$|m|(A) = \sup_{\{A_i\}} S(A; \{A_i\}) \quad \forall A \in \mathcal{F}$$



$$\begin{aligned}
\varphi(\rho_n) &= \varphi(\rho_n^+) - \varphi(\rho_n^-) \\
&= \sum_{i=1}^n \left[ \rho_n^+(A_i) \varphi(\delta_{A_i}) - \rho_n^-(A_i) \varphi(\delta_{A_i}) \right] \\
&= \sum_{i=1}^n \left[ (\rho_n^+(A_i) - \rho_n^-(A_i)) \int_{A_i} T(\omega) d\mu(\omega) \right] \\
&= \sum_{i=1}^n \rho_n(A_i) \int_{A_i} T(\omega) d\mu(\omega) \\
&= \int_{\Omega} T(\omega) \left( \sum_{i=1}^n \rho_n(\omega) \delta_{A_i}(\omega) \right) d\mu(\omega)
\end{aligned}$$

Taking the limit of  $n \rightarrow \infty$ , by the convergence in  $\|\cdot\|_s$  of  $\rho_n$  to  $\rho$ , we obtain

$$\varphi(\rho) = \int_{\Omega} T(\omega) d\rho(\omega)$$

■

**Problem:** Let  $S(\Omega)$  be the space of all finite measure on the measurable space  $(\Omega, \mathcal{F})$  and consider  $\mathcal{P}_+(\Omega) \subset S(\Omega)$  the subspace of probability measures.

Let  $V$  be a vector space and  $\varphi: S(\Omega) \rightarrow V$  a linear map. Given  $p \in V$ , consider

$$S_V = S(\Omega) \cap \varphi^{-1}(p).$$

Show that  $\rho_p = \arg \max_{\rho \in S_V} H(\rho)$  form an exponential family for each  $p \in V$  (almost everywhere in  $\Omega$ )

**Proof:** Let  $p \in V$  and note we have two constraints in this case

$$i) \varphi(\rho) = p$$

$$ii) \rho_p(\Omega) = 1$$

Based on that, we can write the Lagrangian

$$\mathcal{L}(\rho, \lambda, \alpha) = H(\rho) + \langle \lambda, p - \varphi(\rho) \rangle + \alpha(1 - \rho(\Omega)),$$

where  $\lambda \in V^*$  and  $\alpha \in \mathbb{R}$  are called Lagrange multipliers. By the lemma above,

we know  $\exists T: S(\Omega) \rightarrow V$  measurable function such that



