

Lecture 2: Further Discussion on Exponential Family, Minimal Sufficiency, Ancillarity and Completeness

Lecturer: Tony Sit

Scribe: Jianting FENG and Yujia LIU

Disclaimer: These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor.

2.1 Exponential Families

Definition 2.1 The model $\{P_\theta : \theta \in \Theta\}$ forms an s -dimensional exponential family if each p_θ has the density of the form:

$$p(x; \theta) = \exp \left(\sum_{i=1}^s \eta_i(\theta) T_i(x) - B(\theta) \right) h(x),$$

where

- $\eta_i(\theta) \in \mathbb{R}$ are called the natural parameters,
- $T_i(x) \in \mathbb{R}$ are its sufficient statistics (by NFFC),
- $B(\theta)$ is the log-partition function (normalization factor)

$$B(\theta) = \log \left(\int \exp \left(\sum_{i=1}^s \eta_i(\theta) T_i(x) \right) h(x) d\mu(x) \right) \in \mathbb{R}$$

- $h(x) \in \mathbb{R}$ is the base measure.

Example 2.2 Beta distribution $P = \{\text{Beta}(\alpha, \beta) : \alpha, \beta > 0\}$, $\theta = (\alpha, \beta)$. The densities take the form

$$\begin{aligned} p(x, \theta) &= x^{\alpha-1} (1-x)^{\beta-1} I(x \in (0, 1)) \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \\ &= \exp \left((\alpha - 1) \log x + (\beta - 1) \log(1 - x) + \log \left(\frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \right) \right) I(x \in (0, 1)) \\ &= \exp \left(\underbrace{\alpha \log x}_{\eta_1(\theta)=\alpha} + \underbrace{\beta \log(1-x)}_{\eta_2(\theta)=\beta} + \underbrace{\log \left(\frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \right)}_{B(\theta)} \right) \underbrace{\frac{I(x \in (0, 1))}{x(1-x)}}_{h(x)} \end{aligned}$$

Definition 2.3 An exponential family is in canonical form when the density has the form

$$p(x; \eta) = \exp \left(\sum_{i=1}^s \eta_i T_i(x) - A(\eta) \right) h(x)$$

Definition 2.4 The set of all valid natural parameters Θ_{nat} is called the natural parameter space: for each $\eta \in \Theta_{nat}$, there exists a normalizing constant $A(\eta)$ such that $\int p(x; \eta) = 1$, or equivalently,

$$\Theta_{nat} = \left\{ \eta : 0 < \int \exp \left(\sum_{i=1}^s \eta_i T_i(x) \right) h(x) d\mu(x) < \infty \right\}$$

2.1.1 Reducing the Dimension

“Superficial dimension” \Leftarrow be a bit more careful.

Case I: The $T_i(x)$ ’s satisfy an affine equality constraint $\forall x \in \mathcal{X}$.

Example 2.5 $X \sim \text{Exp}(\eta_1, \eta_2)$, $p(x; \eta_1, \eta_2) = \exp(-\eta_1 x - \eta_2 x + \log(\eta_1 + \eta_2)) \times I(x \geq 0)$. Here $T_1(x) = T_2(x) = x$ (i.e. they are linearly dependent).

We can rewrite it as

$$p(x; \eta_1, \eta_2) = \exp(-(\eta_1 + \eta_2)x + \log(\eta_1 + \eta_2)) I(x \geq 0)$$

Definition 2.6 If $\mathcal{P} = \{P_\theta : \theta \in \Omega\}$, then θ is unidentifiable if for two parameters $\theta_1 \neq \theta_2$, $P_{\theta_1} = P_{\theta_2}$.

E.g. In the above example,

$$p(x; \eta_1 + \alpha, \eta_2 - \alpha) = p(x; \eta_1, \eta_2)$$

for any $\alpha < \eta_2$.

Case II: The η_i ’s satisfy an affine equality constraint for all $\eta \in H$.

Example 2.7

$$\begin{aligned} p(x; \eta) &\propto \exp(\eta_1 x + \eta_2 x^2) \text{ for all } (\eta_1, \eta_2) \text{ satisfying } \eta_1 + \eta_2 = 1 \\ &= \exp(\eta_1(x - x^2) + x^2) \end{aligned}$$

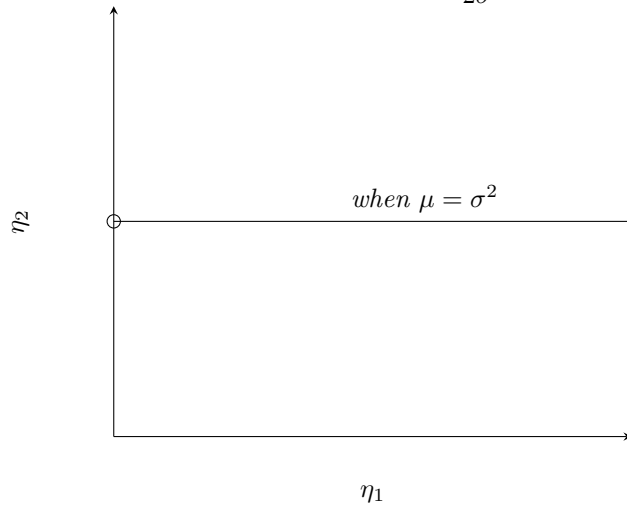
Definition 2.8 A canonical exponential family $\mathcal{P} = \{P_\eta : \eta \in H\}$ is minimal if

- $\sum_{i=1}^s \lambda_i T_i(x) = \lambda_0, \forall x \in \mathcal{X} \Rightarrow \lambda_i = 0, \forall i \in \{0, \dots, s\}$ (no affine T_i equality constraints)
- $\sum_{i=1}^s \lambda_i \eta_i = \lambda_0, \forall \eta \in H \Rightarrow \lambda_i = 0, \forall i \in \{0, \dots, s\}$ (no affine η_i equality constraints).

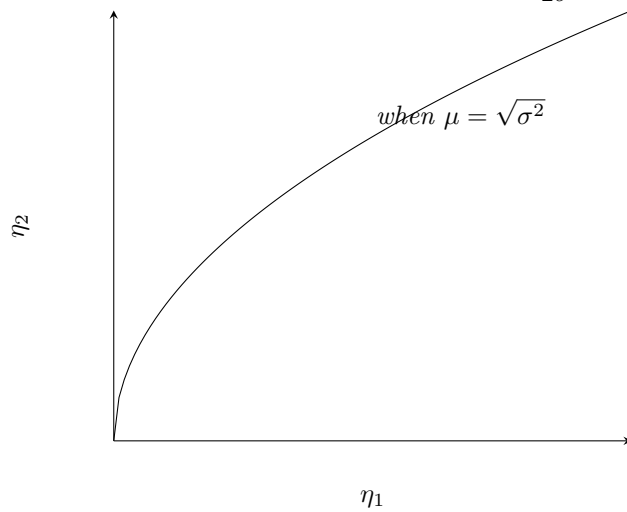
Definition 2.9 Suppose $\mathcal{P} = \{P_\eta : \eta \in H\}$ is an s -dimensional minimal exponential family. If H contains an open s -dimensional rectangle, then \mathcal{P} is called **full-rank**. Otherwise, \mathcal{P} is **curved**. In curved exponential families, the η_i ’s are related in a non-linear way.

Example 2.10 Normal distribution $\mathcal{N}(\mu, \sigma^2)$, $\eta_1 = \frac{1}{2\sigma^2}$, $T_1(x) = -x^2$, $\eta_2 = \frac{\mu}{\sigma^2}$, $T_2(x) = x$

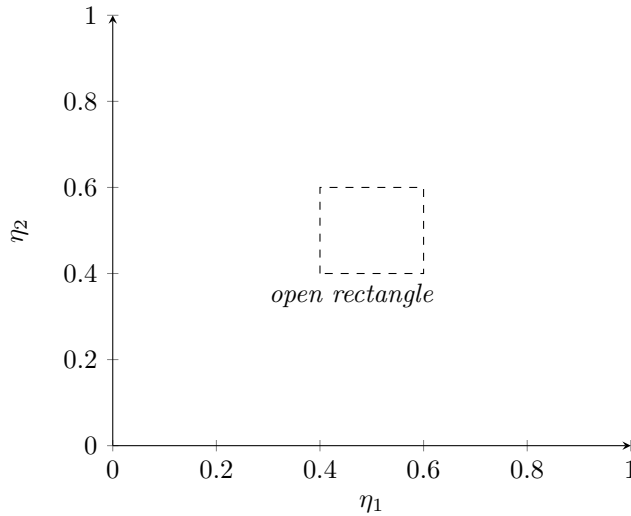
- Non-minimal case: when $\mu = \sigma^2$, $\eta_1 = \frac{1}{2\sigma^2}$, $\eta_2 = 1$



- Minimal and curved: when $\mu = \sqrt{\sigma^2}$, so $\eta_1 = \frac{1}{2\sigma^2}$, $\eta_2 = \frac{1}{\sqrt{\sigma^2}}$, $\eta_2^2 = 2\eta_1$



- Minimal and full-rank (most common): when the natural parameter space is $(0, +\infty) \times \mathbb{R}$



2.1.2 Properties of Exponential Families

Property 1: If $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} p(x; \theta) = \exp(\sum_{i=1}^s \eta_i(\theta) T_i(x) - B(\theta)) h(x)$, then

$$p(x_1, \dots, x_n; \theta) = \exp \left(\sum_{i=1}^s \eta_i(\theta) \sum_{j=1}^n T_i(x_j) - nB(\theta) \right) \prod_{j=1}^n h(x_j)$$

By NFFC, $(\sum_{j=1}^n T_1(x_j), \dots, \sum_{j=1}^n T_s(x_j))$ is a sufficient statistic.

Property 2: If f is interable and $\eta \in \Theta_{\text{nat}}$, then $G(f, \eta) = \int f(x) \exp(\sum_{i=1}^s \eta_i T_i(x)) h(x) d\mu(x)$ is infinitely differentiable with respect to η and the derivatives can be obtained by differentiate under the integral sign. (See Thm 2.7.1 of TSH).

Example 2.11 (Moments of T_i 's) Take $f(x) = 1$, then

$$\begin{aligned} G(f, \eta) &= \int \exp \left(\sum_{j=1}^s \eta_j T_j(x) \right) h(x) d\mu(x) = \exp(A(\eta)) \\ \frac{\partial G(f, \eta)}{\partial \eta_i} &= \int T_i(x) \exp \left(\sum_{j=1}^s \eta_j T_j(x) \right) h(x) d\mu(x) = \frac{\partial A(\eta)}{\partial \eta_i} \times \exp(A(\eta)) \\ \frac{\partial A(\eta)}{\partial \eta_i} &= \int T_i(x) \exp \left(\sum_{j=1}^s \eta_j T_j(x) - A(\eta) \right) h(x) d\mu(x) \\ &= \mathbb{E}_\eta[T_i(X)] \\ \frac{\partial^2 A(\eta)}{\partial \eta_i \partial \eta_j} &= \text{Cov}_\eta(T_i(X), T_j(X)) \end{aligned}$$

Theorem 2.12 (Pitman-Koopman-Darmois) Amongst families of exponential distributions, whose **domain does not depend on/vary with the parameters being estimated**, only in exponential families is there a sufficient statistic whose dimension remains bounded as the sample size increases.

(\Rightarrow Non-exponential families of distribution require non-parametric statistic to fully capture the information in data)

2.2 Minimal Sufficiency

A notion of max achievable lossless data reduction.

Definition 2.13 (Minimal Sufficiency) A sufficient statistic T is minimal if for every sufficient statistic T' and for any $x, y \in \mathcal{X}$, $T(x) = T(y)$ whenever $T'(x) = T'(y)$. In other words, T is a function of T' . (There exists f such that $T(x) = f(T'(x))$ for any $x \in \mathcal{X}$).

Theorem 2.14 Let $\{p(x; \theta), \theta \in \Theta\}$ be a family of densities w.r.t. same measure μ . Suppose that there exists a statistic T such that for every $x, y \in \mathcal{X}$:

$$\frac{p(x; \theta)}{p(y; \theta)} = c_{x,y} \Leftrightarrow p(x; \theta) = c_{x,y} p(y; \theta) \Leftrightarrow T(x) = T(y)$$

for every θ and some $c_{x,y} \in \mathbb{R}$. Then T is a minimal sufficient statistic.

Proof: We first prove that T is sufficient. Start with

$$\begin{aligned} T(\mathcal{X}) &= \{t : t = T(x) \text{ for some } x \in \mathcal{X}\} \\ &= \text{range of } T. \end{aligned}$$

For each $t \in T(\mathcal{X})$, we consider the preimage $A_t = \{x : T(x) = t\}$ and select an arbitrary representative x_t from each A_t . Then, for any $y \in \mathcal{X}$, we have $y \in A_{T(y)}$ and $x_{T(y)} \in A_{T(y)}$. By the definition of A_t , this implies that $T(y) = T(x_{T(y)})$. From the assumption of the theorem,

$$p(y; \theta) = c_{y, x_{T(y)}} p(x_{T(y)}; \theta) = h(y) g_\theta(T(y))$$

which yields sufficiency of T by the NFFC.

Consider another sufficient statistic T' . By NFFC,

$$p(x; \theta) = \tilde{g}_\theta(T'(x)) \tilde{h}(x)$$

Take any x, y such that $T'(x) = T'(y)$, then

$$\begin{aligned} p(x; \theta) &= \tilde{g}_\theta(T'(x)) \tilde{h}(x) \\ &= \tilde{g}_\theta(T'(y)) \tilde{h}(y) \cdot \frac{\tilde{h}(x)}{\tilde{h}(y)} \\ &= p(y; \theta) C_{x,y}. \end{aligned}$$

Hence, $T(x) = T(y)$ by the assumption of the theorem. So $T'(x) = T'(y)$ implies $T(x) = T(y)$ for any sufficient statistic T' and any x and y . As a result, T is a minimal sufficient statistic. \blacksquare

Remark 2.15 (Ex 3.12, Keener) For any minimal s -dimensional exponential family the statistic $(\sum_i T_1(X_i), \dots, \sum_i T_s(X_i))$ is a minimal sufficient statistic.

Example 2.16 (Curved exponential family) Let $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} \mathcal{N}(\sigma, \sigma^2)$, $\theta = \sigma > 0$.

$$\begin{aligned}\frac{p(x; \theta)}{p(y; \theta)} &= \frac{\exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n x_i^2 + \frac{\sigma}{\sigma^2} \sum_{i=1}^n x_i - \frac{n\sigma^2}{2\sigma^2}\right)}{\exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n y_i^2 + \frac{\sigma}{\sigma^2} \sum_{i=1}^n y_i - \frac{n\sigma^2}{2\sigma^2}\right)} \\ &= \exp\left(-\frac{1}{2\sigma^2} \left\{ \sum_{i=1}^n x_i^2 - \sum_{i=1}^n y_i^2 \right\} + \frac{1}{\sigma} \left\{ \sum_{i=1}^n x_i - \sum_{i=1}^n y_i \right\}\right)\end{aligned}$$

Is $T(X) = (T_1(X), T_2(X)) = (\sum_{i=1}^n X_i^2, \sum_{i=1}^n X_i)$ minimal sufficient?

First, if $T(x) = T(y)$ for some $x, y \in \mathcal{X}$, then the ratio above is equal to 1, hence does not depend on θ . Therefore, T is sufficient.

Second, if for some x, y , the ratio is independent of θ , notice that the ratio $\rightarrow 1$ as $\sigma \rightarrow \infty$. Therefore, $C_{x,y} = 1$ and $\log C_{x,y} = 0 = \log \left(\frac{p(x; \theta)}{p(y; \theta)} \right)$. This implies

$$\begin{aligned}\frac{1}{2\sigma^2} (T_1(y) - T_1(x)) + \frac{1}{\sigma} (T_2(x) - T_2(y)) &= 0, \quad \forall \sigma \\ \Leftrightarrow T_1(y) - T_1(x) &= 2\sigma (T_2(y) - T_2(x)) \quad \forall \sigma.\end{aligned}$$

As $\sigma \rightarrow 0$, RHS $\rightarrow 0$. So, $T_2(y) = T_2(x)$. Consequently, T is a minimal sufficient statistic.

Example 2.17 Let $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} U(0, \theta)$ and $T(x) = \max(X_1, \dots, X_n)$. In that case for $x = (x_1, \dots, x_n)$ such that $x_i > 0, i = 1, \dots, n$,

$$p(x; \theta) = \prod_{i=1}^n \frac{1}{\theta} I(x_i < \theta) = \frac{1}{\theta^n} I(T(x) < \theta)$$

If $T(x) = T(y)$, then $p(x; \theta) = \underbrace{1}_{c_{x,y} \perp \theta} \times p(y; \theta) \Rightarrow$ sufficiency. Conversely, if $x, y > 0$ are supported by the same θ 's, then $\{\theta \text{ supporting } x\} = (T(x), \infty) = (T(y), \infty) = \{\theta \text{ supporting } y\}$. Therefore, $T(x) = T(y)$ and T is a minimal sufficient statistic.

2.3 Ancillarity and Completeness

Sufficient statistic/minimal sufficient statistics don't achieve data reduction in a significant way.

Example 2.18 Consider $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} \text{CauchyLoc}(\theta)$, whose density is given by

$$p(x; \theta) = \frac{1}{\pi} \frac{1}{1 + (x - \theta)^2} = f(x - \theta)$$

then $(X_{(1)}, \dots, X_{(n)})$ is minimal sufficient. (See TPE 1.5).

Another example is double exponential location model $p(x; \theta) \propto \exp(-|x - \theta|)$. Amount of ancillary information present in its minimal sufficient statistics.

Definition 2.19 A statistic A is ancillary for $X \sim P_\theta \in \mathcal{P}$ if the distribution of $A(X)$ does not depend on θ .

Example 2.20 Consider again $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} \text{CauchyLoc}(\theta)$. Then $A(X) = X_{(n)} - X_{(1)}$ is ancillary even though (X_1, \dots, X_n) is minimal sufficient. To see this, note that $X_i = Z_i + \theta$ for $Z_i \stackrel{i.i.d.}{\sim} \text{CauchyLoc}(\theta)$, we can see that $X_{(i)} = Z_{(i)} + \theta$ and $A(X) = A(Z) \perp \theta$.

Definition 2.21 A statistic A is first-order ancillary for $X \sim P_\theta \in \mathcal{P}$ if $\mathbb{E}_\theta[A(X)]$ does not depend on θ .

Definition 2.22 A statistic T is complete for $X \sim P_\theta \in \mathcal{P}$ if no non-constant function of T is first-order ancillary. In other words, if $\mathbb{E}_\theta[f(T(X))] = 0$ for all θ , then $f(T(X)) = 0$ with probability 1 for all θ .

Some properties:

1. If T is complete sufficient, then T is minimal sufficient. (Bahadur's theorem)
2. Complete sufficient statistic yields optimal unbiased estimators.

Example 2.23 Let $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} \text{Bernoulli}(\theta)$, $\theta \in (0, 1)$. Then $T(X) = \sum_{i=1}^n X_i$ is sufficient.

Suppose $\mathbb{E}_\theta[f(T(X))] = 0$ for all $\theta \in (0, 1)$,

$$\sum_{j=0}^n f(j) \binom{n}{j} \theta^j (1-\theta)^{n-j} = 0, \quad \forall \theta \in (0, 1)$$

Dividing both sides by θ^n and reparameterizing $\beta = \frac{\theta}{1-\theta}$ we can rewrite it as

$$\sum_{j=0}^n f(j) \binom{n}{j} \beta^j = 0, \quad \forall \beta > 0$$

if f is non-zero, then LHS is a polynomial of degree at most n . However, an n th-degree polynomial has at most n roots. Hence, it is impossible for the LHS to be equal to 0 for every $\beta > 0$ unless $f = 0$. Therefore, T is complete.

Example 2.24 Let $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} \mathcal{N}(\theta, \sigma^2)$ with unknown θ and an known $\sigma^2 > 0$. Is $\bar{X}_n = n^{-1} \sum_{i=1}^n X_i$ complete for this model?

Consider the special case of $n = 1$ and $\sigma = 1$. $T(X) = X \sim \mathcal{N}(\theta, 1)$

$$\mathbb{E}_\theta[f(X)] = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} f(x) \exp\left(-\frac{(x-\theta)^2}{2}\right) dx = 0, \quad \forall \theta \in \mathbb{R}$$

Multiplying both sides by $\frac{\theta^2}{\sqrt{2\pi}e^{\frac{\theta^2}{2}}}$, we have

$$\int_{-\infty}^{\infty} f(x) \exp\left(-\frac{x^2}{2}\right) \exp(\theta x) dx = 0, \quad \forall \theta \in \mathbb{R}.$$

We decompose f into its positive and negative parts as $f(x) = f_+(x) - f_-(x)$, where $f_+(x) = \max(f(x), 0)$, and $f_-(x) = \max(-f(x), 0)$. Note that $f_+ \geq 0$ and $f_- \geq 0$. For all $x \in \mathbb{R}$ $f_+(x) = f_-(x)$ if and only if $f_+(x) = f_-(x) = 0$.

Suppose f_+ and f_- have non-zero components, and we may write

$$\frac{\int_{-\infty}^{\infty} f_+(x) e^{-\frac{x^2}{2}} e^{\theta x} dx}{\int_{-\infty}^{\infty} f_+(x) e^{-\frac{x^2}{2}} dx} = \frac{\int_{-\infty}^{\infty} f_-(x) e^{-\frac{x^2}{2}} e^{\theta x} dx}{\int_{-\infty}^{\infty} f_-(x) e^{-\frac{x^2}{2}} dx}$$

Note that

$$\frac{f_+(x) e^{-\frac{x^2}{2}}}{\int_{-\infty}^{\infty} f_+(x) e^{-\frac{x^2}{2}} dx}$$

defines a probability density.

The equality of the mgfs implies equality of the densities, which in turn implies $f_+(x) = f_-(x)$ a.e.. Then $f_+(x) = f_-(x) = 0$ a.e., or in other words, $f(x) = 0$ a.e.. Hence T is complete.