

## FAMILIES OF DISTRIBUTIONS

## 4.1 Location-Scale Families

**Definition: Location Scale Family**

A **location-scale family** is a family of distributions formed by *translation* and *rescaling* of a standard family member.

Suppose that  $f(x)$  is a pdf. Then if  $\mu$  and  $\sigma > 0$  are constants then

$$f(x|\mu, \sigma) = \frac{1}{\sigma} f((x - \mu)/\sigma)$$

is also a pdf;  $f(x|\mu, \sigma) \geq 0$ , and

$$\int_{-\infty}^{\infty} f(x|\mu, \sigma) dx = \int_{-\infty}^{\infty} \frac{1}{\sigma} f((x - \mu)/\sigma) dx = \int_{-\infty}^{\infty} f(y) dy = 1$$

setting  $y = (x - \mu)/\sigma$  in the penultimate integral.

- $f(x|\mu, \sigma)$  is termed a **location-scale family**
- if  $\sigma = 1$  we have a **location family**:  $f(x|\mu) = f(x - \mu)$
- if  $\mu = 0$  we have a **scale family**:  $f(x|\sigma) = f(x/\sigma)/\sigma$

**Example : Normal distribution family**

$$\begin{aligned} f(x) &= \left(\frac{1}{2\pi}\right)^{1/2} \exp\left\{-\frac{1}{2}x^2\right\} \\ f(x|\mu, \sigma) &= \left(\frac{1}{2\pi\sigma^2}\right)^{1/2} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^2\right\} \end{aligned}$$

**Example : Exponential distribution family**

$$\begin{aligned} f(x) &= e^{-x} & x > 0 \\ f(x|\mu, \sigma) &= \frac{1}{\sigma} e^{-(x-\mu)/\sigma} & x > \mu \\ f(x|\mu) &= e^{-(x-\mu)} & x > \mu \end{aligned}$$

Note that  $X$  is a random variable with pdf  $f_X(x) = f(x|\mu, \sigma)$  (the location-scale family member) **if and only if** there exists another random variable  $Z$  with  $f_Z(z) = f(z)$  (the standard member) such that

$$X = \sigma Z + \mu$$

that is, if  $X$  is a linear (location-scale) transformation of a standard random variable  $Z$ .

## 4.2 Exponential Families

### Definition: Exponential Family

A family of pdfs/pmfs is called an **exponential family** if it can be expressed

$$f(x|\underline{\theta}) = h(x)c(\underline{\theta}) \exp \left\{ \sum_{j=1}^k w_j(\underline{\theta}) t_j(x) \right\} = h(x)c(\underline{\theta}) \exp \left\{ \underline{w}(\underline{\theta})^\top \underline{t}(x) \right\}$$

for all  $x \in \mathbb{R}$ , where  $\underline{\theta} \in \Theta$  is a  $d$ -dimensional parameter vector, and

- $h(x) \geq 0$  is a function that does not depend on  $\underline{\theta}$
- $c(\underline{\theta}) \geq 0$  is a function that does not depend on  $x$
- $\underline{t}(x) = (t_1(x), \dots, t_k(x))^\top$  is a vector of real-valued functions that do not depend on  $\underline{\theta}$
- $\underline{w}(x) = (w_1(\underline{\theta}), \dots, w_k(\underline{\theta}))^\top$  is a vector of real-valued functions that do not depend on  $x$

An exponential family distribution is termed **natural** if  $k = 1$  and  $t_1(x) = x$ . Note that the support of an exponential family distribution  $f(x|\underline{\theta})$  **cannot** depend on  $\underline{\theta}$ .

**Example :**  $\text{Binomial}(n, \theta)$  for  $0 < \theta < 1$

For  $x \in \{0, 1, \dots, n\} \equiv \mathbb{X}$ ,

$$f(x|\theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x} = \binom{n}{x} (1-\theta)^n \left( \frac{\theta}{1-\theta} \right)^x = \binom{n}{x} (1-\theta)^n \exp \left\{ \log \left( \frac{\theta}{1-\theta} \right) x \right\}$$

- $k = 1$
- $h(x) = I_{\mathbb{X}}(x) \binom{n}{x}$ , where  $I_A(x)$  is the **indicator function** for set  $A$

$$I_A(x) = \begin{cases} 1 & x \in A \\ 0 & x \notin A \end{cases}$$

- $c(\underline{\theta}) = (1-\theta)^n$
- $t_1(x) = x$
- $w_1(\underline{\theta}) = \log \left( \frac{\theta}{1-\theta} \right)$

**Example :**  $\text{Normal}(\mu, \sigma^2)$

For  $x \in \mathbb{R}$ ,

$$f(x|\mu, \sigma^2) = \left( \frac{1}{2\pi\sigma^2} \right)^{1/2} \exp \left\{ -\frac{1}{2\sigma^2} (x - \mu)^2 \right\} = \left( \frac{1}{2\pi\sigma^2} \right)^{1/2} \exp \left\{ -\frac{\mu^2}{2\sigma^2} \right\} \exp \left\{ -\frac{x^2}{2\sigma^2} + \frac{\mu x}{\sigma^2} \right\}$$

- $k = 2, \underline{\theta} = (\mu, \sigma^2)^\top$
- $h(x) = 1$
- $c(\underline{\theta}) = c(\mu, \sigma^2) = \left( \frac{1}{2\pi\sigma^2} \right)^{1/2} \exp \left\{ -\frac{\mu^2}{2\sigma^2} \right\}$
- $t_1(x) = -\frac{x^2}{2}, t_2(x) = x$
- $w_1(\underline{\theta}) = \frac{1}{\sigma^2}, w_2(\underline{\theta}) = \frac{\mu}{\sigma^2}$

**Example :** Suppose, for  $\theta > 0$

$$f(x|\theta) = \frac{1}{\theta} \exp \left\{ 1 - \frac{x}{\theta} \right\} \quad x > \theta$$

and zero otherwise. Then

- $k = 1, \underline{\theta} = \theta$
- $h(x) = eI_{[\theta, \infty)}(x)$
- $c(\underline{\theta}) = 1/\theta$
- $t_1(x) = x, w_1(\underline{\theta}) = 1/\theta$

but the support of  $f(x|\theta)$  depends on  $\theta$  so this is **not** an exponential family distribution.

#### 4.2.1 Parameterization

We can **reparameterize** from  $\underline{\theta}$  to  $\underline{\eta} = (\eta_1, \dots, \eta_k)^\top$  by setting  $\eta_j = w_j(\underline{\theta})$  for each  $j$ , and write

$$f(x|\underline{\eta}) = h(x)c^*(\underline{\eta}) \exp \left\{ \sum_{j=1}^k \eta_j t_j(x) \right\} = h(x)c^*(\underline{\eta}) \exp \left\{ \underline{\eta}^\top \underline{t}(x) \right\}.$$

$\underline{\eta}$  is termed the **natural** or **canonical** parameter. Let  $\mathcal{H}$  be the region of  $\mathbb{R}^k$  defined by

$$\mathcal{H} \equiv \left\{ \underline{\eta} : \int_{-\infty}^{\infty} h(x) \exp \left\{ \underline{\eta}^\top \underline{t}(x) \right\} dx < \infty \right\}$$

$\mathcal{H}$  is the **natural parameter space**. For  $\underline{\eta} \in \mathcal{H}$ , we must have

$$c^*(\underline{\eta}) = \left[ \int_{-\infty}^{\infty} h(x) \exp \left\{ \underline{\eta}^\top \underline{t}(x) \right\} dx \right]^{-1}$$

It can be shown that  $\mathcal{H}$  is a **convex** set, that is, for  $0 \leq \lambda \leq 1$ ,

$$\underline{\eta}_1, \underline{\eta}_2 \in \mathcal{H} \implies \lambda \underline{\eta}_1 + (1 - \lambda) \underline{\eta}_2 \in \mathcal{H}.$$

Note that

$$\mathcal{H}_\Theta = \left\{ \underline{w}(\underline{\theta}) = (w_1(\underline{\theta}), \dots, w_k(\underline{\theta}))^\top : \underline{\theta} \in \Theta \right\} \subseteq \mathcal{H}.$$

**Example :** *Binomial*( $n, \theta$ )

$$\eta = \log \left( \frac{\theta}{1 - \theta} \right) \iff \theta = \frac{e^\eta}{1 + e^\eta}$$

so that

$$f(x|\eta) = \left\{ \binom{n}{x} I_{\{0,1,\dots,n\}}(x) \right\} \frac{e^{\eta x}}{(1 + e^\eta)^n}$$

Natural parameter space:

$$\int_{-\infty}^{\infty} h(x) \exp \left\{ \underline{\eta}^\top \underline{t}(x) \right\} dx = \sum_{x=0}^n \binom{n}{x} \exp \{ \eta x \} < \infty \quad \forall \eta \quad \therefore \quad \mathcal{H} \equiv \mathbb{R}.$$

**Example :**  $Normal(\mu, \sigma^2)$

Natural parameters:

$$\underline{\eta} = (\eta_1, \eta_2)^T = (1/\sigma^2, \mu/\sigma^2)^T$$

so that

$$f(x|\underline{\eta}) = \left(\frac{\eta_1}{2\pi}\right)^{1/2} \exp\left\{-\frac{\eta_2^2}{2\eta_1}\right\} \exp\left\{-\frac{\eta_1 x^2}{2} + \eta_2 x\right\}$$

Natural parameter space: this density will be integrable with respect to  $x$  if and only if  $\eta_1 > 0$ , so  $\mathcal{H} \equiv \mathbb{R}^+ \times \mathbb{R}$ .

- An exponential family indexed by parameter  $\underline{\theta}$  (or  $\underline{\eta}$ ) is termed **regular** if
  - I.  $\mathcal{H} \equiv \mathcal{H}_\Theta$ .
  - II. In the natural parameterization, neither the  $\eta_j$  nor the  $t_i(x)$  satisfy linearity constraints.
  - III.  $\mathcal{H}$  is an open set in  $\mathbb{R}^k$ .
- If only I. and II. hold, the exponential family is termed **full**.
- An exponential family indexed by parameter  $\underline{\theta}$  is termed **curved** if

$$\dim(\underline{\theta}) = d < k$$

#### 4.2.2 Expectation and Variance for Exponential Families

##### Definition: Score Function

For pmf/pdf  $f_X$  with  $d$ -dimensional parameter  $\underline{\theta}$ , the **score function**,  $\underline{S}(x; \underline{\theta})$ , is a  $d \times 1$  vector with  $j$ th element equal to

$$S_j(x; \underline{\theta}) = \frac{\partial}{\partial \theta_j} \log f_X(x|\underline{\theta}).$$

The quantity  $\underline{S}(X; \underline{\theta})$  is a  $d$ -dimensional **random variable**.

**Lemma** Under certain regularity conditions

$$\mathbb{E}_{f_X}[\underline{S}(X; \underline{\theta})] = \underline{0}$$

**Proof** In the case  $d = 1$ ; let

$$\dot{f}_X(x|\theta) = \frac{d}{d\theta} f_X(x|\theta)$$

Then

$$\begin{aligned} \mathbb{E}_{f_X}[S(X; \theta)] &= \int S(x; \theta) f_X(x|\theta) dx = \int \left\{ \frac{d}{d\theta} \log f_X(x|\theta) \right\} f_X(x|\theta) dx \\ &= \int \left\{ \frac{\dot{f}_X(x|\theta)}{f_X(x|\theta)} \right\} f_X(x|\theta) dx \\ &= \int \frac{d}{d\theta} f_X(x|\theta) dx = \frac{d}{d\theta} \left\{ \int f_X(x|\theta) dx \right\} = 0 \end{aligned}$$

provided that the order of the differentiation wrt  $\theta$  and the integration wrt  $x$  can be exchanged.

**Definition: Fisher Information**

For pmf/pdf  $f_X$  with  $d$ -dimensional parameter  $\underline{\theta}$ , the **Fisher Information**,  $\mathcal{I}(\underline{\theta})$ , is a  $d \times d$  matrix defined as the variance-covariance matrix of the score random variable  $\underline{S}$ , that is

$$\mathcal{I}(\underline{\theta}) = \text{Var}_{f_X}[\underline{S}(X; \underline{\theta})] = \mathbb{E}_{f_X}[\underline{S}(X; \underline{\theta})\underline{S}(X; \underline{\theta})^\top]$$

with  $(i, j)$ th element equal to

$$\mathbb{E}_{f_X}[S_i(X; \underline{\theta})S_j(X; \underline{\theta})]$$

The Fisher Information is a constant  $d \times d$  matrix in which each of the elements is a function of  $\underline{\theta}$ .

**Lemma** Under certain regularity conditions, if the pmf/pdf is twice partially differentiable with respect to the elements of  $\underline{\theta}$ , then

$$\mathcal{I}(\underline{\theta}) = -\mathbb{E}_{f_X}[\underline{\Psi}(X; \underline{\theta})]$$

where  $\underline{\Psi}(X; \underline{\theta})$  is the  $d \times d$  matrix of second partial derivatives with  $(i, j)$ th element equal to

$$\frac{\partial^2}{\partial \theta_i \partial \theta_j} \log f_X(x|\underline{\theta}).$$

**Proof** In the case  $d = 1$ ; from above

$$\int \left\{ \frac{d}{d\theta} \log f_X(x|\underline{\theta}) \right\} f_X(x|\theta) dx = 0$$

so therefore, differentiating again wrt  $\theta$

$$\int \left[ \left\{ \frac{d^2}{d\theta^2} \log f_X(x|\underline{\theta}) f_X(x|\theta) \right\} + \left\{ \frac{d}{d\theta} \log f_X(x|\underline{\theta}) \frac{d}{d\theta} f_X(x|\theta) \right\} \right] dx = 0 \quad (1)$$

But

$$\frac{d}{d\theta} \log f_X(x|\underline{\theta}) = \frac{\dot{f}_X(x|\underline{\theta})}{f_X(x|\theta)} \quad \therefore \quad \dot{f}_X(x|\theta) = \frac{d}{d\theta} f_X(x|\theta) = f_X(x|\theta) \frac{d}{d\theta} \log f_X(x|\underline{\theta})$$

so therefore

$$\int \frac{d}{d\theta} \log f_X(x|\underline{\theta}) \frac{d}{d\theta} f_X(x|\theta) dx = \int \left\{ \frac{d}{d\theta} \log f_X(x|\underline{\theta}) \right\}^2 f_X(x|\theta) dx$$

and so substituting into equation (1) above, we have

$$\int \left\{ \frac{d^2}{d\theta^2} \log f_X(x|\underline{\theta}) f_X(x|\theta) \right\} dx = - \int \left\{ \frac{d}{d\theta} \log f_X(x|\underline{\theta}) \right\}^2 f_X(x|\theta) dx$$

of equivalently

$$\mathbb{E}_{f_X} \left[ \frac{d^2}{d\theta^2} \log f_X(x|\underline{\theta}) \right] = -\mathbb{E}_{f_X} \left[ \left\{ \frac{d}{d\theta} \log f_X(x|\underline{\theta}) \right\}^2 \right] = \mathbb{E}_{f_X} [S(X; \theta)^2]$$

so that, as  $\mathbb{E}_{f_X} [S(X; \theta)] = 0$ ,

$$\mathbb{E}_{f_X} \left[ \frac{d^2}{d\theta^2} \log f_X(x|\underline{\theta}) \right] = -\text{Var}_{f_X} [S(X; \theta)].$$

Note that if  $\underline{a} = (a_1, \dots, a_d)^\top$ , then

$$\text{Var}_{f_X} [\underline{a}^\top \underline{S}(X; \underline{\theta})] = \underline{a}^\top \mathcal{I}(\underline{\theta}) \underline{a}$$

**Example :** *Binomial*( $n, \theta$ )

$$f(x|\theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x} \quad x \in \{0, 1, \dots, n\}$$

so that

$$S(x; \theta) = \frac{d}{d\theta} \log f_X(x|\theta) = \frac{x}{\theta} - \frac{n-x}{1-\theta} = \frac{x-n\theta}{\theta(1-\theta)}.$$

Hence

$$\mathbb{E}_{f_X}[S(X; \theta)] = \mathbb{E}_{f_X} \left[ \frac{X-n\theta}{\theta(1-\theta)} \right] = \frac{\mathbb{E}_{f_X}[X] - n\theta}{\theta(1-\theta)} = 0$$

as  $X \sim \text{Binomial}(n, \theta)$  yields  $\mathbb{E}_{f_X}[X] = n\theta$ . For the second derivative

$$\frac{d^2}{d\theta^2} \log f_X(x|\theta) = -\frac{x}{\theta^2} - \frac{n-x}{(1-\theta)^2}$$

so that

$$\mathcal{I}(\theta) = -\mathbb{E}_{f_X} \left[ \frac{d^2}{d\theta^2} \log f_X(x|\theta) \right] = \frac{\mathbb{E}_{f_X}[X]}{\theta^2} + \frac{n - \mathbb{E}_{f_X}[X]}{(1-\theta)^2}$$

and as  $\mathbb{E}_{f_X}[X] = n\theta$ , we have

$$\mathcal{I}(\theta) = \frac{n\theta}{\theta^2} + \frac{n-n\theta}{(1-\theta)^2} = \frac{1}{\theta(1-\theta)}$$

**Example :** *Poisson*( $\lambda$ )

$$f(x|\lambda) = \frac{e^{-\lambda} \lambda^x}{x!} \quad x \in \{0, 1, \dots\}$$

so that

$$S(x; \lambda) = \frac{d}{d\lambda} \log f_X(x|\lambda) = \frac{x}{\lambda} - 1$$

Hence

$$\mathbb{E}_{f_X}[S(X; \lambda)] = \mathbb{E}_{f_X} \left[ \frac{X}{\lambda} - 1 \right] = \frac{\mathbb{E}_{f_X}[X]}{\lambda} - 1 = 0$$

as  $X \sim \text{Poisson}(\lambda)$  yields  $\mathbb{E}_{f_X}[X] = \lambda$ . For the second derivative

$$\frac{d^2}{d\lambda^2} \log f_X(x|\lambda) = -\frac{x}{\lambda^2}$$

so that

$$\mathcal{I}(\lambda) = -\mathbb{E}_{f_X} \left[ \frac{d^2}{d\lambda^2} \log f_X(x|\lambda) \right] = \frac{\mathbb{E}_{f_X}[X]}{\lambda^2}$$

and as  $\mathbb{E}_{f_X}[X] = \lambda$ , we have

$$\mathcal{I}(\lambda) = \frac{1}{\lambda}$$

## Results for the Exponential Family

If

$$f_X(x|\underline{\theta}) = h(x)c(\underline{\theta}) \exp \left\{ \sum_{j=1}^k w_j(\underline{\theta})t_j(x) \right\}$$

then, for  $l = 1, \dots, d$ ,

$$S_l(x; \underline{\theta}) = \frac{\partial}{\partial \theta_l} \log f_X(x|\underline{\theta}) = \frac{\partial}{\partial \theta_l} \log c(\underline{\theta}) + \sum_{j=1}^k \dot{w}_{jl}(\underline{\theta})t_j(x) = \frac{\dot{c}_l(\underline{\theta})}{c(\underline{\theta})} + \sum_{j=1}^k \dot{w}_{jl}(\underline{\theta})t_j(x)$$

where

$$\dot{c}_l(\underline{\theta}) = \frac{\partial c(\underline{\theta})}{\partial \theta_l} \quad \dot{w}_{jl}(\underline{\theta}) = \frac{\partial w_j(\underline{\theta})}{\partial \theta_l}.$$

But, for each  $l$ ,  $\mathbb{E}_{f_X}[S_l(X; \underline{\theta})] = 0$ , so therefore, for  $l = 1, \dots, d$ ,

$$\mathbb{E}_{f_X} \left[ \sum_{j=1}^k \dot{w}_{jl}(\underline{\theta})t_j(X) \right] = -\frac{\dot{c}_l(\underline{\theta})}{c(\underline{\theta})} = -\frac{\partial}{\partial \theta_l} \log c(\underline{\theta}).$$

By a similar calculation

$$\text{Var}_{f_X} \left[ \sum_{j=1}^k \dot{w}_{jl}(\underline{\theta})t_j(X) \right] = -\frac{\partial^2}{\partial \theta_l^2} \log c(\underline{\theta}) - \mathbb{E}_{f_X} \left[ \sum_{j=1}^k \ddot{w}_{jl}(\underline{\theta})t_j(X) \right]$$

where

$$\ddot{w}_{jl}(\underline{\theta}) = \frac{\partial^2 w_j(\underline{\theta})}{\partial \theta_l^2}$$

**Example :** *Binomial*( $n, \theta$ )

$$f(x|\theta) = \binom{n}{x} (1-\theta)^n \exp \left\{ \log \left( \frac{\theta}{1-\theta} \right) x \right\}$$

so that

$$w_1(\theta) = \log \left( \frac{\theta}{1-\theta} \right) \quad \log c(\theta) = n \log(1-\theta) \quad S(x; \theta) = -\frac{n}{1-\theta} + \frac{x}{\theta(1-\theta)}.$$

From the result above

$$\mathbb{E}_{f_X} [\dot{w}_{11}(\theta)t_1(X)] = -\frac{\partial}{\partial \theta} \log c(\theta)$$

that is

$$\mathbb{E}_{f_X} \left[ \frac{1}{\theta(1-\theta)} X \right] = \frac{n}{1-\theta} \quad \therefore \quad \mathbb{E}_{f_X}[X] = n\theta.$$

Note that in the **natural** (canonical) parameterization

$$\log f_X(x|\underline{\eta}) = \log h(x) + \log c^*(\underline{\eta}) + \sum_{j=1}^k \eta_j t_j(X)$$

so that, using the arguments above for  $l = 1, \dots, d$ ,

$$\mathbb{E}_{f_X} [t_l(X)] = -\frac{\partial}{\partial \eta_l} \log c^*(\underline{\eta}) \quad \text{Var}_{f_X} [t_l(X)] = -\frac{\partial^2}{\partial \eta_l^2} \log c^*(\underline{\eta})$$

### 4.2.3 Independent random variables from the Exponential Family

Suppose that  $X_1, \dots, X_n$  are independent and identically distributed random variables, with pmf/pdf  $f_X(x|\theta)$  in the Exponential Family. Then the joint pmf/pdf for  $\underline{X} = (X_1, \dots, X_n)^\top$  takes the form

$$f_{\underline{X}}(\underline{x}|\theta) = \prod_{i=1}^n f_X(x_i|\theta) = \prod_{i=1}^n h(x_i)c(\theta) \exp \left\{ \sum_{j=1}^k w_j(\theta)t_j(x_i) \right\} = H(\underline{x})C(\theta) \exp \left\{ \sum_{j=1}^k w_j(\theta)T_j(\underline{x}) \right\}$$

where

$$H(\underline{x}) = \prod_{i=1}^n h(x_i) \quad C(\theta) = \{c(\theta)\}^n \quad T_j(\underline{x}) = \sum_{i=1}^n t_j(x_i).$$

### 4.2.4 Alternative construction of the Exponential Family

Suppose that  $f(x)$  is a pmf/pdf with corresponding mgf  $M(t)$  (presumed to exist in a neighbourhood of zero), so that

$$M(t) = \int e^{tx} f(x) dx = \exp\{K(t)\}$$

and  $K(t) = \log M(t)$  is the **cumulant generating function**. Now suppose that  $f(x) = \exp\{g(x)\}$ . Then

$$\exp\{K(t)\} = M(t) = \int e^{tx} f(x) dx = \int e^{tx} e^{g(x)} dx = \int e^{tx+g(x)} dx.$$

Hence, dividing through by  $\exp\{K(t)\}$ , we have that

$$\int e^{tx+g(x)-K(t)} dx = 1$$

and also that the integrand is non-negative. Thus, for all  $t$  for which  $M(t)$  exists,

$$f(x|t) = \exp\{tx + g(x) - K(t)\} = f(x) \exp\{tx - K(t)\}$$

is a valid pdf. If we set  $t = \eta$ ,  $h(x) = f(x) = \exp\{g(x)\}$  and  $c^*(\eta) = \exp\{-K(\eta)\}$ , then

$$f(x|\eta) = h(x)c^*(\eta) \exp\{\eta x\}$$

and we see that  $f(x|\eta)$  is an exponential family member with natural parameter  $\eta$ . The pmf/pdf  $f(x|t)$  is termed the **exponential tilting** of  $f(x)$ , with expectation

$$-\frac{d}{dt} \log c^*(t) = -\frac{d}{dt} \{-K(t)\} = \dot{K}(t)$$

and variance

$$-\frac{d^2}{dt^2} \log c^*(t) = -\frac{d^2}{dt^2} \{-K(t)\} = \ddot{K}(t).$$

Note further that if

$$f(x|\eta) = h(x)c^*(\eta) \exp\{\eta x\}.$$

then the corresponding mgf is, for  $t$  small enough,

$$M(t) = \int e^{tx} h(x)c^*(\eta) \exp\{\eta x\} dx = c^*(\eta) \int h(x) \exp\{(\eta + t)x\} dx = \frac{c^*(\eta)}{c^*(\eta + t)}.$$



### 4.2.5 The Exponential Dispersion Model

Consider the model

$$f(x|\underline{\theta}, \phi) = \exp \left\{ d(x, \phi) + \frac{\log c(\underline{\theta})}{r(\phi)} + \frac{1}{r(\phi)} \sum_{j=1}^k w_j(\underline{\theta}) t_j(x) \right\} = h(x) c(\underline{\theta}) \exp \left\{ \underline{w}(\underline{\theta})^\top \underline{t}(x) \right\}$$

where  $r(\phi) > 0$  is a function of **dispersion** parameter  $\phi > 0$ .

In this model, using the previous results, we see that the expectation is unchanged compared to the Exponential Family model by the presence of the term  $r(\phi)$ , but the variance is modified by a factor of  $r(\phi)$ .

**Example :** *Binomial*( $n, \theta$ )

$$f_X(x|\theta) = \binom{n}{x} I_{\{0,1,\dots,n\}}(x) \exp \left\{ n \log(1 - \theta) + \log \left( \frac{\theta}{1 - \theta} \right) x \right\}$$

Let  $Y = X/n$ , so that

$$f_Y(y|\theta, \phi) = \binom{1/\phi}{y/\phi} I_{\{0,\phi,2\phi,\dots,1\}}(y/\phi) \exp \left\{ \frac{1}{\phi} \left[ \log(1 - \theta) + y \log \left( \frac{\theta}{1 - \theta} \right) \right] \right\}$$

where  $\phi = 1/n$ . Note that

$$\mathbb{E}_{f_Y}[Y] = \theta = \mu$$

say, and

$$\text{Var}_{f_Y}[Y] = \phi \theta (1 - \theta) = \phi V(\mu)$$

where  $V(\mu) = \mu(1 - \mu)$  is the **variance function**.

Thus the exponential dispersion model allows separate modelling of mean and variance.

### 4.3 Convolution Families

The **convolution** of functions  $g$  and  $h$  is a function written  $g \circ h$ , which is defined by

$$g \circ h(y) = \int_{-\infty}^{\infty} g(x) h(y - x) dx.$$

Now if  $X_1$  and  $X_2$  are independent random variables with marginal pdfs  $f_{X_1}$  and  $f_{X_2}$  respectively, then the random variable  $Y = X_1 + X_2$  has a pdf that can be determined using the multivariate transformation result. If we use dummy variable  $Z = X_1$ , then

$$\left. \begin{array}{l} Z = X_1 \\ Y = X_1 + X_2 \end{array} \right\} \iff \left\{ \begin{array}{l} X_1 = Z \\ X_2 = Y - Z \end{array} \right.$$

which is a transformation with Jacobian 1. Thus

$$f_Y(y) = \int_{-\infty}^{\infty} f_{Z,Y}(z, y) dz = \int_{-\infty}^{\infty} f_{X_1,X_2}(z, y - z) dz = \int_{-\infty}^{\infty} f_{X_1}(x) f_{X_2}(y - x) dx$$

so we can see that the pdf of  $Y$  is computed as the convolution of  $f_{X_1}$  and  $f_{X_2}$ .

A family of distributions,  $\mathcal{F}$ , is **closed under convolution** if

$$f_1, f_2 \in \mathcal{F} \implies f_1 \circ f_2 \in \mathcal{F}$$

For independent random variables  $X_1$  and  $X_2$  with pdfs  $f_1$  and  $f_2$  in a family  $\mathcal{F}$ , closure under convolution implies that the random variable  $Y = X_1 + X_2$  also has a pdf in  $\mathcal{F}$ .

This concept is closely related to the idea of **infinite divisibility**, **decomposability**, and **self-decomposability**.

- **Infinite Divisibility** : A probability distribution for rv  $X$  is *infinitely divisible* if, for all positive integers  $n$ , there exists a **sequence of independent and identically distributed** rvs  $Z_{n1}, \dots, Z_{nn}$  such that  $X$  and

$$Z_n = \sum_{j=1}^n Z_{nj}$$

have the same distribution, that is, the characteristic function of  $X$  can be written

$$C_X(t) = \{C_Z(t)\}^n$$

for some characteristic function  $C_Z$ .

- **Decomposability** : A probability distribution for rv  $X$  is *decomposable* if

$$C_X(t) = C_{X_1}(t)C_{X_2}(t)$$

for two characteristic functions  $C_{X_1}$  and  $C_{X_2}$  so that

$$X = X_1 + X_2$$

where  $X_1$  and  $X_2$  are **independent** rvs with characteristic functions  $C_{X_1}$  and  $C_{X_2}$ .

- **Self-Decomposability** : A probability distribution for rv  $X$  is *self-decomposable* if

$$C_X(t) = \{C_{X_1}(t)\}^2$$

for characteristic function  $C_{X_1}$  so that

$$X = X_1 + X_2$$

where  $X_1$  and  $X_2$  are **independent identically distributed** rvs with characteristic function  $C_{X_1}$ .

#### 4.4 Hierarchical Models

A hierarchical model is a model constructed by considering a series of distributions at different levels of a “hierarchy” that together, after marginalization, combine to yield the distribution of the observable quantities.

##### Example : A three-level model

Consider the **three-level** hierarchical model:

LEVEL 3 :	$\lambda > 0$	Fixed parameter
LEVEL 2 :	$N \sim \text{Poisson}(\lambda)$	
LEVEL 1 :	$X N = n, \theta \sim \text{Binomial}(n, \theta)$	

Then the marginal pmf for  $X$  is given by

$$f_X(x|\theta, \lambda) = \sum_{n=0}^{\infty} f_{X|N}(x|n, \theta, \lambda) f_N(n|\lambda).$$

By elementary calculation, we see that  $X \sim \text{Poisson}(\lambda\theta)$

$$f_X(x|\theta, \lambda) = \frac{(\lambda\theta)^x e^{-\lambda\theta}}{x!} \quad x = 0, 1, \dots$$

### Example : A three-level model

Consider the **three-level** hierarchical model:

$$\begin{array}{lll} \text{LEVEL 3 :} & \alpha, \beta > 0 & \text{Fixed parameters} \\ \text{LEVEL 2 :} & Y \sim \text{Gamma}(\alpha, \beta) & \\ \text{LEVEL 1 :} & X|Y = y \sim \text{Poisson}(y) & \end{array}$$

Then the marginal pdf for  $X$  is given by

$$f_X(x|\alpha, \beta) = \int_0^\infty f_{X|Y}(x|y) f_Y(y|\alpha, \beta) dy.$$

A general  $K$ -level hierarchical model can be specified in terms of  $K$  vector random variables:

$$\begin{array}{ll} \text{LEVEL } K : & \underline{X}_K = (X_{K1}, \dots, X_{Kn_K})^\top \\ & \vdots \\ \text{LEVEL } 2 : & \underline{X}_2 = (X_{21}, \dots, X_{2n_2})^\top \\ \text{LEVEL } 1 : & \underline{X}_1 = (X_{11}, \dots, X_{1n_1})^\top \end{array}$$

The hierarchical model specifies the joint distribution via a series of **conditional independence** assumptions, so that

$$f_{\underline{X}_1, \dots, \underline{X}_K}(\underline{x}_1, \dots, \underline{x}_K) = f_{\underline{X}_K}(\underline{x}_K) \prod_{k=1}^{K-1} f_{\underline{X}_k | \underline{X}_{k+1}}(\underline{x}_k | \underline{x}_{k+1})$$

where

$$f_{\underline{X}_k | \underline{X}_{k+1}}(\underline{x}_k | \underline{x}_{k+1}) = \prod_{j=1}^{n_k} f_k(x_{kj} | \underline{x}_{k+1})$$

that is, at level  $k$  in the hierarchy, the random variables are taken to be **conditionally independent** given the values of variables at level  $k + 1$ .

The uppermost level, Level  $K$ , can be taken to be a degenerate model, with mass function equal to 1 at a set of fixed values.

### Example : A three-level model

Consider the **three-level** hierarchical model:

$$\begin{array}{lll} \text{LEVEL 3 :} & \theta, \tau^2 > 0 & \text{Fixed parameters} \\ \text{LEVEL 2 :} & M_1, \dots, M_L \sim \text{Normal}(\theta, \tau^2) & \text{Independent} \\ \text{LEVEL 1 :} & \text{For } l = 1, \dots, L : & \\ & X_{l1}, \dots, X_{ln_l} | M_l = m_l \sim \text{Normal}(m_l, 1) & \\ & \text{where all the } X_{lj} \text{ are conditionally independent given } M_1, \dots, M_L & \end{array}$$

For random variables  $X, Y$  and  $Z$ , we write  $X \perp Y | Z$  if  $X$  and  $Y$  are conditionally independent given  $Z$ , so that in the above model

$$X_{l_1 j_1} \perp X_{l_2 j_2} | M_1, \dots, M_L$$

for all  $l_1, j_1, l_2, j_2$ .

## Special Cases of Hierarchical Models

### 1. Finite Mixture Models

LEVEL 3 :  $L \geq 1$  (integer),  $\pi_1, \dots, \pi_L$  with  $0 \leq \pi_l \leq 1$  and  $\sum_{l=1}^L \pi_l = 1$ , and  $\theta_1, \dots, \theta_L$

LEVEL 2 :  $X \sim f_X(x|\underline{\pi}, L)$  with  $\mathbb{X} \equiv \{1, 2, \dots, L\}$  such that  $P[X = l] = \pi_l$

LEVEL 1 :  $Y|X = l \sim f_l(y|\theta_l)$

where  $f_l$  is some pmf or pdf with parameters  $\theta_l$ . Then

$$f_Y(y|\underline{\pi}, \underline{\theta}, L) = \sum_{l=1}^L f_{Y|X}(y|x) f_X(x) = \sum_{l=1}^L f_l(y|\theta_l) \pi_l$$

This is a **finite mixture distribution**: the observed  $Y$  are drawn from  $L$  distinct sub-populations characterized by pmf/pdf  $f_1, \dots, f_L$  and parameters  $\theta_1, \dots, \theta_L$ , with sub-population proportions  $\pi_1, \dots, \pi_L$ .

Note that if  $M_1, \dots, M_L$  are the mgfs corresponding to  $f_1, \dots, f_L$ , then

$$M_Y(t) = \sum_{l=1}^L \pi_l M_l(t)$$

### 2. Random Sums

LEVEL 3 :  $\underline{\theta}, \phi$  (fixed parameters)

LEVEL 2 :  $X \sim f_X(x|\phi)$  with  $\mathbb{X} \equiv \{0, 1, 2, \dots\}$

LEVEL 1 :  $Y_1, \dots, Y_n|X = x \sim f_Y(y|\underline{\theta})$  (independent), and  $S = \sum_{i=1}^x Y_i$

Then, by the law of iterated expectation,

$$\begin{aligned} M_S(t) = \mathbb{E}_{f_S} [e^{tS}] &= \mathbb{E}_{f_X} \left[ \mathbb{E}_{f_{S|X}} [e^{tS} | X = x] \right] \\ &= \mathbb{E}_{f_X} \left[ \mathbb{E}_{f_{Y|X}} \left[ \exp \left\{ t \sum_{i=1}^x Y_i \right\} | X = x \right] \right] \\ &= \mathbb{E}_{f_X} [\{M_Y(t)\}^X] \\ &= G_X(M_Y(t)) \end{aligned}$$

where  $G_X$  is the factorial mgf (or pgf) for  $X$ . By a similar calculation,

$$G_S(t) = G_X(G_Y(t)).$$

For example, if  $X \sim \text{Poisson}(\phi)$ , then

$$G_S(t) = \exp \{ \phi (G_Y(t) - 1) \}$$

is the pgf of  $S$ . Expanding the pgf as a power series in  $t$  yields the pmf of  $S$ .

### Example : Branching Process

Consider a sequence of generations of an organism; let  $S_i$  be the total number of individuals in the  $i$ th generation, for  $i = 0, 1, 2, \dots$ . Suppose that  $f_X$  is a pmf with support  $\mathbb{X} \equiv \{0, 1, 2, \dots\}$ .

- **Generation 0** :  $S_0 \sim f_X(x|\phi)$
- **Generation 1** : Given  $S_0 = s_0$ , let

$$S_{11}, \dots, S_{1s_0} | S_0 = s_0 \quad \text{such that} \quad S_{1j} \sim f_X(x|\phi), \text{ with } S_{1j_1} \perp S_{1j_2} \text{ for all } j_1, j_2$$

and set

$$S_1 = \sum_{j=1}^{s_0} S_{1j}$$

is the total number of individuals in the 1st generation.  $S_{1j}$  is the number of offspring of the  $j$ th individual in the zeroth generation.

- **Generation i** : Given  $S_{i-1} = s_{i-1}$ , let

$$S_{i1}, \dots, S_{is_{i-1}} | S_{i-1} = s_{i-1} \quad \text{such that} \quad S_{ij} \sim f_X(x|\phi) \text{ (independent)}$$

and set

$$S_i = \sum_{j=1}^{s_{i-1}} S_{ij}$$

Let  $G_i$  be the pgf of  $S_i$ . Then, by recursion, we have

$$G_i(t) = G_{i-1}(G_X(t)) = G_{i-2}(G_X(G_X(t))) = \dots = G_X(G_X(\dots G_X(G_X(t)) \dots))$$

that is, an  $i + 1$ -fold iterated calculation.

### 3. Location-Scale Mixtures

LEVEL 3 :	$\theta$	Fixed parameters
LEVEL 2 :	$M, V \sim f_{M,V}(m, v \theta)$	
LEVEL 1 :	$Y M = m, V = v \sim f_{Y M,V}(y m, v)$	

where

$$f_{Y|M,V}(y|m, v) = \frac{1}{v} f\left(\frac{y-m}{v}\right)$$

that is a location-scale family distribution, mixed over different location and scale parameters with *mixing distribution*  $f_{M,V}$ .

### Example : Scale Mixtures of Normal Distributions

LEVEL 3 :	$\theta$	
LEVEL 2 :	$V \sim f_V(v \theta)$	
LEVEL 1 :	$Y V = v \sim f_{Y V}(y v) \equiv \text{Normal}(0, g(v))$	

for some positive function  $g$ .

For example, if

$$Y|V = v \sim \text{Normal}(0, v^{-1}) \quad V \sim \text{Gamma}\left(\frac{1}{2}, \frac{1}{2}\right)$$

then by elementary calculations, we find that

$$f_Y(y) = \frac{1}{\pi} \frac{1}{1+y^2} \quad y \in \mathbb{R} \quad \therefore \quad Y \sim \text{Cauchy}.$$

The scale mixture of normal distributions family includes the *Student*, *Double Exponential* and *Logistic* as special cases.

Moments of location-scale mixtures can be computed using the law of iterated expectation. The location-scale mixture construction allows the modelling of

- **skewness** through the mixture over different *locations*
- **kurtosis** through the mixture over different *scales*

### Example : Location-Scale Mixtures of Normal Distributions

Suppose  $M$  and  $V$  are independent, with

$$M \sim \text{Exponential}(1/2) \quad V \sim \text{Gamma}(2, 1/2)$$

and

$$Y|M = m, V = v \sim \text{Normal}(m, 1/v)$$

Then the marginal distribution of  $Y$  is given by

$$f_Y(y) = \int_0^\infty \int_0^\infty f_{Y|M,V}(y|m, v) f_M(m) f_V(v) dm dv$$

which can most readily be examined by simulation. The figure below depicts a histogram of 10000 values simulated from the model, and demonstrates the skewness of the marginal of  $Y$ .

