

# COMS 4771 Lecture 21

## 1. Exponential families

# EXPONENTIAL FAMILIES

# EXPONENTIAL FAMILIES

We saw that solutions to the maximum entropy optimization problem are distributions of the form

$$P_{\boldsymbol{\eta}}(x) = \exp\left\{\langle \boldsymbol{\eta}, \boldsymbol{T}(x) \rangle - G(\boldsymbol{\eta})\right\} \cdot \pi(x) \quad \forall x \in \mathcal{X}. \quad (\star)$$

- ▶  $\pi$ : base distribution over domain  $\mathcal{X}$  (okay if unnormalized).
- ▶  $\boldsymbol{T}: \mathcal{X} \rightarrow \mathbb{R}^k$ : vector-valued feature function.
- ▶  $G(\boldsymbol{\eta})$ : log of normalizer that makes  $P_{\boldsymbol{\eta}}$  a probability distribution over  $\mathcal{X}$ .

# EXPONENTIAL FAMILIES

We saw that solutions to the maximum entropy optimization problem are distributions of the form

$$P_{\boldsymbol{\eta}}(x) = \exp\left\{\langle \boldsymbol{\eta}, \boldsymbol{T}(x) \rangle - G(\boldsymbol{\eta})\right\} \cdot \pi(x) \quad \forall x \in \mathcal{X}. \quad (\star)$$

- ▶  $\pi$ : base distribution over domain  $\mathcal{X}$  (okay if unnormalized).
- ▶  $\boldsymbol{T}: \mathcal{X} \rightarrow \mathbb{R}^k$ : vector-valued feature function.
- ▶  $G(\boldsymbol{\eta})$ : log of normalizer that makes  $P_{\boldsymbol{\eta}}$  a probability distribution over  $\mathcal{X}$ .

Each choice of  $\pi$  and  $\boldsymbol{T}$  leads to a **family of probability distributions**

$$\left\{ P_{\boldsymbol{\eta}} \text{ as in } (\star) : \boldsymbol{\eta} \in \mathbb{R}^k \right\}^{\dagger}$$

We call such a family an **exponential family**.

# EXPONENTIAL FAMILIES

We saw that solutions to the maximum entropy optimization problem are distributions of the form

$$P_{\boldsymbol{\eta}}(x) = \exp\left\{\langle \boldsymbol{\eta}, \boldsymbol{T}(x) \rangle - G(\boldsymbol{\eta})\right\} \cdot \pi(x) \quad \forall x \in \mathcal{X}. \quad (\star)$$

- ▶  $\pi$ : base distribution over domain  $\mathcal{X}$  (okay if unnormalized).
- ▶  $\boldsymbol{T}: \mathcal{X} \rightarrow \mathbb{R}^k$ : vector-valued feature function.
- ▶  $G(\boldsymbol{\eta})$ : log of normalizer that makes  $P_{\boldsymbol{\eta}}$  a probability distribution over  $\mathcal{X}$ .

Each choice of  $\pi$  and  $\boldsymbol{T}$  leads to a **family of probability distributions**

$$\left\{ P_{\boldsymbol{\eta}} \text{ as in } (\star) : \boldsymbol{\eta} \in \mathbb{R}^k \right\}^{\ddagger}$$

We call such a family an **exponential family**.

<sup>‡</sup> **This is not entirely accurate, due to an important technicality.**

# IMPORTANT TECHNICALITY

The log of the normalizer (which is also called the **log partition function**),

$$G(\boldsymbol{\eta}) = \ln Z(\boldsymbol{\eta}) = \ln \left( \sum_{x \in \mathcal{X}} \exp \left\{ \langle \boldsymbol{\eta}, \mathbf{T}(x) \rangle \right\} \pi(x) \right),$$

**could be infinite** for some values of  $\boldsymbol{\eta} \in \mathbb{R}^k$ .

# IMPORTANT TECHNICALITY

The log of the normalizer (which is also called the **log partition function**),

$$G(\boldsymbol{\eta}) = \ln Z(\boldsymbol{\eta}) = \ln \left( \sum_{x \in \mathcal{X}} \exp \left\{ \langle \boldsymbol{\eta}, \mathbf{T}(x) \rangle \right\} \pi(x) \right),$$

**could be infinite** for some values of  $\boldsymbol{\eta} \in \mathbb{R}^k$ .

**Example:**  $\mathcal{X} = \mathbb{N}$ ,  $T_1(x) = \ln(x)$ ,  $\pi(x) = 1/x^2$ , so

$$G(\eta) = \sum_{x \in \mathbb{N}} \exp \left\{ \eta \ln(x) \right\} \cdot \frac{1}{x^2}$$

# IMPORTANT TECHNICALITY

The log of the normalizer (which is also called the **log partition function**),

$$G(\boldsymbol{\eta}) = \ln Z(\boldsymbol{\eta}) = \ln \left( \sum_{x \in \mathcal{X}} \exp \left\{ \langle \boldsymbol{\eta}, \mathbf{T}(x) \rangle \right\} \pi(x) \right),$$

**could be infinite** for some values of  $\boldsymbol{\eta} \in \mathbb{R}^k$ .

**Example:**  $\mathcal{X} = \mathbb{N}$ ,  $T_1(x) = \ln(x)$ ,  $\pi(x) = 1/x^2$ , so

$$G(\eta) = \sum_{x \in \mathbb{N}} \exp \left\{ \eta \ln(x) \right\} \cdot \frac{1}{x^2} = \sum_{x \in \mathbb{N}} \frac{1}{x^{2-\eta}}$$

which finite if and only if  $\eta < 1$ .



# IMPORTANT TECHNICALITY

The **log of the normalizer** (which is also called the **log partition function**),

$$G(\boldsymbol{\eta}) = \ln Z(\boldsymbol{\eta}) = \ln \left( \sum_{x \in \mathcal{X}} \exp \left\{ \langle \boldsymbol{\eta}, \mathbf{T}(x) \rangle \right\} \pi(x) \right),$$

**could be infinite** for some values of  $\boldsymbol{\eta} \in \mathbb{R}^k$ .

**Example:**  $\mathcal{X} = \mathbb{N}$ ,  $T_1(x) = \ln(x)$ ,  $\pi(x) = 1/x^2$ , so

$$G(\eta) = \sum_{x \in \mathbb{N}} \exp \left\{ \eta \ln(x) \right\} \cdot \frac{1}{x^2} = \sum_{x \in \mathbb{N}} \frac{1}{x^{2-\eta}}$$

which finite if and only if  $\eta < 1$ .

**Parameter values  $\eta \geq 1$  cause  $G(\eta)$  to be infinite, and hence do not yield valid probability distributions.**

# REVISED DEFINITION

## Revised definition

The **exponential family** corresponding to  $\mathbf{T}$  and  $\pi$  is

$$\{P_{\boldsymbol{\eta}} \text{ as in } (\star) : \boldsymbol{\eta} \in \mathcal{N}\}$$

where

$$\mathcal{N} = \left\{ \boldsymbol{\eta} \in \mathbb{R}^k : G(\boldsymbol{\eta}) < \infty \right\}$$

is the **natural parameter space** for this exponential family.

$$\begin{aligned} P_{\boldsymbol{\eta}}(x) &= \exp\left\{ \langle \boldsymbol{\eta}, \mathbf{T}(x) \rangle - G(\boldsymbol{\eta}) \right\} \cdot \pi(x) \quad \forall x \in \mathcal{X}, \\ G(\boldsymbol{\eta}) &= \ln \left( \sum_{x \in \mathcal{X}} \exp\left\{ \langle \boldsymbol{\eta}, \mathbf{T}(x) \rangle \right\} \pi(x) \right). \end{aligned} \tag{\star}$$

## EXAMPLE #1

$$\mathcal{X} = \{0, 1\}, \quad T_1(x) = \mathbb{1}\{x = 1\} = x, \quad \pi(x) = 1.$$

# EXAMPLE #1

equal meaning

$$\mathcal{X} = \{0, 1\}, \quad T_1(x) = \mathbb{1}\{x = 1\} = x, \quad \pi(x) = 1.$$

What is the natural parameter space  $\mathcal{N} = \{\eta \in \mathbb{R} : G(\eta) < \infty\}$ ?

# EXAMPLE #1

$$\mathcal{X} = \{0, 1\}, \quad T_1(x) = \mathbb{1}\{x = 1\} = x, \quad \pi(x) = 1.$$

What is the natural parameter space  $\mathcal{N} = \{\eta \in \mathbb{R} : G(\eta) < \infty\}$ ?

$$G(\eta) = \ln \left( \sum_{x \in \{0, 1\}} \exp\{\eta x\} \right) = \ln(1 + \exp(\eta)) < \infty \quad \text{for all } \eta \in \mathbb{R},$$

so  $\mathcal{N} = \mathbb{R}$ .

# EXAMPLE #1

$$\mathcal{X} = \{0, 1\}, \quad T_1(x) = \mathbb{1}\{x = 1\} = x, \quad \pi(x) = 1.$$

What is the natural parameter space  $\mathcal{N} = \{\eta \in \mathbb{R} : G(\eta) < \infty\}$ ?

$$G(\eta) = \ln \left( \sum_{x \in \{0,1\}} \exp\{\eta x\} \right) = \ln(1 + \exp(\eta)) < \infty \quad \text{for all } \eta \in \mathbb{R},$$

so  $\mathcal{N} = \mathbb{R}$ .

Probability distributions have the form

$$P_\eta(x) = \exp\left\{\eta x - \ln(1 + \exp(\eta))\right\}$$

# EXAMPLE #1

$$\mathcal{X} = \{0, 1\}, \quad T_1(x) = \mathbb{1}\{x = 1\} = x, \quad \pi(x) = 1.$$

What is the natural parameter space  $\mathcal{N} = \{\eta \in \mathbb{R} : G(\eta) < \infty\}$ ?

$$G(\eta) = \ln \left( \sum_{x \in \{0,1\}} \exp\{\eta x\} \right) = \ln(1 + \exp(\eta)) < \infty \quad \text{for all } \eta \in \mathbb{R},$$

so  $\mathcal{N} = \mathbb{R}$ .

Probability distributions have the form

$$P_\eta(x) = \exp\left\{\eta x - \ln(1 + \exp(\eta))\right\} = \frac{\exp(\eta x)}{1 + \exp(\eta)}$$

## EXAMPLE #1

$$\mathcal{X} = \{0, 1\}, \quad T_1(x) = \mathbb{1}\{x = 1\} = x, \quad \pi(x) = 1.$$

What is the natural parameter space  $\mathcal{N} = \{\eta \in \mathbb{R} : G(\eta) < \infty\}$ ?

$$G(\eta) = \ln \left( \sum_{x \in \{0,1\}} \exp\{\eta x\} \right) = \ln(1 + \exp(\eta)) < \infty \quad \text{for all } \eta \in \mathbb{R},$$

so  $\mathcal{N} = \mathbb{R}$ .

Probability distributions have the form

$$P_\eta(x) = \exp\left\{\eta x - \ln(1 + \exp(\eta))\right\} = \frac{\exp(\eta x)}{1 + \exp(\eta)} = \begin{cases} \frac{1}{1 + e^\eta} & \text{if } x = 0; \\ \frac{e^\eta}{1 + e^\eta} & \text{if } x = 1. \end{cases}$$



# EXAMPLE #1

$\pi(x)$  is a constant function.

$\mathcal{X} = \{0, 1\}$ ,  $T_1(x) = \mathbb{1}\{x = 1\} = x$ ,  $\pi(x) = 1$ . equal meaning

What is the natural parameter space  $\mathcal{N} = \{\eta \in \mathbb{R} : G(\eta) < \infty\}$ ?

$$G(\eta) = \ln \left( \sum_{x \in \{0,1\}} \exp\{\eta x\} \right) = \ln(1 + \exp(\eta)) < \infty \quad \text{for all } \eta \in \mathbb{R},$$

so  $\mathcal{N} = \mathbb{R}$ .

Probability distributions have the form

$$P_\eta(x) = \exp\left\{\eta x - \ln(1 + \exp(\eta))\right\} = \frac{\exp(\eta x)}{1 + \exp(\eta)} = \begin{cases} \frac{1}{1 + e^\eta} & \text{if } x = 0; \\ \frac{e^\eta}{1 + e^\eta} & \text{if } x = 1. \end{cases}$$

**Bernoulli distributions**  $\text{Bern}\left(\frac{e^\eta}{1+e^\eta}\right)$

## EXAMPLE #1.5

$$\mathcal{X} = \{0, 1, \dots, d-1\}, \quad T_i(x) = \mathbb{1}\{x = i\} \text{ for } i \in \{1, \dots, d-1\}, \quad \pi(x) = 1.$$

## EXAMPLE #1.5

$$\mathcal{X} = \{0, 1, \dots, d-1\}, \quad T_i(x) = \mathbb{1}\{x = i\} \text{ for } i \in \{1, \dots, d-1\}, \quad \pi(x) = 1.$$

$$G(\boldsymbol{\eta}) = \ln \left( \sum_{x=0}^{d-1} \exp \left\{ \langle \boldsymbol{\eta}, \mathbf{T}(x) \rangle \right\} \right) = \ln \left( 1 + \sum_{i=1}^{d-1} e^{\eta_i} \right)$$

$$\text{so } \mathcal{N} = \mathbb{R}^{d-1}.$$

## EXAMPLE #1.5

$\mathcal{X} = \{0, 1, \dots, d-1\}$ ,  $T_i(x) = \mathbb{1}\{x = i\}$  for  $i \in \{1, \dots, d-1\}$ ,  $\pi(x) = 1$ .

$$G(\boldsymbol{\eta}) = \ln \left( \sum_{x=0}^{d-1} \exp \left\{ \langle \boldsymbol{\eta}, \mathbf{T}(x) \rangle \right\} \right) = \ln \left( 1 + \sum_{i=1}^{d-1} e^{\eta_i} \right)$$

so  $\mathcal{N} = \mathbb{R}^{d-1}$ .

Probability distributions have the form

$$P_{\boldsymbol{\eta}}(x) = \exp \left\{ \langle \boldsymbol{\eta}, \mathbf{T}(x) \rangle - \ln \left( 1 + \sum_{i=1}^{d-1} e^{\eta_i} \right) \right\}$$

## EXAMPLE #1.5

$\mathcal{X} = \{0, 1, \dots, d-1\}$ ,  $T_i(x) = \mathbb{1}\{x = i\}$  for  $i \in \{1, \dots, d-1\}$ ,  $\pi(x) = 1$ .

$$G(\boldsymbol{\eta}) = \ln \left( \sum_{x=0}^{d-1} \exp \left\{ \langle \boldsymbol{\eta}, \mathbf{T}(x) \rangle \right\} \right) = \ln \left( 1 + \sum_{i=1}^{d-1} e^{\eta_i} \right)$$

so  $\mathcal{N} = \mathbb{R}^{d-1}$ .

Probability distributions have the form

$$P_{\boldsymbol{\eta}}(x) = \exp \left\{ \langle \boldsymbol{\eta}, \mathbf{T}(x) \rangle - \ln \left( 1 + \sum_{i=1}^{d-1} e^{\eta_i} \right) \right\} = \begin{cases} \frac{1}{1 + \sum_{i=1}^{d-1} e^{\eta_i}} & \text{if } x = 0; \\ \frac{e^{\eta_i}}{1 + \sum_{i=1}^{d-1} e^{\eta_i}} & \text{if } x = i \neq 0. \end{cases}$$

## EXAMPLE #1.5

$\mathcal{X} = \{0, 1, \dots, d-1\}$ ,  $T_i(x) = \mathbb{1}\{x = i\}$  for  $i \in \{1, \dots, d-1\}$ ,  $\pi(x) = 1$ .

$$G(\boldsymbol{\eta}) = \ln \left( \sum_{x=0}^{d-1} \exp \left\{ \langle \boldsymbol{\eta}, \mathbf{T}(x) \rangle \right\} \right) = \ln \left( 1 + \sum_{i=1}^{d-1} e^{\eta_i} \right)$$

so  $\mathcal{N} = \mathbb{R}^{d-1}$ .

Probability distributions have the form

$$P_{\boldsymbol{\eta}}(x) = \exp \left\{ \langle \boldsymbol{\eta}, \mathbf{T}(x) \rangle - \ln \left( 1 + \sum_{i=1}^{d-1} e^{\eta_i} \right) \right\} = \begin{cases} \frac{1}{1 + \sum_{i=1}^{d-1} e^{\eta_i}} & \text{if } x = 0; \\ \frac{e^{\eta_i}}{1 + \sum_{i=1}^{d-1} e^{\eta_i}} & \text{if } x = i \neq 0. \end{cases}$$

**Categorical distributions** (generalizes Bernoulli)

For  $X \sim P_{\boldsymbol{\eta}}$ ,  $\Pr[X = i] \propto e^{\eta_i}$  for  $i \neq 0$ .

## EXAMPLE #1.5

no need to include 0, redundant! can get the information from other  $T_i$

$$\mathcal{X} = \{0, 1, \dots, d-1\}, \quad T_i(x) = \mathbb{1}\{x = i\} \text{ for } i \in \{1, \dots, d-1\}, \quad \pi(x) = 1.$$

$$G(\boldsymbol{\eta}) = \ln \left( \sum_{x=0}^{d-1} \exp \left\{ \langle \boldsymbol{\eta}, \mathbf{T}(x) \rangle \right\} \right) = \ln \left( 1 + \sum_{i=1}^{d-1} e^{\eta_i} \right)$$

so  $\mathcal{N} = \mathbb{R}^{d-1}$ .

Probability distributions have the form

$$P_{\boldsymbol{\eta}}(x) = \exp \left\{ \langle \boldsymbol{\eta}, \mathbf{T}(x) \rangle - \ln \left( 1 + \sum_{i=1}^{d-1} e^{\eta_i} \right) \right\} = \begin{cases} \frac{1}{1 + \sum_{i=1}^{d-1} e^{\eta_i}} & \text{if } x = 0; \\ \frac{e^{\eta_i}}{1 + \sum_{i=1}^{d-1} e^{\eta_i}} & \text{if } x = i \neq 0. \end{cases}$$

**Categorical distributions** (generalizes Bernoulli)

For  $X \sim P_{\boldsymbol{\eta}}$ ,  $\Pr[X = i] \propto e^{\eta_i}$  for  $i \neq 0$ .

**Note:** keep track of  $d-1$  features—almost one for every  $x \in \mathcal{X}$ .

## EXAMPLE #2

$\mathcal{X}$  = non-negative integers,  $T_1(x) = x$ ,  $\pi(x) = 1/x!$ .



## EXAMPLE #2

$\mathcal{X}$  = non-negative integers,  $T_1(x) = x$ ,  $\pi(x) = 1/x!$ .

Using the fact that  $\ln\left(\sum_{j=0}^{\infty} z^j/j!\right) = \ln \exp(z) = z$ ,

$$G(\eta) = \ln\left(\sum_{x=0}^{\infty} \exp\{\eta x\} \cdot \frac{1}{x!}\right) = \exp(\eta),$$

so  $\mathcal{N} = \mathbb{R}$ .

## EXAMPLE #2

$\mathcal{X}$  = non-negative integers,  $T_1(x) = x$ ,  $\pi(x) = 1/x!$ .

Using the fact that  $\ln\left(\sum_{j=0}^{\infty} z^j/j!\right) = \ln \exp(z) = z$ ,

$$G(\eta) = \ln\left(\sum_{x=0}^{\infty} \exp\{\eta x\} \cdot \frac{1}{x!}\right) = \exp(\eta),$$

so  $\mathcal{N} = \mathbb{R}$ .

Probability distributions have the form

$$P_{\eta}(x) = \exp\{\eta x - \exp(\eta)\} \cdot \frac{1}{x!}$$

## EXAMPLE #2

$\mathcal{X}$  = non-negative integers,  $T_1(x) = x$ ,  $\pi(x) = 1/x!$ .

Using the fact that  $\ln\left(\sum_{j=0}^{\infty} z^j/j!\right) = \ln \exp(z) = z$ ,

$$G(\eta) = \ln\left(\sum_{x=0}^{\infty} \exp\{\eta x\} \cdot \frac{1}{x!}\right) = \exp(\eta),$$

so  $\mathcal{N} = \mathbb{R}$ .

Probability distributions have the form

$$P_{\eta}(x) = \exp\{\eta x - \exp(\eta)\} \cdot \frac{1}{x!} = \frac{(e^{\eta})^x \cdot \exp(-e^{\eta})}{x!}.$$

## EXAMPLE #2

$\mathcal{X}$  = non-negative integers,  $T_1(x) = x$ ,  $\pi(x) = 1/x!$ .

Using the fact that  $\ln\left(\sum_{j=0}^{\infty} z^j/j!\right) = \ln \exp(z) = z$ ,

$$G(\eta) = \ln\left(\sum_{x=0}^{\infty} \exp\{\eta x\} \cdot \frac{1}{x!}\right) = \exp(\eta), \quad \text{always converges}$$

so  $\mathcal{N} = \mathbb{R}$ .

Probability distributions have the form

$$P_{\eta}(x) = \exp\{\eta x - \exp(\eta)\} \cdot \frac{1}{x!} = \frac{(e^{\eta})^x \cdot \exp(-e^{\eta})}{x!}.$$

**Poisson distributions**  $\text{Pois}(e^{\eta})$

# CONTINUOUS DOMAINS

If  $\mathcal{X}$  is a continuous domain (e.g.,  $\mathbb{R}^d$ ), we can also obtain exponential families comprised of probability densities having the form

$$\begin{aligned} p_{\boldsymbol{\eta}}(x) &= \exp\left\{\langle \boldsymbol{\eta}, \mathbf{T}(x) \rangle - G(\boldsymbol{\eta})\right\} \cdot \pi(x) \quad \forall x \in \mathcal{X}, \\ G(\boldsymbol{\eta}) &= \ln\left(\int_{\mathcal{X}} \exp\left\{\langle \boldsymbol{\eta}, \mathbf{T}(x) \rangle\right\} \cdot \pi(x) \, dx\right). \end{aligned} \tag{*}$$

where  $\pi$  is a “base measure” (often a density) on  $\mathcal{X}$ .

# CONTINUOUS DOMAINS

If  $\mathcal{X}$  is a continuous domain (e.g.,  $\mathbb{R}^d$ ), we can also obtain exponential families comprised of probability densities having the form

$$\begin{aligned} p_{\boldsymbol{\eta}}(x) &= \exp\left\{\langle \boldsymbol{\eta}, \mathbf{T}(x) \rangle - G(\boldsymbol{\eta})\right\} \cdot \pi(x) \quad \forall x \in \mathcal{X}, \\ G(\boldsymbol{\eta}) &= \ln\left(\int_{\mathcal{X}} \exp\left\{\langle \boldsymbol{\eta}, \mathbf{T}(x) \rangle\right\} \cdot \pi(x) \, dx\right). \end{aligned} \tag{*}$$

where  $\pi$  is a “base measure” (often a density) on  $\mathcal{X}$ .

The exponential family corresponding to  $\mathbf{T}$  and  $\pi$  is

$$\{P_{\boldsymbol{\eta}} \text{ as in } (*) : \boldsymbol{\eta} \in \mathcal{N}\} \quad \text{where } \mathcal{N} = \left\{\boldsymbol{\eta} \in \mathbb{R}^k : G(\boldsymbol{\eta}) < \infty\right\}.$$

# CONTINUOUS DOMAINS

If  $\mathcal{X}$  is a continuous domain (e.g.,  $\mathbb{R}^d$ ), we can also obtain exponential families comprised of probability densities having the form

$$\begin{aligned} p_{\boldsymbol{\eta}}(x) &= \exp\left\{\langle \boldsymbol{\eta}, \mathbf{T}(x) \rangle - G(\boldsymbol{\eta})\right\} \cdot \pi(x) \quad \forall x \in \mathcal{X}, \\ G(\boldsymbol{\eta}) &= \ln\left(\int_{\mathcal{X}} \exp\left\{\langle \boldsymbol{\eta}, \mathbf{T}(x) \rangle\right\} \cdot \pi(x) \, dx\right). \end{aligned} \tag{*}$$

where  $\pi$  is a “base measure” (often a density) on  $\mathcal{X}$ .

The exponential family corresponding to  $\mathbf{T}$  and  $\pi$  is

$$\{P_{\boldsymbol{\eta}} \text{ as in } (*) : \boldsymbol{\eta} \in \mathcal{N}\} \quad \text{where } \mathcal{N} = \left\{\boldsymbol{\eta} \in \mathbb{R}^k : G(\boldsymbol{\eta}) < \infty\right\}.$$

(A more general treatment requires a little bit of measure theory, which we'll forego.)

## EXAMPLE #3

$$\mathcal{X} = \mathbb{R}, \quad T_1(x) = x, \quad \pi(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}.$$



## EXAMPLE #3

$$\mathcal{X} = \mathbb{R}, \quad T_1(x) = x, \quad \pi(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}.$$

$$G(\eta) = \ln \left( \int_{\mathbb{R}} \exp\{\eta x\} \cdot \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx \right) = \frac{\eta^2}{2},$$

so  $\mathcal{N} = \mathbb{R}$ .

## EXAMPLE #3

$$\mathcal{X} = \mathbb{R}, \quad T_1(x) = x, \quad \pi(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}.$$

$$G(\eta) = \ln \left( \int_{\mathbb{R}} \exp\{\eta x\} \cdot \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx \right) = \frac{\eta^2}{2},$$

so  $\mathcal{N} = \mathbb{R}$ .

$$p_{\eta}(x) = \exp\left\{\eta x - \eta^2/2\right\} \cdot \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$

## EXAMPLE #3

$$\mathcal{X} = \mathbb{R}, \quad T_1(x) = x, \quad \pi(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}.$$

$$G(\eta) = \ln \left( \int_{\mathbb{R}} \exp\{\eta x\} \cdot \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx \right) = \frac{\eta^2}{2},$$

so  $\mathcal{N} = \mathbb{R}$ .

$$p_{\eta}(x) = \exp\left\{\eta x - \eta^2/2\right\} \cdot \frac{1}{\sqrt{2\pi}} e^{-x^2/2} = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{(x - \eta)^2}{2}\right\}.$$

## EXAMPLE #3

$$\mathcal{X} = \mathbb{R}, \quad T_1(x) = x, \quad \pi(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}.$$

$$G(\eta) = \ln \left( \int_{\mathbb{R}} \exp\{\eta x\} \cdot \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx \right) = \frac{\eta^2}{2},$$

so  $\mathcal{N} = \mathbb{R}$ .

$$p_{\eta}(x) = \exp\left\{\eta x - \eta^2/2\right\} \cdot \frac{1}{\sqrt{2\pi}} e^{-x^2/2} = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{(x - \eta)^2}{2}\right\}.$$

**Gaussian distributions with unit variance  $N(\eta, 1)$**

## EXAMPLE #4

$$\mathcal{X} = \mathbb{R}, \quad T_1(x) = x, \quad T_2(x) = x^2, \quad \pi(x) = \frac{1}{\sqrt{2\pi}} \text{ (this is actually okay).}$$

## EXAMPLE #4

$\mathcal{X} = \mathbb{R}$ ,  $T_1(x) = x$ ,  $T_2(x) = x^2$ ,  $\pi(x) = \frac{1}{\sqrt{2\pi}}$  (this is actually okay).

$$G(\eta) = \ln \left( \int_{\mathbb{R}} \exp \left\{ \eta_1 x + \eta_2 x^2 \right\} \cdot \frac{1}{\sqrt{2\pi}} dx \right)$$

## EXAMPLE #4

$\mathcal{X} = \mathbb{R}$ ,  $T_1(x) = x$ ,  $T_2(x) = x^2$ ,  $\pi(x) = \frac{1}{\sqrt{2\pi}}$  (this is actually okay).

$$\begin{aligned} G(\eta) &= \ln \left( \int_{\mathbb{R}} \exp \left\{ \eta_1 x + \eta_2 x^2 \right\} \cdot \frac{1}{\sqrt{2\pi}} dx \right) \\ &= \ln \left( \int_{\mathbb{R}} \exp \left\{ -\frac{1}{2} \cdot \frac{x^2 + \eta_1 x / \eta_2}{-1/(2\eta_2)} \right\} \cdot \frac{1}{\sqrt{2\pi}} dx \right) \end{aligned}$$

$$(\text{completing the square } \dots) = -\ln \sqrt{-2\eta_2} - \frac{\eta_1^2}{4\eta_2}.$$

## EXAMPLE #4

$\mathcal{X} = \mathbb{R}$ ,  $T_1(x) = x$ ,  $T_2(x) = x^2$ ,  $\pi(x) = \frac{1}{\sqrt{2\pi}}$  (this is actually okay).

$$\begin{aligned} G(\eta) &= \ln \left( \int_{\mathbb{R}} \exp \left\{ \eta_1 x + \eta_2 x^2 \right\} \cdot \frac{1}{\sqrt{2\pi}} dx \right) \\ &= \ln \left( \int_{\mathbb{R}} \exp \left\{ -\frac{1}{2} \cdot \frac{x^2 + \eta_1 x / \eta_2}{-1/(2\eta_2)} \right\} \cdot \frac{1}{\sqrt{2\pi}} dx \right) \end{aligned}$$

$$(\text{completing the square } \dots) = -\ln \sqrt{-2\eta_2} - \frac{\eta_1^2}{4\eta_2}.$$

The “completing the square” step assumes that  $-1/(2\eta_2) > 0$ , or else the integral would diverge. So  $\mathcal{N} = \mathbb{R} \times \mathbb{R}_{--}$  (i.e.,  $\eta_2$  must be negative).



## EXAMPLE #4

$\mathcal{X} = \mathbb{R}$ ,  $T_1(x) = x$ ,  $T_2(x) = x^2$ ,  $\pi(x) = \frac{1}{\sqrt{2\pi}}$  (this is actually okay).

$$\begin{aligned} G(\eta) &= \ln \left( \int_{\mathbb{R}} \exp \left\{ \eta_1 x + \eta_2 x^2 \right\} \cdot \frac{1}{\sqrt{2\pi}} dx \right) \\ &= \ln \left( \int_{\mathbb{R}} \exp \left\{ -\frac{1}{2} \cdot \frac{x^2 + \eta_1 x / \eta_2}{-1/(2\eta_2)} \right\} \cdot \frac{1}{\sqrt{2\pi}} dx \right) \end{aligned}$$

$$(\text{completing the square } \dots) = -\ln \sqrt{-2\eta_2} - \frac{\eta_1^2}{4\eta_2}.$$

The “completing the square” step assumes that  $-1/(2\eta_2) > 0$ , or else the integral would diverge. So  $\mathcal{N} = \mathbb{R} \times \mathbb{R}_{--}$  (i.e.,  $\eta_2$  must be negative).

$$p_{\eta}(x) = \exp \left\{ \eta_1 x + \eta_2 x^2 + \ln \sqrt{-2\eta_2} + \frac{\eta_1^2}{4\eta_2} \right\} \cdot \frac{1}{\sqrt{2\pi}}$$

## EXAMPLE #4

$\mathcal{X} = \mathbb{R}$ ,  $T_1(x) = x$ ,  $T_2(x) = x^2$ ,  $\pi(x) = \frac{1}{\sqrt{2\pi}}$  (this is actually okay).

$$\begin{aligned} G(\eta) &= \ln \left( \int_{\mathbb{R}} \exp \left\{ \eta_1 x + \eta_2 x^2 \right\} \cdot \frac{1}{\sqrt{2\pi}} dx \right) \\ &= \ln \left( \int_{\mathbb{R}} \exp \left\{ -\frac{1}{2} \cdot \frac{x^2 + \eta_1 x / \eta_2}{-1/(2\eta_2)} \right\} \cdot \frac{1}{\sqrt{2\pi}} dx \right) \end{aligned}$$

$$(\text{completing the square } \dots) = -\ln \sqrt{-2\eta_2} - \frac{\eta_1^2}{4\eta_2}.$$

The “completing the square” step assumes that  $-1/(2\eta_2) > 0$ , or else the integral would diverge. So  $\mathcal{N} = \mathbb{R} \times \mathbb{R}_{--}$  (i.e.,  $\eta_2$  must be negative).

$$p_{\eta}(x) = \exp \left\{ \eta_1 x + \eta_2 x^2 + \ln \sqrt{-2\eta_2} + \frac{\eta_1^2}{4\eta_2} \right\} \cdot \frac{1}{\sqrt{2\pi}} = \dots$$

$$\text{Gaussian distributions } \mathcal{N} \left( -\frac{\eta_1}{2\eta_2}, -\frac{1}{2\eta_2} \right)$$

# PARAMETERIZATION

# NATURAL PARAMETERS

In most of the examples so far, the “natural parameters”  $\boldsymbol{\eta}$  are not the parameters we typically use for the distribution families.

# NATURAL PARAMETERS

In most of the examples so far, the “natural parameters”  $\boldsymbol{\eta}$  are not the parameters we typically use for the distribution families.

**Example #1:** Bernoulli distribution  $\text{Bern}(p)$  for  $p = \frac{e^\eta}{1+e^\eta}$ .

$$X \sim \text{Bern}(p) \implies \mathbb{E}[X] = p.$$

# NATURAL PARAMETERS

In most of the examples so far, the “natural parameters”  $\boldsymbol{\eta}$  are not the parameters we typically use for the distribution families.

**Example #1:** Bernoulli distribution  $\text{Bern}(p)$  for  $p = \frac{e^\eta}{1+e^\eta}$ .

$$X \sim \text{Bern}(p) \implies \mathbb{E}[X] = p.$$

**Example #2:** Poisson distribution  $\text{Pois}(\lambda)$  for  $\lambda = e^\eta$ .

$$X \sim \text{Pois}(\lambda) \implies \mathbb{E}[X] = \lambda.$$

# NATURAL PARAMETERS

In most of the examples so far, the “natural parameters”  $\eta$  are not the parameters we typically use for the distribution families.

**Example #1:** Bernoulli distribution  $\text{Bern}(p)$  for  $p = \frac{e^\eta}{1+e^\eta}$ .

$$X \sim \text{Bern}(p) \implies \mathbb{E}[X] = p.$$

**Example #2:** Poisson distribution  $\text{Pois}(\lambda)$  for  $\lambda = e^\eta$ .

$$X \sim \text{Pois}(\lambda) \implies \mathbb{E}[X] = \lambda.$$

**Example #4:** Gaussian distribution  $N(\mu, \sigma^2)$  for  $\mu = -\frac{\eta_1}{2\eta_2}$  and  $\sigma^2 = -\frac{1}{2\eta_2}$ .

$$X \sim N(\mu, \sigma^2) \implies \mathbb{E}[X] = \mu, \quad \mathbb{E}[X^2] = \sigma^2 + \mu^2.$$

# NATURAL PARAMETERS

In most of the examples so far, the “natural parameters”  $\eta$  are not the parameters we typically use for the distribution families.

**Example #1:** Bernoulli distribution  $\text{Bern}(p)$  for  $p = \frac{e^\eta}{1+e^\eta}$ .

$$X \sim \text{Bern}(p) \implies \mathbb{E}[X] = p.$$

usually use the form  
of  $\mathbb{E}(X)$

**Example #2:** Poisson distribution  $\text{Pois}(\lambda)$  for  $\lambda = e^\eta$ .

$$X \sim \text{Pois}(\lambda) \implies \mathbb{E}[X] = \lambda.$$

**Example #4:** Gaussian distribution  $N(\mu, \sigma^2)$  for  $\mu = -\frac{\eta_1}{2\eta_2}$  and  $\sigma^2 = -\frac{1}{2\eta_2}$ .

$$X \sim N(\mu, \sigma^2) \implies \mathbb{E}[X] = \mu, \quad \mathbb{E}[X^2] = \sigma^2 + \mu^2.$$

Is there a more general relationship between these expectations and the natural parameters?



# EXPECTATIONS

Let  $\{P_{\boldsymbol{\eta}} : \boldsymbol{\eta} \in \mathcal{N}\}$  be the exponential family corresponding to  $\boldsymbol{T}$  and  $\pi$ .

**Fact:** For any  $\boldsymbol{\eta} \in \mathcal{N}$  (except on the boundary), all derivatives of the log partition function  $G$  exist at  $\boldsymbol{\eta}$ , and

$$X \sim P_{\boldsymbol{\eta}} \implies \nabla G(\boldsymbol{\eta}) = \mathbb{E}[\boldsymbol{T}(X)].$$

# EXPECTATIONS

Let  $\{P_{\boldsymbol{\eta}} : \boldsymbol{\eta} \in \mathcal{N}\}$  be the exponential family corresponding to  $\boldsymbol{T}$  and  $\pi$ .

**Fact:** For any  $\boldsymbol{\eta} \in \mathcal{N}$  (except on the boundary), all derivatives of the log partition function  $G$  exist at  $\boldsymbol{\eta}$ , and

$$X \sim P_{\boldsymbol{\eta}} \implies \nabla G(\boldsymbol{\eta}) = \mathbb{E}[\boldsymbol{T}(X)].$$

**Proof in discrete case:**

Recall that  $G(\boldsymbol{\eta}) = \ln\left(\sum_{x \in \mathcal{X}} \exp\{\langle \boldsymbol{\eta}, \boldsymbol{T}(x) \rangle\} \cdot \pi(x)\right)$ , so its gradient is

$$\nabla G(\boldsymbol{\eta}) = \frac{\sum_{x \in \mathcal{X}} \exp\{\langle \boldsymbol{\eta}, \boldsymbol{T}(x) \rangle\} \cdot \pi(x) \cdot \boldsymbol{T}(x)}{\sum_{x \in \mathcal{X}} \exp\{\langle \boldsymbol{\eta}, \boldsymbol{T}(x) \rangle\} \cdot \pi(x)}$$

# EXPECTATIONS

Let  $\{P_{\boldsymbol{\eta}} : \boldsymbol{\eta} \in \mathcal{N}\}$  be the exponential family corresponding to  $\boldsymbol{T}$  and  $\pi$ .

**Fact:** For any  $\boldsymbol{\eta} \in \mathcal{N}$  (except on the boundary), all derivatives of the log partition function  $G$  exist at  $\boldsymbol{\eta}$ , and

$$X \sim P_{\boldsymbol{\eta}} \implies \nabla G(\boldsymbol{\eta}) = \mathbb{E}[\boldsymbol{T}(X)].$$

**Proof in discrete case:**

Recall that  $G(\boldsymbol{\eta}) = \ln\left(\sum_{x \in \mathcal{X}} \exp\{\langle \boldsymbol{\eta}, \boldsymbol{T}(x) \rangle\} \cdot \pi(x)\right)$ , so its gradient is

$$\nabla G(\boldsymbol{\eta}) = \frac{\sum_{x \in \mathcal{X}} \exp\{\langle \boldsymbol{\eta}, \boldsymbol{T}(x) \rangle\} \cdot \pi(x) \cdot \boldsymbol{T}(x)}{\sum_{x \in \mathcal{X}} \exp\{\langle \boldsymbol{\eta}, \boldsymbol{T}(x) \rangle\} \cdot \pi(x)} = \sum_{x \in \mathcal{X}} P_{\boldsymbol{\eta}}(x) \cdot \boldsymbol{T}(x) = \mathbb{E}[\boldsymbol{T}(X)].$$

□

# EXPECTATIONS

Let  $\{P_{\boldsymbol{\eta}} : \boldsymbol{\eta} \in \mathcal{N}\}$  be the exponential family corresponding to  $\boldsymbol{T}$  and  $\pi$ .

**Fact:** For any  $\boldsymbol{\eta} \in \mathcal{N}$  (except on the boundary), all derivatives of the log partition function  $G$  exist at  $\boldsymbol{\eta}$ , and

$$X \sim P_{\boldsymbol{\eta}} \implies \nabla G(\boldsymbol{\eta}) = \mathbb{E}[\boldsymbol{T}(X)].$$

$\lambda$  is the target for gradient

$X$  is also a random variable

**Proof in discrete case:**

Recall that  $G(\boldsymbol{\eta}) = \ln\left(\sum_{x \in \mathcal{X}} \exp\{\langle \boldsymbol{\eta}, \boldsymbol{T}(x) \rangle\} \cdot \pi(x)\right)$ , so its gradient is

$$\nabla G(\boldsymbol{\eta}) = \frac{\sum_{x \in \mathcal{X}} \exp\{\langle \boldsymbol{\eta}, \boldsymbol{T}(x) \rangle\} \cdot \pi(x) \cdot \boldsymbol{T}(x)}{\sum_{x \in \mathcal{X}} \exp\{\langle \boldsymbol{\eta}, \boldsymbol{T}(x) \rangle\} \cdot \pi(x)} = \sum_{x \in \mathcal{X}} P_{\boldsymbol{\eta}}(x) \cdot \boldsymbol{T}(x) = \mathbb{E}[\boldsymbol{T}(X)].$$

the gradient gives the



In fact,  $\nabla^2 G(\boldsymbol{\eta}) = \text{cov}(\boldsymbol{T}(X))$ , and higher-order derivatives correspond to higher-order moments of  $\boldsymbol{T}(X)$ .

†

# EXPECTATIONS: EXAMPLES

**Example #1:** Bernoulli distribution  $\text{Bern}(p)$  for  $p = \frac{e^\eta}{1+e^\eta}$ .

$$G(\eta) = \ln(1 + \exp(\eta))$$

# EXPECTATIONS: EXAMPLES

**Example #1:** Bernoulli distribution  $\text{Bern}(p)$  for  $p = \frac{e^\eta}{1+e^\eta}$ .

$$G(\eta) = \ln(1 + \exp(\eta)), \quad \frac{\partial G(\eta)}{\partial \eta} = \frac{e^\eta}{1 + e^\eta} = p = \mathbb{E}[X].$$

# EXPECTATIONS: EXAMPLES

**Example #1:** Bernoulli distribution  $\text{Bern}(p)$  for  $p = \frac{e^\eta}{1+e^\eta}$ .

$$G(\eta) = \ln(1 + \exp(\eta)), \quad \frac{\partial G(\eta)}{\partial \eta} = \frac{e^\eta}{1 + e^\eta} = p = \mathbb{E}[X].$$

**Example #2:** Poisson distribution  $\text{Pois}(\lambda)$  for  $\lambda = e^\eta$ .

$$G(\eta) = \exp(\eta)$$

# EXPECTATIONS: EXAMPLES

**Example #1:** Bernoulli distribution  $\text{Bern}(p)$  for  $p = \frac{e^\eta}{1+e^\eta}$ .

$$G(\eta) = \ln(1 + \exp(\eta)), \quad \frac{\partial G(\eta)}{\partial \eta} = \frac{e^\eta}{1 + e^\eta} = p = \mathbb{E}[X].$$

**Example #2:** Poisson distribution  $\text{Pois}(\lambda)$  for  $\lambda = e^\eta$ .

$$G(\eta) = \exp(\eta), \quad \frac{\partial G(\eta)}{\partial \eta} = e^\eta = \lambda = \mathbb{E}[X].$$



# EXPECTATIONS: EXAMPLES

**Example #1:** Bernoulli distribution  $\text{Bern}(p)$  for  $p = \frac{e^\eta}{1+e^\eta}$ .

$$G(\eta) = \ln(1 + \exp(\eta)), \quad \frac{\partial G(\eta)}{\partial \eta} = \frac{e^\eta}{1 + e^\eta} = p = \mathbb{E}[X].$$

**Example #2:** Poisson distribution  $\text{Pois}(\lambda)$  for  $\lambda = e^\eta$ .

$$G(\eta) = \exp(\eta), \quad \frac{\partial G(\eta)}{\partial \eta} = e^\eta = \lambda = \mathbb{E}[X].$$

**Example #4:** Gaussian distribution  $N(\mu, \sigma^2)$  for  $\mu = -\frac{\eta_1}{2\eta_2}$  and  $\sigma^2 = -\frac{1}{2\eta_2}$ .

$$G(\eta_1, \eta_2) = -\ln \sqrt{-2\eta_2} - \frac{\eta_1^2}{4\eta_2}$$

# EXPECTATIONS: EXAMPLES

**Example #1:** Bernoulli distribution  $\text{Bern}(p)$  for  $p = \frac{e^\eta}{1+e^\eta}$ .

$$G(\eta) = \ln(1 + \exp(\eta)), \quad \frac{\partial G(\eta)}{\partial \eta} = \frac{e^\eta}{1 + e^\eta} = p = \mathbb{E}[X].$$

**Example #2:** Poisson distribution  $\text{Pois}(\lambda)$  for  $\lambda = e^\eta$ .

$$G(\eta) = \exp(\eta), \quad \frac{\partial G(\eta)}{\partial \eta} = e^\eta = \lambda = \mathbb{E}[X].$$

**Example #4:** Gaussian distribution  $N(\mu, \sigma^2)$  for  $\mu = -\frac{\eta_1}{2\eta_2}$  and  $\sigma^2 = -\frac{1}{2\eta_2}$ .

$$G(\eta_1, \eta_2) = -\ln \sqrt{-2\eta_2} - \frac{\eta_1^2}{4\eta_2}, \quad \frac{\partial G(\eta_1, \eta_2)}{\partial \eta_1} = -\frac{\eta_1}{2\eta_2} = \mu = \mathbb{E}[X]$$

# EXPECTATIONS: EXAMPLES

**Example #1:** Bernoulli distribution  $\text{Bern}(p)$  for  $p = \frac{e^\eta}{1+e^\eta}$ .

$$G(\eta) = \ln(1 + \exp(\eta)), \quad \frac{\partial G(\eta)}{\partial \eta} = \frac{e^\eta}{1 + e^\eta} = p = \mathbb{E}[X].$$

**Example #2:** Poisson distribution  $\text{Pois}(\lambda)$  for  $\lambda = e^\eta$ .

$$G(\eta) = \exp(\eta), \quad \frac{\partial G(\eta)}{\partial \eta} = e^\eta = \lambda = \mathbb{E}[X].$$

**Example #4:** Gaussian distribution  $N(\mu, \sigma^2)$  for  $\mu = -\frac{\eta_1}{2\eta_2}$  and  $\sigma^2 = -\frac{1}{2\eta_2}$ .

$$G(\eta_1, \eta_2) = -\ln \sqrt{-2\eta_2} - \frac{\eta_1^2}{4\eta_2}, \quad \frac{\partial G(\eta_1, \eta_2)}{\partial \eta_1} = -\frac{\eta_1}{2\eta_2} = \mu = \mathbb{E}[X],$$

$$\frac{\partial G(\eta_1, \eta_2)}{\partial \eta_2} = -\frac{1}{2\eta_2} + \frac{\eta_1^2}{4\eta_2^2} = \sigma^2 + \mu^2 = \mathbb{E}[X^2].$$

# EXPECTATIONS: EXAMPLES

**Example #1:** Bernoulli distribution  $\text{Bern}(p)$  for  $p = \frac{e^\eta}{1+e^\eta}$ .

$$G(\eta) = \ln(1 + \exp(\eta)), \quad \frac{\partial G(\eta)}{\partial \eta} = \frac{e^\eta}{1 + e^\eta} = p = \mathbb{E}[X].$$

**Example #2:** Poisson distribution  $\text{Pois}(\lambda)$  for  $\lambda = e^\eta$ .

$$G(\eta) = \exp(\eta), \quad \frac{\partial G(\eta)}{\partial \eta} = e^\eta = \lambda = \mathbb{E}[X].$$

**Example #4:** Gaussian distribution  $N(\mu, \sigma^2)$  for  $\mu = -\frac{\eta_1}{2\eta_2}$  and  $\sigma^2 = -\frac{1}{2\eta_2}$ .

$$G(\eta_1, \eta_2) = -\ln \sqrt{-2\eta_2} - \frac{\eta_1^2}{4\eta_2}, \quad \frac{\partial G(\eta_1, \eta_2)}{\partial \eta_1} = -\frac{\eta_1}{2\eta_2} = \mu = \mathbb{E}[X],$$

$$\frac{\partial G(\eta_1, \eta_2)}{\partial \eta_2} = -\frac{1}{2\eta_2} + \frac{\eta_1^2}{4\eta_2^2} = \sigma^2 + \mu^2 = \mathbb{E}[X^2].$$

**In each of these cases, the distributions are also parameterized by these expectations. Is this true for all exponential families?**

# EXPECTATIONS: EXAMPLES

we don't care about  $\lambda$ , we want to get  $p$

**Example #1:** Bernoulli distribution  $\text{Bern}(p)$  for  $p = \frac{e^\eta}{1+e^\eta}$ .

$$G(\eta) = \ln(1 + \exp(\eta)), \quad \frac{\partial G(\eta)}{\partial \eta} = \frac{e^\eta}{1 + e^\eta} = p = \mathbb{E}[X].$$

**Example #2:** Poisson distribution  $\text{Pois}(\lambda)$  for  $\lambda = e^\eta$ .

$$G(\eta) = \exp(\eta), \quad \frac{\partial G(\eta)}{\partial \eta} = e^\eta = \lambda = \mathbb{E}[X].$$

**Example #4:** Gaussian distribution  $N(\mu, \sigma^2)$  for  $\mu = -\frac{\eta_1}{2\eta_2}$  and  $\sigma^2 = -\frac{1}{2\eta_2}$ .

$$G(\eta_1, \eta_2) = -\ln \sqrt{-2\eta_2} - \frac{\eta_1^2}{4\eta_2}, \quad \frac{\partial G(\eta_1, \eta_2)}{\partial \eta_1} = -\frac{\eta_1}{2\eta_2} = \mu = \mathbb{E}[X],$$

$$\frac{\partial G(\eta_1, \eta_2)}{\partial \eta_2} = -\frac{1}{2\eta_2} + \frac{\eta_1^2}{4\eta_2^2} = \sigma^2 + \mu^2 = \mathbb{E}[X^2].$$

In each of these cases, the distributions are also parameterized by these expectations. Is this true for all exponential families? **Yes!**

# EXPECTATION PARAMETERS

$g := \nabla G$  is an invertible map between natural parameters and expectations:

$$\mu = g(\eta) \quad \Longleftrightarrow \quad \eta = g^{-1}(\mu).$$

# EXPECTATION PARAMETERS

$g := \nabla G$  is an invertible map between natural parameters and expectations:

$$\mu = g(\eta) \iff \eta = g^{-1}(\mu).$$

**Fact:**  $G$  is a strictly convex function.<sup>‡</sup> (Proof by generalized Cauchy-Schwarz ineq.)

<sup>‡</sup>Strict convexity, as opposed to just convexity, actually requires that the exponential family be *minimal*: feature functions should be linearly independent.

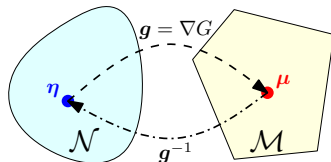
# EXPECTATION PARAMETERS

$g := \nabla G$  is an invertible map between natural parameters and expectations:

$$\mu = g(\eta) \iff \eta = g^{-1}(\mu).$$

**Fact:**  $G$  is a strictly convex function.<sup>‡</sup> (Proof by generalized Cauchy-Schwarz ineq.)

**Corollary:**  $g$  is a 1-to-1 map.



$$\mathcal{M} := \{g(\eta) : \eta \in \mathcal{N}\} = \text{expectation parameter space.}$$

<sup>‡</sup>Strict convexity, as opposed to just convexity, actually requires that the exponential family be *minimal*: feature functions should be linearly independent.



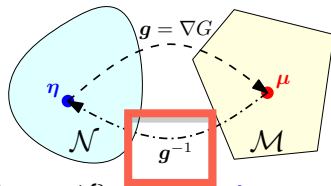
# EXPECTATION PARAMETERS

$g := \nabla G$  is an invertible map between natural parameters and expectations:

$$\mu = g(\eta) \iff \eta = g^{-1}(\mu).$$

**Fact:**  $G$  is a strictly convex function.<sup>‡</sup> (Proof by generalized Cauchy-Schwarz ineq.)

**Corollary:**  $g$  is a 1-to-1 map.



$\mathcal{M} := \{g(\eta) : \eta \in \mathcal{N}\} =$  **expectation parameter space.**

Therefore, exponential families can equivalently be parameterized by the **expectation (mean) parameters**  $g(\eta) = \mathbb{E}_{X \sim P_\eta}[T(X)]$ .

<sup>‡</sup>Strict convexity, as opposed to just convexity, actually requires that the exponential family be *minimal*: feature functions should be linearly independent.

# SPECIAL CASE: GAUSSIANS WITH UNIT COVARIANCE

**Example:**

$$\mathcal{X} = \mathbb{R}^d, \quad \mathbf{T}(\mathbf{x}) = \mathbf{x}, \quad \pi(\mathbf{x}) = \frac{1}{(2\pi)^{d/2}} e^{-\|\mathbf{x}\|_2^2/2}.$$

$$G(\boldsymbol{\eta}) = \ln \left( \int_{\mathbb{R}^d} \exp\{\langle \boldsymbol{\eta}, \mathbf{x} \rangle\} \cdot \frac{1}{(2\pi)^{d/2}} e^{-\|\mathbf{x}\|_2^2/2} d\mathbf{x} \right) = \frac{1}{2} \|\boldsymbol{\eta}\|_2^2.$$

# SPECIAL CASE: GAUSSIANS WITH UNIT COVARIANCE

**Example:**

$$\mathcal{X} = \mathbb{R}^d, \quad \mathbf{T}(\mathbf{x}) = \mathbf{x}, \quad \pi(\mathbf{x}) = \frac{1}{(2\pi)^{d/2}} e^{-\|\mathbf{x}\|_2^2/2}.$$

$$G(\boldsymbol{\eta}) = \ln \left( \int_{\mathbb{R}^d} \exp\{\langle \boldsymbol{\eta}, \mathbf{x} \rangle\} \cdot \frac{1}{(2\pi)^{d/2}} e^{-\|\mathbf{x}\|_2^2/2} d\mathbf{x} \right) = \frac{1}{2} \|\boldsymbol{\eta}\|_2^2.$$

In this case,

$$\nabla G(\boldsymbol{\eta}) = \boldsymbol{\eta}.$$

**Natural parameters and expectation parameters coincide ( $\mathcal{N} = \mathcal{M}$ ).**

# PARAMETER ESTIMATION

Let  $\mathcal{P} = \{P_{\boldsymbol{\eta}} : \boldsymbol{\eta} \in \mathcal{N}\}$  be the exponential family corresponding to  $\boldsymbol{T}$  and  $\pi$ .

# PARAMETER ESTIMATION

Let  $\mathcal{P} = \{P_{\boldsymbol{\eta}} : \boldsymbol{\eta} \in \mathcal{N}\}$  be the exponential family corresponding to  $\boldsymbol{T}$  and  $\pi$ .

**Estimation:** Given data  $x_1, x_2, \dots, x_n \in \mathcal{X}$ , how do we pick  $P_{\boldsymbol{\eta}} \in \mathcal{P}$ ?

# PARAMETER ESTIMATION

Let  $\mathcal{P} = \{P_{\boldsymbol{\eta}} : \boldsymbol{\eta} \in \mathcal{N}\}$  be the exponential family corresponding to  $\boldsymbol{T}$  and  $\pi$ .

**Estimation:** Given data  $x_1, x_2, \dots, x_n \in \mathcal{X}$ , how do we pick  $P_{\boldsymbol{\eta}} \in \mathcal{P}$ ?

We'll look at two approaches to estimation:

1. Maximum entropy principle.
2. Maximum likelihood principle.

# PARAMETER ESTIMATION

Let  $\mathcal{P} = \{P_{\eta} : \eta \in \mathcal{N}\}$  be the exponential family corresponding to  $\mathbf{T}$  and  $\pi$ .

**Estimation:** Given data  $x_1, x_2, \dots, x_n \in \mathcal{X}$ , how do we pick  $P_{\eta} \in \mathcal{P}$ ?

We'll look at two approaches to estimation:

1. Maximum entropy principle.
2. Maximum likelihood principle.

In both cases, we'll use the fact that

$$\mathbb{E}_{X \sim P_{\eta}}[\mathbf{T}(X)] = \mathbf{g}(\eta) \quad \forall \eta \in \mathcal{N}$$

where  $\mathbf{g} : \mathcal{N} \rightarrow \mathbb{R}^k$  is gradient of log partition function  $G$ .



# MAXIMUM ENTROPY PRINCIPLE

**Maximum entropy principle** (for discrete  $\mathcal{X}$ ): given data  $x_1, x_2, \dots, x_n \in \mathcal{X}$ , pick the distribution that solves the optimization problem

$$\begin{aligned} \min_{P \in \Delta(\mathcal{X})} \quad & \text{RE}(P \| \pi) \\ \text{s.t.} \quad & \mathbb{E}_{X \sim P}[\mathbf{T}(X)] = \hat{\mathbb{E}}[\mathbf{T}(X)] \left( = \frac{1}{n} \sum_{i=1}^n \mathbf{T}(x_i) \right). \end{aligned}$$

# MAXIMUM ENTROPY PRINCIPLE

**Maximum entropy principle** (for discrete  $\mathcal{X}$ ): given data  $x_1, x_2, \dots, x_n \in \mathcal{X}$ , pick the distribution that solves the optimization problem

$$\begin{aligned} \min_{P \in \Delta(\mathcal{X})} \quad & \text{RE}(P \| \pi) \\ \text{s.t.} \quad & \mathbb{E}_{X \sim P}[\mathbf{T}(X)] = \hat{\mathbb{E}}[\mathbf{T}(X)] \left( = \frac{1}{n} \sum_{i=1}^n \mathbf{T}(x_i) \right). \end{aligned}$$

The solution  $P_\star$  has the form  $P_\star = P_{\eta_\star} \in \mathcal{P}$ , and therefore

$$\mathbf{g}(\eta_\star) = \mathbb{E}_{X \sim P_{\eta_\star}}[\mathbf{T}(X)] = \hat{\mathbb{E}}[\mathbf{T}(X)]$$

and

$$\eta_\star = \mathbf{g}^{-1}(\hat{\mathbb{E}}[\mathbf{T}(X)]).$$

# MAXIMUM ENTROPY PRINCIPLE

**Maximum entropy principle** (for discrete  $\mathcal{X}$ ): given data  $x_1, x_2, \dots, x_n \in \mathcal{X}$ , pick the distribution that solves the optimization problem

$$\begin{aligned} \min_{P \in \Delta(\mathcal{X})} \quad & \text{RE}(P \| \pi) \\ \text{s.t.} \quad & \mathbb{E}_{X \sim P}[\mathbf{T}(X)] = \hat{\mathbb{E}}[\mathbf{T}(X)] \left( = \frac{1}{n} \sum_{i=1}^n \mathbf{T}(x_i) \right). \end{aligned}$$

The solution  $P_\star$  has the form  $P_\star = P_{\eta_\star} \in \mathcal{P}$ , and therefore

it's interesting!

$$\mathbf{g}(\eta_\star) = \mathbb{E}_{X \sim P_{\eta_\star}}[\mathbf{T}(X)] = \hat{\mathbb{E}}[\mathbf{T}(X)]$$

and

$$\eta_\star = \mathbf{g}^{-1}(\hat{\mathbb{E}}[\mathbf{T}(X)]).$$

Pick the distribution  $P_\eta \in \mathcal{P}$  whose **expectation parameters**  $\mathbf{g}(\eta) = \mu$  are **equal to the corresponding empirical expectations**.

# MAXIMUM LIKELIHOOD PRINCIPLE

**Maximum likelihood principle:** given data  $x_1, x_2, \dots, x_n \in \mathcal{X}$ , pick  $\boldsymbol{\eta} \in \mathcal{N}$  that maximizes the log-likelihood function  $\mathcal{L}(\boldsymbol{\eta}) := \sum_{i=1}^n \ln P_{\boldsymbol{\eta}}(x_i)$ .

# MAXIMUM LIKELIHOOD PRINCIPLE

**Maximum likelihood principle:** given data  $x_1, x_2, \dots, x_n \in \mathcal{X}$ , pick  $\boldsymbol{\eta} \in \mathcal{N}$  that maximizes the log-likelihood function  $\mathcal{L}(\boldsymbol{\eta}) := \sum_{i=1}^n \ln P_{\boldsymbol{\eta}}(x_i)$ .

$$\mathcal{L}(\boldsymbol{\eta}) = \sum_{i=1}^n \left( \langle \boldsymbol{\eta}, \boldsymbol{T}(x_i) \rangle - G(\boldsymbol{\eta}) + \ln \pi(x_i) \right)$$

which is a *concave* function of  $\boldsymbol{\eta}$ .

# MAXIMUM LIKELIHOOD PRINCIPLE

**Maximum likelihood principle:** given data  $x_1, x_2, \dots, x_n \in \mathcal{X}$ , pick  $\eta \in \mathcal{N}$  that maximizes the log-likelihood function  $\mathcal{L}(\eta) := \sum_{i=1}^n \ln P_{\eta}(x_i)$ .

$$\mathcal{L}(\eta) = \sum_{i=1}^n \left( \langle \eta, T(x_i) \rangle - G(\eta) + \ln \pi(x_i) \right)$$

which is a *concave* function of  $\eta$ .

Pick  $\eta_{\text{ML}} \in \mathcal{N}$  so that gradient of  $\mathcal{L}$  at  $\eta_{\text{ML}}$  is zero:

$$\nabla \mathcal{L}(\eta_{\text{ML}}) = \sum_{i=1}^n \left( T(x_i) - \nabla G(\eta_{\text{ML}}) \right) = \mathbf{0}$$

# MAXIMUM LIKELIHOOD PRINCIPLE

**Maximum likelihood principle:** given data  $x_1, x_2, \dots, x_n \in \mathcal{X}$ , pick  $\boldsymbol{\eta} \in \mathcal{N}$  that maximizes the log-likelihood function  $\mathcal{L}(\boldsymbol{\eta}) := \sum_{i=1}^n \ln P_{\boldsymbol{\eta}}(x_i)$ .

$$\mathcal{L}(\boldsymbol{\eta}) = \sum_{i=1}^n \left( \langle \boldsymbol{\eta}, \mathbf{T}(x_i) \rangle - G(\boldsymbol{\eta}) + \ln \pi(x_i) \right)$$

which is a *concave* function of  $\boldsymbol{\eta}$ .

Pick  $\boldsymbol{\eta}_{\text{ML}} \in \mathcal{N}$  so that gradient of  $\mathcal{L}$  at  $\boldsymbol{\eta}_{\text{ML}}$  is zero:

$$\nabla \mathcal{L}(\boldsymbol{\eta}_{\text{ML}}) = \sum_{i=1}^n \left( \mathbf{T}(x_i) - \nabla G(\boldsymbol{\eta}_{\text{ML}}) \right) = \mathbf{0},$$

i.e.,

$$\nabla G(\boldsymbol{\eta}_{\text{ML}}) = \frac{1}{n} \sum_{i=1}^n \mathbf{T}(x_i) = \widehat{\mathbb{E}}[\mathbf{T}(X)].$$

# MAXIMUM LIKELIHOOD PRINCIPLE

**Maximum likelihood principle:** given data  $x_1, x_2, \dots, x_n \in \mathcal{X}$ , pick  $\boldsymbol{\eta} \in \mathcal{N}$  that maximizes the log-likelihood function  $\mathcal{L}(\boldsymbol{\eta}) := \sum_{i=1}^n \ln P_{\boldsymbol{\eta}}(x_i)$ .

$$\mathcal{L}(\boldsymbol{\eta}) = \sum_{i=1}^n \left( \langle \boldsymbol{\eta}, \mathbf{T}(x_i) \rangle - G(\boldsymbol{\eta}) + \ln \pi(x_i) \right)$$

which is a *concave* function of  $\boldsymbol{\eta}$ .

Pick  $\boldsymbol{\eta}_{\text{ML}} \in \mathcal{N}$  so that gradient of  $\mathcal{L}$  at  $\boldsymbol{\eta}_{\text{ML}}$  is zero:

$$\nabla \mathcal{L}(\boldsymbol{\eta}_{\text{ML}}) = \sum_{i=1}^n \left( \mathbf{T}(x_i) - \nabla G(\boldsymbol{\eta}_{\text{ML}}) \right) = \mathbf{0},$$

i.e.,

$$\nabla G(\boldsymbol{\eta}_{\text{ML}}) = \frac{1}{n} \sum_{i=1}^n \mathbf{T}(x_i) = \widehat{\mathbb{E}}[\mathbf{T}(X)].$$

For exponential families, the maximum likelihood principle is the  
**same as the maximum entropy principle!**

use maximum  
likelihood  
method.



# EMPIRICAL EXPECTATIONS

Maximum entropy / maximum likelihood for exponential families is reduced to **computing empirical expectations of feature functions**

$$\mathbf{b} := \hat{\mathbb{E}}[\mathbf{T}(X)] = \frac{1}{n} \sum_{j=1}^n \mathbf{T}(x_j),$$

and then applying the inverse of  $\nabla G$  to obtain the natural parameter estimate.

# EMPIRICAL EXPECTATIONS

Maximum entropy / maximum likelihood for exponential families is reduced to **computing empirical expectations of feature functions**

$$\mathbf{b} := \hat{\mathbb{E}}[\mathbf{T}(X)] = \frac{1}{n} \sum_{j=1}^n \mathbf{T}(x_j),$$

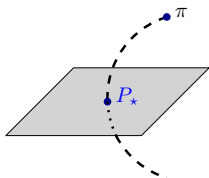
reduce to  
compute feature  
function

and then applying the inverse of  $\nabla G$  to obtain the natural parameter estimate.

Unfortunately, except for some simple exponential families (like the well-known distribution families), computing the inverse of  $\nabla G$  can be difficult.

# ITERATIVE PROJECTION ALGORITHM

## Csiszar's iterative projection algorithm

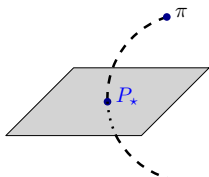


Let  $H_i := \{P : \mathbb{E}_{X \sim P}[T_i(X)] = b_i\}$  for each  $i \in [k]$ , so

$$H := \{P : \mathbb{E}_{X \sim P}[\mathbf{T}(X)] = \mathbf{b}\} = \bigcap_{i=1}^k H_i.$$

# ITERATIVE PROJECTION ALGORITHM

## Csiszar's iterative projection algorithm



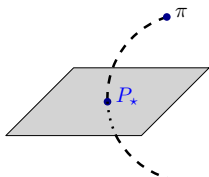
Let  $H_i := \{P : \mathbb{E}_{X \sim P}[T_i(X)] = b_i\}$  for each  $i \in [k]$ , so

$$H := \{P : \mathbb{E}_{X \sim P}[\mathbf{T}(X)] = \mathbf{b}\} = \bigcap_{i=1}^k H_i.$$

**Goal:** compute **entropy projection**  $P_*$  of  $\pi$  onto  $H$ .

# ITERATIVE PROJECTION ALGORITHM

## Csiszar's iterative projection algorithm



Let  $H_i := \{P : \mathbb{E}_{X \sim P}[T_i(X)] = b_i\}$  for each  $i \in [k]$ , so

$$H := \{P : \mathbb{E}_{X \sim P}[\mathbf{T}(X)] = \mathbf{b}\} = \bigcap_{i=1}^k H_i.$$

**Goal:** compute **entropy projection**  $P_*$  of  $\pi$  onto  $H$ .

Start with  $P^{(0)} := \pi$ .

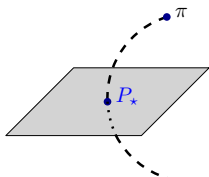
For  $t = 1, 2, \dots$ :

$$P^{(t)} := \arg \min_{P \in \Delta(\mathcal{X})} \text{RE}(P \| P^{(t-1)}) \quad \text{subject to } P \in H_{i_t}.$$

Here,  $(i_1, i_2, \dots)$  is a sequence like  $(1, 2, \dots, k, 1, 2, \dots, k, \dots)$ .

# ITERATIVE PROJECTION ALGORITHM

## Csiszar's iterative projection algorithm



Let  $H_i := \{P : \mathbb{E}_{X \sim P}[T_i(X)] = b_i\}$  for each  $i \in [k]$ , so

$$H := \{P : \mathbb{E}_{X \sim P}[\mathbf{T}(X)] = \mathbf{b}\} = \bigcap_{i=1}^k H_i.$$

**Goal:** compute **entropy projection**  $P_\star$  of  $\pi$  onto  $H$ .

Start with  $P^{(0)} := \pi$ .

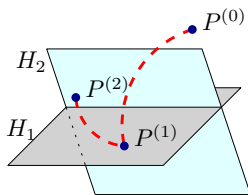
For  $t = 1, 2, \dots$ :

$$P^{(t)} := \arg \min_{P \in \Delta(\mathcal{X})} \text{RE}(P \| P^{(t-1)}) \quad \text{subject to } P \in H_{i_t}.$$

Here,  $(i_1, i_2, \dots)$  is a sequence like  $(1, 2, \dots, k, 1, 2, \dots, k, \dots)$ .

**This converges to the entropy projection  $P_\star$  of  $\pi$  onto  $H$ .**

# ITERATIVE PROJECTION ALGORITHM

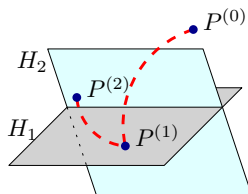


Start with  $P^{(0)} := \pi$ .

For  $t = 1, 2, \dots$ :

$$P^{(t)} := \arg \min_{P \in \Delta(\mathcal{X})} \text{RE}(P \| P^{(t-1)}) \quad \text{s.t.} \quad P \in H_{i_t}.$$

# ITERATIVE PROJECTION ALGORITHM



Start with  $P^{(0)} := \pi$ .

For  $t = 1, 2, \dots$ :

$$P^{(t)} := \arg \min_{P \in \Delta(\mathcal{X})} \text{RE}(P \| P^{(t-1)}) \quad \text{s.t.} \quad P \in H_{i_t}.$$

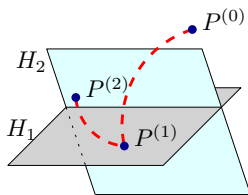
**Entropy projection of  $Q$  onto  $H_i$ :** Find  $\eta_i$  such that

$$G_i(\eta) := \ln \left( \sum_{x \in \mathcal{X}} \exp\{\eta T_i(x)\} \cdot Q(x) \right)$$

has derivative at  $\eta_i$  equal to  $b_i$ . **Can be solved using a line search.**



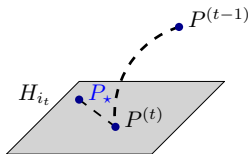
# ITERATIVE PROJECTION ALGORITHM



Start with  $P^{(0)} := \pi$ .

For  $t = 1, 2, \dots$ :

$$P^{(t)} := \arg \min_{P \in \Delta(\mathcal{X})} \text{RE}(P \| P^{(t-1)}) \quad \text{s.t.} \quad P \in H_{i_t}.$$



**Convergence:** since  $P^{(t)}$  is the entropy projection of  $P^{(t-1)}$  onto  $H_{i_t}$ , and  $P_\star \in H_{i_t}$ ,

$$\text{RE}(P_\star \| P^{(t)}) = \text{RE}(P_\star \| P^{(t-1)}) - \text{RE}(P^{(t)} \| P^{(t-1)}).$$

no negative

## EXAMPLE: SPECIES DISTRIBUTION MODELING

Phillips, Dudík, and Schapire (2004): model the geographic distribution of particular animal species.

# EXAMPLE: SPECIES DISTRIBUTION MODELING

Phillips, Dudík, and Schapire (2004): model the geographic distribution of particular animal species.

- ▶  $\mathcal{X}$  = set of geographic locations.

# EXAMPLE: SPECIES DISTRIBUTION MODELING

Phillips, Dudík, and Schapire (2004): model the geographic distribution of particular animal species.

- ▶  $\mathcal{X}$  = set of geographic locations.
- ▶ Sample  $x_1, x_2, \dots, x_n \in \mathcal{X}$ , corresponding to locations where a particular animal species was observed.

# EXAMPLE: SPECIES DISTRIBUTION MODELING

Phillips, Dudík, and Schapire (2004): model the geographic distribution of particular animal species.

- ▶  $\mathcal{X}$  = set of geographic locations.
- ▶ Sample  $x_1, x_2, \dots, x_n \in \mathcal{X}$ , corresponding to locations where a particular animal species was observed.
- ▶ Feature functions  $T_1, T_2, \dots$ , of the form:

# EXAMPLE: SPECIES DISTRIBUTION MODELING

Phillips, Dudík, and Schapire (2004): model the geographic distribution of particular animal species.

- ▶  $\mathcal{X}$  = set of geographic locations.
- ▶ Sample  $x_1, x_2, \dots, x_n \in \mathcal{X}$ , corresponding to locations where a particular animal species was observed.
- ▶ Feature functions  $T_1, T_2, \dots$ , of the form:
  - ▶ environmental features of geographic location, e.g.,

$T_i(x)$  = temperature at  $x$  in Celsius;

# EXAMPLE: SPECIES DISTRIBUTION MODELING

Phillips, Dudík, and Schapire (2004): model the geographic distribution of particular animal species.

- ▶  $\mathcal{X}$  = set of geographic locations.
- ▶ Sample  $x_1, x_2, \dots, x_n \in \mathcal{X}$ , corresponding to locations where a particular animal species was observed.
- ▶ Feature functions  $T_1, T_2, \dots$ , of the form:
  - ▶ environmental features of geographic location, e.g.,

$T_i(x)$  = temperature at  $x$  in Celsius;

- ▶ products of environmental features;

# EXAMPLE: SPECIES DISTRIBUTION MODELING

Phillips, Dudík, and Schapire (2004): model the geographic distribution of particular animal species.

- ▶  $\mathcal{X}$  = set of geographic locations.
- ▶ Sample  $x_1, x_2, \dots, x_n \in \mathcal{X}$ , corresponding to locations where a particular animal species was observed.
- ▶ Feature functions  $T_1, T_2, \dots$ , of the form:
  - ▶ environmental features of geographic location, e.g.,

$T_i(x)$  = temperature at  $x$  in Celsius;

- ▶ products of environmental features;
- ▶ indicator of whether environmental variable exceeds threshold, e.g.,

$T_i(x) = \mathbb{1}\{\text{temperature at } x \text{ exceeds } 20 \text{ degrees}\}.$



# EXAMPLE: SPECIES DISTRIBUTION MODELING

Phillips, Dudík, and Schapire (2004): model the geographic distribution of particular animal species.

- ▶  $\mathcal{X}$  = set of geographic locations.
- ▶ Sample  $x_1, x_2, \dots, x_n \in \mathcal{X}$ , corresponding to locations where a particular animal species was observed.
- ▶ Feature functions  $T_1, T_2, \dots$ , of the form:
  - ▶ environmental features of geographic location, e.g.,

$$T_i(x) = \text{temperature at } x \text{ in Celsius;}$$

- ▶ products of environmental features;
- ▶ indicator of whether environmental variable exceeds threshold, e.g.,

$$T_i(x) = \mathbb{1}\{\text{temperature at } x \text{ exceeds } 20 \text{ degrees}\}.$$

- ▶ Base distribution  $\pi = \text{uniform}$ .

# EXAMPLE: SPECIES DISTRIBUTION MODELING

Phillips, Dudík, and Schapire (2004): model the geographic distribution of particular animal species.

- ▶  $\mathcal{X}$  = set of geographic locations.
- ▶ Sample  $x_1, x_2, \dots, x_n \in \mathcal{X}$ , corresponding to locations where a particular animal species was observed.
- ▶ Feature functions  $T_1, T_2, \dots$ , of the form:
  - ▶ environmental features of geographic location, e.g.,

$T_i(x)$  = temperature at  $x$  in Celsius;

- ▶ products of environmental features;
- ▶ indicator of whether environmental variable exceeds threshold, e.g.,

$T_i(x) = \mathbb{1}\{\text{temperature at } x \text{ exceeds } 20 \text{ degrees}\}.$

- ▶ Base distribution  $\pi$  = uniform.

Lots of feature functions, but  $n$  (number of sightings of an animal) is small.

# EXAMPLE: SPECIES DISTRIBUTION MODELING

More details:

- ▶ Used a variant of maximum entropy where only require

$$\left| \mathbb{E}_{X \sim P}[T_i(X)] - \widehat{\mathbb{E}}[T_i(X)] \right| \leq \beta$$

for some *bound parameter*  $\beta > 0$  (determined using some other method based on the sample size).

# EXAMPLE: SPECIES DISTRIBUTION MODELING

More details:

- ▶ Used a variant of maximum entropy where only require

$$\left| \mathbb{E}_{X \sim P}[T_i(X)] - \widehat{\mathbb{E}}[T_i(X)] \right| \leq \beta$$

for some *bound parameter*  $\beta > 0$  (determined using some other method based on the sample size).

- ▶ Corresponds to *penalized maximum likelihood*: pick  $\boldsymbol{\eta} \in \mathcal{N}$  to maximize

$$\mathcal{L}(\boldsymbol{\eta}) - \beta \|\boldsymbol{\eta}\|_1$$

(similar to Lasso).

Better able to cope with many feature functions despite small sample size.

# EXAMPLE: SPECIES DISTRIBUTION MODELING

More details:

- ▶ Used a variant of maximum entropy where only require

$$\left| \mathbb{E}_{X \sim P}[T_i(X)] - \hat{\mathbb{E}}[T_i(X)] \right| \leq \beta$$

for some *bound parameter*  $\beta > 0$  (determined using some other method based on the sample size).

- ▶ Corresponds to *penalized maximum likelihood*: pick  $\boldsymbol{\eta} \in \mathcal{N}$  to maximize

$$\mathcal{L}(\boldsymbol{\eta}) - \beta \|\boldsymbol{\eta}\|_1$$

(similar to Lasso).

Better able to cope with many feature functions despite small sample size.

- ▶ Used resulting  $P_{\boldsymbol{\eta}}$  to rank locations by *habitability* for animal species.

# RECAP

- ▶ **Exponential families:** naturally arise as solution to maximum entropy optimization problem.

# RECAP

- ▶ **Exponential families:** naturally arise as solution to maximum entropy optimization problem.
- ▶ Many commonly used distribution families are exponential families.

# RECAP

- ▶ **Exponential families:** naturally arise as solution to maximum entropy optimization problem.
- ▶ Many commonly used distribution families are exponential families.
- ▶ Log partition function (and its gradient) connects the *natural parameters* and the *expectation parameters*.



- ▶ **Exponential families:** naturally arise as solution to maximum entropy optimization problem.
- ▶ Many commonly used distribution families are exponential families.
- ▶ Log partition function (and its gradient) connects the *natural parameters* and the *expectation parameters*.
- ▶ Maximum entropy and maximum likelihood estimation are the same for exponential families:

Pick the distribution whose expectation parameters are equal to the corresponding empirical expectations.

- ▶ **Exponential families:** naturally arise as solution to maximum entropy optimization problem.
- ▶ Many commonly used distribution families are exponential families.
- ▶ Log partition function (and its gradient) connects the *natural parameters* and the *expectation parameters*.
- ▶ Maximum entropy and maximum likelihood estimation are the same for exponential families:  
  
Pick the distribution whose expectation parameters are equal to the corresponding empirical expectations.
- ▶ Many variants of maximum entropy (see, e.g., Phillips, Dudík, and Schapire, 2004).

- ▶ **Exponential families:** naturally arise as solution to maximum entropy optimization problem.
- ▶ Many commonly used distribution families are exponential families.
- ▶ Log partition function (and its gradient) connects the *natural parameters* and the *expectation parameters*.
- ▶ Maximum entropy and maximum likelihood estimation are the same for exponential families:  

Pick the distribution whose expectation parameters are equal to the corresponding empirical expectations.
- ▶ Many variants of maximum entropy (see, e.g., Phillips, Dudík, and Schapire, 2004).
- ▶ Exponential families also important in / related to other subjects:
  - ▶ Generalized linear models (e.g., linear regression, logistic regression).
  - ▶ Bayesian inference (due to convenience of *conjugate priors*).
  - ▶ Convex optimization (via *Bregman divergences*).
  - ▶ ...