

# **Machine Learning 2**

## **Exponential Families - Basics**

Patrick Forré

# Exponential Families

- Exponential Families are sets of probability distributions that have a certain convenient structure relating the (data) arguments to the parameters.
- Most families of distributions we have seen so far can be written in exponential-family-form.
- Examples are Gaussian, Gamma, Poisson, Multinomial (with fixed number of trials), Bernoulli, Categorical, Dirichlet distributions, etc. etc.
- Counter examples are Student-t, Cauchy and Uniform distributions (with parameterized bounds).

# Exponential Families - Definition

- A set of probability distributions  $\{p(x | \theta) \mid \theta \in \Theta\}$  forms an Exponential Family, if it can be written in the form:

$$p(x | \theta) = h(x) \cdot \exp[\eta^\top \cdot T(x) - A(\eta)]$$

where:

- $h(x)$  is the **base measure**,
- $\eta = (\eta_1(\theta), \dots, \eta_s(\theta))^\top$  are the **natural parameters**,
- $T(x) = (T_1(x), \dots, T_s(x))^\top$  are the **sufficient statistics**,
- $A(\eta)$  is **log-partition function**, aka log-normalizer, aka cumulant function.

# Exponential Families - Remarks

- We usually aim at a representation with a **minimal number  $s$**  of natural parameters:  $\eta = (\eta_1(\theta), \dots, \eta_s(\theta))$ .
- $A(\eta) = \log \left( \int \exp(\eta^\top T(x)) \cdot h(x) dx \right)$  as **log-normalizer**, as  $Z(\eta) := \exp(A(\eta))$  is the normalizing constant.
- **Exponential Families** are NOT the same as **exponential distributions**, but:

# Example: Exponential distributions

- $p(x|\theta) = \theta \cdot e^{-\theta \cdot x}$  for  $x \geq 0$ , otherwise: 0. Parameters:  $\theta \in \mathbb{R}_{>0}$

$$= \underbrace{\mathbb{1}_{[0,\infty)}(x)}_{h(x)} \cdot \exp(\underbrace{-\theta \cdot x}_{T(x)} + \underbrace{\log(\theta)}_{A(\eta)})$$

- $h(x) = \mathbb{1}_{[0,\infty)}(x)$

$$p(x|y) = \mathbb{1}_{[0,\infty)}(x) \cdot \exp(y \cdot x - (-\log(-y)))$$

- $\eta(\theta) = -\theta$

$$\theta(\eta) = -y$$

- $T(x) = x$

- $A(\eta) = -\log(-y)$

# How to find an exponential family structure?

- Let the set of probability distributions  $\{p(x | \theta) \mid \theta \in \Theta\}$  be given.
- Be clear what are **parameters**  $\theta$  and what possible fixed **hyper-parameters**  $\alpha$ .
- Find **support**:  $S = \{x \mid p(x | \theta) > 0 \forall \theta\}$  and try to write for  $x \in S$ :  
$$\log p(x | \theta) = \eta_1(\theta) \cdot T_1(x) + \dots + \eta_s(\theta) \cdot T_s(x) + B(x) + C(\theta),$$
by **factoring out** everything and **separating**  $\theta$ -terms and  $x$ -terms, and  
**gathering** them into a minimal number of functions:  $\eta_i(\theta)$ ,  $T_i(x)$ ,  $B(x)$ ,  $C(\theta)$ .
- Put  $h(x) := \mathbb{1}_S(x) \cdot \exp(B(x))$ .
- **Express**  $\theta$  as a function  $\theta(\eta)$  of  $\eta$  and put  $A(\eta) := -C(\theta(\eta))$ .

# Counter-Example: Student-t distributions

$$\bullet p(x | \mu, \Sigma, \nu) = \frac{\Gamma(\frac{D}{2} + \frac{\nu}{2})}{\Gamma(\frac{\nu}{2})} \frac{1}{(\pi \nu)^{D/2} (\det \Sigma)^{\nu/2}} \left[ 1 + \nu^{-1} (x - \mu)^T \Sigma^{-1} (x - \mu) \right]^{-\frac{D}{2} - \frac{\nu}{2}}$$

NOT an exponential family

- can't be written as exponential family (not so easy to prove).
- looks like Gaussian, but with heavier tails, polynomial decay (instead of exponential decay)

# Machine Learning 2

## Exponential Families

- Example: Multinomial distributions (fixed K, M)

Patrick Forré

# Example: Multinomial distribution (fixed M, K)

$K+1 := \# \text{classes}$

$M := \# \text{trial}$

$$x = (x_0, \dots, x_K) \in \mathbb{N}^{K+1}$$

$$\sum_{k=0}^K x_k = M$$

$$\Delta^K := \left\{ \pi = (\pi_0, \dots, \pi_K) \in \mathbb{R}^{K+1} \mid \begin{array}{l} \pi_k \geq 0 \\ \sum_{k=0}^K \pi_k = 1 \end{array} \right\}$$

$$\begin{aligned} p(x|\pi) &= \binom{M}{x} \cdot \pi^x \\ &= \begin{cases} \frac{M!}{x_0! \cdots x_K!} & \text{if } \sum_{k=0}^K x_k = M \\ 0 & \text{o. otherwise} \end{cases} \cdot \pi_0^{x_0} \cdots \pi_K^{x_K} \end{aligned}$$

$$S = \left\{ x \in \mathbb{N}^{K+1} \mid \sum_{k=0}^K x_k = M \right\}$$

# Example: Multinomial distribution (fixed M, K)

$$\begin{aligned}
 \log p(x|\pi) &= \log \binom{M}{x} + \log \pi^x = \log \binom{M}{x} + \sum_{u=0}^K x_u \cdot \log \pi_u \\
 &= \log \binom{M}{x} + x_0 \log \pi_0 + \sum_{u=1}^K x_u \log \pi_u \\
 &= \log \binom{M}{x} + (M - \sum_{u=1}^K x_u) \cdot \log \left(1 - \sum_{u=1}^K \pi_u\right) + \sum_{u=1}^K x_u \log(\pi_u) \\
 &= \underbrace{\log \binom{M}{x}}_{B(x)} + \underbrace{M \cdot \log \left(1 - \sum_{u=1}^K \pi_u\right)}_{C(\pi)} + \underbrace{\sum_{u=0}^K x_u \cdot \log \left( \frac{\pi_u}{1 - \sum_{j=1}^u \pi_j} \right)}_{T_h(x)} \eta_u
 \end{aligned}$$

$x_0 = M - \sum_{u=1}^K x_u$   
 $\pi_0 = 1 - \sum_{u=1}^K \pi_u$

# Example: Multinomial distribution (fixed M, K)

$$y_u = \log\left(\frac{\pi_u}{1 - \sum_{j=1}^u \pi_j}\right)$$

$$u = 1, \dots, k$$

$$1 + \sum_{u=1}^k \exp(y_u) = 1 + \frac{\sum_{k=1}^u \pi_k}{1 - \sum_{j=1}^u \pi_j} = \frac{1 - \sum_{j=1}^u \pi_j + \sum_{u=1}^k \pi_u}{1 - \sum_{j=1}^u \pi_j} = \frac{\left(1 - \sum_{j=1}^u \pi_j\right)^{-1}}{1 - \sum_{j=1}^u \pi_j}$$

$$\Rightarrow \exp(y_u) = \frac{\pi_u}{1 - \sum_{j=1}^u \pi_j} = \pi_u \left(1 + \sum_{j=1}^u \exp(y_j)\right)$$

$$\Rightarrow \pi_u = \frac{\exp(y_u)}{1 + \sum_{j=1}^u \exp(y_j)} = \text{softmax}_u(y) = g_u(y)$$

# Example: Multinomial distribution (fixed M, K)

$$A(y) = -C(\pi(y)) = -C(g(y)) = -M \underbrace{\log \left( 1 - \sum_{j=1}^K \pi_j \right)}_{\left( 1 + \sum_{j=1}^K \exp(y_j) \right)^{-1}}$$
$$= +M \cdot \log \left( 1 + \sum_{j=1}^K \exp(y_j) \right)$$

$$p(x|y) = \binom{M}{x} \cdot \exp \left[ \sum_{k=1}^K y_k \cdot x_k - \left( M \cdot \log \left( 1 + \sum_{j=1}^K \exp(y_j) \right) \right) \right]$$

y exponential family with minimum number of natural parameters  $y_i$ : sufficient statistics  $T_i$ .

# Example: Binomial distribution (fixed N)

$$P(x|\pi) = \binom{N}{x} \pi^x (1-\pi)^{N-x}$$
$$= \text{multinomial}(M=N, K=1)$$

L, exponential family, same quantities as before

# Example: Bernoulli distributions

- $p(x | \pi) = \text{Multinomial } (k=1, N=1)$ 
  - ↳ exponential family  
same quantities as before
- $h(x) =$
- $\eta(\theta) =$        $\theta(\eta) =$
- $T(x) =$
- $A(\eta) =$

# **Machine Learning 2**

**Exponential Families**

**- Example: Gaussian distributions**

Patrick Forré

# Example: D-dim Gaussian distributions

- $p(x|\mu, \Sigma) = (2\pi)^{-D/2} \det(\Sigma)^{-1/2} \exp\left[-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right]$

$$\log p(x|\mu, \Sigma) = \underbrace{-\frac{D}{2} \log(2\pi)}_{B(x)} + \underbrace{\mu^T \Sigma^{-1} x}_{y_1^T T_1(x)} - \underbrace{\frac{1}{2} x^T \Sigma^{-1} x}_{? ?} - \underbrace{\frac{1}{2} \mu^T \Sigma^{-1} \mu}_{C(\mu, \Sigma)} - \frac{1}{2} \log \det(\Sigma)$$

$$\begin{aligned}
 -\frac{1}{2} x^T \Sigma^{-1} x &= \text{Tr}\left(-\frac{1}{2} x^T \Sigma^{-1} x\right) = \text{Tr}\left(\underbrace{-\frac{1}{2} \Sigma^{-1}}_{h_2} \underbrace{x x^T}_{T_2(x)}\right) \\
 &= \text{vec}\left(-\frac{1}{2} \Sigma^{-1}\right)^T \text{vec}(x x^T)
 \end{aligned}$$

# Example: D-dim Gaussian distributions

$$\bullet \cancel{p(x|\mu, \Sigma)} =$$

$$\Rightarrow \Sigma = -\frac{1}{2} y_2^{-1} \quad (\Sigma^{-1} = -2y_2)$$

$$y_1^T = \mu^T \Sigma^{-1}$$

$$y_2 = -\frac{1}{2} \Sigma^{-1}$$

$$\mu^T = y_1^T \Sigma = -\frac{1}{2} y_1^T y_2^{-1}$$

$$\mu = -\frac{1}{2} y_2^{-T} y_1$$

$$\begin{aligned} A(y) &= -C(\mu, \Sigma) = \frac{1}{2} \mu^T \Sigma^{-1} \mu + \frac{1}{2} \log \det(\Sigma) \\ &= \frac{1}{2} y_1^T \left( -\frac{1}{2} y_2^{-T} y_1 \right) + \frac{1}{2} \log \det \left( -\frac{1}{2} y_2^{-1} \right) \\ &= -\frac{1}{4} y_1^T y_2^{-T} y_1 - \frac{1}{2} \log \det(-2y_2) \end{aligned}$$

# Example: D-dim Gaussian distributions

- ~~$p(x | \mu, \Sigma)$~~  =

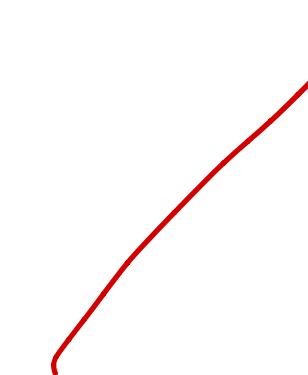
$$p(x|\gamma) = (2\pi)^{-D/2} \exp \left[ \gamma_1^T x + \text{Tr}(\gamma_2 \cdot (\mathbf{x}\mathbf{x}^T)) - A(\gamma) \right]$$

$$A(\gamma) = -\frac{1}{4} \gamma_1^T \gamma_2^{-T} \gamma_1 - \frac{1}{2} \log \det(-2\gamma_2)$$

# Example: 1-dim Gaussian distributions

- $p(x | \mu, \sigma^2) =$ 
  - $\rightsquigarrow D = 1$
  - $\rightsquigarrow \Sigma = \sigma^2$
  - $\rightsquigarrow$  same equations hold

- $h(x) =$
- $\eta(\mu, \sigma^2) =$
- $T(x) =$
- $A(\eta) =$



# **Machine Learning 2**

**Exponential Families**  
**- Conjugate Prior, Bayesian Update**

Patrick Forré

# Conjugate prior - Definition

- Let  $M = \{p(x|\theta) \mid \theta \in \Theta\}$  be a statistical model and  $D$  some data.
- A **conjugate prior** w.r.t.  $M$  is a family of priors  $F = \{p(\theta|\alpha) \mid \alpha \in A\}$  such that all possible posteriors  $p(\theta|\alpha, D)$  also lie in  $F$ .

- Remember:

$$\text{posterior} - p(\theta|\alpha, D) = \frac{\underset{\text{evidence}}{\cancel{p(D|\alpha)}} \cdot \underset{\substack{\text{likelihood} \\ p(D|\theta)}}{p(D|\theta)} \cdot \underset{\text{prior}}{p(\theta|\alpha)}}{\underset{\text{evidence}}{p(D|\alpha)}} \propto \underset{\substack{\text{likelihood} \\ p(D|\theta)}}{p(D|\theta)} \cdot \underset{\text{prior}}{p(\theta|\alpha)}$$
$$p(D|\alpha) = \int p(D|\theta) p(\theta|\alpha) d\theta$$

# Conjugate prior for Exponential Families

- Consider an **Exponential Family** in a natural parameterization:

$$\{ p(x|\eta) = h(x) \cdot \exp(\eta^\top T(x) - A(\eta)) \quad | \quad \eta \in \mathcal{N} \}$$

- Then a **conjugate prior** is given by:

$$p(\eta|\tau, \nu) \propto \exp[\gamma^\top \tau - \nu \cdot A(\gamma)]$$

$$p(\gamma | \tau, \nu) \propto \exp(\gamma^\top \tau - \nu \cdot A(\gamma))$$

$$p(\gamma | \tau, \nu, X) \propto \prod_{n=1}^N p(x^{(n)} | \gamma) \cdot p(\gamma | \tau, \nu)$$

$$\propto \prod_{n=1}^N \exp[\gamma^\top T(x^{(n)}) - A(\gamma)] \cdot \exp[\gamma^\top \tau - \nu \cdot A(\gamma)]$$

$$= \exp\left[\gamma^\top \left(\underbrace{\sum_{n=1}^N T(x^{(n)})}_{\tilde{\tau}} + \tau\right) - \underbrace{(N+\nu) \cdot A(\gamma)}_{\tilde{\nu}}\right]$$

$\propto p(\gamma | \tilde{\tau}, \tilde{\nu}) \Rightarrow \text{conjugate.}$

# Bayesian update rule

$$\begin{array}{ccc} \text{prior} & & \text{posterior} \\ \tau & \xrightarrow{\quad} & \tau + \sum_{n=1}^N T(x^{(n)}) \\ \nu & \xrightarrow{\quad} & \nu + N \end{array}$$

# Example: Multinomial distributions (fixed M) $K$

$$p(x|\pi) = \binom{M}{x} \cdot \pi^x$$

$$P(x|y) = \binom{M}{x} \exp \left[ \sum_{u=0}^K y_u \cdot x_u \right]$$

$$p(y|\tau) \propto \exp \left[ \sum_{u=0}^K y_u \cdot \tau_u \right]$$

$$y_u = \log \pi_u \quad , u=0, \dots, K$$

$$A(y) = 0$$

$$\tau_u \quad , u=0, \dots, K$$

↳ go back to original parametrization

$$\alpha_u := \tau_u + 1$$

$$\sim p(\pi|\alpha) \propto \prod_{u=0}^K \pi_u^{\alpha_u - 1} \quad \text{Dirichlet distribution}$$

# Example: Multinomial distributions (fixed M)

natural  
param.

Prior

$$\tau_0 \\ \vdots \\ \tau_K$$



posterior

$$\tau_0 + \sum_{n=1}^N x_0^{(n)} \\ \vdots \\ \tau_K + \sum_{n=1}^N x_K^{(n)}$$

Bayesian update  
rule:

original  
param.

$$\alpha_0 \\ \vdots \\ \alpha_K$$



$$\alpha_0 + \sum_{n=1}^N x_0^{(n)} \\ \vdots \\ \alpha_K + \sum_{n=1}^N x_K^{(n)}$$

Dirichlet ( $\alpha$ )

# **Machine Learning 2**

## **Exponential Families - Deriving Moments**

Patrick Forré

# Computing Moments for Exponential Families

$$A(\gamma)$$

(minimal number of natural parameters)

$$1.) \nabla_{\gamma} A(\gamma) = \bar{E}_{\gamma}[T(x)]$$

$$2.) \frac{\partial^2}{\partial \gamma_i \partial \gamma_j} A(\gamma) = \text{cov}_{\gamma}(T_i(x), T_j(x))$$

$$A(\gamma) = \log \int h(x) \exp(\gamma^T T(x)) dx$$

Proof:  $\frac{\partial}{\partial \gamma_k} A(\gamma) = \frac{1}{\int h(x) \cdot \exp(\gamma^T T(x)) dx} \cdot \int T_k(x) h(x) \exp(\gamma^T T(x)) dx$

$$= \int T_k(x) \underbrace{\exp(-A(\gamma)) - h(x) \exp(\gamma^T T(x))}_{P(x|\gamma)} dx = \bar{E}_{\gamma}[T_k]$$

$$\begin{aligned}
\frac{\partial^2}{\partial \gamma_i \partial \gamma_j} A(\gamma) &= \frac{\partial}{\partial \gamma_i} \left[ \int T_j(x) \cdot h(x) \exp [y^T T(x) - A(y)] dx \right] \\
&= \int T_j(x) \cdot (T_i(x) - \mathbb{E}_y[T_i(x)]) \cdot \underbrace{h(x) \cdot \exp (y^T T(x) - A(y))}_{p(x|y)} dx \\
&= \mathbb{E}_y[T_j(x) \cdot (T_i(x) - \mathbb{E}_y[T_i(x)])] \\
&= \mathbb{E}_y[T_j(x) \cdot T_i(x)] - \mathbb{E}_y[T_j(x)] \cdot \mathbb{E}_y[T_i(x)] \\
&= \text{cov}_y(T_i(x), T_j(x))
\end{aligned}$$

# Example: Exponential distributions

$$p(x|\lambda) = \mathbb{1}_{[0,\infty)}(x) \cdot \lambda \cdot \exp(-\lambda \cdot x)$$

$$p(x|y) = \mathbb{1}_{[0,\infty)}(x) \exp(y \cdot x - (-\log(-y)))$$

$$\begin{aligned} y &= -\lambda, \quad T(x) = x \\ A[y] &= -\log(-y) \end{aligned}$$

$$\mathbb{E}[T(x)] = A'(y) = +\frac{1}{-y} = \frac{1}{\lambda} \quad (\checkmark)$$

$$\text{Var}(T(x)) = A''(y) = \frac{1}{y^2} = \frac{1}{\lambda^2} \quad (\checkmark)$$

# Example: Gaussian distributions

$$N(x|\mu, \sigma^2) \quad , \quad A(y) = -\frac{y_1^2}{4y_2} - \frac{1}{2} \log(-2y_2) \quad \begin{aligned} y_1 &= \frac{\mu}{\sigma^2} \\ y_2 &= -\frac{1}{2\sigma^2} \end{aligned}$$

$$T_1(x) = x \quad , \quad T_2(x) = x^2$$

- $\mathbb{E}[X] = \mathbb{E}[T_1(x)] = \partial_1 A(y) = -\frac{y_1}{2y_2} (= \left(\frac{\mu}{\sigma^2}\right) \cdot (\sigma^2) = \mu \quad \checkmark)$
- $\text{Var}(X) = \text{Var}(T_1(x)) = \partial_{1,1} A(y) = -\frac{1}{2y_2} (= \sigma^2 \quad \checkmark)$
- $\mathbb{E}[X^2] = \mathbb{E}[T_2(x)] = \partial_2 A(y) = \frac{y_1^2}{4y_2^2} - \frac{1}{2y_2} (= +\mu^2 + \sigma^2 \quad \checkmark)$
- $\text{Var}(X^2) = \text{Var}(T_2(x)) = \partial_{2,2} A(y) = (-2) \frac{y_1^2}{4y_2^3} + \frac{1}{2y_2^2} (= 4\mu^2 + 2\sigma^2 ?)$

# **Machine Learning 2**

## **Exponential Families - Summary**

Patrick Forré

# Summary

- Exponential Family:  $p(x | \eta) = h(x) \cdot \exp(\eta^\top T(x) - A(\eta))$
- How to transform typical distributions into exponential-family-form.
- Conjugate prior:  $p(\eta | \tau, \nu) \propto \exp(\eta^\top \tau - \nu \cdot A(\eta))$
- Bayesian update:  $(\tau, \nu) \mapsto \left( \tau + \sum_{n=1}^N T(x^{(n)}), \nu + N \right)$
- Moments:  $\mathbb{E}[T_k(X)] = \frac{\partial}{\partial \eta_k} A(\eta), \quad \text{cov}(T_i(X), T_j(X)) = \frac{\partial^2}{\partial \eta_i \partial \eta_j} A(\eta)$
- Counter-example: Student-t distributions