

# Deep Exponential Families

Rajesh Ranganath   Linpeng Tang   Laurent Charlin   David Blei

Paper Presentation by  
Evangelos Chatzipantazis  
for CIS620

October 2020



# Table of Contents

- 1 Exponential Families
  - Examples
  - Counter Examples
  - Sufficiency
  - Conjugacy
- 2 Black-box VI
- 3 Going Deep with Exponential Families
- 4 Discussion

# Table of Contents

- 1 Exponential Families
  - Examples
  - Counter Examples
  - Sufficiency
  - Conjugacy
- 2 Black-box VI
- 3 Going Deep with Exponential Families
- 4 Discussion

# Exponential Families

- A family  $\mathcal{P} = \{P_\theta : \theta \in \Theta \subset \mathbb{R}^k\}$  of probability measures on  $\mathcal{X} \subset \mathbb{R}^d$

# Exponential Families

- A family  $\mathcal{P} = \{P_\theta : \theta \in \Theta \subset \mathbb{R}^k\}$  of probability measures on  $\mathcal{X} \subset \mathbb{R}^d$  that have densities  $p_\theta(x)$  of the form:

$$\begin{aligned} p_\theta(x) &= h(x) \exp\left(\sum_{i=1}^m \eta_i(\theta) T_i(x) - A(\theta)\right) \\ &= \frac{h(x)}{Z(\theta)} \exp(\eta(\theta)^T T(x)) \end{aligned}$$

forms an m-parametric exponential family.

# Exponential Families

- A family  $\mathcal{P} = \{P_\theta : \theta \in \Theta \subset \mathbb{R}^k\}$  of probability measures on  $\mathcal{X} \subset \mathbb{R}^d$  that have densities  $p_\theta(x)$  of the form:

$$\begin{aligned} p_\theta(x) &= h(x) \exp\left(\sum_{i=1}^m \eta_i(\theta) T_i(x) - A(\theta)\right) \\ &= \frac{h(x)}{Z(\theta)} \exp(\eta(\theta)^T T(x)) \end{aligned}$$

forms an m-parametric exponential family.

- Disclaimer: The definition above can be extended to the Lebesgue measure (eg. the counting measure for pmfs) or any other ( $\sigma$ -finite) measure, but we stick to the case where  $\mathcal{X} \subset \mathbb{R}^d$  and a pdf exists for clarity

# Exponential Families

- A family  $\mathcal{P} = \{P_\theta : \theta \in \Theta \subset \mathbb{R}^k\}$  of distributions on  $\mathcal{X} \subset \mathbb{R}^d$  that have densities  $p_\theta(x)$  of the form:

$$\begin{aligned} p_\theta(x) &= h(x) \exp \left( \sum_{i=1}^m \eta_i(\theta) T_i(x) - A(\theta) \right) \\ &= \frac{h(x)}{Z(\theta)} \exp (\eta(\theta)^T T(x)) \end{aligned}$$

forms an m-parametric exponential family.

$h : \mathbb{R}^d \rightarrow \mathbb{R}^+$  (support or base measure)

$\eta : \Theta \rightarrow \mathbb{R}^m$  (natural parameter)

$T : \mathbb{R}^d \rightarrow \mathbb{R}^m$  (sufficient statistic)

$A : \Theta \rightarrow \mathbb{R}$  (log partition)

# Exponential Families

- (natural parametrization)  $\mathcal{P} = \{p_\eta : \eta \in \eta(\Theta)\}$

$$p_\eta(x) = h(x) \exp(\eta^\top T(x) - A(\eta))$$



# Exponential Families

- (natural parametrization)  $\mathcal{P} = \{p_\eta : \eta \in \eta(\Theta)\}$

$$p_\eta(x) = h(x) \exp(\eta^\top T(x) - A(\eta))$$

- $\eta = (\eta_1, \dots, \eta_m) \in \mathcal{H}$  (natural parameter space)

# Exponential Families

- (natural parametrization)  $\mathcal{P} = \{p_\eta : \eta \in \eta(\Theta)\}$

$$p_\eta(x) = h(x) \exp(\eta^\top T(x) - A(\eta))$$

- $\eta = (\eta_1, \dots, \eta_m) \in \mathcal{H}$  (natural parameter space)
- $\theta$ 's indicate members of the family.  $\eta$ 's may be redundant.

# Exponential Families

- (natural parametrization)  $\mathcal{P} = \{p_\eta : \eta \in \eta(\Theta)\}$

$$p_\eta(x) = h(x) \exp(\eta^\top T(x) - A(\eta))$$

- $\eta = (\eta_1, \dots, \eta_m) \in \mathcal{H}$  (natural parameter space)
- $\theta$ 's indicate members of the family.  $\eta$ 's may be redundant.
- $A(\eta) = \ln \int_{\mathbb{R}^d} h(x) \exp(\eta^\top T(x)) dx \in \mathbb{R}$ , if  $\eta \in \mathcal{H}$

# Exponential Families

- (natural parametrization)  $\mathcal{P} = \{p_\eta : \eta \in \eta(\Theta)\}$

$$p_\eta(x) = h(x) \exp(\eta^\top T(x) - A(\eta))$$

- $\eta = (\eta_1, \dots, \eta_m) \in \mathcal{H}$  (natural parameter space)
- $\theta$ 's indicate members of the family.  $\eta$ 's may be redundant.
- $A(\eta) = \ln \int_{\mathbb{R}^d} h(x) \exp(\eta^\top T(x)) dx \in \mathbb{R}$ , if  $\eta \in \mathcal{H}$
- for a specific distribution  $\eta$ ,  $T$ ,  $h$  are not uniquely defined.

# Exponential Families

- (natural parametrization)  $\mathcal{P} = \{p_\eta : \eta \in \eta(\Theta)\}$

$$p_\eta(x) = h(x) \exp(\eta^\top T(x) - A(\eta))$$

- $\eta = (\eta_1, \dots, \eta_m) \in \mathcal{H}$  (natural parameter space)
- $\theta$ 's indicate members of the family.  $\eta$ 's may be redundant.
- $A(\eta) = \ln \int_{\mathbb{R}^d} h(x) \exp(\eta^\top T(x)) dx \in \mathbb{R}$ , if  $\eta \in \mathcal{H}$
- for a specific distribution  $\eta$ ,  $T$ ,  $h$  are not uniquely defined.
- $T(x)$  is the  $m$ -dimensional sufficient statistic. Important to find **minimal**  $m$  ( $\Leftarrow$  linearly independent  $T$ ).

# Exponential Families

- (natural parametrization)  $\mathcal{P} = \{p_\eta : \eta \in \eta(\Theta)\}$

$$p_\eta(x) = h(x) \exp(\eta^\top T(x) - A(\eta))$$

- $\eta = (\eta_1, \dots, \eta_m) \in \mathcal{H}$  (natural parameter space)
- $\theta$ 's indicate members of the family.  $\eta$ 's may be redundant.
- $A(\eta) = \ln \int_{\mathbb{R}^d} h(x) \exp(\eta^\top T(x)) dx \in \mathbb{R}$ , if  $\eta \in \mathcal{H}$
- for a specific distribution  $\eta$ ,  $T$ ,  $h$  are not uniquely defined.
- $T(x)$  is the  $m$ -dimensional sufficient statistic. Important to find **minimal**  $m$  ( $\Leftarrow$  linearly independent  $T$ ).
- (minimal)  $k < (\text{minimal}) m \implies \text{curved family} \implies k$  independent parameters embedded in a  $m$  dimensional parameter space.
  - For at least 1 member of a family (at least 1  $\theta$ )

# Exponential Families

- (natural parametrization)  $\mathcal{P} = \{p_\eta : \eta \in \eta(\Theta)\}$

$$p_\eta(x) = h(x) \exp(\eta^\top T(x) - A(\eta))$$

- Misconception:
  - There is no such thing as THE exponential family.
  - $T(, h)$  fixed and  $\eta$  varies  $\implies$  *same family* of distributions.
  - For example,  $\mathcal{N}(x|\mu, \sigma)$  is a *family* of gaussian distributions. This family is actually AN exponential family.

# How to detect an exponential family

- ▶ Takeaway: Observations and parameters must **factorize**
  - ▶ Either directly.
  - ▶ Or both in base and the exponent

$$p(x) = Z(\theta)^{-1} h(x) \exp(\eta(\theta)^\top T(x))$$



# Examples of Exponential Families

$$p(x) = Z(\theta)^{-1} h(x) \exp(\eta(\theta)^\top T(x))$$

1) Exponential (1 parametric):

$$p(x; \theta) = \theta \exp(-\theta x) I(x \geq 0), \quad \Theta = \mathbb{R}^{*+}$$

$\eta(\theta) = \theta$  "canonical/natural form"

# Examples of Exponential Families

$$p(x) = Z(\theta)^{-1} h(x) \exp(\eta(\theta)^\top T(x))$$

1) Exponential (1 parametric):

$$p(x; \theta) = \theta \exp(-\theta x) I(x \geq 0), \quad \Theta = \mathbb{R}^{*+}$$

$$\eta(\theta) = \theta \text{ "canonical/natural form"}$$

2) Binomial (n fixed) (1 parametric):

•

$$\begin{aligned} p(k; \theta) &= \binom{n}{k} \theta^k (1 - \theta)^{(n-k)} \\ &= (1 - \theta)^n \binom{n}{k} I(k \in [n]) \exp\left(\ln \frac{\theta}{1 - \theta} k\right) \end{aligned}$$

•  $\Theta = (0, 1)$

•  $\eta = \ln \frac{\theta}{1 - \theta}$  (logits)  $\implies A(\eta) = n \log(e^\eta + 1)$ ,  $\mathcal{H} = \mathbb{R}$

# Examples of Exponential Families

$$p(x) = Z(\theta)^{-1} h(x) \exp(\eta(\theta)^\top T(x))$$

## 3) Gaussian (2-parametric)

$$\begin{aligned} p(x; \mu, \sigma) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) \\ &= \frac{1}{\sqrt{2\pi}} \frac{1}{\sigma} \exp\left(-\frac{\mu^2}{2\sigma^2}\right) \exp\left(-\left(\frac{1}{2\sigma^2}, \frac{-\mu}{\sigma^2}\right)^\top (x^2, x)\right), \quad \Theta = \mathbb{R} \times \mathbb{R}^{*+} \end{aligned}$$

# Examples of Exponential Families

$$p(x) = Z(\theta)^{-1} h(x) \exp(\eta(\theta)^\top T(x))$$

## 3) Gaussian (2-parametric)

$$\begin{aligned} p(x; \mu, \sigma) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) \\ &= \frac{1}{\sqrt{2\pi}} \frac{1}{\sigma} \exp\left(-\frac{\mu^2}{2\sigma^2}\right) \exp\left(-\left(\frac{1}{2\sigma^2}, \frac{-\mu}{\sigma^2}\right)^\top (x^2, x)\right), \quad \Theta = \mathbb{R} \times \mathbb{R}^{*+} \end{aligned}$$

- $\mathcal{N}(x|\mu, \mu^2)$  is a *curved family* ( $\dim(\theta) < \dim(\eta)$ ).

Statistics are again 2 dimensional ( $T_1(x), T_2(x)$ ) =  $(x^2, x)$ )

- Check  $\text{Cov}_\theta(T(X)) \succ 0$  for some  $\theta$  (for minimality) ✓
- Strict 2-dimensional family (Stonger). ( $\Leftarrow$ )  
 $\{1, T_1(x), T_2(x)\}$  linearly independent for all  $\theta$  (w.r.t the measure).  
And  $\{1, \eta_1(\theta), \eta_2(\theta)\}$  linearly independent.
- $\eta$ 's could/should still have non-linear (curved) dependencies!



# How to get away with Curved

- Multinomial:  $M$  independent trials,  $K$  events,  $X_i = \#(\text{event } i \text{ occurred in } M \text{ trials})$ ,  $\pi_i = \mathbb{P}[X^{(n)} = i]$ ,  $X = (X_1, \dots, X_K)$

# How to get away with Curved

- Multinomial:  $M$  independent trials,  $K$  events,  $X_i = \#(\text{event } i \text{ occurred in } M \text{ trials})$ ,  $\pi_i = \mathbb{P}[X^{(n)} = i]$ ,  $X = (X_1, \dots, X_K)$

- 

$$p(x; \{\pi_i\}_{i=1}^K) = \frac{M!}{x_1! \cdots x_K!} \pi_1^{x_1} \cdots \pi_K^{x_K} = \frac{M!}{x_1! \cdots x_K!} e^{(\sum_{i=1}^K x_i \ln \pi_i)}$$

# How to get away with Curved

- Multinomial:  $M$  independent trials,  $K$  events,  $X_i = \#(\text{event } i \text{ occurred in } M \text{ trials})$ ,  $\pi_i = \mathbb{P}[X^{(n)} = i]$ ,  $X = (X_1, \dots, X_K)$

- 

$$p(x; \{\pi_i\}_{i=1}^K) = \frac{M!}{x_1! \dots x_K!} \pi_1^{x_1} \dots \pi_K^{x_K} = \frac{M!}{x_1! \dots x_K!} e^{(\sum_{i=1}^K x_i \ln \pi_i)}$$

- $h(x) = \frac{M!}{x_1! \dots x_K!} I(x_i \in \mathbb{Z}, \sum_{i=1}^K x_i = M)$  (absorbed in dom. measure)

# How to get away with Curved

- Multinomial:  $M$  independent trials,  $K$  events,  $X_i = \#(\text{event } i \text{ occurred in } M \text{ trials})$ ,  $\pi_i = \mathbb{P}[X^{(n)} = i]$ ,  $X = (X_1, \dots, X_K)$

- 

$$p(x; \{\pi_i\}_{i=1}^K) = \frac{M!}{x_1! \cdots x_K!} \pi_1^{x_1} \cdots \pi_K^{x_K} = \frac{M!}{x_1! \cdots x_K!} e^{(\sum_{i=1}^K x_i \ln \pi_i)}$$

- $h(x) = \frac{M!}{x_1! \cdots x_K!} I(x_i \in \mathbb{Z}, \sum_{i=1}^K x_i = M)$  (absorbed in dom. measure)
- $T(X) = X$ ,  $\eta = \{\ln \pi_i\}_{i=1}^K$ ,  $A(\eta) = 0$  (?).



# How to get away with Curved

- Multinomial:  $M$  independent trials,  $K$  events,  $X_i = \#(\text{event } i \text{ occurred in } M \text{ trials})$ ,  $\pi_i = \mathbb{P}[X^{(n)} = i]$ ,  $X = (X_1, \dots, X_K)$

- 

$$p(x; \{\pi_i\}_{i=1}^K) = \frac{M!}{x_1! \cdots x_K!} \pi_1^{x_1} \cdots \pi_K^{x_K} = \frac{M!}{x_1! \cdots x_K!} e^{(\sum_{i=1}^K x_i \ln \pi_i)}$$

- $h(x) = \frac{M!}{x_1! \cdots x_K!} I(x_i \in \mathbb{Z}, \sum_{i=1}^K x_i = M)$  (absorbed in dom. measure)
- $T(X) = X$ ,  $\eta = \{\ln \pi_i\}_{i=1}^K$ ,  $A(\eta) = 0$  (?).
- Is this family curved? or full-rank? Is it *strictly*  $K$ -dimensional?

# How to get away with Curved

- Intuitively there are only  $K - 1$  independent parameters but  $\eta$  has dimension  $K$ .  $\Theta = \{\pi : \pi_i \in (0, 1), \sum_{i=1}^K \pi_i = 1\}$
- $\mathcal{H} = \mathbb{R}^K$  and not only the subspace  $\mathcal{H}^- = \{\eta : \sum_{i=1}^K e_i^\eta = 1\}$ .
  - $A(\eta) = 0$  if  $\eta \in \mathcal{H}^-$
  - $A(\eta) \neq 0$  if  $\eta \in \mathcal{H}$ !
  - Inconvenient to work on the ambient space. Redundant representation!
- $\mathcal{X} = \{x_j \in \mathbb{N}, \sum_{i=1}^K x_i = M\} \implies \sum_{i=1}^K T_i(x) = M \implies \{1, T_1(x), \dots, T_K(x)\}$  not linear independent  $\implies$  not minimal representation! (More than sufficient)

# Multinomial Version 2

- $x_K = M - x_{K-1} - \dots - x_1$
- $\pi_K = 1 - \sum_{i=1}^{K-1} \pi_i$

$$\begin{aligned} p(x; \pi) &= \binom{M}{x_1, \dots, x_K} e^{(\sum_{i=1}^{K-1} x_i \ln \pi_i + (M - \sum_{i=1}^{K-1} x_i) \ln \pi_K)} \\ &= \binom{M}{x_1, \dots, x_K} \pi_K^M e^{\sum_{i=1}^{K-1} x_i \ln \frac{\pi_i}{\pi_K}} \end{aligned}$$

- Representation is minimal. Not curved, full-rank family, of order  $K-1$ .
- $\pi_k = \frac{e^{\eta_k}}{\sum_{i=1}^{K-1} e^{\eta_i} + 1}$ ,  $k \in [K-1]$  (softmax)
- $A(\eta) = M \ln (\sum_{i=1}^{K-1} e^{\eta_i} + 1)$
- $\mathbb{E}[T_i(X)] = \mathbb{E}[X_i] = M\pi_i = \nabla_{\eta_i} A(\eta)$
- $\text{Cov}[T_i(X), T_j(X)] = \text{Cov}(X_i, X_j) = -M\pi_i\pi_j = \frac{\partial^2 A(\eta)}{\partial \eta_i \partial \eta_j}$

# More exponential families

Name	sufficient stats	domain	use case
Bernoulli	$\phi(x) = [x]$	$\mathbb{X} = \{0; 1\}$	coin toss
Poisson	$\phi(x) = [x]$	$\mathbb{X} = \mathbb{R}_+$	emails per day
Laplace	$\phi(x) = [1, x]^\top$	$\mathbb{X} = \mathbb{R}$	floods
Helmert ( $\chi^2$ )	$\phi(x) = [x, -\log x]$	$\mathbb{X} = \mathbb{R}$	variances
Dirichlet	$\phi(x) = [\log x]$	$\mathbb{X} = \mathbb{R}_+$	class probabilities
Euler ( $\Gamma$ )	$\phi(x) = [x, \log x]$	$\mathbb{X} = \mathbb{R}_+$	variances
Wishart	$\phi(X) = [X, \log  X ]$	$\mathbb{X} = \{X \in \mathbb{R}^{N \times N} \mid v^\top X v \geq 0 \forall v \in \mathbb{R}^N\}$	covariances
Gauss	$\phi(X) = [X, XX^\top]$	$\mathbb{X} = \mathbb{R}^N$	functions
Boltzmann	$\phi(X) = [X, \text{triag}(XX^\top)]$	$\mathbb{X} = \{0; 1\}^N$	thermodynamics

Figure: More Exponential Families. Figure taken from P.Hennig

# Counter Examples

- As we saw many exponential families: Beta, Gamma, Poisson, Laplace, Chi-squared, Wishart.

# Counter Examples

- As we saw many exponential families: Beta, Gamma, Poisson, Laplace, Chi-squared, Wishart.
- Do we know any parametric families that are not exponential?

# Counter Examples

- As we saw many exponential families: Beta, Gamma, Poisson, Laplace, Chi-squared, Wishart.
- Do we know any parametric families that are not exponential?
- Remember the parameters and observations must *factorize*.

# Counter Examples

- As we saw many exponential families: Beta, Gamma, Poisson, Laplace, Chi-squared, Wishart.
- Do we know any parametric families that are not exponential?
- Remember the parameters and observations must *factorize*.
- A non Example:  $p_{\theta}(x) = U(0, \theta) = \frac{1}{\theta} I(x \in [0, \theta])$



# Counter Examples

- As we saw many exponential families: Beta, Gamma, Poisson, Laplace, Chi-squared, Wishart.
- Do we know any parametric families that are not exponential?
- Remember the parameters and observations must *factorize*.
- A non Example:  $p_{\theta}(x) = U(0, \theta) = \frac{1}{\theta} I(x \in [0, \theta])$
- A note: Of course, for any fixed  $\theta$  we get a uniform distribution, which is a trivial exponential family. We cannot gather all those trivial exponential families into a parametric family that is still exponential!

# More Counter Examples

- Terms  $1 + f(x)g(\theta)$  do not factorize:
  - Cauchy:  $\frac{1}{\pi\theta} \frac{1}{(x-x_0)^2/\theta^2+1}$
  - Student't, etc.
- Note:  $h(x)$  is fixed for all  $\theta$ . In exponential families the support is fixed.
  - A non example:  $p(x; \{\theta, n\}) = \text{Bin}(x; \{\theta, n\})$
- Most *mixtures* are not exponential families: Mixture of Gaussians.
- Most *Compound Probability Distributions* are not exponential!
  - Except Conjugacy!
  - Relevant to today's paper.

# Sufficient Statistic

- A *Statistic* is a function of the sample (only).

$$T : x \in \mathcal{X} \rightarrow T(x) = t \in \mathcal{T}$$

Since  $X$  is a random variable, so is  $T(X)$ .

# Sufficient Statistic

- A *Statistic* is a function of the sample (only).

$$T : x \in \mathcal{X} \rightarrow T(x) = t \in \mathcal{T}$$

Since  $X$  is a random variable, so is  $T(X)$ .

## Sufficient Statistic

A Statistic  $T$  is *Sufficient* for the statistical model  $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$  of  $X$  if the conditional distribution of  $X$  given  $T = t$  is independent of  $\theta$ .

- ▶ How frequentists think about it:  $p(x|T(x), \theta) = p(x|T(x))$   
▶ How bayesians think about it:  $p(\theta|x, T(x)) = p(\theta|T(x))$

# Sufficient Statistic

- A *Statistic* is a function of the sample (only).

$$T : x \in \mathcal{X} \rightarrow T(x) = t \in \mathcal{T}$$

Since  $X$  is a random variable, so is  $T(X)$ .

## Sufficient Statistic

A Statistic  $T$  is *Sufficient* for the statistical model  $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$  of  $X$  if the conditional distribution of  $X$  given  $T = t$  is independent of  $\theta$ .

- - ▶ How frequentists think about it:  $p(x|T(x), \theta) = p(x|T(x))$
  - ▶ How bayesians think about it:  $p(\theta|x, T(x)) = p(\theta|T(x))$
  - ▶ Creates an *information bottleneck* between our data and our parameters.

# Sufficient Statistic

- A *Statistic* is a function of the sample (only).

$$T : x \in \mathcal{X} \rightarrow T(x) = t \in \mathcal{T}$$

Since  $X$  is a random variable, so is  $T(X)$ .

## Sufficient Statistic

A Statistic  $T$  is *Sufficient* for the statistical model  $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$  of  $X$  if the conditional distribution of  $X$  given  $T = t$  is independent of  $\theta$ .

- - ▶ How frequentists think about it:  $p(x|T(x), \theta) = p(x|T(x))$
  - ▶ How bayesians think about it:  $p(\theta|x, T(x)) = p(\theta|T(x))$
  - ▶ Creates an **information bottleneck** between our data and our parameters.
  - ▶ **Data reduction** technique. You can throw away the data as long as the statistics are known. No more information for inference on  $\theta$

# Sufficiency

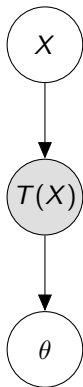


Figure: a) Bayesian

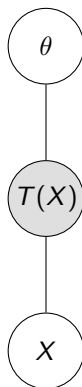


Figure: c) Neyman Factorization

$$p(x, T(x), \theta) = \psi_1(T(x), \theta) \psi_2(T(x), x)$$

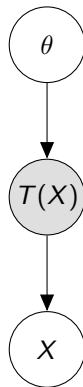


Figure: b) Frequentist

# Factorization Theorem

## Theorem (Neyman Factorization Theorem)

Let  $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$  be a statistical model, with probability function  $p(\cdot; \theta)$ . A statistic  $T$  is sufficient for  $\mathcal{P}$  if and only if there exist non-negative functions  $g(\cdot; \theta)$  and  $h$  such that the probability function satisfies:

$$p(x; \theta) = g(T(x); \theta)h(x)$$

- ▶  $\theta$  is connected with  $X$  only through  $T(X)$ .
- ▶ We can see why  $T(x)$  is a sufficient statistic in exponential families.



# Factorization Theorem

## Theorem (Neyman Factorization Theorem)

Let  $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$  be a statistical model, with probability function  $p(\cdot; \theta)$ . A statistic  $T$  is sufficient for  $\mathcal{P}$  if and only if there exist non-negative functions  $g(\cdot; \theta)$  and  $h$  such that the probability function satisfies:

$$p(x; \theta) = g(T(x); \theta)h(x)$$

- ▶  $\theta$  is connected with  $X$  only through  $T(X)$ .
- ▶ We can see why  $T(x)$  is a sufficient statistic in exponential families.

## Example ( $X \sim \mathcal{N}(\mu, \sigma)$ Draw $n$ i.i.d. samples)

$$p(x; \theta) = \left(\frac{1}{2\pi\sigma^2}\right)^{n/2} \exp\left(-\frac{1}{2\sigma^2}\left(\sum_{i=1}^n x_i^2 - 2\mu \sum_{i=1}^n x_i + n\mu^2\right)\right) \Rightarrow$$

$$T(x) = \left(\sum_{i=1}^n x_i, \sum_{i=1}^n x_i^2\right) \text{ are sufficient statistics. } h=1, g \text{ rest.}$$

- ▶ With  $n$  i.i.d gaussian samples, the most we can infer about  $\theta$  is in the empirical mean and the empirical variance

- ▶ A Sufficient Statistic *partitions* the sample space.

# Minimal Sufficient Statistics

- ▶ A Sufficient Statistic *partitions* the sample space.
- ▶ We would like to find the *coarsest* partition possible. We call that a *minimal sufficient* statistic. (Assume strictly m-parametric family)

# Minimal Sufficient Statistics

- ▶ A Sufficient Statistic *partitions* the sample space.
- ▶ We would like to find the *coarsest* partition possible. We call that a *minimal sufficient* statistic. (Assume strictly m-parametric family)
- ▶ It is essentially a function of any other sufficient statistic  
 $(T'(x) = T'(y) \implies T(x) = T(y), \forall x, y \in \mathcal{X})$

# Minimal Sufficient Statistics

- ▶ A Sufficient Statistic *partitions* the sample space.
- ▶ We would like to find the *coarsest* partition possible. We call that a *minimal sufficient* statistic. (Assume strictly m-parametric family)
- ▶ It is essentially a function of any other sufficient statistic  
 $(T'(x) = T'(y) \implies T(x) = T(y), \forall x, y \in \mathcal{X})$
- ▶ Criterion:  $\frac{p(x;\theta)}{p(y;\theta)}$  independent of  $\theta, \forall x, y \in \mathcal{X} \implies T(x) = T(y)$

# Minimal Sufficient Statistics

- ▶ A Sufficient Statistic *partitions* the sample space.
- ▶ We would like to find the *coarsest* partition possible. We call that a *minimal sufficient* statistic. (Assume strictly m-parametric family)
- ▶ It is essentially a function of any other sufficient statistic  
 $(T'(x) = T'(y) \implies T(x) = T(y), \forall x, y \in \mathcal{X})$
- ▶ Criterion:  $\frac{p(x;\theta)}{p(y;\theta)}$  independent of  $\theta, \forall x, y \in \mathcal{X} \implies T(x) = T(y)$
- ▶ For the example above  $T(x) = (\sum_{i=1}^n x_i, \sum_{i=1}^n x_i^2)$  is a minimal sufficient statistic.

# Minimal Sufficient Statistics

- ▶ A Sufficient Statistic *partitions* the sample space.
- ▶ We would like to find the *coarsest* partition possible. We call that a *minimal sufficient* statistic. (Assume strictly m-parametric family)
- ▶ It is essentially a function of any other sufficient statistic  
 $(T'(x) = T'(y) \implies T(x) = T(y), \forall x, y \in \mathcal{X})$
- ▶ Criterion:  $\frac{p(x;\theta)}{p(y;\theta)}$  independent of  $\theta, \forall x, y \in \mathcal{X} \implies T(x) = T(y)$
- ▶ For the example above  $T(x) = (\sum_{i=1}^n x_i, \sum_{i=1}^n x_i^2)$  is a minimal sufficient statistic.
- ▶ Note: Different than the *strictly m-parametric* condition or the minimality in m, we discussed above.

# Minimal Sufficiency in Exponential Family

- For a sample, of i.i.d. random variables from a strictly m-parametric exponential family it holds:

$$T_{(n)}(x) = \left( \sum_{i=1}^n T_1(x_i), \dots, \sum_{i=1}^n T_m(x_i) \right)$$

is a minimal sufficient statistic.

- Note: Dimensions of  $T$  fixed; independent of  $n$ !
- Sufficiency still holds for non-strict exponential families, minimality not. So curved families still have sufficient statistics.
- In some sense, exponential family is the *only* family with minimal, finite-dimensional, sufficient statistics ([Pitman–Koopman–Darmois theorem](#))
  - $U(0, \theta)(?)$  : non fixed domain!



# Conjugacy

- Conjugate Priors: Important for Bayesian Statistics, inference.

# Conjugacy

- Conjugate Priors: Important for Bayesian Statistics, inference.
- For a likelihood function in exponential *the conjugate prior is again an exponential family*.

# Conjugacy

- Conjugate Priors: Important for Bayesian Statistics, inference.
- For a likelihood function in exponential *the conjugate prior is again an exponential family*.
- Conjugate prior  $p_{\pi}$  of the parameter  $\eta$  of an exponential family.

- Consider the exponential family:

$$p(x|\eta) = h(x) \exp(\eta^\top T(x) - \ln A(\eta))$$

# Conjugacy

- Consider the exponential family:

$$p(x|\eta) = h(x) \exp(\eta^\top T(x) - \ln A(\eta))$$

- Its conjugate is the exponential family:

$$p_\pi(\eta|\alpha, \nu) = \exp\left(\begin{pmatrix} \alpha \\ \nu \end{pmatrix}^\top \begin{pmatrix} \eta \\ A(\eta) \end{pmatrix} - \ln F(\alpha, \nu)\right)$$

# Conjugacy

- Consider the exponential family:

$$p(x|\eta) = h(x) \exp(\eta^\top T(x) - \ln A(\eta))$$

- Its conjugate is the exponential family:

$$p_\pi(\eta|\alpha, \nu) = \exp\left(\begin{pmatrix} \alpha \\ \nu \end{pmatrix}^\top \begin{pmatrix} \eta \\ A(\eta) \end{pmatrix} - \ln F(\alpha, \nu)\right)$$

- ▶  $F(\alpha, \nu) = \int_{\mathcal{H}} \exp\left(\begin{pmatrix} \alpha \\ \nu \end{pmatrix}^\top \begin{pmatrix} \eta \\ A(\eta) \end{pmatrix}\right) d\eta$
- ▶  $\nu$  : pseudo-observations from prior.
- ▶  $\alpha$ : effective amount these pseudo-observations contribute to the sufficient statistic (vector).

# Posterior and Bayesian Inference

- ▶ Posterior (n iid samples):

$$p(\eta|X, \alpha, \nu) = p_\pi(\eta|\alpha + \sum_{i=1}^n T(x_i), \nu + n)$$

- ▶ Same family as the prior (T, h)
- ▶ Predictive (1 sample):

$$p(x) = \int_{\eta} p(x|\eta) p_\pi(\eta|\alpha, \nu) d\eta = h(x) \frac{F(T(x) + \alpha, \nu + 1)}{F(\alpha, \nu)}$$

- ▶ Calculation of  $F(\alpha, \nu)$ , often intractable! Even if  $A(\eta)$  is known.

# Parameter Estimation in Exponential Families

- for  $N$  iid datapoints, parameter estimation in  $\mathcal{O}(N)$  from sufficient statistics.
- $\mathcal{H}$  is a convex set.
- As we saw in the examples above (for full rank families):
  - $\mathbb{E}[T(X)] = \nabla_{\eta} A(\eta)$
  - $[Cov(T_i(X), T_j(X))] = Hessian[A(\eta)]$
  - Thus  $A(\eta)$  is convex (strictly for minimal families).
- $l(\eta) = \ln(\prod_{n=1}^N h(x_n)) + \eta^{\top}(\sum_{n=1}^N T(x_n)) - NA(\eta)$ 
  - $\nabla_{\eta} l(\eta) = 0 \iff \mathbb{E}[T(X)] = \nabla_{\eta} A(\eta) = \sum_{n=1}^N T(x_n)$
  - $Hessian[l(\eta)] = Hessian[A(\eta)] = [Cov(T_i, T_j)] \succcurlyeq 0$
  - Data appear only through the sufficient stat. for  $\eta$  estimation.
  - Not always closed form! But iterative algorithms converge!
- Note:  $l(\eta)$  concave (in  $\eta$ ), does not mean pdf is unimodal (in  $x$ ). See Beta(1/2,1/2).



- Unbiased MLE estimator:  $\mathbb{E}[\sum_{n=1}^N T(x_n)] = \mathbb{E}[T(X)]$  (not for  $\eta$ , but there exists an equivalent *mean parametrization*)

# Cramer-Rao, Rao-Blackwell and more

- Unbiased MLE estimator:  $\mathbb{E}[\sum_{n=1}^N T(x_n)] = \mathbb{E}[T(X)]$  (not for  $\eta$ , but there exists an equivalent *mean parametrization*)
- Attains **Cramer-Rao** bound:  $I(\eta) = -\mathbb{E}[\frac{\partial^2 \ln p(X|\eta)}{d\eta^2}] = \text{Cov}(T(X))$  (optimal)

# Cramer-Rao, Rao-Blackwell and more

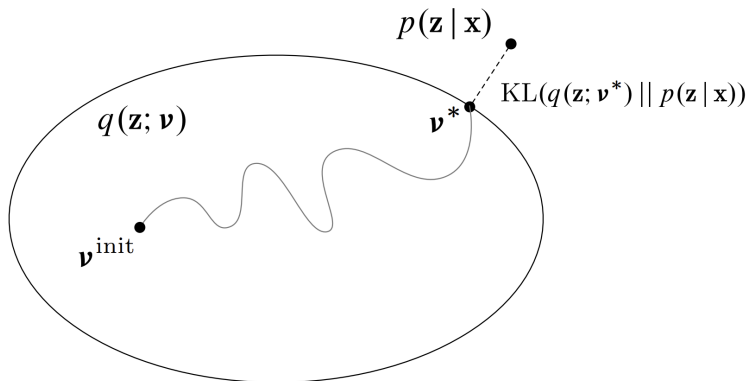
- Unbiased MLE estimator:  $\mathbb{E}[\sum_{n=1}^N T(x_n)] = \mathbb{E}[T(X)]$  (not for  $\eta$ , but there exists an equivalent *mean parametrization*)
- Attains **Cramer-Rao** bound:  $I(\eta) = -\mathbb{E}[\frac{\partial^2 \ln p(X|\eta)}{d\eta^2}] = \text{Cov}(T(X))$  (optimal)
- Other interesting properties:
  - Solutions to **Maximum Entropy problems**, consistent with constraints in expected values. For example, Normal is the maximum entropy distribution, out of all those with a bounded variance.
  - **Bregman Divergence** between parameters (wrt log partition) is the KL divergence between the distributions with those parameters.
  - **Method of Moments** converge to the maximum likelihood estimator.
  - Full rank exponential family  $\implies$  Complete and Sufficient stat.  $\implies$  UMVUE for mean-value param. (**Rao-Blackwell**, **Lehmann-Scheffe**)
  - **Conditional** Conjugacy

# Table of Contents

- 1 Exponential Families
  - Examples
  - Counter Examples
  - Sufficiency
  - Conjugacy
- 2 Black-box VI
- 3 Going Deep with Exponential Families
- 4 Discussion

# Black-box VI [Ranganath, Gerrish, and Blei 2014]

- Variational Inference:
  - Recasts posterior inference as an optimization problem.
  - Approximate only if the variational family does not include the posterior; or the optimization algorithm cannot find it.
  - Always strive for richer variational families. No notion of "overfitting".



- Variational Family does not need to be parametric. Calculus of Variations is functional optimization! Still waiting for mathematicians to step up...
- For now, we posit a parametric posterior  $q(z; \nu) \in \mathcal{Q}$  over the latent variables and optimize over  $\nu$ :
  - (Typical) mean-field family:  $q(z) = \prod_{i=1}^K q_i(z_i)$
  - Amortization:  $q(z_i|x_i; \nu)$ , shared  $\nu$
  - (Typical) Algorithms: SGD performs VI. [Chaudhari and Soatto 2018]
  - (Typical) Objective:  $KL(q(z; \nu) || p(z|x)) \rightarrow \min_{\nu}$ 
    - KL divergence is problematic ( $\infty$  if no overlap; in high dim. spaces, no overlap  $\implies$  no gradient.
    - Mode seeking behavior of  $KL(q||p)$  plus mean-field assumptions underestimate the variance!

- Optimizing  $\nu$  is one task, we also need to optimize  $\theta$ , the parameters of our model!
- ELBO in Reverse:
  - $ELBO(\theta, q) = \mathbb{E}_{x_i \in p_{data}} \mathbb{E}_{z_i \sim q_i(z_i|x_i)} [\ln(p(z_i)p(x_i|z_i; \theta)) - \ln q_i(z_i|x_i)]$
  - $ELBO(\theta, q) = -\mathbb{E}_{x_i \sim p_{data}} [KL(q_i(z_i|x_i) || p(z_i|x_i; \theta))] + \ln p(D; \theta)$
  - Maximization:  $\max_{\theta} \max_{q \in \mathcal{Q}} ELBO(\theta, q) \leq \max_{\theta} \ln p(D; \theta)$
  - " = " if  $p(z|x) \in \mathcal{Q}$

## Motivation behind BBVI:

- In need of **scalable** variational inference (massive data). In need of **generic** variational inference, no model specific bounds.
- Follow noisy, unbiased estimates of the gradient! Double stochasticity in sampling dataset and approximate posterior.
- Reformulate ELBO's gradient as an expectation. Use Monte Carlo sampling.



- How to optimize functions of the form:

$$\mathcal{L}(\nu, \theta) = \mathbb{E}_{q_\nu(z)}[f(z; \theta)]$$

- "Easy" w.r.t  $\theta$ , Difficult w.r.t  $\nu$  (Stochastic Optimization!)
- **REINFORCE** Estimator (or log-derivative trick, or score-function estimator)

$$\nabla_\nu \mathcal{L}(\nu, \theta) = \mathbb{E}_{q_\nu(z)}[\nabla_\nu \ln q_\nu(z) f(z)]$$

- Note:  $E_q[\nabla_\nu \ln q_\nu(z)] = 0 \implies$  oscillates. around 0.

About the log-derivative trick:

- The estimator has HUGE variance, not much better than random search ([Mania, Guy, and Recht 2018]). HUGE variance in stochastic optimization theory  $\implies$  slow convergence.
- Pretends to be a 1-st order estimator, it is actually a 0-th order estimator (no  $f(z; \theta)$  gradient)
- Control Variate Techniques can reduce it but not enough to be practical (if not careful, can magnify it).
- We have seen another estimator ([Reparametrization trick](#)) ([Kingma and Welling 2014]).
  - Very small variance when it is applicable (location-scale distributions;  $z = T(\epsilon, \nu)$ ).
  - Latent has to be continuous! Otherwise no differentiable reparametrization exists! (Gumbel-Softmax relaxation (?)).

## Black-box Variational Inference Necessary Criteria:

- Sample  $q(z; \nu)$  easily.
- Evaluate  $\ln p(x, z), \ln q(z)$  easily.
- Compute  $\nabla_{\nu} \ln q_{\nu}(z)$  easily.

# Table of Contents

- 1 Exponential Families
  - Examples
  - Counter Examples
  - Sufficiency
  - Conjugacy
- 2 Black-box VI
- 3 Going Deep with Exponential Families
- 4 Discussion

# Deep Exponential Families

Model Specification:

$$p(t_i) = \mathcal{N}(0, I)$$

$$p(x_i | t_i) = \mathcal{N}(Wt_i + b, \Sigma)$$

$$\implies p(x_i) = (?)$$

# Deep Exponential Families

- Model Specification :

$$\begin{aligned}p(t_i) &= \mathcal{N}(0, I) \\p(x_i|t_i) &= \mathcal{N}(Wt_i + b, \Sigma) \\ \implies p(x_i) &= \mathcal{N}(b, WW^\top + \Sigma)\end{aligned}$$

# Deep Exponential Families

- Model Specification :

$$\begin{aligned}p(t_i) &= \mathcal{N}(0, I) \\p(x_i|t_i) &= \mathcal{N}(Wt_i + b, \Sigma) \\ \implies p(x_i) &= \mathcal{N}(b, WW^\top + \Sigma)\end{aligned}$$

- Model Specification:

$$\begin{aligned}p(t_i) &= \mathcal{N}(0, I) \\p(x_i|t_i) &= \mathcal{N}(\text{ResNet}(t_i), \Sigma) \\ \implies p(x_i) &= (?)\end{aligned}$$

# Deep Exponential Families

- Model Specification :

$$\begin{aligned}p(t_i) &= \mathcal{N}(0, I) \\p(x_i|t_i) &= \mathcal{N}(Wt_i + b, \Sigma) \\ \implies p(x_i) &= \mathcal{N}(b, WW^\top + \Sigma)\end{aligned}$$

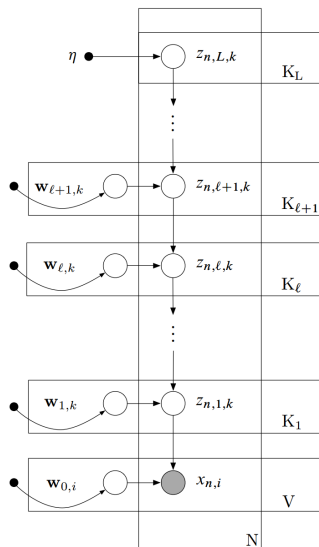
- Model Specification:

$$\begin{aligned}p(t_i) &= \mathcal{N}(0, I) \\p(x_i|t_i) &= \mathcal{N}(\text{ResNet}(t_i), \Sigma) \\ \implies p(x_i) &= (?)\end{aligned}$$

- Indicates that Deep Exponential Families will be more expressive than shallow ones.

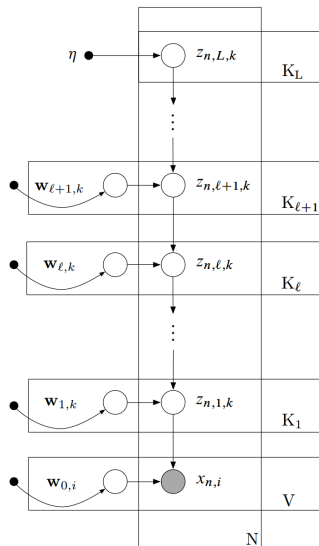


# Generative Model: The story of our data



- Chain Exponential Family in a hierarchy.
- Each layer's draw is input to the natural parameters of the next.
- Plate notation.
  - $\eta$ , hyperparameters of  $w_{l \in L}$ :  
Shared (for data and z-nodes)
  - $w_{l,k}$  shared across data, local for each z-node.

# Generative Model: The story of our data



- Model specification (for each  $x_n$ ,  $n \in [N]$ ):

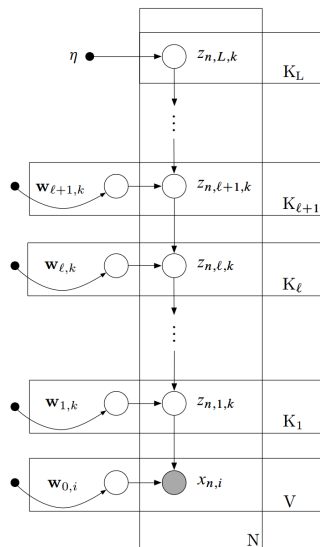
$$p(z_{n,L,k}) = \text{ExpFam}_L(\eta) \quad \forall k \in [K_L]$$

$$p(z_{n,l,k} | z_{n,l+1,k}, \mathbf{w}_{l,k}) = \text{ExpFam}_l(g_l(z_{n,l+1,k}^\top \mathbf{w}_{l,k}))$$

$$p(x_{n,i} | z_{n,1,k}, \mathbf{w}_{0,i}) = \text{ExpFam}_0(g_0(z_{n,1,k}^\top \mathbf{w}_{0,i}))$$

$$W_l \sim p(W_l)$$

# Generative Model: The story of our data



- Model specification (for each  $x_n$ ,  $n \in [N]$ ):

$$p(z_{n,L,k}) = \text{ExpFam}_L(\eta) \quad \forall k \in [K_L]$$

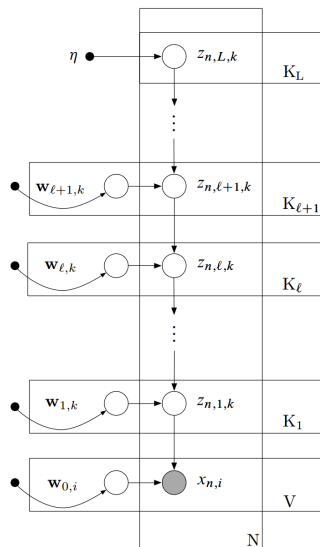
$$p(z_{n,l,k} | z_{n,l+1}, \mathbf{w}_{l,k}) = \text{ExpFam}_l(g_l(z_{n,l+1}^\top \mathbf{w}_{l,k}))$$

$$p(x_{n,i} | z_{n,1}, \mathbf{w}_{0,i}) = \text{ExpFam}_0(g_0(z_{n,1}^\top \mathbf{w}_{0,i}))$$

$$\mathbf{W}_l \sim p(\mathbf{W}_l)$$

- $z_{n,l,k}$ : scalar.  
 $z_{n,l+1}$ :  $K_{l+1}$ -vector.  
 $\mathbf{w}_{l,k}$ :  $K_{l+1}$ -vector.  
 $K_l$  vectors like  $\mathbf{w}_{l,k}$  in layer  $l$ .  
 $\mathbf{W}_l \in \mathbb{R}^{K_l \times K_{l+1}}$

# Generative Model: The story of our data



- Model specification (for each  $x_n$ ,  $n \in [N]$ ):

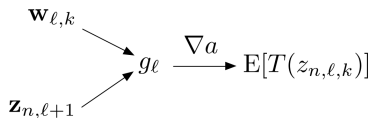
$$p(z_{n,L,k}) = \text{ExpFam}_L(\eta) \quad \forall k \in [K_L]$$

$$p(z_{n,l,k} | z_{n,l+1,k}, \mathbf{w}_{l,k}) = \text{ExpFam}_l(g_l(z_{n,l+1,k}^\top \mathbf{w}_{l,k}))$$

$$p(x_{n,i} | z_{n,1}, \mathbf{w}_{0,i}) = \text{ExpFam}_0(g_0(z_{n,1}^\top \mathbf{w}_{0,i}))$$

$$\mathbf{W}_l \sim p(\mathbf{W}_l)$$

- $z_{n,l,k}$ : scalar.  
 $z_{n,l+1}$ :  $K_{l+1}$ -vector.  
 $\mathbf{w}_{l,k}$ :  $K_{l+1}$ -vector.  
 $K_l$  vectors like  $\mathbf{w}_{l,k}$  in layer  $l$ .  
 $\mathbf{W}_l \in \mathbb{R}^{K_l \times K_{l+1}}$
- $\text{ExpFam}_l()$  can be different across layers. But it has to be defined a priori.



As we have already seen:

$$\mathbb{E}_{z_{n,l,k}}[T(z_{n,l,k})] = \nabla_{\eta} a(\eta) \Big|_{g_l(\mathbf{w}_{l,k}^{\top} \mathbf{z}_{n,l+1})}$$

- Misconception: Sufficient Statistics are fixed a priori. Expected Sufficient Statistics are functions of the parameters.
- Two kinds of non-linearities (from  $\mathbf{w}$  to  $z$ ): 1) link function  $g(\mathbf{x})$ , 2) log-normalizer's gradient  $\nabla_{\eta} a(\eta)$ .

# Example Link Function

- We reformulated the Binomial into its natural parametrization, above.
- If we do the same for Bernoulli we get:

$$p(z_{n,l,k}|\eta) = \exp(\eta z_{n,l,k} - \ln(1 + \exp \eta))$$

where  $z_{n,l,k} \in \{0, 1\}$ ,  $\eta = \text{a.s.} g_l(\mathbf{z}_{n,l+1}^\top \mathbf{w}_{l,k})$

- With identity link function:

$$\mathbb{E}[z_{n,l,k}|\eta] = \frac{1}{1 + \exp(-\eta)} = \sigma(\mathbf{z}_{n,l+1}^\top \mathbf{w}_{l,k})$$

Sigmoid Belief Network!

- Gamma PDF:

$$Z \sim \text{Gamma}(\alpha, \beta) \iff p(z) = \frac{1}{\Gamma(\alpha)} \beta^\alpha z^{\alpha-1} e^{-\beta z}, \quad z, \alpha, \beta > 0$$

# Sparse Gamma DEF

- Gamma PDF:

$$Z \sim \text{Gamma}(\alpha, \beta) \iff p(z) = \frac{1}{\Gamma(\alpha)} \beta^\alpha z^{\alpha-1} e^{-\beta z}, \quad z, \alpha, \beta > 0$$

- Shape parameter  $\alpha$ , Rate parameter  $\beta$ .

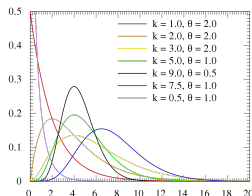


# Sparse Gamma DEF

- Gamma PDF:

$$Z \sim \text{Gamma}(\alpha, \beta) \iff p(z) = \frac{1}{\Gamma(a)} \beta^\alpha z^{\alpha-1} e^{-\beta z}, \quad z, \alpha, \beta > 0$$

- Shape parameter  $\alpha$ , Rate parameter  $\beta$ .

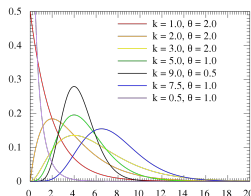


# Sparse Gamma DEF

- Gamma PDF:

$$Z \sim \text{Gamma}(\alpha, \beta) \iff p(z) = \frac{1}{\Gamma(a)} \beta^\alpha z^{\alpha-1} e^{-\beta z}, \quad z, \alpha, \beta > 0$$

- Shape parameter  $\alpha$ , Rate parameter  $\beta$ .



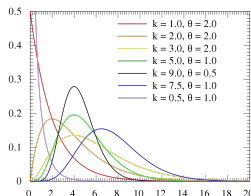
- $\mathbb{E}[z] = \alpha/\beta$ ,  $\text{Var}[z] = \alpha/\beta^2$

# Sparse Gamma DEF

- Gamma PDF:

$$Z \sim \text{Gamma}(\alpha, \beta) \iff p(z) = \frac{1}{\Gamma(\alpha)} \beta^\alpha z^{\alpha-1} e^{-\beta z}, \quad z, \alpha, \beta > 0$$

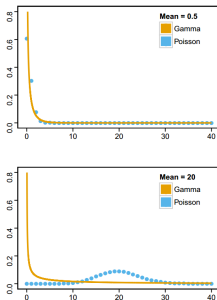
- Shape parameter  $\alpha$ , Rate parameter  $\beta$ .



- $\mathbb{E}[z] = \alpha/\beta, \quad \text{Var}[z] = \alpha/\beta^2$
- $p(z) = z^{-1} \exp((\alpha, \beta) \cdot (\ln z, -z) - \ln \Gamma(\alpha) - \alpha \ln \beta)$

# Sparse Gamma DEF

- $p(W_l) \sim \text{Gamma}(\alpha_{W_l}, \beta_{W_l})$  to ensure  $\mathbf{z}_{l+1}^\top \mathbf{w}_{l,k} > 0$
- Set  $\alpha_l, \alpha_{W_l} < 1 \implies$  soft spike-slab prior (**sparse gamma**).
- Sparse Gamma DEFs for documents means that an observable does not need to express every "super-topic" in the "concept" it expresses.
- $g_\alpha(\cdot) = \alpha_l = \text{const}$ ,  $g_\beta(\cdot) = \frac{\alpha_l}{\mathbf{z}_{n,l+1}^\top \mathbf{w}_{l,k}} \implies$
- $\mathbb{E}[z_{n,l,k}] = \mathbf{z}_{n,l+1}^\top \mathbf{w}_{l,k}$



z-Dist	$\mathbf{z}_{\ell+1}$	W-dist	$\mathbf{w}_{\ell,k}$	$g_\ell$	$E[T(z_{\ell,k})]$
Gamma	$R_+^{K_{\ell+1}}$	Gamma	$R_+^{K_{\ell+1}}$	[constant; inverse]	$[z_{\ell+1}^\top \mathbf{w}_{\ell,k}; \Psi(\alpha_\ell) - \log(\alpha) + \log(z_{\ell+1}^\top \mathbf{w}_{\ell,k})]$
Bernoulli	$\{0, 1\}^{K_{\ell+1}}$	Normal	$R_+^{K_{\ell+1}}$	identity	$\sigma(z_{\ell+1}^\top \mathbf{w}_{\ell,k})$
Poisson	$N^{K_{\ell+1}}$	Gamma	$R_+^{K_{\ell+1}}$	log	$z_{\ell+1}^\top \mathbf{w}_{\ell,k}$
Poisson	$N^{K_{\ell+1}}$	Normal	$R_+^{K_{\ell+1}}$	log-softmax	$\log(1 + \exp(z_{\ell+1}^\top \mathbf{w}_{\ell,k}))$

Table 1: A summary of all the DEFs we present in terms of their layer distributions, weight distributions, and link functions.

- Compounding Exponential Families is not an Exponential Family (Except Conjugacy!)
- Black-box Variational Inference to train DEFs.
  - Monte Carlo approximation of variational objective (and gradient).
  - Stochastic Optimization routine, follows noisy unbiased gradients.
- Mean-field family:

$$q(z, W) = q(\mathbf{W}_0) \prod_{l=1}^L q(\mathbf{W}_l) \prod_{n=1}^N q(z_{n,l})$$

- $q(\mathbf{W}_l | \xi_l)$ ,  $q(z_{n,l})$  fully factorized. Same family as the model distribution  $p$ . In actual model,  $\mathbf{z}_{n,l} \not\perp \mathbf{z}_{m,l} | \{x_m, x_n\}$  (common cause  $W$ 's)
- $z_{n,l,k} \sim \text{ExpFam}_l(\lambda_{n,l,k})$

- $\mathcal{L}(\lambda, \xi) = \mathbb{E}_{q(z, W; \xi, \lambda)} [\ln p(x, z, W) - \ln q(z, W)]$

# Inference BBVI

- $\mathcal{L}(\lambda, \xi) = \mathbb{E}_{q(z, W; \xi, \lambda)} [\ln p(x, z, W) - \ln q(z, W)]$
- Intractable to compute estimation! But we "only" need to compute derivative!



# Inference BBVI

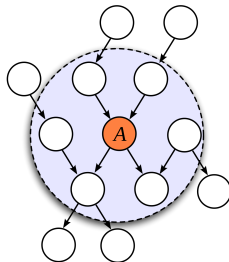
- $\mathcal{L}(\lambda, \xi) = \mathbb{E}_{q(z, W; \xi, \lambda)}[\ln p(x, z, W) - \ln q(z, W)]$
- Intractable to compute estimation! But we "only" need to compute derivative!
- $\nabla_{\lambda, \xi} \mathcal{L}(\lambda, \xi) = \mathbb{E}_q[\nabla_{\lambda, \xi} \ln q (\ln p(x, W, z) - \ln q(W, z))]$   
 $\approx \frac{1}{S} \sum_{i=1}^S \nabla_{\lambda, \xi} \ln q(z_s, W_s | \lambda, \xi) (\ln p(x, W_s, z_s) - \ln q(W_s, z_s | \lambda, \xi))$   
 $z_s, W_s \sim q(z, W | \lambda, \xi)$

# Inference BBVI

- $\mathcal{L}(\lambda, \xi) = \mathbb{E}_{q(z, W; \xi, \lambda)} [\ln p(x, z, W) - \ln q(z, W)]$
- Intractable to compute estimation! But we "only" need to compute derivative!
- $\nabla_{\lambda, \xi} \mathcal{L}(\lambda, \xi) = \mathbb{E}_q [\nabla_{\lambda, \xi} \ln q (\ln p(x, W, z) - \ln q(W, z))]$   
 $\approx \frac{1}{S} \sum_{i=1}^S \nabla_{\lambda, \xi} \ln q(z_s, W_s | \lambda, \xi) (\ln p(x, W_s, z_s) - \ln q(W_s, z_s | \lambda, \xi))$   
 $z_s, W_s \sim q(z, W | \lambda, \xi)$
- Each row in the jacobian is actually sparse (for this graphical model).  
The dependency on  $z_{n,l,k}$  is only through its markov blanket.  
(Computational graph resolves it).

# Inference BBVI

- $\mathcal{L}(\lambda, \xi) = \mathbb{E}_{q(z, W; \xi, \lambda)} [\ln p(x, z, W) - \ln q(z, W)]$
- Intractable to compute estimation! But we "only" need to compute derivative!
- $\nabla_{\lambda, \xi} \mathcal{L}(\lambda, \xi) = \mathbb{E}_q [\nabla_{\lambda, \xi} \ln q (\ln p(x, W, z) - \ln q(W, z))]$   
 $\approx \frac{1}{S} \sum_{i=1}^S \nabla_{\lambda, \xi} \ln q(z_s, W_s | \lambda, \xi) (\ln p(x, W_s, z_s) - \ln q(W_s, z_s | \lambda, \xi))$   
 $z_s, W_s \sim q(z, W | \lambda, \xi)$
- Each row in the jacobian is actually sparse (for this graphical model).  
The dependency on  $z_{n,l,k}$  is only through its markov blanket.  
(Computational graph resolves it).



- Coordinate gradients:

$$\nabla_{\lambda_{n,l,k}} \mathcal{L} = \mathbb{E}_q[\nabla_{\lambda_{n,l,k}} \ln q(z_{n,l,k})(\ln p_{n,l,k}(x, z, W) - \ln q(z_{n,l,k}))]$$

- Markov Blanket of  $z_{n,l,k}$ :

$$\ln p_{n,l,k}(x, z, W) = \ln p(z_{n,l,k} | z_{n,l+1}, W_{l,k}) + \ln p(z_{n,l-1} | z_{n,l}, W_{l-1})$$

---

**Algorithm 1** BBVI for DEFs

---

**Input:** data  $X$ , model  $p$ ,  $L$  layers.

**Initialize**  $\lambda, \xi$  randomly,  $t = 1$ .

**repeat**

    Sample a datapoint  $x$

**for**  $s = 1$  to  $S$  **do**

$z_x[s], W[s] \sim q$

$p[s] = \log p(z_x[s], W[s], x)$

$q[s] = \log q(z_x[s], W[s])$

$g[s] = \nabla \log q(z_x[s], W[s])$

**end for**

    Compute gradient using BBVI

    Update variational parameters for  $z$  and  $W$

**until** change in validation likelihood is small

---

- Hierarchical Clustering of words into topics, groups of topics etc.

# Text Modelling

- Hierarchical Clustering of words into topics, groups of topics etc.
- N Documents  $\{x_n\}_{n=1}^N$ .  
V-dimensional  $x_{n,i} = \#$  term i in doc n. (Observable)

# Text Modelling

- Hierarchical Clustering of words into topics, groups of topics etc.
- N Documents  $\{x_n\}_{n=1}^N$ .  
V-dimensional  $x_{n,i} = \#$  term i in doc n. (Observable)
- Observation Likelihood:  $p(x_{n,i} | \mathbf{z}_{n,1}, \mathbf{w}_{0,i}) = \text{Poisson}(g(\mathbf{z}_{n,1}^\top \mathbf{w}_{0,i}))$



# Text Modelling

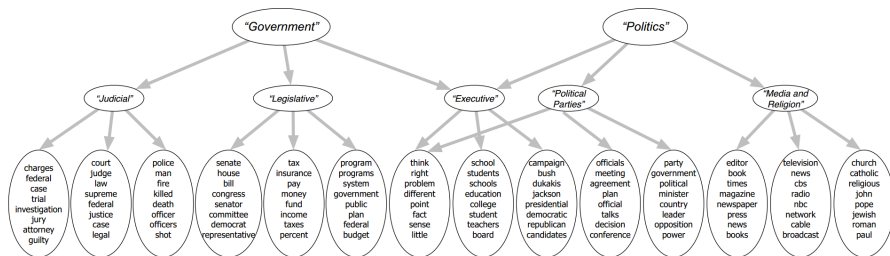
- Hierarchical Clustering of words into topics, groups of topics etc.
- N Documents  $\{x_n\}_{n=1}^N$ .  
V-dimensional  $x_{n,i} = \#$  term i in doc n. (Observable)
- Observation Likelihood:  $p(x_{n,i} | \mathbf{z}_{n,1}, \mathbf{w}_{0,i}) = \text{Poisson}(g(\mathbf{z}_{n,1}^\top \mathbf{w}_{0,i}))$
- $[\mathbf{W}_0]_{i,j} \sim \text{Gamma}(\alpha, \beta) \implies \mathbf{W}_0$  puts positive mass on groups of terms : "topics"!

# Text Modelling

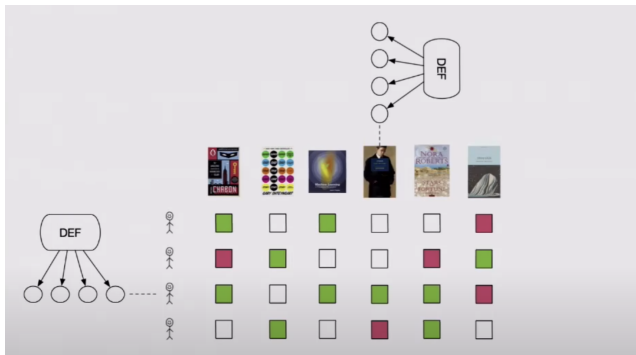
- Hierarchical Clustering of words into topics, groups of topics etc.
- N Documents  $\{x_n\}_{n=1}^N$ .  
V-dimensional  $x_{n,i} = \#$  term i in doc n. (Observable)
- Observation Likelihood:  $p(x_{n,i} | \mathbf{z}_{n,1}, \mathbf{w}_{0,i}) = \text{Poisson}(g(\mathbf{z}_{n,1}^\top \mathbf{w}_{0,i}))$
- $[\mathbf{W}_0]_{i,j} \sim \text{Gamma}(\alpha, \beta) \implies \mathbf{W}_0$  puts positive mass on groups of terms : "topics"!
- "topics"  $\mathbf{z}_{n,1,k} = \#$  topic k in document n
- "super topics"  $\mathbf{z}_{n,2,k} = \#$  super topic k in document n.
- "concepts"  $\mathbf{z}_{n,3,k}$ , etc.

- Hierarchical Clustering of words into topics, groups of topics etc.
- N Documents  $\{x_n\}_{n=1}^N$ .  
V-dimensional  $x_{n,i} = \#$  term i in doc n. (Observable)
- Observation Likelihood:  $p(x_{n,i} | \mathbf{z}_{n,1}, \mathbf{w}_{0,i}) = \text{Poisson}(g(\mathbf{z}_{n,1}^\top \mathbf{w}_{0,i}))$
- $[\mathbf{W}_0]_{i,j} \sim \text{Gamma}(\alpha, \beta) \implies \mathbf{W}_0$  puts positive mass on groups of terms : "topics"!
- "topics"  $z_{n,1,k} = \#$  topic k in document n
- "super topics"  $z_{n,2,k} = \#$  super topic k in document n.
- "concepts"  $z_{n,3,k}$ , etc.
- $p(\mathbf{z}_{n,1} | \mathbf{z}_{n,2}, \mathbf{W}_1) =$  "distribution of topics in a document given the super-topics in the same document". Bernoulli, Sparse Gamma, Poisson etc.

# Text Modelling



# Double DEF



- Draw user preferences:  $\theta_i \sim DEF()$
- Draw item attributes:  $\beta_j \sim DEF()$
- Draw rating:  $y_{i,j} = f(\theta_i^\top \beta_j)$

# Table of Contents

- 1 Exponential Families
  - Examples
  - Counter Examples
  - Sufficiency
  - Conjugacy
- 2 Black-box VI
- 3 Going Deep with Exponential Families
- 4 Discussion

# Discussion

- Differences with Deep GPs
- BBVI in Probabilistic Programming.
- Directed versus Undirected Models (Explaining away)
- Survival Analysis

# Bibliography



Pratik Chaudhari and Stefano Soatto. *Stochastic gradient descent performs variational inference, converges to limit cycles for deep networks*. 2018. arXiv: 1710.11029 [cs.LG].



Diederik P Kingma and Max Welling. *Auto-Encoding Variational Bayes*. 2014. arXiv: 1312.6114 [stat.ML].



Horia Mania, Aurelia Guy, and Benjamin Recht. *Simple random search provides a competitive approach to reinforcement learning*. 2018. arXiv: 1803.07055 [cs.LG].



Rajesh Ranganath, Sean Gerrish, and David Blei. “Black Box Variational Inference”. In: ed. by Samuel Kaski and Jukka Corander. Vol. 33. *Proceedings of Machine Learning Research*. Reykjavik, Iceland: PMLR, 22–25 Apr 2014, pp. 814–822. URL: <http://proceedings.mlr.press/v33/ranganath14.html>

