

# Variational inference

## Scalable solutions

Helge Langseth and Thomas Dyhre Nielsen

# Introduction

## Day 1: Bayesian networks – Definition and inference

- Definition of Bayesian networks: Syntax and semantics
- Exact inference
- Approximate inference using MCMC

## Day 2: Variational inference – Introduction and basis

- Approximate inference through the *Kullback-Leibler divergence*
- *Variational Bayes*
- The *mean-field* approach to Variational Bayes

## Day 3: Variational Bayes – cont'd

- Solving the VB equations
- Introducing Exponential families

## Day 4: Scalable Variational Bayes

- Variational message passing
- Stochastic gradient ascent
- Stochastic variational inference

## Day 5: Current approaches and extensions

- Variational Auto Encoders
- Black Box variational inference
- Probabilistic Programming Languages

## Algorithm:

- We have observed  $\mathbf{X} = \mathbf{x}$ , and have access to the full joint  $p(\mathbf{z}, \mathbf{x})$ .
- We posit a *variational family* of distributions  $q_j(\cdot \mid \boldsymbol{\lambda}_j)$ , i.e., we choose the distributional form, while wanting to optimize the parameterization  $\boldsymbol{\lambda}_j$ .
- The posterior approximation is assumed to factorize according to the mean-field assumption, and we use the KL ( $q(\mathbf{z}) \parallel p(\mathbf{z} \mid \mathbf{x})$ ) as our objective.

## Algorithm:

Repeat until negligible improvement in terms of  $\mathcal{L}(q)$ :

- For each  $j$ :
  - Calculate  $\mathbb{E}_{q_{-j}} [\log p(\mathbf{z}, \mathbf{x})]$  using current estimates for  $q_i(\cdot \mid \boldsymbol{\lambda}_i)$ ,  $i \neq j$ .
  - Choose  $\boldsymbol{\lambda}_j$  so that  $q_j(z_j \mid \boldsymbol{\lambda}_j) \propto \exp(\mathbb{E}_{q_{-j}} [\log p(\mathbf{z}, \mathbf{x})])$ .
- Calculate the new  $\mathcal{L}(q)$ .

As we realized last time, calculations of  $\mathbb{E}_{q_{-j}} [\log p(\mathbf{z}, \mathbf{x})]$  and  $\mathcal{L}(q)$  are quite tedious – and apparently must be done separately for each model we make.

This **harms the applicability** of variational inference, even under the **quite restrictive** mean field assumption.

# The Exponential Family

## Definition of ExpFam models

Consider a distribution  $f_{\mathbf{x}}(\mathbf{x} \mid \boldsymbol{\theta})$ , and assume it can be written as:

$$f_{\mathbf{x}}(\mathbf{x} \mid \boldsymbol{\eta}) = \exp \left( h(\mathbf{x}) + \boldsymbol{\eta}^T \mathbf{t}(\mathbf{x}) - A(\boldsymbol{\eta}) \right)$$

Here  $h(\mathbf{x})$  is the log base measure,  $\boldsymbol{\eta}$  the natural parameters,  $\mathbf{t}(\mathbf{x})$  the sufficient statistics, and  $A(\boldsymbol{\eta})$  the log partition function.

## Positives with the Exponential Family:

- It is the only family of distributions with finite-sized sufficient statistics\*;
- It is the only family of distributions that has conjugate priors;
- It simplifies the operations of variational inference;
- It has simple mathematical procedures for calculating moments, MLEs, Bayesian posteriors, ...

\* Under certain regularity conditions...

## Examples

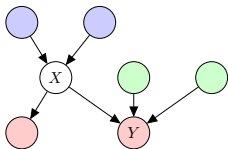
Distributions that are in the ExpFam family of distributions, and therefore have a shared overarching theory include Bernoulli, Gamma, Normal, Poisson,  $\chi^2$ , Beta, Dirichlet, Categorical, Multinomial, and Wishart.

## Variational message passing

# The Conjugate Exponential Family (abridged version)

Referring to the conditioning parameters as parents, consider a conditional distribution in exponential form

$$\log f(x \mid \text{pa}(x)) = h_x(x) + \boldsymbol{\eta}_X(\text{pa}(x))^\top \mathbf{t}_x(x) - A_x(\text{pa}(x))$$



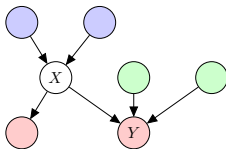
# The Conjugate Exponential Family (abridged version)

Referring to the conditioning parameters as parents, consider a conditional distribution in exponential form

$$\log f(x \mid \text{pa}(x)) = h_x(x) + \boldsymbol{\eta}_X(\text{pa}(x))^{\top} \mathbf{t}_x(x) - A_x(\text{pa}(x))$$

with a child  $y$  also in exponential form:

$$\log f(y \mid x, \text{cp}(x)) = h_y(y) + \boldsymbol{\eta}_y(x, \text{cp}(x))^{\top} \mathbf{t}_y(y) - A_y(x, \text{cp}(x))$$





# The Conjugate Exponential Family (abridged version)

Referring to the conditioning parameters as parents, consider a conditional distribution in exponential form

$$\log f(x \mid \text{pa}(x)) = h_x(x) + \boldsymbol{\eta}_X(\text{pa}(x))^\top \mathbf{t}_x(x) - A_x(\text{pa}(x))$$

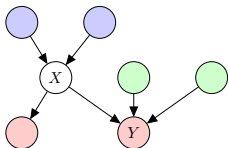
with a child  $y$  also in exponential form:

$$\log f(y \mid x, \text{cp}(x)) = h_y(y) + \boldsymbol{\eta}_y(x, \text{cp}(x))^\top \mathbf{t}_y(y) - A_y(x, \text{cp}(x))$$

## Conjugacy

Conjugacy requires that  $\log f(x \mid \text{pa}(x))$  and  $\log f(y \mid x, \text{cp}(x))$  have the same functional form wrt.  $x$  and so the latter can be written as:

$$\log f(y \mid x, \text{cp}(x)) = \boldsymbol{\eta}_{xy}(y, \text{cp}(x))^\top \mathbf{t}_x(x) - A_{xy}(y, \text{cp}(x))$$



# The Conjugate Exponential Family (abridged version)

Referring to the conditioning parameters as parents, consider a conditional distribution in exponential form

$$\log f(x \mid \text{pa}(x)) = h_x(x) + \boldsymbol{\eta}_X(\text{pa}(x))^\top \mathbf{t}_x(x) - A_x(\text{pa}(x))$$

with a child  $y$  also in exponential form:

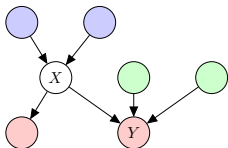
$$\log f(y \mid x, \text{cp}(x)) = h_y(y) + \boldsymbol{\eta}_y(x, \text{cp}(x))^\top \mathbf{t}_y(y) - A_y(x, \text{cp}(x))$$

## Conjugacy

Conjugacy requires that  $\log f(x \mid \text{pa}(x))$  and  $\log f(y \mid x, \text{cp}(x))$  have the same functional form wrt.  $x$  and so the latter can be written as:

$$\log f(y \mid x, \text{cp}(x)) = \boldsymbol{\eta}_{xy}(y, \text{cp}(x))^\top \mathbf{t}_x(x) - A_{xy}(y, \text{cp}(x))$$

- $\log f(y \mid x, \text{cp}(x))$  is linear in  $\mathbf{t}_x(x)$  and  $\mathbf{t}_y(y)$  (and in  $\mathbf{t}_z(z)$  for any  $z \in \text{cp}(x)$ ).
- $\rightsquigarrow \log f(y \mid x, \text{cp}(x))$  is a multi-linear function



Assume that  $X$  is normal distributed with mean  $m$  and precision  $b$ :

$$\log f(x | m, b) = \underbrace{-\frac{1}{2} \log(2\pi)}_{h(x)} + \underbrace{\begin{bmatrix} b \cdot m \\ -\frac{b}{2} \end{bmatrix}}_{\boldsymbol{\eta}(m,b)^\top} \underbrace{\begin{bmatrix} x \\ x^2 \end{bmatrix}}_{\mathbf{t}(x)} - \underbrace{\left( \frac{b \cdot m^2}{2} - \frac{1}{2} \log b \right)}_{A(\boldsymbol{\eta}(m,b))}$$

Assume that  $X$  is normal distributed with mean  $m$  and precision  $b$ :

$$\log f(x | m, b) = \underbrace{-\frac{1}{2} \log(2\pi)}_{h(x)} + \underbrace{\begin{bmatrix} b \cdot m \\ -\frac{b}{2} \end{bmatrix}}_{\boldsymbol{\eta}(m,b)^\top} \underbrace{\begin{bmatrix} x \\ x^2 \end{bmatrix}}_{\mathbf{t}(x)} - \underbrace{\left( \frac{b \cdot m^2}{2} - \frac{1}{2} \log b \right)}_{A(\boldsymbol{\eta}(m,b))}$$

## Normal prior for mean

We can then express  $f(x | m, b)$  in terms of  $m$  as

$$\log f(x | m, b) = -\frac{1}{2} \log(2\pi) + \begin{bmatrix} b \cdot x \\ -\frac{b}{2} \end{bmatrix}^\top \begin{bmatrix} m \\ m^2 \end{bmatrix} - \left( \frac{b \cdot x^2}{2} - \frac{1}{2} \log b \right)$$

Thus, conjugacy implies that the prior distribution over  $m$  is a normal distribution.

Assume that  $X$  is normal distributed with mean  $m$  and precision  $b$ :

$$\log f(x | m, b) = \underbrace{-\frac{1}{2} \log(2\pi)}_{h(x)} + \underbrace{\begin{bmatrix} b \cdot m \\ -\frac{b}{2} \end{bmatrix}}_{\boldsymbol{\eta}(m,b)^\top} \underbrace{\begin{bmatrix} x \\ x^2 \end{bmatrix}}_{\mathbf{t}(x)} - \underbrace{\left( \frac{b \cdot m^2}{2} - \frac{1}{2} \log b \right)}_{A(\boldsymbol{\eta}(m,b))}$$

## Normal prior for mean

We can then express  $f(x | m, b)$  in terms of  $m$  as

$$\log f(x | m, b) = -\frac{1}{2} \log(2\pi) + \begin{bmatrix} b \cdot x \\ -\frac{b}{2} \end{bmatrix}^\top \begin{bmatrix} m \\ m^2 \end{bmatrix} - \left( \frac{b \cdot x^2}{2} - \frac{1}{2} \log b \right)$$

Thus, conjugacy implies that the prior distribution over  $m$  is a normal distribution.

## Gamma prior for precision

We can also express  $f(x | m, b)$  in terms of  $b$ :

$$\log f(x | m, b) = -\frac{1}{2} \log(2\pi) + \begin{bmatrix} -\frac{1}{2}(x - m)^2 \\ \frac{1}{2} \end{bmatrix}^\top \begin{bmatrix} b \\ \log(b) \end{bmatrix} - 0$$

Thus, implying that  $b$  follows a gamma distribution.

## Variational message passing

Variational message passing is a variational-based inference algorithm for Bayesian networks:

- Applies to conjugate-exponential models (can be generalized to non-conjugate models)
- Works by sending local messages between nodes in the graphical model.

### Advantage

No need for the tedious model-specific variational derivations that we saw during the last lecture.

Given a Bayesian network that factorizes according to

$$p(x_1, \dots, x_n) = \prod_{i=1}^N p(x_i \mid \text{pa}(x_i))$$

we saw last time that the mean-field approximation for a variable  $X_j$  is given by

$$\begin{aligned} \log q(x_j) &= \mathbb{E} \left[ \sum_{i=1}^N \log p(x_i \mid \text{pa}(x_i)) \right] + c \\ &= \mathbb{E} [\log p(x_j \mid \text{pa}(x_j))] + \mathbb{E} \left[ \sum_{y \in \text{ch}(x_j)} \log p(y \mid x_j, \text{cp}(x_j)) \right] + c' \end{aligned}$$

Given a Bayesian network that factorizes according to

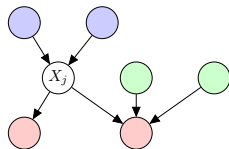
$$p(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i \mid \text{pa}(x_i))$$

we saw last time that the mean-field approximation for a variable  $X_j$  is given by

$$\log q(x_j) = \mathbb{E}[\log p(x_j \mid \text{pa}(x_j))] + \mathbb{E}\left[\sum_{y \in \text{ch}(x_j)} \log p(y \mid x_j, \text{cp}(x_j))\right] + c'$$

Therefore, the only contributions to  $q(x_j)$  come from

- $X_j$ 's **parents** through the term  $\mathbb{E}[\log p(x_j \mid \text{pa}(x_j))]$ .
- $X_j$ 's **children** through the term  $\mathbb{E}\left[\sum_{y \in \text{ch}(x_j)} \log p(y \mid \text{pa}(y))\right]$ .
- $X_j$ 's **co-parents** through the term  $\mathbb{E}\left[\sum_{y \in \text{ch}(x_i)} \log p(y \mid x_j, \text{cp}(x_j))\right]$ .



This is exactly **the Markov Blanket** for  $X_j$ , which we also exploited during Gibbs sampling and our previous VI derivations.



The distributions involved in the mean-field approximation for variable  $X_j$

$$\log q(x_j) = \mathbb{E}[\log p(x_j \mid \text{pa}(x_j))] + \mathbb{E}\left[\sum_{y \in \text{ch}(x_j)} \log p(y \mid x_j, \text{cp}(x_j))\right] + c'$$

can be expressed in exponential form:

- $\log p(x_j \mid \text{pa}(x_j)) = h_{X_j}(x_j) + \boldsymbol{\eta}_{X_j}(\text{pa}(x_j))^T \mathbf{t}(x_j) - A_{X_j}(\boldsymbol{\eta}_{X_j}(\text{pa}(x_j)))$
- $\log p(y \mid x_j, \text{cp}(x_j)) = h_Y(y) + \boldsymbol{\eta}_Y(x_j, \text{cp}(x_j))^T \mathbf{t}(y) - A_Y(\boldsymbol{\eta}_Y(x_j, \text{cp}(x_j)))$

The distributions involved in the mean-field approximation for variable  $X_j$

$$\log q(x_j) = \mathbb{E}[\log p(x_j \mid \text{pa}(x_j))] + \mathbb{E}\left[\sum_{\mathbf{y} \in \text{ch}(x_j)} \log p(\mathbf{y} \mid x_j, \text{cp}(x_j))\right] + c'$$

can be expressed in exponential form:

- $\log p(x_j \mid \text{pa}(x_j)) = h_{X_j}(x_j) + \boldsymbol{\eta}_{X_j}(\text{pa}(x_j))^T \mathbf{t}(x_j) - A_{X_j}(\boldsymbol{\eta}_{X_j}(\text{pa}(x_j)))$
- $\log p(\mathbf{y} \mid x_j, \text{cp}(x_j)) = h_Y(\mathbf{y}) + \boldsymbol{\eta}_Y(x_j, \text{cp}(x_j))^T \mathbf{t}(\mathbf{y}) - A_Y(\boldsymbol{\eta}_{X_j}(x_j, \text{cp}(x_j)))$

Due to conjugacy, we can rewrite  $\log p(\mathbf{y} \mid x_j, \text{cp}(x_j))$  in terms of  $\mathbf{t}(x_j)$ :

$$\log p(\mathbf{y} \mid x_j, \text{cp}(x_j)) = \boldsymbol{\eta}_{X_j, Y}(\mathbf{y}, \text{cp}(x_j))^T \mathbf{t}(x_j) - A_{X_j, Y}(\boldsymbol{\eta}_{X_j, Y}(\mathbf{y}, \text{cp}(x_j)))$$

Thus, we end up with:

$$\begin{aligned} \log q(x_j) = & \mathbb{E}\left[h_{X_j}(x_j) + \boldsymbol{\eta}_{X_j}(\text{pa}(x_j))^T \mathbf{t}(x_j) - A_{X_j}(\boldsymbol{\eta}_{X_j}(\text{pa}(x_j)))\right] \\ & + \mathbb{E}\left[\sum_{\mathbf{Y} \in \text{ch}(x_j)} \boldsymbol{\eta}_{X_j, Y}(\mathbf{Y}, \text{cp}(x_j))^T \mathbf{t}(x_j) - A_{X_j, Y}(\boldsymbol{\eta}_{X_j, Y}(\mathbf{Y}, \text{cp}(x_j)))\right] + c' \end{aligned}$$

The distributions involved in the mean-field approximation for variable  $X_j$

$$\log q(x_j) = \mathbb{E}[\log p(x_j \mid \text{pa}(x_j))] + \mathbb{E}\left[\sum_{\mathbf{y} \in \text{ch}(x_j)} \log p(\mathbf{y} \mid x_j, \text{cp}(x_j))\right] + c'$$

can be expressed in exponential form:

- $\log p(x_j \mid \text{pa}(x_j)) = h_{X_j}(x_j) + \boldsymbol{\eta}_{X_j}(\text{pa}(x_j))^T \mathbf{t}(x_j) - A_{X_j}(\boldsymbol{\eta}_{X_j}(\text{pa}(x_j)))$
- $\log p(\mathbf{y} \mid x_j, \text{cp}(x_j)) = h_Y(\mathbf{y}) + \boldsymbol{\eta}_Y(x_j, \text{cp}(x_j))^T \mathbf{t}(\mathbf{y}) - A_Y(\boldsymbol{\eta}_{X_j}(x_j, \text{cp}(x_j)))$

Due to conjugacy, we can rewrite  $\log p(\mathbf{y} \mid x_j, \text{cp}(x_j))$  in terms of  $\mathbf{t}(x_j)$ :

$$\log p(\mathbf{y} \mid x_j, \text{cp}(x_j)) = \boldsymbol{\eta}_{X_j, Y}(\mathbf{y}, \text{cp}(x_j))^T \mathbf{t}(x_j) - A_{X_j, Y}(\boldsymbol{\eta}_{X_j, Y}(\mathbf{y}, \text{cp}(x_j)))$$

Rearranging and absorbing into constant  $c$  gives:

$$\log q(x_j) = \left[ \mathbb{E}(\boldsymbol{\eta}_{X_j}(\text{pa}(x_j))) + \sum_{\mathbf{Y} \in \text{ch}(x_j)} \mathbb{E}(\boldsymbol{\eta}_{X_j, Y}(\mathbf{Y}, \text{cp}(x_j))) \right]^T \mathbf{t}(x_j) + h_{X_j}(x_j) + c$$

We have:

$$\log q(x_j) = \left[ \mathbb{E}(\boldsymbol{\eta}_{X_j}(\text{pa}(x_j))) + \sum_{Y \in \text{ch}(x_j)} \mathbb{E}(\boldsymbol{\eta}_{X_j, Y}(\textcolor{red}{Y}, \text{cp}(x_j))) \right]^T \mathbf{t}(x_j) + h_{X_j}(x_j) + c$$

As seen before, both  $\boldsymbol{\eta}_{X_j}(\text{pa}(x_j))$  and each  $\boldsymbol{\eta}_{X_j, Y}(\textcolor{red}{y}, \text{cp}(x_j))$  are multi-linear functions of the natural statistics vectors of their dependent variables, hence:

$$\begin{aligned} \mathbb{E}(\boldsymbol{\eta}_{X_j}(\text{pa}(x_j))) &= \tilde{\boldsymbol{\eta}}_{X_j}(\{\mathbb{E}(\mathbf{t}(\textcolor{blue}{X}_i))\}_{X_i \in \text{pa}(X_j)}) \\ \mathbb{E}(\boldsymbol{\eta}_{X_j, Y}(\textcolor{red}{y}, \text{cp}(x_j))) &= \tilde{\boldsymbol{\eta}}_{X_j, Y}(\mathbb{E}(\mathbf{t}(\textcolor{red}{Y})), \{\mathbb{E}(\mathbf{t}(\textcolor{green}{X}_i))\}_{X_i \in \text{cp}(X_j)}). \end{aligned}$$

We have:

$$\log q(x_j) = \left[ \mathbb{E}(\boldsymbol{\eta}_{X_j}(\text{pa}(x_j))) + \sum_{Y \in \text{ch}(x_j)} \mathbb{E}(\boldsymbol{\eta}_{X_j, Y}(\textcolor{red}{Y}, \text{cp}(x_j))) \right]^T \mathbf{t}(x_j) + h_{X_j}(x_j) + c$$

As seen before, both  $\boldsymbol{\eta}_{X_j}(\text{pa}(x_j))$  and each  $\boldsymbol{\eta}_{X_j, Y}(\textcolor{red}{y}, \text{cp}(x_j))$  are multi-linear functions of the natural statistics vectors of their dependent variables, hence:

$$\begin{aligned} \mathbb{E}(\boldsymbol{\eta}_{X_j}(\text{pa}(x_j))) &= \tilde{\boldsymbol{\eta}}_{X_j}(\{\mathbb{E}(\mathbf{t}(\textcolor{blue}{X}_i))\}_{X_i \in \text{pa}(X_j)}) \\ \mathbb{E}(\boldsymbol{\eta}_{X_j, Y}(\textcolor{red}{y}, \text{cp}(x_j))) &= \tilde{\boldsymbol{\eta}}_{X_j, Y}(\mathbb{E}(\mathbf{t}(\textcolor{red}{Y})), \{\mathbb{E}(\mathbf{t}(\textcolor{green}{X}_i))\}_{X_i \in \text{cp}(X_j)}). \end{aligned}$$

The expectations over the individual natural statistics can be found using the log-normalizer trick:

$$\nabla A_X(\boldsymbol{\eta}) = \mathbb{E}(\mathbf{t}(X)).$$

# Calculating the expectations in VMP

We have:

$$\log q(x_j) = \left[ \mathbb{E}(\boldsymbol{\eta}_{X_j}(\text{pa}(x_j))) + \sum_{Y \in \text{ch}(x_j)} \mathbb{E}(\boldsymbol{\eta}_{X_j,Y}(\mathbf{Y}, \text{cp}(x_j))) \right]^\top \mathbf{t}(x_j) + h_{X_j}(x_j) + c$$

As seen before, both  $\boldsymbol{\eta}_{X_j}(\text{pa}(x_j))$  and each  $\boldsymbol{\eta}_{X_j,Y}(\mathbf{y}, \text{cp}(x_j))$  are multi-linear functions of the natural statistics vectors of their dependent variables, hence:

$$\begin{aligned} \mathbb{E}(\boldsymbol{\eta}_{X_j}(\text{pa}(x_j))) &= \tilde{\boldsymbol{\eta}}_{X_j}(\{\mathbb{E}(\mathbf{t}(X_i))\}_{X_i \in \text{pa}(X_j)}) \\ \mathbb{E}(\boldsymbol{\eta}_{X_j,Y}(\mathbf{y}, \text{cp}(x_j))) &= \tilde{\boldsymbol{\eta}}_{X_j,Y}(\mathbb{E}(\mathbf{t}(\mathbf{Y})), \{\mathbb{E}(\mathbf{t}(X_i))\}_{X_i \in \text{cp}(X_j)}). \end{aligned}$$

The expectations over the individual natural statistics can be found using the log-normalizer trick:

$$\nabla A_X(\boldsymbol{\eta}) = \mathbb{E}(\mathbf{t}(X)).$$

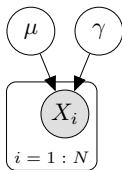
## New natural parameter vector and implied message passing scheme

$$\boldsymbol{\eta}_{X_j}^* = \tilde{\boldsymbol{\eta}}_{X_j}(\underbrace{\{\mathbb{E}(\mathbf{t}(X_i))\}_{X_i \in \text{pa}(X_j)}}_{\text{Messages from parents: } \mathbf{m}_{X_i \rightarrow X_j}}) + \sum_{Y \in \text{ch}(X_j)} \underbrace{\tilde{\boldsymbol{\eta}}_{X_j,Y}(\mathbb{E}(\mathbf{t}(\mathbf{Y})), \{\mathbb{E}(\mathbf{t}(X_i))\}_{X_i \in \text{cp}(X_j)})}_{\text{Messages from children: } \mathbf{m}_{Y \rightarrow X_j}}$$

## VMP algorithm

- ➊ Initialize each variational distribution  $q(x_i)$  by its moment vector  $\mathbb{E}(\mathbf{t}(x_i))$ .
- ➋ For each variable  $X_i$ :
  - ➊ Retrieve messages from all parents and children.
  - ➋ Computed new natural parameter vector  $\boldsymbol{\eta}_{x_i}^*$ .
  - ➌ Computed new moment parameters  $\mathbb{E}(\mathbf{t}(x_i))$ .
- ➌ Calculate ELBO if needed (not described here!)
- ➍ Repeat from step 1 unless termination criteria reached.

## Model specification



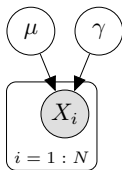
$$\log p(\gamma \mid \alpha, \beta) = \begin{bmatrix} -\beta \\ \alpha - 1 \end{bmatrix}^{\top} \begin{bmatrix} \gamma \\ \log(\gamma) \end{bmatrix} + c_1$$

$$\log p(\mu \mid m, b) = \begin{bmatrix} b \cdot m \\ -\frac{b}{2} \end{bmatrix}^{\top} \begin{bmatrix} \mu \\ \mu^2 \end{bmatrix} + c_2$$

$$\begin{aligned} \log p(x_i \mid \mu, \gamma) &= \begin{bmatrix} \gamma\mu \\ -\frac{\gamma}{2} \end{bmatrix}^{\top} \begin{bmatrix} x_i \\ x_i^2 \end{bmatrix} + c_3 \\ &= \begin{bmatrix} \gamma x_i \\ -\frac{\gamma}{2} \end{bmatrix}^{\top} \begin{bmatrix} \mu \\ \mu^2 \end{bmatrix} + c_4 \\ &= \begin{bmatrix} -\frac{1}{2}x_i^2 + x_i\mu - \frac{1}{2}\mu^2 \\ \frac{1}{2} \end{bmatrix}^{\top} \begin{bmatrix} \gamma \\ \log(\gamma) \end{bmatrix} + c_5 \end{aligned}$$



## Model specification

Updating  $q(\mu)$ 

- 1 Calc. message from co-parent ( $\gamma$ ):

$$\mathbf{m}_{\gamma \rightarrow x_i} = \begin{bmatrix} \mathbb{E}(\gamma) \\ \mathbb{E}(\log(\gamma)) \end{bmatrix}$$

$$\log p(\gamma | \alpha, \beta) = \begin{bmatrix} -\beta \\ \alpha - 1 \end{bmatrix}^\top \begin{bmatrix} \gamma \\ \log(\gamma) \end{bmatrix} + c_1$$

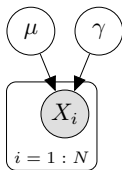
$$\log p(\mu | m, b) = \begin{bmatrix} b \cdot m \\ -\frac{b}{2} \end{bmatrix}^\top \begin{bmatrix} \mu \\ \mu^2 \end{bmatrix} + c_2$$

$$\log p(x_i | \mu, \gamma) = \begin{bmatrix} \gamma\mu \\ -\frac{\gamma}{2} \end{bmatrix}^\top \begin{bmatrix} x_i \\ x_i^2 \end{bmatrix} + c_3$$

$$= \begin{bmatrix} \gamma x_i \\ -\frac{\gamma}{2} \end{bmatrix}^\top \begin{bmatrix} \mu \\ \mu^2 \end{bmatrix} + c_4$$

$$= \begin{bmatrix} -\frac{1}{2}x_i^2 + x_i\mu - \frac{1}{2}\mu_i^2 \\ \frac{1}{2} \end{bmatrix}^\top \begin{bmatrix} \gamma \\ \log(\gamma) \end{bmatrix} + c_5$$

## Model specification



$$\log p(\gamma \mid \alpha, \beta) = \begin{bmatrix} -\beta \\ \alpha - 1 \end{bmatrix}^\top \begin{bmatrix} \gamma \\ \log(\gamma) \end{bmatrix} + c_1$$

$$\log p(\mu \mid m, b) = \begin{bmatrix} b \cdot m \\ -\frac{b}{2} \end{bmatrix}^\top \begin{bmatrix} \mu \\ \mu^2 \end{bmatrix} + c_2$$

$$\begin{aligned} \log p(x_i \mid \mu, \gamma) &= \begin{bmatrix} \gamma\mu \\ -\frac{\gamma}{2} \end{bmatrix}^\top \begin{bmatrix} x_i \\ x_i^2 \end{bmatrix} + c_3 \\ &= \begin{bmatrix} \gamma x_i \\ -\frac{\gamma}{2} \end{bmatrix}^\top \begin{bmatrix} \mu \\ \mu^2 \end{bmatrix} + c_4 \\ &= \begin{bmatrix} -\frac{1}{2}x_i^2 + x_i\mu - \frac{1}{2}\mu_i^2 \\ \frac{1}{2} \end{bmatrix}^\top \begin{bmatrix} \gamma \\ \log(\gamma) \end{bmatrix} + c_5 \end{aligned}$$

Updating  $q(\mu)$ 

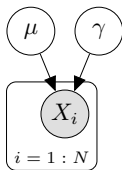
- 1 Calc. message from co-parent ( $\gamma$ ):

$$\mathbf{m}_{\gamma \rightarrow x_i} = \begin{bmatrix} \mathbb{E}(\gamma) \\ \mathbb{E}(\log(\gamma)) \end{bmatrix}$$

- 2 Send messages from all  $x_i$ :

$$\mathbf{m}_{x_i \rightarrow \mu} = \begin{bmatrix} \mathbb{E}(\gamma)x_i \\ -\frac{1}{2} \mathbb{E}(\gamma) \end{bmatrix}$$

## Model specification



$$\log p(\gamma \mid \alpha, \beta) = \begin{bmatrix} -\beta \\ \alpha - 1 \end{bmatrix}^\top \begin{bmatrix} \gamma \\ \log(\gamma) \end{bmatrix} + c_1$$

$$\log p(\mu \mid m, b) = \begin{bmatrix} b \cdot m \\ -\frac{b}{2} \end{bmatrix}^\top \begin{bmatrix} \mu \\ \mu^2 \end{bmatrix} + c_2$$

$$\begin{aligned} \log p(x_i \mid \mu, \gamma) &= \begin{bmatrix} \gamma \mu \\ -\frac{\gamma}{2} \end{bmatrix}^\top \begin{bmatrix} x_i \\ x_i^2 \end{bmatrix} + c_3 \\ &= \begin{bmatrix} \gamma x_i \\ -\frac{\gamma}{2} \end{bmatrix}^\top \begin{bmatrix} \mu \\ \mu^2 \end{bmatrix} + c_4 \\ &= \begin{bmatrix} -\frac{1}{2}x_i^2 + x_i\mu - \frac{1}{2}\mu^2 \\ \frac{1}{2} \end{bmatrix}^\top \begin{bmatrix} \gamma \\ \log(\gamma) \end{bmatrix} + c_5 \end{aligned}$$

Updating  $q(\mu)$ 

- 1 Calc. message from co-parent ( $\gamma$ ):

$$\mathbf{m}_{\gamma \rightarrow x_i} = \begin{bmatrix} \mathbb{E}(\gamma) \\ \mathbb{E}(\log(\gamma)) \end{bmatrix}$$

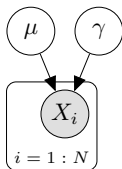
- 2 Send messages from all  $x_i$ :

$$\mathbf{m}_{x_i \rightarrow \mu} = \begin{bmatrix} \mathbb{E}(\gamma)x_i \\ -\frac{1}{2} \mathbb{E}(\gamma) \end{bmatrix}$$

- 3 Update natural parameter of  $q(\mu)$ :

$$\boldsymbol{\eta}_\mu^* = \begin{bmatrix} b \cdot m \\ -\frac{b}{2} \end{bmatrix} + \sum_{i=1}^N \mathbf{m}_{x_i \rightarrow \mu}$$

## Model specification

Updating  $q(\gamma)$ 

- 1 Calc. message from co-parent ( $\mu$ ):

$$\mathbf{m}_{\mu \rightarrow x_i} = \begin{bmatrix} \mathbb{E}(\mu) \\ \mathbb{E}(\mu^2) \end{bmatrix}$$

(based on updated parameters.)

$$\log p(\gamma | \alpha, \beta) = \begin{bmatrix} -\beta \\ \alpha - 1 \end{bmatrix}^\top \begin{bmatrix} \gamma \\ \log(\gamma) \end{bmatrix} + c_1$$

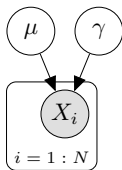
$$\log p(\mu | m, b) = \begin{bmatrix} b \cdot m \\ -\frac{b}{2} \end{bmatrix}^\top \begin{bmatrix} \mu \\ \mu^2 \end{bmatrix} + c_2$$

$$\log p(x_i | \mu, \gamma) = \begin{bmatrix} \gamma \mu \\ -\frac{\gamma}{2} \end{bmatrix}^\top \begin{bmatrix} x_i \\ x_i^2 \end{bmatrix} + c_3$$

$$= \begin{bmatrix} \gamma x_i \\ -\frac{\gamma}{2} \end{bmatrix}^\top \begin{bmatrix} \mu \\ \mu^2 \end{bmatrix} + c_4$$

$$= \begin{bmatrix} -\frac{1}{2}x_i^2 + x_i\mu - \frac{1}{2}\mu_i^2 \\ \frac{1}{2} \end{bmatrix}^\top \begin{bmatrix} \gamma \\ \log(\gamma) \end{bmatrix} + c_5$$

## Model specification



$$\log p(\gamma \mid \alpha, \beta) = \begin{bmatrix} -\beta \\ \alpha - 1 \end{bmatrix}^\top \begin{bmatrix} \gamma \\ \log(\gamma) \end{bmatrix} + c_1$$

$$\log p(\mu \mid m, b) = \begin{bmatrix} b \cdot m \\ -\frac{b}{2} \end{bmatrix}^\top \begin{bmatrix} \mu \\ \mu^2 \end{bmatrix} + c_2$$

$$\begin{aligned} \log p(x_i \mid \mu, \gamma) &= \begin{bmatrix} \gamma\mu \\ -\frac{\gamma}{2} \end{bmatrix}^\top \begin{bmatrix} x_i \\ x_i^2 \end{bmatrix} + c_3 \\ &= \begin{bmatrix} \gamma x_i \\ -\frac{\gamma}{2} \end{bmatrix}^\top \begin{bmatrix} \mu \\ \mu^2 \end{bmatrix} + c_4 \\ &= \begin{bmatrix} -\frac{1}{2}x_i^2 + x_i\mu - \frac{1}{2}\mu_i^2 \\ \frac{1}{2} \end{bmatrix}^\top \begin{bmatrix} \gamma \\ \log(\gamma) \end{bmatrix} + c_5 \end{aligned}$$

Updating  $q(\gamma)$ 

- 1 Calc. message from co-parent ( $\mu$ ):

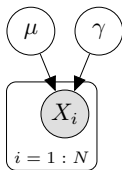
$$\mathbf{m}_{\mu \rightarrow x_i} = \begin{bmatrix} \mathbb{E}(\mu) \\ \mathbb{E}(\mu^2) \end{bmatrix}$$

(based on updated parameters.)

- 2 Send messages from all  $x_i$ :

$$\mathbf{m}_{x_i \rightarrow \gamma} = \begin{bmatrix} -\frac{1}{2}x_i^2 + x_i \mathbb{E}(\mu) - \frac{1}{2} \mathbb{E}(\mu^2) \\ \frac{1}{2} \end{bmatrix}$$

## Model specification



$$\log p(\gamma | \alpha, \beta) = \begin{bmatrix} -\beta \\ \alpha - 1 \end{bmatrix}^\top \begin{bmatrix} \gamma \\ \log(\gamma) \end{bmatrix} + c_1$$

$$\log p(\mu | m, b) = \begin{bmatrix} b \cdot m \\ -\frac{b}{2} \end{bmatrix}^\top \begin{bmatrix} \mu \\ \mu^2 \end{bmatrix} + c_2$$

$$\log p(x_i | \mu, \gamma) = \begin{bmatrix} \gamma \mu \\ -\frac{\gamma}{2} \end{bmatrix}^\top \begin{bmatrix} x_i \\ x_i^2 \end{bmatrix} + c_3$$

$$= \begin{bmatrix} \gamma x_i \\ -\frac{\gamma}{2} \end{bmatrix}^\top \begin{bmatrix} \mu \\ \mu^2 \end{bmatrix} + c_4$$

$$= \begin{bmatrix} -\frac{1}{2}x_i^2 + x_i\mu - \frac{1}{2}\mu^2 \\ \frac{1}{2} \end{bmatrix}^\top \begin{bmatrix} \gamma \\ \log(\gamma) \end{bmatrix} + c_5$$

Updating  $q(\gamma)$ 

- 1 Calc. message from co-parent ( $\mu$ ):

$$\mathbf{m}_{\mu \rightarrow x_i} = \begin{bmatrix} \mathbb{E}(\mu) \\ \mathbb{E}(\mu^2) \end{bmatrix}$$

(based on updated parameters.)

- 2 Send messages from all  $x_i$ :

$$\mathbf{m}_{x_i \rightarrow \gamma} = \begin{bmatrix} -\frac{1}{2}x_i^2 + x_i \mathbb{E}(\mu) - \frac{1}{2} \mathbb{E}(\mu^2) \\ \frac{1}{2} \end{bmatrix}$$

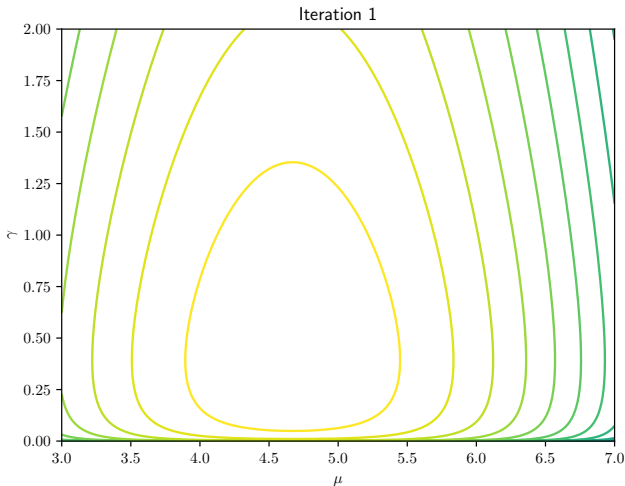
- 3 Update natural parameter of  $q(\gamma)$ :

$$\boldsymbol{\eta}_\gamma^* = \begin{bmatrix} -\beta \\ \alpha - 1 \end{bmatrix} + \sum_{i=1}^N \mathbf{m}_{x_i \rightarrow \gamma}$$

## Data

Four data points sampled from a normal distribution with mean 5 and variance 1.

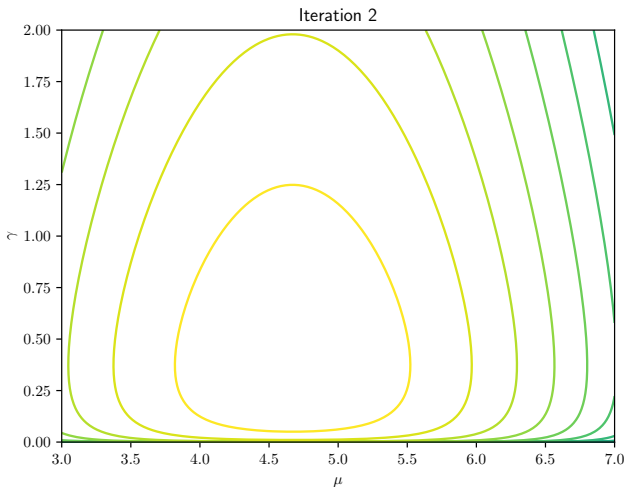
## Posteriors



## Data

Four data points sampled from a normal distribution with mean 5 and variance 1.

## Posteriors

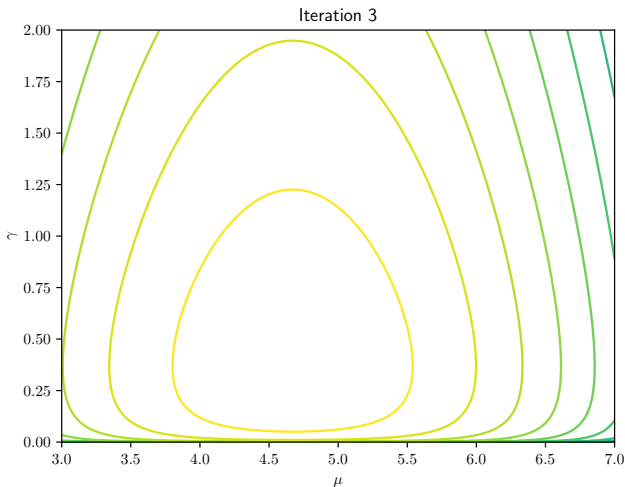




## Data

Four data points sampled from a normal distribution with mean 5 and variance 1.

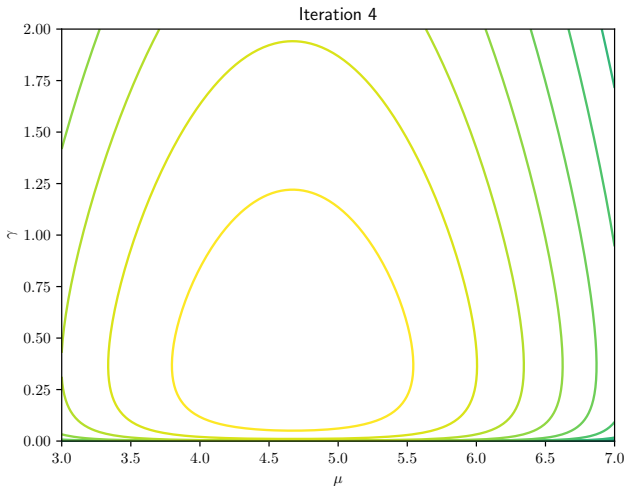
## Posteriors



## Data

Four data points sampled from a normal distribution with mean 5 and variance 1.

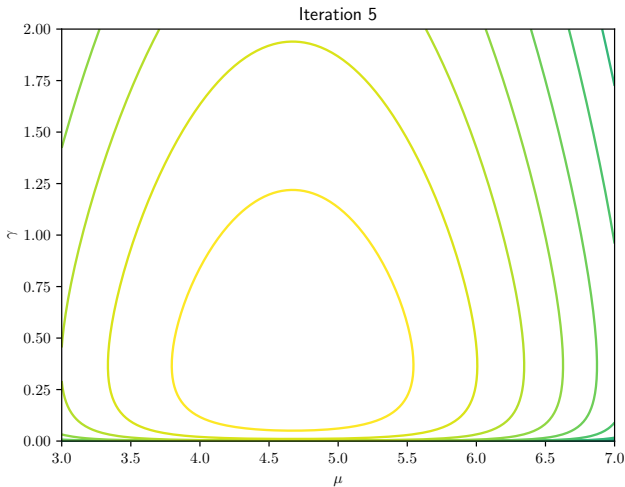
## Posteriors



## Data

Four data points sampled from a normal distribution with mean 5 and variance 1.

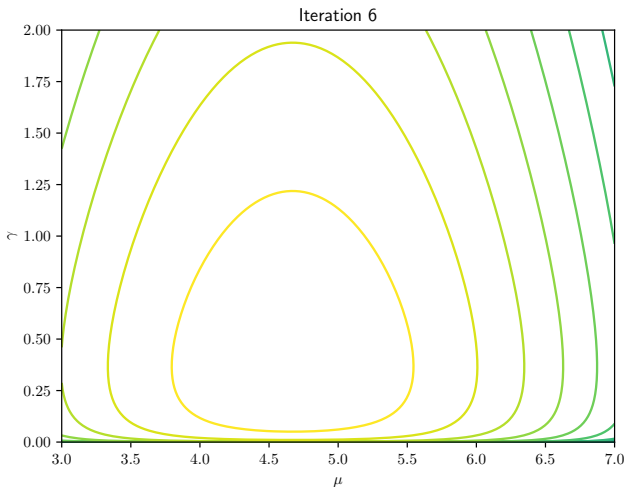
## Posteriors



## Data

Four data points sampled from a normal distribution with mean 5 and variance 1.

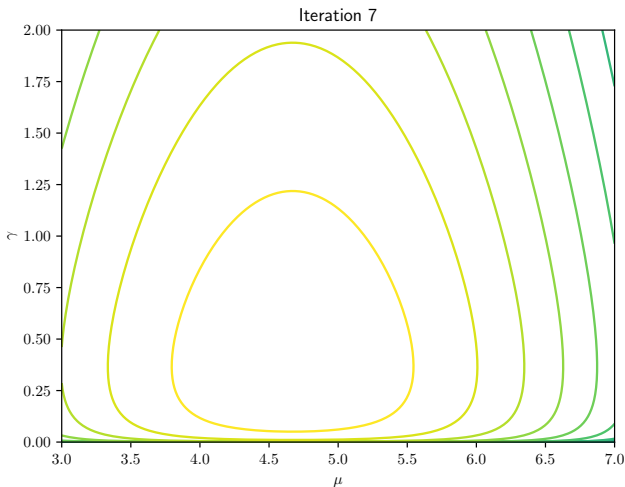
## Posteriors



## Data

Four data points sampled from a normal distribution with mean 5 and variance 1.

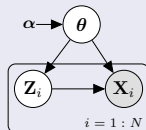
## Posteriors



# Stochastic Variational Inference

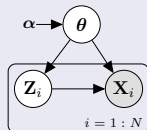
## Model of interest

- $\theta$  are *global* hidden variables, using  $\alpha$  as (hyper-)parameters.
- $\mathbf{Z}_i$  is a vector of latent variables *local* to  $\mathbf{X}_i$ 
  - $\mathbf{Z}_i$  describes the internal structure of  $\mathbf{X}_i$  (like in a factor analysis model).
  - Notice that  $\{\mathbf{X}_i, \mathbf{Z}_i\} \perp\!\!\!\perp \{\mathbf{X}_j, \mathbf{Z}_j\} \mid \theta$  for  $i \neq j$ .
- All distributions belong to the (conjugate) Exponential Family.
- $\mathbf{X}_i$  is observed, hence we have a data-set  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ .



## Model of interest

- $\theta$  are *global* hidden variables, using  $\alpha$  as (hyper-)parameters.
- $\mathbf{Z}_i$  is a vector of latent variables *local* to  $\mathbf{X}_i$ 
  - $\mathbf{Z}_i$  describes the internal structure of  $\mathbf{X}_i$  (like in a factor analysis model).
  - Notice that  $\{\mathbf{X}_i, \mathbf{Z}_i\} \perp\!\!\!\perp \{\mathbf{X}_j, \mathbf{Z}_j\} \mid \theta$  for  $i \neq j$ .
- All distributions belong to the (conjugate) Exponential Family.
- $\mathbf{X}_i$  is observed, hence we have a data-set  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ .



## Example of use:

- $\theta$  represents *topics* for text document and which words are used for each topic.
- $\mathbf{X}_i$  is a text document represented by a bag-of-words.
- $\mathbf{Z}_i$  encodes which topics  $\mathbf{X}_i$  discusses.

## Goals:

- Infer the local latent representation  $\mathbf{Z}_i$  for each observation  $\mathbf{x}_i$
- Infer the global representation  $\theta$ . **Typically this is the most important goal.**



## Algorithm

- 1 Initialize all variational parameters randomly.
- 2 Repeat
  - (a) For each local variational parameter-vector  $\eta_{z_j}$ :

$$\eta_{z_j} \leftarrow \tilde{\eta}_{z_j}(\mathbf{m}_{\theta \rightarrow z_j}) + \mathbf{m}_{x_i \rightarrow z_i}.$$

- (b) Update the variational parameters for global parameter  $\theta$ :

$$\eta_{\theta} \leftarrow \eta_{\theta}(\alpha) + \sum_{i=1}^N (\mathbf{m}_{z_i \rightarrow \theta} + \mathbf{m}_{x_i \rightarrow \theta}).$$

- 3 Until we converge wrt.  $\mathcal{L}(q)$ .

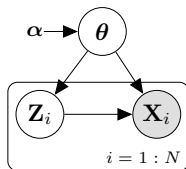
**Computational problem:**

- In the first iteration, each  $\eta_{z_j}$  is using the random initialization of  $\eta_{\theta}$ .
  - This may jeopardize the results of  $\eta_{z_j}$ , and this waste computation.
  - We do  $N$  local updates using  $\eta_{\theta}$ .
  - If  $N$  is large this can be a considerable loss of computational time.
- In turn,  $\eta_{\theta}$  will be updated based on poorly adjusted  $\eta_{z_i}$  values, and the whole process converges slowly.

Recall the VMP architecture

$\eta_\theta$  is updates according to

$$\eta_\theta(\alpha) + \sum_{i=1}^N (\mathbf{m}_{\mathbf{z}_i \rightarrow \theta} + \mathbf{m}_{\mathbf{x}_i \rightarrow \theta})$$



## Opportunity for speed-up through parallelization

- If we distribute the dataset  $\mathcal{D} = \{\mathbf{x}_i, \dots, \mathbf{x}_N\}$  on several computational nodes, the node using data-partition  $\mathcal{D}_j$  will calculate

$$\sum_{i: \mathbf{x}_i \in \mathcal{D}_j} (\mathbf{m}_{\mathbf{z}_i \rightarrow \theta} + \mathbf{m}_{\mathbf{x}_i \rightarrow \theta})$$

- In a map-reduce organization, the master-node sums the messages from the slaves, updates  $\eta_\theta$ , and distributes the  $\theta$ -message back to the slaves.
- Each slave updates its local latent variables  $\{\eta_{\mathbf{z}_i}\}_{i: \mathbf{x}_i \in \mathcal{D}_j}$ .

## “Crazy” idea: Subsampling:

Instead of distributing the dataset we just **subsample** a dataset (“minibatch”) from  $\mathcal{D}$ , say a single observation  $\mathbf{x}_i$  and use that subsample to update  $\lambda_\theta$ :

❶ Initialize all variational parameters randomly to  $\lambda$ .

❷ Repeat forever

(a) Select one  $\mathbf{x}_i$  randomly from  $\mathcal{D}$ .

(b) Update  $\eta_{\mathbf{z}_i}$ :

$$\eta_{\mathbf{z}_i} \leftarrow \tilde{\eta}_{\mathbf{z}_i}(\mathbf{m}_{\theta \rightarrow \mathbf{z}_i}) + \mathbf{m}_{\mathbf{x}_i \rightarrow \mathbf{z}_i}$$

(c) Update the variational parameters for  $\theta$ :

$$\eta_\theta \leftarrow \eta_\theta(\alpha) + \mathbf{m}_{\mathbf{z}_i \rightarrow \theta} + \mathbf{m}_{\mathbf{x}_i \rightarrow \theta}$$

## “Crazy” idea: Subsampling:

Instead of distributing the dataset we just **subsample** a dataset (“minibatch”) from  $\mathcal{D}$ , say a single observation  $\mathbf{x}_i$  and use that subsample to update  $\lambda_\theta$ :

❶ Initialize all variational parameters randomly to  $\lambda$ .

❷ Repeat forever

(a) Select one  $\mathbf{x}_i$  randomly from  $\mathcal{D}$ .

(b) Update  $\eta_{\mathbf{z}_i}$ :

$$\eta_{\mathbf{z}_i} \leftarrow \tilde{\eta}_{\mathbf{z}_i}(\mathbf{m}_{\theta \rightarrow \mathbf{z}_i}) + \mathbf{m}_{\mathbf{x}_i \rightarrow \mathbf{z}_i}$$

(c) Update the variational parameters for  $\theta$ :

$$\eta_\theta \leftarrow \eta_\theta(\alpha) + \mathbf{m}_{\mathbf{z}_i \rightarrow \theta} + \mathbf{m}_{\mathbf{x}_i \rightarrow \theta}$$

Bad news – and some good:

**Bad news:** This does not work!

However, the **good news** is that we can fix it!

# A small side-step: Gradient Ascent

Gradient ascent algorithm for maximizing a function  $f(\lambda)$ :

❶ Initialize  $\lambda^{(0)}$  randomly.

❷ For  $t = 1, \dots$  :

$$\lambda^{(t+1)} \leftarrow \lambda^{(t)} + \rho \cdot \nabla_{\lambda} f(\lambda^{(t)})$$

$\lambda^{(t)}$  converges to a (local) optimum of  $f(\cdot)$  if:

- $f$  is “sufficiently nice”;
- The learning-rate  $\rho$  is “sufficiently small”.

## Stochastic gradient ascent algorithm for maximizing a function $f(\lambda)$ :

If we have access to  $g(\lambda)$  – an **unbiased estimate** of the gradient – it still works!

- 1 Initialize all variational parameters randomly to  $\lambda^{(0)}$ .
- 2 For  $t = 1, \dots$ :

$$\lambda^{(t)} \leftarrow \lambda^{(t-1)} + \rho_t \cdot g\left(\lambda^{(t-1)}\right)$$

$\lambda_t$  converges to a (local) optimum of  $f(\cdot)$  if:

- $f$  is “sufficiently nice”;
- $g(\lambda)$  is a random variable with  $\mathbb{E}[g(\lambda)] = \nabla_{\lambda} f(\lambda)$  and finite variance.
- The learning-rates  $\{\rho_t\}$  is a Robbins-Monro – sequence:
  - $\sum_t \rho_t = \infty$
  - $\sum_t \rho_t^2 < \infty$

The (Euclidian) gradient points in the direction of the solution of

$$\arg \max_{d\lambda} f(\lambda + d\lambda) \quad \text{subject to } \|d\lambda\|_2 < \epsilon$$

- This, however, fails to recognize that we work with probability distributions:
  - The distributions  $\mathcal{N}(\mu = 0, \tau = 10^{-6})$  and  $\mathcal{N}(\mu = 10, \tau = 10^{-6})$  are “close” as both are virtually uniform on  $\mathbb{R}$ , but have distance 10 parameter space.
  - $\mathcal{N}(\mu = 0, \tau = 10^6)$  and  $\mathcal{N}(\mu = 1, \tau = 10^6)$  are “separated”, even though their distance in  $\lambda$ -space is only 1.

The (Euclidian) gradient points in the direction of the solution of

$$\arg \max_{d\lambda} f(\lambda + d\lambda) \quad \text{subject to } \|d\lambda\|_2 < \epsilon$$

- This, however, fails to recognize that we work with probability distributions:
  - The distributions  $\mathcal{N}(\mu = 0, \tau = 10^{-6})$  and  $\mathcal{N}(\mu = 10, \tau = 10^{-6})$  are “close” as both are virtually uniform on  $\mathbb{R}$ , but have distance 10 parameter space.
  - $\mathcal{N}(\mu = 0, \tau = 10^6)$  and  $\mathcal{N}(\mu = 1, \tau = 10^6)$  are “separated”, even though their distance in  $\lambda$ -space is only 1.
- **Natural gradients** take the information geometry into account.
- The natural gradients are found by pre-multiplying with the inverse Fisher matrix:

$$\tilde{\nabla}_{\lambda} f(\lambda) = \mathbf{H}_{\lambda}^{-1} \nabla_{\lambda} f(\lambda)$$

where  $\mathbf{H}_{\lambda}$  is defined as  $\mathbf{H}_{\lambda} = -\mathbb{E}_X [\nabla_{\lambda}^2 \log f(X | \lambda) | \lambda]$ .



The (Euclidian) gradient points in the direction of the solution of

$$\arg \max_{d\lambda} f(\lambda + d\lambda) \quad \text{subject to } \|d\lambda\|_2 < \epsilon$$

- This, however, fails to recognize that we work with probability distributions:
  - The distributions  $\mathcal{N}(\mu = 0, \tau = 10^{-6})$  and  $\mathcal{N}(\mu = 10, \tau = 10^{-6})$  are “close” as both are virtually uniform on  $\mathbb{R}$ , but have distance 10 parameter space.
  - $\mathcal{N}(\mu = 0, \tau = 10^6)$  and  $\mathcal{N}(\mu = 1, \tau = 10^6)$  are “separated”, even though their distance in  $\lambda$ -space is only 1.
- **Natural gradients** take the information geometry into account.
- The natural gradients are found by pre-multiplying with the inverse Fisher matrix:

$$\tilde{\nabla}_{\lambda} f(\lambda) = \mathbf{H}_{\lambda}^{-1} \nabla_{\lambda} f(\lambda)$$

where  $\mathbf{H}_{\lambda}$  is defined as  $\mathbf{H}_{\lambda} = -\mathbb{E}_X [\nabla_{\lambda}^2 \log f(X | \lambda) | \lambda]$ .

- The same operation can obviously also be done in a sub-sample setting, with the same  $\mathbf{H}_{\lambda}$  – since  $\mathbf{H}_{\lambda}$  is data-independent.

**Code Task: Maximum Likelihood in a simple Gaussian model**

We are looking for the ML estimators for them mean  $\mu$  and precision  $\tau$  in a simple Gaussian model.

Recall that

$$f(\mu, \tau) = \sum_{i=1}^N \log p(x_i | \mu, \tau) = -\frac{N}{2} \log(2\pi) + \frac{N}{2} \log \tau - \frac{\tau}{2} \sum_{i=1}^N (x_i - \mu)^2$$

This should give us easy access to both  $\nabla_{(\mu, \tau)} f(\mu, \tau)$  as well as the Fisher information matrix  $-\mathbb{E}_X [\nabla_{(\mu, \tau)}^2 \log p(X | \mu, \tau) | \mu, \tau]$ .

You are asked to fix `calculate_gradient` in the file  
`students_ML_via_SGD.ipynb`.

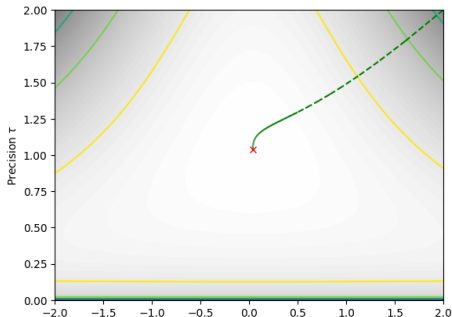
Examine the effect of

- Changing the batch-size (all data or just a single observation per step)
- Switching between natural and Euclidian gradients.

## Example: Maximum log likelihood in a Gaussian model

We have access to  $N = 1000$  observations from a Gaussian distribution with unknown mean  $\mu$  and precision  $\tau$ . Use  $\lambda = [\mu, \tau]^\top$ .

$$f(\lambda) = \sum_{i=1}^N \log p(x_i | \lambda) = \frac{N}{2} \log \tau - \frac{N}{2} \log(2\pi) - \frac{\tau}{2} \sum_{i=1}^N (x_i - \mu)^2$$
$$\nabla_{\lambda} = \begin{bmatrix} -N\tau\mu + \tau \sum_{i=1}^N x_i \\ \frac{N}{2\tau} - \frac{1}{2} \sum_{i=1}^N (x_i - \mu)^2 \end{bmatrix} \quad \text{Cost of calculation: } O(N)$$



## Example: Maximum log likelihood in a Gaussian model

We consider the same maximum likelihood problem, but instead of the gradient based on the full sample, we only have a **mini-batch of a single example**  $x_t$  at iteration  $t$ :

$$\mathbf{g}(\boldsymbol{\lambda} | x_t) = N \cdot \begin{bmatrix} -\tau\mu + \tau x_t \\ \frac{1}{2\tau} - \frac{1}{2} (x_t - \mu)^2 \end{bmatrix}$$

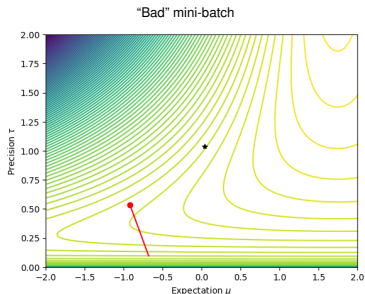
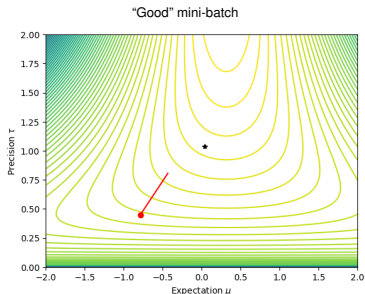
Randomness in  $\mathbf{g}$  is a consequence of the random data selection process, and it follows that  $\mathbb{E}[\mathbf{g}(\boldsymbol{\lambda})] = \nabla_{\boldsymbol{\lambda}} f(\boldsymbol{\lambda})$  because we re-scaled through a multiplication of  $N$ .

## Example: Maximum log likelihood in a Gaussian model

We consider the same maximum likelihood problem, but instead of the gradient based on the full sample, we only have a **mini-batch of a single example**  $x_t$  at iteration  $t$ :

$$\mathbf{g}(\boldsymbol{\lambda} | x_t) = N \cdot \begin{bmatrix} -\tau\mu + \tau x_t \\ \frac{1}{2\tau} - \frac{1}{2} (x_t - \mu)^2 \end{bmatrix}$$

Randomness in  $\mathbf{g}$  is a consequence of the random data selection process, and it follows that  $\mathbb{E}[\mathbf{g}(\boldsymbol{\lambda})] = \nabla_{\boldsymbol{\lambda}} f(\boldsymbol{\lambda})$  because we re-scaled through a multiplication of  $N$ .

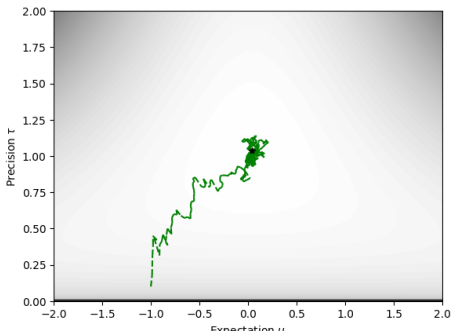


## Example: Maximum log likelihood in a Gaussian model

We consider the same maximum likelihood problem, but instead of the gradient based on the full sample, we only have a **mini-batch of a single example**  $x_t$  at iteration  $t$ :

$$\mathbf{g}(\boldsymbol{\lambda} | x_t) = N \cdot \begin{bmatrix} -\tau\mu + \tau x_t \\ \frac{1}{2\tau} - \frac{1}{2} (x_t - \mu)^2 \end{bmatrix}$$

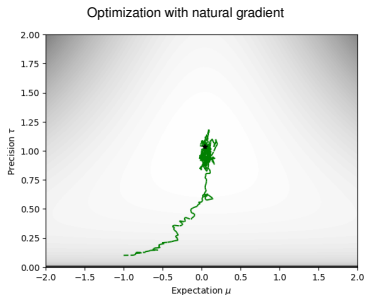
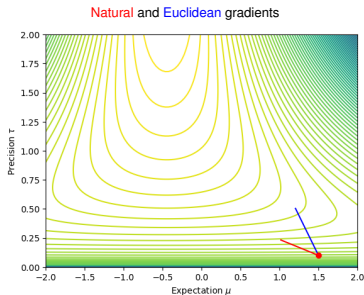
Randomness in  $\mathbf{g}$  is a consequence of the random data selection process, and it follows that  $\mathbb{E}[\mathbf{g}(\boldsymbol{\lambda})] = \nabla_{\boldsymbol{\lambda}} f(\boldsymbol{\lambda})$  because we re-scaled through a multiplication of  $N$ .



## Example: Maximum log likelihood in a Gaussian model

We consider the same maximum likelihood problem. The Fisher information matrix is given by

$$\begin{aligned}\mathbf{H}_{\lambda} &= -\mathbb{E}_X [\nabla_{\lambda}^2 \log p(X | \lambda) | \lambda] \\ &= -\mathbb{E}_X \left( \begin{bmatrix} -\tau & X - \mu \\ X - \mu & -\frac{1}{2\tau^2} \end{bmatrix} \right) = \begin{bmatrix} \tau & 0 \\ 0 & \frac{1}{2\tau^2} \end{bmatrix}\end{aligned}$$



## The Stochastic Variational Inference (SVI) algorithm

- 1 Initialize  $\lambda^{(0)}$ ; Set  $t = 1$ ; Let  $\{\rho_t\}_{t=1}^{\infty}$  be a Robbins-Moreno – sequence.
- 2 Repeat forever
  - (a) Select a single observation  $\mathbf{x}_i$  from  $\mathcal{D}$ .
  - (b) Compute its local variational parameter-vector  $\eta_i$ :

$$\eta_{\mathbf{z}_j} \leftarrow \tilde{\eta}_{\mathbf{z}_j}(\mathbf{m}_{\theta \rightarrow \mathbf{z}_j}) + \mathbf{m}_{\mathbf{x}_i \rightarrow \mathbf{z}_j}$$

- (c) Update the variational parameters for  $\theta$ :

$$\lambda_{\theta}^{(t)} \leftarrow (1 - \rho_t) \lambda_{\theta}^{(t-1)} + \rho_t (\eta_{\theta}(\alpha) + N \cdot (\mathbf{m}_{\mathbf{z}_i \rightarrow \theta} + \mathbf{m}_{\mathbf{x}_i \rightarrow \theta}))$$

- (d)  $t \leftarrow t + 1$

- The algorithm holds **theoretical guarantees** of convergence – it can be seen as a **stochastic gradient ascent** algorithm over  $\mathcal{L}(q)$ .
- We work in natural parameter space – this **simplifies the calculations** considerably, and **improves performance**.
- SVI offers **substantive improvements** over VMP on massive data-sets –  $\lambda_{\theta}$  can often converge before all data is read once.



## The last exponential toppings

Remember the definition of the Exponential Family:

$$f_X(x | \boldsymbol{\eta}) = \exp \left( h(\mathbf{x}) + \boldsymbol{\eta}^\top \mathbf{t}(\mathbf{x}) - A(\boldsymbol{\eta}) \right)$$

Maximum likelihood estimator for  $\boldsymbol{\eta}$  from a dataset  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ :

- First notice that

$$\log f(\mathcal{D} | \boldsymbol{\eta}) = \sum_i h(\mathbf{x}_i) + \boldsymbol{\eta}^\top \sum_i \mathbf{t}(\mathbf{x}_i) - N \cdot A(\boldsymbol{\eta}).$$

- We thus look for

$$\boldsymbol{\eta}^* = \arg \max_{\boldsymbol{\eta}} \sum_i \boldsymbol{\eta}^\top \mathbf{t}(\mathbf{x}_i) - N \cdot A(\boldsymbol{\eta})$$

- Since  $\frac{dA(\boldsymbol{\eta})}{d\boldsymbol{\eta}} = \mathbb{E}[\mathbf{t}(\mathbf{X})]$ , it follows that  $\boldsymbol{\eta}^*$  ensures *moment matching* for  $\mathbf{t}(\mathbf{X})$ :

$$\mathbb{E}[\mathbf{t}(\mathbf{X})] = \frac{1}{N} \sum_i \mathbf{t}(\mathbf{x}_i)$$

- ... and convexity of  $A(\boldsymbol{\eta})$  ensures that this is the unique global optimum.

Remember the definition or the Exponential Family:

$$f_{\mathbf{X}}(x \mid \boldsymbol{\eta}) = \exp \left( h(\mathbf{x}) + \boldsymbol{\eta}^{\top} \mathbf{t}(\mathbf{x}) - A(\boldsymbol{\eta}) \right)$$

Maximum likelihood estimator for  $\boldsymbol{\eta}$  from a dataset  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ :

- Example: The multivariate Gaussian has  $t(\mathbf{X}) = [\mathbf{X}, \mathbf{X}\mathbf{X}^{\top}]$ .
- Remember that  $\mathbb{E}[\mathbf{t}(\mathbf{X})] = [\mathbb{E}(\mathbf{X}), \mathbb{E}(\mathbf{X}\mathbf{X}^{\top})] = [-\frac{1}{2}\boldsymbol{\eta}_2^{-1}\boldsymbol{\eta}_1, -\frac{1}{2}\boldsymbol{\eta}_2^{-1}]$ .
- This gives us

$$\boldsymbol{\eta}_1^* = \left( \sum_i \mathbf{x}_i \right) \cdot \left( \sum_i \mathbf{x}_i \mathbf{x}_i^{\top} \right)^{-1} \text{ and } \boldsymbol{\eta}_2^* = -\frac{N}{2} \left( \sum_i \mathbf{x}_i \mathbf{x}_i^{\top} \right)^{-1}.$$

- While harder to interpret than the *moment parameter* MLEs,

$$\boldsymbol{\mu}^* = \frac{1}{N} \sum_i \mathbf{x}_i \text{ and } \boldsymbol{\Sigma}^* = \frac{1}{N} \sum_i \mathbf{x}_i \mathbf{x}_i^{\top},$$

the *natural parameter* MLEs are found with a unifying theory.

## Conjugacy plays a crucial role in Bayesian inference:

- Assume we observe variables  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  from the model  $f(\mathbf{x} | \boldsymbol{\eta})$ .
- Prior knowledge about  $\boldsymbol{\eta}$  is encoded in  $f(\boldsymbol{\eta} | \boldsymbol{\nu})$ , we seek the posterior  $f(\boldsymbol{\eta} | \mathcal{D}, \boldsymbol{\nu}) \propto f(\boldsymbol{\eta} | \boldsymbol{\nu}) \prod_i f(\mathbf{x}_i | \boldsymbol{\eta})$ .
- $f(\boldsymbol{\eta} | \boldsymbol{\nu})$  is a conjugate for  $f(\mathbf{x} | \boldsymbol{\eta})$  if  $f(\boldsymbol{\eta} | \mathcal{D}, \boldsymbol{\nu})$  has the same form as  $f(\boldsymbol{\eta} | \boldsymbol{\nu})$  when both are seen as functions of  $\boldsymbol{\eta}$ .
- Then inference amounts to “updating”  $\boldsymbol{\nu}$  to a new value, defined by  $\boldsymbol{\nu}$  and  $\mathcal{D}$ .

## Examples:

- Prior Multinomial + Multinomial likelihood .
- Prior  $\text{Beta}(p)$  + Likelihood  $\text{Bernoulli}(x | p)$  – also for Binomial for fixed  $n$ .
- Prior  $\text{Normal}(\mu)$  + Likelihood  $\text{Normal}(x | \mu)$  with mean  $\mu$  and known variance.
- Prior  $\text{Gamma}(\tau)$  + Likelihood  $\text{Normal}(x | \tau)$  with known mean and variance  $\tau^{-1}$ .

## Conjugacy plays a crucial role in Bayesian inference:

- Assume we observe variables  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  from the model  $f(\mathbf{x} | \boldsymbol{\eta})$ .
- Prior knowledge about  $\boldsymbol{\eta}$  is encoded in  $f(\boldsymbol{\eta} | \boldsymbol{\nu})$ , we seek the posterior  $f(\boldsymbol{\eta} | \mathcal{D}, \boldsymbol{\nu}) \propto f(\boldsymbol{\eta} | \boldsymbol{\nu}) \prod_i f(\mathbf{x}_i | \boldsymbol{\eta})$ .
- $f(\boldsymbol{\eta} | \boldsymbol{\nu})$  is a conjugate for  $f(\mathbf{x} | \boldsymbol{\eta})$  if  $f(\boldsymbol{\eta} | \mathcal{D}, \boldsymbol{\nu})$  has the same form as  $f(\boldsymbol{\eta} | \boldsymbol{\nu})$  when both are seen as functions of  $\boldsymbol{\eta}$ .
- Then inference amounts to “updating”  $\boldsymbol{\nu}$  to a new value, defined by  $\boldsymbol{\nu}$  and  $\mathcal{D}$ .

## The Exponential Family:

- Assume a likelihood-model  $f_x(\mathbf{x} | \boldsymbol{\eta}) = \exp(h(\mathbf{x}) + \boldsymbol{\eta}^\top \mathbf{t}(\mathbf{x}) - A_x(\boldsymbol{\eta}))$ .
- Then the conjugate prior is  $f_\eta(\boldsymbol{\eta} | \boldsymbol{\tau}, \nu) = \exp(h_\eta(\boldsymbol{\eta}) + \boldsymbol{\tau}^\top \boldsymbol{\eta} - \nu A_x(\boldsymbol{\eta}) - A_\eta(\boldsymbol{\tau}))$ .
- Notice how the prior is not in ExpFam form; can be done by setting

$$\tilde{\boldsymbol{\tau}} = [\boldsymbol{\tau}^\top, -\nu]^\top, \quad \tilde{\mathbf{t}}_\eta(\boldsymbol{\eta}) = [\boldsymbol{\eta}^\top, A_x(\boldsymbol{\eta})]^\top, \quad \tilde{A}_\eta(\tilde{\boldsymbol{\tau}}) = A_\eta(\boldsymbol{\tau}).$$

from which we obtain  $f_\eta(\boldsymbol{\eta} | \tilde{\boldsymbol{\tau}}) = \exp(h_\eta(\boldsymbol{\eta}) + \tilde{\boldsymbol{\tau}}^\top \tilde{\mathbf{t}}_\eta(\boldsymbol{\eta}) - \tilde{A}_\eta(\tilde{\boldsymbol{\tau}}))$ .

↪ specific form can be exploited during, e.g., SVI for calculating the gradients.

## Model formulation

Assume a likelihood-model  $f_x(\mathbf{x} | \boldsymbol{\eta}) = \exp(h_x(\mathbf{x}) + \boldsymbol{\eta}^\top \mathbf{t}(\mathbf{x}) - A_x(\boldsymbol{\eta}))$ , from which we have observed the dataset  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ .

The prior is  $f_\eta(\boldsymbol{\eta} | \boldsymbol{\tau}, \nu) = \exp(h_\eta(\boldsymbol{\eta}) + \boldsymbol{\tau}^\top \boldsymbol{\eta} - \nu A_x(\boldsymbol{\eta}) - A_\eta(\boldsymbol{\tau}))$ .

$$\begin{aligned}\log f_\eta(\boldsymbol{\eta} | \mathcal{D}, \boldsymbol{\tau}, \nu) &\propto \log f_x(\mathcal{D} | \boldsymbol{\eta}) + \log f_\eta(\boldsymbol{\eta} | \boldsymbol{\tau}, \nu) \\&= \sum_i h_x(\mathbf{x}_i) + \boldsymbol{\eta}^\top \sum_i \mathbf{t}(\mathbf{x}_i) - N \cdot A_x(\boldsymbol{\eta}) \\&\quad + h_\eta(\boldsymbol{\eta}) + \boldsymbol{\tau}^\top \boldsymbol{\eta} - \nu \cdot A_x(\boldsymbol{\eta}) - A_\eta(\boldsymbol{\tau}) \\&\propto h_\eta(\boldsymbol{\eta}) + \boldsymbol{\eta}^\top \left\{ \boldsymbol{\tau} + \sum_i \mathbf{t}(\mathbf{x}_i) \right\} - \{N + \nu\} \cdot A_x(\boldsymbol{\eta}) - A_\eta(\boldsymbol{\tau}) \\ \log f_\eta(\boldsymbol{\eta} | \mathcal{D}, \boldsymbol{\tau}, \nu) &= \log f_x(\boldsymbol{\eta} | \boldsymbol{\tau} + \sum_i \mathbf{t}(\mathbf{x}_i), \nu + N)\end{aligned}$$

## Interpretation:

Posterior is found directly by collecting the summed sufficient statistics and no. observations. This **always** works for **every** conjugate exponential family.