

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053

---

# Exponential Family

---

## Contents

<b>1</b>	<b>The Exponential Family</b>	<b>2</b>
1.1	Definition . . . . .	2
1.2	Example: Dirichlet distribution . . . . .	2
1.3	Example: Truncated normal distribution . . . . .	3
1.4	i.i.d samples from an exponential family . . . . .	3
<b>2</b>	<b>Exponential Family: Maximum Likelihood Estimation</b>	<b>3</b>
<b>3</b>	<b>Exponential Family: Expectation Maximization</b>	<b>4</b>
<b>4</b>	<b>Exponential Family: Conjugate Priors</b>	<b>4</b>
4.1	Example: Normal prior on mean of multivariate Gaussian with known covariance .	5
4.2	Example: Inverse Wishart prior on covariance matrix of multivariate Gaussian with known mean . . . . .	5
4.3	Example: Bayesian linear regression with normal prior on regression weights . . .	6
4.4	General formalism . . . . .	7
<b>5</b>	<b>Exponential Family: Conditional Conjugacy</b>	<b>8</b>
5.1	Example: Bayesian normal model with conditionally conjugate prior . . . . .	8
<b>A</b>	<b>EF representation of Multivariate Gaussian in message passing</b>	<b>10</b>
<b>B</b>	<b>More on Bayesian multivariate linear regression</b>	<b>11</b>
<b>C</b>	<b>Multivariate completing the square</b>	<b>12</b>

# 1 The Exponential Family

We are interested in the exponential family primarily because it makes inference easier. When a problem can be cast within the exponential family framework, inference can be tied to general principles, and parameter updates often have nice interpretations. This is true regardless of whether we're doing maximum likelihood, expectation maximization, variational inference, or Gibbs sampling.

## 1.1 Definition

We define an *exponential family* of probability distributions as those distributions whose density has the following form

$$p(x | \theta) = h(x) \exp\{\eta(\theta)^T t(x) - a(\theta)\} \quad (1.1.1)$$

where we refer to  $h$  as the base measure,  $\eta$  as the natural parameter,  $t$  as the sufficient statistics, and  $a$  as the log normalizer.<sup>1</sup>

**Remark 1.1.1.** (*Non-uniqueness of natural parameter*) Note from (1.1.1) that natural parameters are not unique since, for example,  $\eta$  could be multiplied by a non-zero constant  $c$  if  $t(x)$  is divided by  $c$ .<sup>2</sup> Thus, we should speak of *a* natural parameter, rather than *the* natural parameter.

## 1.2 Example: Dirichlet distribution

**Example 1.2.1.** (Dirichlet Distribution) We can write the density of the Dirichlet distribution in exponential form:

$$\begin{aligned} p(\pi | \alpha) &= \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \pi_1^{\alpha_1-1} \dots \pi_K^{\alpha_K-1} \\ &= \exp \left\{ \sum_{k=1}^K (\alpha_k - 1) \log \pi_k - \left[ \sum_k \log \Gamma(\alpha_k) - \log \Gamma(\sum \alpha_k) \right] \right\} \end{aligned}$$

with natural parameter  $\eta(\alpha) = [\alpha_1 - 1, \dots, \alpha_K - 1]^T$ , sufficient statistics  $t(\pi) = \log \pi = [\log \pi_1, \dots, \log \pi_K]^T$ , base measure  $h(\pi) = 1$ , and log normalizer  $a(\alpha) = \sum_k \log \Gamma(\alpha_k) - \log \Gamma(\sum_k \alpha_k)$ .  $\square$

For an example of how the natural parametrization can help provide insight into message passing, see Section A.<sup>3</sup>

<sup>1</sup>TODO: It would be helpful to get more solid on integrating against probability measure here, so that I can set this up in a more precise way, as Jordan does. He remarks on this somewhere in his exponential family lecture notes. What also may be helpful is this beautiful excerpt from pp.38 of [1]: "[...] we represent the probability distribution as a density  $p$  absolutely continuous with respect to some measure  $\eta$ . This base measure  $\eta$  might be the counting measure on  $\{0, 1, \dots, r-1\}$ , in which case  $p$  is a probability mass function; alternatively, for a continuous random vector, the base measure  $\eta$  could be the ordinary Lebesgue measure on  $\mathbb{R}$ ."

<sup>2</sup>Are they unique up to scalar multiplication?

<sup>3</sup>**Remark** The exponential family representation of the Dirichlet, as given in Example 1.2.1, is useful when we want to compute the expectation of a log probability from a Dirichlet distributed probability vector (as happens in the derivation of LDA with variational inference; see my notes on variational inference).

In those notes, we see

$$\begin{aligned} \mathbb{E}[\log \pi_k] &= \mathbb{E}[t_k(p)] \stackrel{(1)}{=} \frac{\partial}{\partial \eta_k} a(\eta) \\ &= \Psi(\alpha_k) - \Psi\left(\sum_k \alpha_k\right) \end{aligned} \quad (1.2.1)$$

where (1) uses a well-known exponential family property and where  $\Psi(\cdot)$  is the first derivative of the log  $\Gamma$  function. It is known as the *digamma function*.  $\square$

TODO: Add multivariate Gaussian example, showing that the natural parameters are the precision  $\Sigma^{-1}$  and precision-weighted mean  $\Sigma^{-1}\mu$ , as we use this in Section 5.1 combined with the exponential family formalism to derive the updates to the mean for a Bayesian normal model with conditionally conjugate prior.

### 1.3 Example: Truncated normal distribution

**Example 1.3.1.** (Truncated normal distribution) The univariate truncated normal distribution  $\mathcal{TN}(\mu, \sigma^2, \Omega)$  results when a normal distribution  $\mathcal{N}(\mu, \sigma^2)$  is truncated to some set  $\Omega \in \mathbb{R}$ .<sup>4</sup> Note that the parameters  $\mu, \sigma^2$  denote the mean and variance of the *parent* normal distribution; i.e. if  $X \sim \mathcal{TN}(\mu, \sigma^2, \Omega)$  then  $\mathbb{E}[X] \neq \mu$  (unless  $\Omega = \mathbb{R}$ ).

If we assume that the truncation set is an interval  $\Omega = (a, b)$  for  $a, b \in \mathbb{R}$ , then the distribution  $\mathcal{TN}(\mu, \sigma^2, (a, b))$  has p.d.f.

$$f(x; \mu, \sigma^2, a, b) = \frac{\phi_{\mu, \sigma^2}(x)}{\Phi_{\mu, \sigma^2}(b) - \Phi_{\mu, \sigma^2}(a)} 1_{a \leq x \leq b} \quad (1.3.1)$$

where  $\phi_{\mu, \sigma^2}$  and  $\Phi_{\mu, \sigma^2}$  denote the pdf and cdf, respectively, of a univariate normal distribution with mean  $\mu$  and variance  $\sigma^2$ .

If we write

$$\begin{aligned} f(x; \mu, \sigma^2, a, b) &= K \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2} \frac{(x - \mu)^2}{\sigma^2}\right) 1_{a \leq x \leq b} \\ &= K \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} x^2 + \frac{\mu}{\sigma^2} x + \frac{\mu^2}{\sigma^2} - \log \sigma\right) 1_{a \leq x \leq b} \end{aligned}$$

where  $K := \Phi_{\mu, \sigma^2}(b) - \Phi_{\mu, \sigma^2}(a)$ , then we see that  $\mathcal{TN}(\mu, \sigma^2, (a, b))$  belongs to the exponential family (1.1.1) where, in this case, we have natural parameter  $\eta = (-\frac{1}{2\sigma^2}, \frac{\mu}{\sigma^2})^T$ , sufficient statistics function  $t(x) = (x^2, x)^T$ , and base measure  $h(x) = K \frac{1}{\sqrt{2\pi}} 1_{a \leq x \leq b}$ .  $\square$

### 1.4 i.i.d samples from an exponential family

If  $\mathbf{x} = (x_1, \dots, x_n)$  are  $n$  independent samples from the same exponential family, then

$$p(\mathbf{x} \mid \theta) = \prod_{i=1}^n h(x_i) \exp\left\{\eta(\theta)^T \sum_{i=1}^n t(x_i) - n a(\eta(\theta))\right\} \quad (1.4.1)$$

## 2 Exponential Family: Maximum Likelihood Estimation

The goal for maximum likelihood is to determine the parameter

$$\theta_{ML} = \operatorname{argmax}_{\theta} \log p(\mathbf{x} \mid \theta) \quad (2.0.1)$$

Let us assume that  $\mathbf{x} = (x_1, \dots, x_n)$  are i.i.d observations from a fixed exponential family, so that the likelihood has form (1.4.1). Let us compute the gradient with respect to the natural parameter  $\eta$  of  $\ell(\eta) := \log p(\mathbf{x} \mid \eta)$

$$\nabla_{\eta} \ell(\eta) = \sum_{i=1}^n t(x_i) - n \nabla_{\eta} a(\eta)$$

Setting the gradient to zero, we obtain

<sup>4</sup>For more information on the truncated normal, see e.g. [https://people.sc.fsu.edu/~jburkardt/presentations/truncated\\_normal.pdf](https://people.sc.fsu.edu/~jburkardt/presentations/truncated_normal.pdf) or <http://parker.ad.siu.edu/Olive/ch4.pdf>.

$$\nabla_{\eta} a(\eta) = \frac{1}{n} \sum_{i=1}^n t(x_i)$$

But  $\nabla_{\eta} a(\eta) = \mathbb{E}[t(X)]$  [2]. Thus, we should set  $\theta_{ML}$  such that

$$\mu(\theta_{ML}) = \frac{1}{n} \sum_{i=1}^n t(x_i)$$

where  $\mu := \mathbb{E}[t(x)]$  refers to the mean parametrization of the likelihood.<sup>5</sup>

### 3 Exponential Family: Expectation Maximization

Some models have latent variables associated with each observation, and so maximum likelihood is not possible. Let us see how expectation maximization looks when the complete data likelihood is in the exponential family.

The expectation maximization algorithm is

$$\theta^{(t+1)} = \operatorname{argmax}_{\theta} \mathbb{E}_{p(z | x, \theta^{(t)})} \left[ \ln p(x, z | \theta) \right] \quad (3.0.1)$$

We see how this plays out in the exponential family by following the logic of Section 2. Let us assume that  $(x, z) = ((x_1, z_1), \dots, (x_n, z_n))$  are  $n$  independent samples from the same exponential family, where  $x$  is observed data and  $z$  is unobserved data. Moreover, let us assume that the complete data likelihood is in the exponential family

$$p(x, z | \theta) = \prod_{i=1}^n h(x_i, z_i) \exp \left\{ \eta(\theta)^T \sum_{i=1}^n t(x_i, z_i) - n a(\eta(\theta)) \right\} \quad (3.0.2)$$

Here we want to find  $\theta$  to optimize

$$f(\theta) = \mathbb{E}_{p(z | x, \theta^{(t)})} \left[ \ln p(x, z | \theta) \right]$$

Following the logic of Section 2, we determine that we should select  $\theta^{(t+1)}$  such that

$$\mu(\theta^{(t+1)}) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{p(z | x, \theta^{(t)})} t(x_i, z_i)$$

where  $\mu := \mathbb{E}[t(x_1, z_1)]$  refers to the mean parametrization of the likelihood.

This is why an EM iteration is often described and/or implemented as performing maximum likelihood with the expected sufficient statistics.

TODO: But is EM *always* equivalent to performing ML with ESS's? Or is this *ONLY* true if I'm working within the exponential family? I need to read up some more on EM theory.

TODO: Check this section, especially with respect to the fact that I am dealing with three parametrizations here -  $\mu, \theta, \nu$ ; that is, mean, arbitrary, and natural, respectively. Really the core problem is that it's not sufficiently clear in how head how and when reparametrizations affect things.

### 4 Exponential Family: Conjugate Priors

TODO: State what a conjugate prior is without using the formalism of Section 4.4

<sup>5</sup>TODO: This switching of parameterization should be handled much more explicitly.

#### 216 4.1 Example: Normal prior on mean of multivariate Gaussian with known covariance 217

218 **TODO.** Note that in Section ??, I've already written up a self-enclosed argument for Bayesian linear  
219 regression; some of that argument can likely be factored out to here.  
220

#### 221 4.2 Example: Inverse Wishart prior on covariance matrix of multivariate Gaussian with 222 known mean

223 Here we will show that the Inverse Wishart is a conjugate prior for the covariance of a multivariate  
224 normally distributed random variable with known mean.  
225

226 This situation comes up

227 **Example 4.2.1.** (*Inverse Wishart prior on the covariance of a Multivariate Normal sampling model*  
228 *with known mean*)  
229

230 Consider the sampling model for  $\mathbf{y} := (\mathbf{y}_1, \dots, \mathbf{y}_n) \stackrel{\text{iid}}{\sim} \mathcal{N}_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

$$\begin{aligned} p(\mathbf{y} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) &\propto |\boldsymbol{\Sigma}|^{-n/2} \exp \left[ -\frac{1}{2} \sum_{i=1}^n (\mathbf{y}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i - \boldsymbol{\mu}) \right] \\ &= |\boldsymbol{\Sigma}|^{-n/2} \exp \left[ -\frac{1}{2} \text{tr}(\boldsymbol{\Sigma}^{-1} \mathbf{S}_\mu) \right] \end{aligned} \quad (4.2.1)$$

231 where  $\mathbf{S}_\mu := \sum_{i=1}^n (\mathbf{y}_i - \boldsymbol{\mu})(\mathbf{y}_i - \boldsymbol{\mu})^T$  is the sum of pairwise deviation products, and where the  
232 equality in (4.2.1) is justified in Remark 4.2.2.  
233

234 Let us take the mean  $\boldsymbol{\mu}$  to be known, and let us take the prior on the covariance  $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$  to be  
235 given by  $\boldsymbol{\Sigma} \sim \mathcal{W}^{-1}(\boldsymbol{\Psi}, \nu)$ , i.e.  
236

$$p(\boldsymbol{\Sigma}) \propto |\boldsymbol{\Sigma}|^{-(\nu+d+1)/2} \exp \left[ -\frac{1}{2} \text{tr}(\boldsymbol{\Sigma}^{-1} \boldsymbol{\Psi}) \right] \quad (4.2.2)$$

237 where  $\boldsymbol{\Sigma} \succ 0$  and  $\nu > d - 1$  to have a proper prior. Note that  $\mathbb{E}[\boldsymbol{\Sigma}] = \frac{\boldsymbol{\Psi}}{\nu - d - 1}$ .  
238

239 It is easy to see from the forms of the likelihood (4.2.1) and prior (4.2.2) that the Inverse Wishart is  
240 a conjugate prior in this context. In particular  
241

$$p(\boldsymbol{\Sigma} \mid \boldsymbol{\mu}, \mathbf{y}) \propto |\boldsymbol{\Sigma}|^{-(\nu+n+d+1)/2} \exp \left[ -\frac{1}{2} \text{tr}(\boldsymbol{\Sigma}^{-1}(\boldsymbol{\Psi} + \mathbf{S}_\mu)) \right] \quad (4.2.3)$$

242 where  $\mathbf{S}_\mu$  was defined above. Thus, we have  
243

$$\boldsymbol{\Sigma} \mid \boldsymbol{\mu}, \mathbf{y} \sim \mathcal{W}^{-1} \left( \boldsymbol{\Psi} + \sum_{i=1}^n (\mathbf{y}_i - \boldsymbol{\mu})(\mathbf{y}_i - \boldsymbol{\mu})^T, \nu + n \right)$$

244 And so the conjugate updates are given by  
245

$$\nu' \leftarrow \nu + n \quad (4.2.4)$$

$$\boldsymbol{\Psi}' \leftarrow \boldsymbol{\Psi} + \sum_{i=1}^n (\mathbf{y}_i - \boldsymbol{\mu})(\mathbf{y}_i - \boldsymbol{\mu})^T \quad (4.2.5)$$

246 **Remark 4.2.1.** (*Interpreting the hyperparameters of the Inverse Wishart*) Note that the hyperpa-  
247 rameters of the Inverse Wishart can be interpreted (as per conjugacy) in the following way: the  
248 covariance was estimated from  $\nu$  observations with a sum of pairwise deviation products  $\boldsymbol{\Psi}$ .<sup>6</sup>  
249

250 **Remark 4.2.2.** (*Expressing the Multivariate Gaussian density in a nice form for the Inverse Wishart*  
251 *prior on the Covariance Matrix*)  
252

<sup>6</sup>This interpretation also makes the formula for  $\mathbb{E}[\boldsymbol{\Sigma}]$  more intuitive.

Here we justify the equality of (4.2.1).

We will show that  $\sum_{i=1}^n \mathbf{x}_i^T \mathbf{A} \mathbf{x}_i = \text{tr}(\mathbf{A} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T)$  for  $\mathbf{x} \in \mathbb{R}^d$ , and  $\mathbf{A} \in \mathbb{R}^{d \times d}$  symmetric.

$$\begin{aligned} \sum_{i=1}^n \mathbf{x}_i^T \mathbf{A} \mathbf{x}_i &= \sum_{i=1}^n \sum_{j,k=1}^n a_{jk} x_{ij} x_{ik} \\ &= \sum_{j,k=1}^n \left( \mathbf{A} \circ \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right)_{jk} \\ &\stackrel{(*)}{=} \text{tr}(\mathbf{A} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T) \end{aligned}$$

where  $\circ$  is the Hadamard, also called the elementwise, operator, and where  $(*)$  holds by properties of the  $\text{tr}$  operator

$$\text{tr}(\mathbf{A} \mathbf{B}) = \sum_{i,j} (\mathbf{A}^T \circ \mathbf{B})_{ij} \stackrel{\mathbf{A} \text{ symmetric}}{=} \sum_{i,j} (\mathbf{A} \circ \mathbf{B})_{ij}$$

### 4.3 Example: Bayesian linear regression with normal prior on regression weights

In this section, we will show that the normal prior on  $\beta$  is a conjugate prior for the regression weights  $\beta$  of a Bayesian multiple regression model with known observation noise  $\sigma^2$ . That is, the posterior on  $\beta$  given  $\mathbf{y} = (y_1, \dots, y_n)^T$  for such a model is also Gaussian.

**Proposition 4.3.1.** *Consider the Bayesian linear multiple regression model with known observation noise  $\sigma^2$*

$$\begin{aligned} \beta &\sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \\ y_i \mid \beta &\stackrel{\text{ind}}{\sim} \mathcal{N}(\mathbf{x}_i^T \beta, \sigma^2), \quad i = 1, \dots, n \end{aligned} \tag{4.3.1}$$

where  $\mathbf{x}_i$  designates the  $i$ -th row of the design matrix  $\mathbf{X} \in \mathbb{R}^{n \times p}$ .

The posterior distribution for (4.3.1) is given by

$$p(\beta \mid \mathbf{y}) \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

where

$$\begin{aligned} \boldsymbol{\Sigma} &= \left( \boldsymbol{\Sigma}_0^{-1} + \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} \right)^{-1} \\ \boldsymbol{\mu} &= \boldsymbol{\Sigma} \left( \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0 + \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{y} \right) \end{aligned}$$

*Proof.* By Bayes rule,

$$p(\beta \mid \mathbf{y}) \propto p(\beta) \exp \left\{ \sum_{i=1}^N -\frac{1}{2\sigma^2} \left( y_i - \mathbf{x}_i^T \beta \right)^2 \right\}$$

and defining  $\boldsymbol{\Omega} \in \mathbb{R}^{n \times n} : \boldsymbol{\Omega} = \text{diag}(\frac{1}{\sigma^2}, \dots, \frac{1}{\sigma^2})$ , we have

$$\begin{aligned} &\stackrel{1}{\propto} p(\beta) \exp \left\{ -\frac{1}{2} (\mathbf{y} - \mathbf{X} \beta)^T \boldsymbol{\Omega} (\mathbf{y} - \mathbf{X} \beta) \right\} \\ &\stackrel{2}{\propto} p(\beta) \exp \left\{ -\frac{1}{2} (\mathbf{X}^+ \mathbf{y} - \beta)^T \mathbf{X}^T \boldsymbol{\Omega} \mathbf{X} (\mathbf{X}^+ \mathbf{y} - \beta) \right\} \end{aligned}$$

where (1) writes the weighted sum of squares in matrix notation, and (2) isolates  $\beta$ , using  $\mathbf{X}^+$ , the Moore-Penrose psuedo-inverse of  $\mathbf{X}$ .<sup>7</sup>

Thus, we see that  $p(\beta \mid \mathbf{y})$  is proportional to the product of two multivariate Gaussians:  $p(\beta)$ , which has mean  $\mu_0$  and covariance  $\Sigma_0$ , and another Gaussian, which has mean  $\mathbf{X}^+\mathbf{y}$  and covariance  $(\mathbf{X}^T\Omega\mathbf{X})^{-1}$ . We know from the exponential family representation of the Gaussian that the resulting distribution can be obtained by summing at the scale of natural parameters – which for the Gaussian are the precision and precision-weighted mean.<sup>8</sup> Using this, we obtain

$$p(\beta \mid \mathbf{y}) \sim \mathcal{N}(\mu, \Sigma)$$

where

$$\begin{aligned}\Sigma &= \left( \Sigma_0^{-1} + \mathbf{X}^T\Omega\mathbf{X} \right)^{-1} \\ \mu &= \Sigma \left( \Sigma_0^{-1}\mu_0 + \mathbf{X}^T\Omega\mathbf{X}\mathbf{X}^+\mathbf{y} \right) \\ &= \Sigma \left( \Sigma_0^{-1}\mu_0 + \mathbf{X}^T\Omega\mathbf{y} \right)\end{aligned}$$

recalling that we defined  $\Omega = \text{diag}(\frac{1}{\sigma^2}, \dots, \frac{1}{\sigma^2})$  completes the proof.  $\square$

**Remark 4.3.1.** For a nice conceptual overview of Bayesian linear regression., see [3] or [4]. Among other things, these resources demonstrate how Bayesian regression makes predictions using an infinite collection of regression models (whose contributions are weighted by their posterior probabilities). They also show how the linear model is less restrictive than it might first seem; it can be used to model nonlinear functional relationships by using nonlinear basis functions.

#### 4.4 General formalism

Here we provide some notes, following [2], about conjugate priors for exponential family data models.

Writing the exponential family density in canonical form, we have

$$p(x \mid \eta) = h(x) \exp\{\eta^T T(x) - A(\eta)\}$$

where  $\eta$  is the canonical parameter,  $T(x)$  are the sufficient statistics,  $h(x)$  is the base measure, and  $A(\eta)$  is the log normalizer (and so is *not* a degree of freedom).

The natural parameter space is

$$\left\{ \eta : \int h(x) \exp\{\eta^T T(x) - A(\eta)\} < \infty \right\}$$

Given a random sample,  $\mathbf{x} = (x_1, x_2, \dots, x_N)$ , we obtain:

$$p(\mathbf{x} \mid \eta) = \left( \prod_{i=1}^N h(x_i) \right) \exp \left\{ \eta^T \sum_{i=1}^N T(x_i) - NA(\eta) \right\}$$

<sup>7</sup>Specifically, since  $\mathbf{X}\mathbf{X}^+ = \mathbf{I}$ , we use

$$\begin{aligned}(\mathbf{y} - \mathbf{X}\beta)^T \Omega (\mathbf{y} - \mathbf{X}\beta) &= (\mathbf{X}\beta - \mathbf{y})^T \Omega (\mathbf{X}\beta - \mathbf{y}) \\ &= \left( \mathbf{X}(\beta - \mathbf{X}^+\mathbf{y}) \right)^T \Omega \left( \mathbf{X}(\beta - \mathbf{X}^+\mathbf{y}) \right) \\ &= (\beta - \mathbf{X}^+\mathbf{y})^T \mathbf{X}^T \Omega \mathbf{X} (\beta - \mathbf{X}^+\mathbf{y})\end{aligned}$$

<sup>8</sup>See, for reference, Section A.

as the likelihood function.

A conjugate prior can be obtained by mimicking the likelihood

$$p(\eta \mid \tau, \eta_0) = H(\tau, \eta_0) \exp\{\tau^T \eta - \eta_0 A(\eta)\} \quad (4.4.1)$$

where now  $H(\tau, \eta_0)$  is the normalizing factor. (For conditions on normalizability, see [?]). Note that  $\tau$  has the dimensionality of the canonical parameter  $\eta$  and  $n_0$  is a scalar.

To verify conjugacy, we compute the posterior density

$$p(\eta \mid \mathbf{x}, \tau, \eta_0) \propto \exp\left\{\left(\tau + \sum_{n=1}^N T(x_n)\right)^T \eta - (n_0 + N)A(\eta)\right\}$$

which retains the form of (4.4.1)

Thus, the prior-to-posterior conversion can be summarized with the following update rules

$$\begin{aligned} \tau &\rightarrow \tau + \sum_{n=1}^N T(x_n) \\ n_0 &\rightarrow n_0 + N \end{aligned} \quad (4.4.2)$$

For conjugate Bayesian models, the predictive posterior distribution,  $p(x_{\text{new}} \mid \mathbf{x})$  is always tractable, because it has the same form (integrating a likelihood against the parameter distribution) as does the evidence term in Bayes law. For exponential family models, the predictive posterior takes the form of a ratio of normalizing factors

$$p(x_{\text{new}} \mid \mathbf{x}) = \frac{H(\tau_{\text{post}}, n_0 + N)}{H(\tau_{\text{post}} + T(x_{\text{new}}), n_0 + N + 1)} \quad (4.4.3)$$

TODO: Redo some of the examples using the exponential family conjugate prior formalism. A possibly useful resource in the giant table at [https://en.wikipedia.org/wiki/Exponential\\_family](https://en.wikipedia.org/wiki/Exponential_family).

## 5 Exponential Family: Conditional Conjugacy

A family of prior distributions for a parameter is called conditionally conjugate if the conditional posterior distribution (often called the *complete conditional*), given the data and all other parameters in the model, is also in that class [5]. The posterior distribution for conditionally conjugate models is easily approximated with Gibbs sampling or Mean Field Variational Inference – the former samples from the complete conditional, whereas the latter takes variational expectations with respect to the natural parameter of the complete conditional.

Below we give perhaps the simplest example.

### 5.1 Example: Bayesian normal model with conditionally conjugate prior

Consider the following model with a normal sampling distribution and conditionally conjugate prior<sup>9</sup>:

$$\begin{aligned} \boldsymbol{\mu} &\sim \mathcal{N}_d(\mathbf{m}_0, \mathbf{V}_0) \\ \boldsymbol{\Sigma} &\sim \mathcal{W}^{-1}(\nu_0, \boldsymbol{\Psi}_0) \\ \mathbf{x}_i \mid \boldsymbol{\mu}, \boldsymbol{\Sigma} &\stackrel{\text{iid}}{\sim} \mathcal{N}_d(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \quad i = 1, \dots, N \end{aligned}$$

We define  $\mathbf{x} := (\mathbf{x}_1, \dots, \mathbf{x}_N)$ , where each  $\mathbf{x}_i \in \mathbb{R}^d$ .

<sup>9</sup>TODO: Prove that the prior, although conditionally conjugate, is not conjugate. (I believe this is true, based on context clues from experience, but I am not currently certain about it.)



The complete conditionals are well-known. In particular

$$\boldsymbol{\mu} \mid \boldsymbol{\Sigma}, \mathbf{x} \sim \mathcal{N}_d(\mathbf{m}, \mathbf{V}) \quad (5.1.1)$$

where

$$\begin{aligned} \mathbf{m} &= \left( \mathbf{V}_0^{-1} + N\boldsymbol{\Sigma}^{-1} \right)^{-1} \left( \mathbf{V}_0^{-1}\mathbf{m}_0 + N\boldsymbol{\Sigma}^{-1}\bar{\mathbf{x}} \right) \\ \mathbf{V} &= \left( \mathbf{V}_0^{-1} + N\boldsymbol{\Sigma}^{-1} \right)^{-1} \end{aligned}$$

and

$$\boldsymbol{\Sigma} \mid \boldsymbol{\mu}, \mathbf{x} \sim \mathcal{W}^{-1}(\nu, \boldsymbol{\Psi}) \quad (5.1.2)$$

where

$$\begin{aligned} \nu &= \nu_0 + N \\ \boldsymbol{\Psi} &= \boldsymbol{\Psi}_0 + \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T \end{aligned}$$

Indeed, we derived (5.1.2) in Section 4.2.<sup>10</sup>

Note that the model is different than the model fully conjugate (Normal-Inverse-Wishart) prior on the pair  $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . The conditionally conjugate prior lacks closed-form posterior updating, but is also more expressive.<sup>11</sup>

These conjugate posterior updates have nice interpretations:

- **Hyperparameter updates for  $(\boldsymbol{\mu} \mid \boldsymbol{\Sigma}, \mathbf{x})$ :** On the precision scale,  $\mathbf{V}$  is the sum of the prior precision matrix  $\mathbf{V}_0^{-1}$  and  $N$  copies of the precision for each observation,  $\boldsymbol{\Sigma}^{-1}$ . Similarly,  $\mathbf{m}$  is the precision-weighted convex combination of  $\mathbf{m}_0$ , the prior mean and the empirical average  $\bar{\mathbf{x}}$ .
- **Hyperparameter updates for  $(\boldsymbol{\Sigma} \mid \boldsymbol{\mu}, \mathbf{x})$ :** The covariance was estimated from  $\nu$  observations with a sum of pairwise deviation products  $\boldsymbol{\Psi}$ .

<sup>10</sup>We still need to add a derivation for (5.1.1) **TODO**, but the birds' eye view for one approach is to use the general formalism for conjugacy updates in the exponential family (4.4.2), noting that the natural parameters for a multivariate Gaussian are its precision and precision-weighted mean.

<sup>11</sup>Is it also more expressive once we move to a variational approximation? i.e., can we get more expressive marginals this way?

## A EF representation of Multivariate Gaussian in message passing

In a dissertation on Gaussian Belief Propagation [6], referred to in [7], a multivariate Gaussian is considered as a Markov Random Field.

In particular, consider the Markov Random field

$$p(x) = \frac{1}{Z} \left( \prod_{i=1}^n \phi(x_i) \prod_{i,j} \psi(x_i, x_j) \right) \quad (\text{A.0.1})$$

Now note that a multivariate Gaussian has a joint distribution which can be expressed as

$$p(x) \propto \exp \left\{ -\frac{1}{2} x^T A x + b^T x \right\}$$

as this is just the exponential family form of a Gaussian (e.g., see [8]), where the natural parameters are given in terms of the *precision*  $\Sigma^{-1}$

$$A = \Sigma^{-1}$$

$$b = \Sigma^{-1} \mu$$

Thus, the multivariate Gaussian is a MRF where the potentials in (A.0.1) are given by

$$\begin{aligned} \psi_{ij}(x_i, x_j) &:= \exp \left\{ -\frac{1}{2} x_i A_{ij} x_j \right\} \\ \phi_i(x_i) &:= \exp \left\{ -\frac{1}{2} A_{ii} x_i^2 + b_i x_i \right\} \end{aligned}$$

This seems to be useful in inference for state space models, where one multiplies multiple “messages” that are different Gaussian densities *over the same variable*. For example, see the equations for  $\mu_t$  and  $\sigma_t^2$  in Section 4 of [7], where messages from the past and the future of a time series model are combined to get a posterior distribution on the state  $z_t$ . The combined parameters have an expression which may at first be puzzling:

$$\mu_t = \frac{\mu_1 \sigma_2^2 + \mu_2 \sigma_1^2}{\sigma_1^2 + \sigma_2^2}, \quad \sigma_t^2 = \frac{\sigma_1^2 \sigma_2^2}{\sigma_1^2 + \sigma_2^2}$$

However, these messages have an intuitive form when considered in terms of the natural parameterizations: the combined mean is a weighted combination of the original means, with the weights given by the precisions. The combined precision (inverse covariance) is given simply by the sum of the original precisions. Very nice!

See Figure 1 for the general expression, which explains the formula in [7]. This is an example of where the natural parametrization provides more insight than the standard parametrization.

**Lemma 12.** *Let  $f_1(x)$  and  $f_2(x)$  be the probability density functions of a Gaussian random variable with two possible densities  $\mathcal{N}(\mu_1, P_1^{-1})$  and  $\mathcal{N}(\mu_2, P_2^{-1})$ , respectively. Then their product,  $f(x) = f_1(x)f_2(x)$  is, up to a constant factor, the probability density function of a Gaussian random variable with distribution  $\mathcal{N}(\mu, P^{-1})$ , where*

$$\mu = P^{-1}(P_1 \mu_1 + P_2 \mu_2), \quad (2.9)$$

$$P^{-1} = (P_1 + P_2)^{-1}. \quad (2.10)$$

Figure 1: Lemma 12 of [6]

## B More on Bayesian multivariate linear regression

Below we give an alternate proof for the posterior of the regression weights in multivariate linear regression, compared to what was given in Proposition 4.3.1. This alternate proof may be of interest. Whereas the proof of proposition 4.3.1 given in Section 4.3 refers to exponential family properties, the proof below does not, and instead uses multivariate completing the square to do the heavy lifting.

Note that proposition B.0.1 below, as stated, is slightly more restrictive in that it assumes the prior mean is zero. This additional restriction is not necessary; the proposition could be rewritten to match Proposition 4.3.1 exactly, and the proof could be adjusted accordingly to match the additional generality. The difference in statements is just an unnecessary presentational blemish.<sup>12</sup>

**Proposition B.0.1.** *Consider the Bayesian linear multiple regression model*

$$\beta \sim \mathcal{N}(\mathbf{0}, \mathbf{V})$$

$$y_i \mid \beta \stackrel{\text{ind}}{\sim} \mathcal{N}(\mathbf{x}_i^T \beta, \sigma^2), \quad i = 1, \dots, n$$

where  $\mathbf{x}_i$  designates the  $i$ -th row of the design matrix  $\mathbf{X} \in \mathbb{R}^{n \times p}$ .

The posterior distribution for this model is given by

$$p(\beta \mid \mathbf{y}) \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

where

$$\begin{aligned} \boldsymbol{\mu} &= \frac{1}{\sigma^2} \boldsymbol{\Sigma} \mathbf{X}^T \mathbf{y} \\ \boldsymbol{\Sigma} &= \left( \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \mathbf{V}^{-1} \right)^{-1} \end{aligned}$$

*Proof.* The posterior on  $\beta$  given  $\mathbf{y} = (y_1, \dots, y_n)^T$  is Gaussian, since

$$\begin{aligned} \ln p(\beta \mid \mathbf{y}) &= \sum_{i=1}^n \ln p(y_i \mid \beta) + \ln p(\beta) \\ &= -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \beta)^2 - \frac{1}{2} \beta^T \mathbf{V}^{-1} \beta + \text{constant} \\ &\stackrel{1}{=} -\frac{1}{2\sigma^2} \left( \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X}^T \beta + \beta^T \mathbf{X}^T \mathbf{X} \beta \right) - \frac{1}{2} \beta^T \mathbf{V}^{-1} \beta + \text{constant} \\ &\stackrel{2}{=} -\frac{1}{2} (\beta - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\beta - \boldsymbol{\mu}) + \text{constant} \end{aligned}$$

where

$$\begin{aligned} \boldsymbol{\mu} &= \frac{1}{\sigma^2} \boldsymbol{\Sigma} \mathbf{X}^T \mathbf{y} \\ \boldsymbol{\Sigma} &= \left( \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \mathbf{V}^{-1} \right)^{-1} \end{aligned}$$

Equality (1) is obtained by noting  $\sum_{i=1}^n (y_i - \mathbf{x}_i^T \beta)^2 = (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta)$ , FOIL-ing, and observing that the cross-products are scalars. Equality (2) is obtained by completing the square, where  $\beta$  plays the role of  $\mathbf{x}$  in (C.0.1), and where in that notation we have  $\mathbf{M} = \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \mathbf{V}^{-1}$  and  $\mathbf{b}^T = \frac{1}{\sigma^2} \mathbf{y}^T \mathbf{X}^T$   $\square$

<sup>12</sup>TODO: fix up the unnecessary presentational blemish – assuming that we don't end up sacrificing too much pedagogical clarity for the sake of generality

## C Multivariate completing the square

A nice overview of multivariate completing the square is given by [9].

Let  $\mathbf{x}, \mathbf{b}$  be  $d$ -dimensional vectors, and let  $\mathbf{M} \in \mathbb{R}^{d \times d}$  be a symmetric invertible matrix. Then

$$\mathbf{x}^T \mathbf{M} \mathbf{x} - 2\mathbf{b}^T \mathbf{x} = (\mathbf{x} - \mathbf{M}^{-1}\mathbf{b})^T \mathbf{M} (\mathbf{x} - \mathbf{M}^{-1}\mathbf{b}) - \mathbf{b}^T \mathbf{M}^{-1} \mathbf{b} \quad (\text{C.0.1})$$

## References

- [1] Martin J Wainwright and Michael Irwin Jordan. *Graphical models, exponential families, and variational inference*. Now Publishers Inc, 2008.
- [2] Michael Jordan. *The exponential family: Conjugate priors*, (accessed September 11, 2020).
- [3] Roger Grosse et al. *Bayesian Linear Regression*. Available at [https://www.cs.toronto.edu/~rgrosse/courses/csc411\\_f18/slides/lec19-slides.pdf](https://www.cs.toronto.edu/~rgrosse/courses/csc411_f18/slides/lec19-slides.pdf).
- [4] Christopher M Bishop. *Pattern recognition and machine learning*. springer, 2006.
- [5] Andrew Gelman et al. Prior distributions for variance parameters in hierarchical models (comment on article by browne and draper). *Bayesian analysis*, 1(3):515–534, 2006.
- [6] Danny Bickson. Gaussian belief propagation: Theory and application. *arXiv preprint arXiv:0811.2518*, 2008.
- [7] Rahul G Krishnan, Uri Shalit, and David Sontag. Structured inference networks for nonlinear state space models. *arXiv preprint arXiv:1609.09869*, 2016.
- [8] Barbara Englehardt. *Gaussian Models*, (accessed November 22, 2020).
- [9] Gregory Gundersen. *Completing the square*. Available at <http://gregorygundersen.com/blog/2019/09/18/completing-the-square/>.