

Minimising the Kullback–Leibler Divergence

Ralf Herbrich

7th August 2005

Abstract

In this note we show that minimising the Kullback–Leibler divergence over a family in the class of exponential distributions is achieved by matching the *expected natural statistic*. We will also give an explicit update formula for distributions with only one likelihood term.

1 Notation

We use $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ to denote a Gaussian density at \mathbf{x} with a mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$,

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) := (2\pi)^{-\frac{n}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right). \quad (1.1)$$

When dealing one dimensional Gaussians the vectors and matrices are replaced by scalars. If p is a density over \mathbf{x} , we will write $\langle g(\mathbf{x}) \rangle_{p(\mathbf{x})}$ as a shorthand notation for the expectation of g over \mathbf{x} , $\int g(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}$. Finally, the *Kullback–Leibler* divergence between two densities p and q is defined by

$$\text{KL}(p||q) := \left\langle \log\left(\frac{p(\mathbf{x})}{q(\mathbf{x})}\right) \right\rangle_{p(\mathbf{x})}. \quad (1.2)$$

2 Minimising in the Exponential Family

A set of distributions over \mathbb{R}^N is in the exponential family if its densities can be written as

$$p_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{1}{Z(\boldsymbol{\theta})} \exp\left(\boldsymbol{\theta}^T \boldsymbol{\phi}(\mathbf{x})\right),$$

where $\boldsymbol{\phi}(\mathbf{x})$ is known as the *natural statistic* of \mathbf{x} and $Z(\boldsymbol{\theta}) := \int \exp(\boldsymbol{\theta}^T \boldsymbol{\phi}(\mathbf{x})) d\mathbf{x}$ ensures normalisation. The exponential family includes many known families of distributions including the Gaussian distribution. For example, in the Gaussian case, the natural statistic $\boldsymbol{\phi}(\mathbf{x})$ is simply the vector of all first and second moments, $\boldsymbol{\phi}(\mathbf{x}) = (x_1, \dots, x_N, x_1^2, x_1 x_2, \dots, x_N x_{N-1}, x_N^2)$. Note that the expected natural statistic of $p_{\boldsymbol{\theta}}(\mathbf{x})$ is given in terms of the gradient of $\log(Z(\boldsymbol{\theta}))$ w.r.t. $\boldsymbol{\theta}$, that is,

$$\nabla_{\boldsymbol{\theta}} \log(Z(\boldsymbol{\theta})) = \frac{\int [\nabla_{\boldsymbol{\theta}} \exp(\boldsymbol{\theta}^T \boldsymbol{\phi}(\mathbf{x}))] d\mathbf{x}}{Z(\boldsymbol{\theta})} = \langle \boldsymbol{\phi}(\mathbf{x}) \rangle_{p_{\boldsymbol{\theta}}(\mathbf{x})}. \quad (2.1)$$

Theorem 1. For any distribution p , the distribution $p_{\boldsymbol{\theta}^*}$ which minimises the Kullback–Leibler divergence, $\text{KL}(p||p_{\boldsymbol{\theta}^*})$, over the exponential family with natural statistic $\boldsymbol{\phi}$ is implicitly given by

$$\langle \boldsymbol{\phi}(\mathbf{x}) \rangle_{p_{\boldsymbol{\theta}^*}(\mathbf{x})} = \langle \boldsymbol{\phi}(\mathbf{x}) \rangle_{p(\mathbf{x})}. \quad (2.2)$$

Proof. Let us recall the Kullback-Leibler divergence from (1.2) and consider it as a function f of the parameters θ ,

$$\begin{aligned} f(\theta) &= \text{KL}(p||p_\theta) = \left\langle \log \left(\frac{p(\mathbf{x})}{p_\theta(\mathbf{x})} \right) \right\rangle_{p(\mathbf{x})} \\ &= \langle \log(p(\mathbf{x})) \rangle_{p(\mathbf{x})} + \langle \log(Z(\theta)) \rangle_{p(\mathbf{x})} - \langle \theta^T \phi(\mathbf{x}) \rangle_{p(\mathbf{x})} \\ &= \langle \log(p(\mathbf{x})) \rangle_{p(\mathbf{x})} + \log(Z(\theta)) - \theta^T \langle \phi(\mathbf{x}) \rangle_{p(\mathbf{x})} . \end{aligned}$$

Recall that a necessary condition for the minimum θ^* is $\nabla_\theta f(\theta^*) = \mathbf{0}$. From (2.1) we have

$$\nabla_\theta f(\theta) = \langle \phi(\mathbf{x}) \rangle_{p_\theta(\mathbf{x})} - \langle \phi(\mathbf{x}) \rangle_{p(\mathbf{x})} .$$

It remains to show that θ^* such that $\langle \phi(\mathbf{x}) \rangle_{p_{\theta^*}(\mathbf{x})} = \langle \phi(\mathbf{x}) \rangle_{p(\mathbf{x})}$ is a minimum. To this end, consider the matrix of second derivatives,

$$\begin{aligned} [\nabla \nabla_\theta f(\theta)]_{i,j} &= \frac{\partial^2 \log(Z(\theta))}{\partial \theta_i \partial \theta_j} = \frac{\partial}{\partial \theta_j} \frac{\int \phi_i(\mathbf{x}) \exp(\theta^T \phi(\mathbf{x})) d\mathbf{x}}{Z(\theta)} \\ &= \langle \phi_i(\mathbf{x}) \phi_j(\mathbf{x}) \rangle_{p_\theta(\mathbf{x})} - \langle \phi_i(\mathbf{x}) \rangle_{p_\theta(\mathbf{x})} \langle \phi_j(\mathbf{x}) \rangle_{p_\theta(\mathbf{x})} . \end{aligned}$$

At the solution θ^* , this is the covariance matrix of the natural statistic $\phi(\mathbf{x})$ over the distribution p_{θ^*} . By definition, this is positive semi-definite matrix (in fact, for every distribution p_θ) and thus we have proven the theorem. \square

Remark. In the case of the Gaussian family, $\{\mathcal{N}(\cdot; \mu, \Sigma)\}$, Theorem 1 reduces to matching the mean and covariance (which are related in a one-to-one way to the first and second moments),

$$\mu^* = \langle \mathbf{x} \rangle_{p(\mathbf{x})} , \quad (2.3)$$

$$\Sigma^* = \left\langle \mathbf{x} \mathbf{x}^T \right\rangle_{p(\mathbf{x})} - \langle \mathbf{x} \rangle_{p(\mathbf{x})} \langle \mathbf{x} \rangle_{p(\mathbf{x})}^T . \quad (2.4)$$

3 Matching the Bayesian Posterior

We will now derive an explicit update formula for matching the expected natural statistic if $p(\mathbf{x})$ has the simple form

$$p(\mathbf{x}) = \frac{1}{\tilde{Z}(\theta)} \cdot t(\mathbf{x}) p_\theta(\mathbf{x}) ,$$

where $\tilde{Z}(\theta) := \int t(\mathbf{x}) p_\theta(\mathbf{x}) d\mathbf{x}$ ensures normalisation¹. In fact, similar to (2.1), the expected natural statistic under $p(\mathbf{x})$ can again be expressed solely in terms of the gradient of $\tilde{Z}(\theta)$ w.r.t. θ . In order to see this, note that

$$\begin{aligned} \nabla_\theta p_\theta(\mathbf{x}) &= \left[\nabla_\theta \frac{1}{Z(\theta)} \right] \exp(\theta^T \phi(\mathbf{x})) + \frac{1}{Z(\theta)} \left[\nabla_\theta \exp(\theta^T \phi(\mathbf{x})) \right] \\ &= -\frac{[\nabla_\theta Z(\theta)]}{Z(\theta)} p_\theta(\mathbf{x}) + \phi(\mathbf{x}) p_\theta(\mathbf{x}) \\ &= -\langle \phi(\mathbf{x}) \rangle_{p_\theta(\mathbf{x})} \cdot p_\theta(\mathbf{x}) + \phi(\mathbf{x}) p_\theta(\mathbf{x}) . \end{aligned}$$

Multiplying both sides by $\tilde{Z}^{-1}(\theta)t(\mathbf{x})$, integrating over \mathbf{x} and re-arranging terms we get

$$\begin{aligned} \tilde{Z}^{-1}(\theta) \nabla_\theta \tilde{Z}(\theta) &= -\langle \phi(\mathbf{x}) \rangle_{p_\theta(\mathbf{x})} + \langle \phi(\mathbf{x}) \rangle_{p(\mathbf{x})} \\ \langle \phi(\mathbf{x}) \rangle_{p(\mathbf{x})} &= \nabla_\theta \log(\tilde{Z}(\theta)) + \langle \phi(\mathbf{x}) \rangle_{p_\theta(\mathbf{x})} . \end{aligned} \quad (3.1)$$

¹Please note that the normalisation constant $\tilde{Z}(\theta)$ should not be confused with the normalisation constant $Z(\theta)$.

Finally, using Theorem 1 and (2.1) we obtain

$$\nabla_{\theta} \log (Z(\theta^*)) = \nabla_{\theta} \log (\tilde{Z}(\theta)) + \nabla_{\theta} \log (Z(\theta)) .$$

All that is required to solve the above equation for a given exponential family is to know the analytical solution of the gradient equation of $\log(Z(\theta))$ and $\log(\tilde{Z}(\theta))$. These two equations only depend on the particular natural statistic function ϕ and the function t . This is applicable, for example, for Gamma densities.

However, some exponential families are usually not parameterised in terms of θ but rather in terms of $\tau(\theta) := \langle \phi(\mathbf{x}) \rangle_{p_{\theta}(\mathbf{x})}$ —a parameterisation also known as the *moment representation*. This representation has particular advantages when minimising the KL divergence as Theorem 1 directly specifies the update equation for the parameters. In this case, (3.1) can still be used together with the chain rule of differentiation to obtain the update equation for a particular class of exponential densities if the mapping to $\tau \mapsto \theta$ is easy to differentiate. We can also follow the above argument simply in the new parameterisation. In the next section we give a detailed derivation for the Gaussian family (which is represented in terms of its moments).

4 Matching the Bayesian Posterior in the Gaussian Family

We consider a family of Gaussians parameterised in terms of its mean, μ , and covariance, Σ ,

$$q(\mathbf{x}) := q(\mathbf{x}; \mu, \Sigma) := \mathcal{N}(\mathbf{x}; \mu, \Sigma) .$$

Our ability to compute (2.3) and (2.4) when $p(\mathbf{x}) \propto t(\mathbf{x})q(\mathbf{x})$ depends only on the tractability of the normalisation constant,

$$\tilde{Z} := \tilde{Z}(\mu, \Sigma) := \int t(\mathbf{x}) q(\mathbf{x}; \mu, \Sigma) d\mathbf{x} .$$

Matching the Mean We will consider the mean of \mathbf{x} under $t(\mathbf{x})q(\mathbf{x})$. First note that

$$\nabla_{\mu} q(\mathbf{x}) = \Sigma^{-1}(\mathbf{x} - \mu) q(\mathbf{x}) ,$$

which can be re-expressed in terms of $\mathbf{x}q(\mathbf{x})$,

$$\mathbf{x}q(\mathbf{x}) = \mu q(\mathbf{x}) + \Sigma \nabla_{\mu} q(\mathbf{x}) .$$

Now multiplying both sides by $\tilde{Z}^{-1}t(\mathbf{x})$, integrating over \mathbf{x} , and exploiting the linearity of the gradient operator gives

$$\begin{aligned} \langle \mathbf{x} \rangle_{p(\mathbf{x})} &= \mu + \tilde{Z}^{-1} \cdot \Sigma \left[\nabla_{\mu} \int t(\mathbf{x}) q(\mathbf{x}) d\mathbf{x} \right] \\ &= \mu + \tilde{Z}^{-1}(\mu, \Sigma) \cdot \Sigma \nabla_{\mu} \tilde{Z}(\mu, \Sigma) \\ &= \mu + \Sigma \nabla_{\mu} \log(\tilde{Z}(\mu, \Sigma)) \\ &= \mu + \Sigma \mathbf{g} , \end{aligned} \tag{4.1}$$

where we have defined $\mathbf{g} := \nabla_{\mu} \log(\tilde{Z}(\mu, \Sigma))$.

The Second Moment Matrix Once again we take gradients² of $q(\mathbf{x})$, but this time with respect to the covariance matrix Σ ,

$$\nabla_{\Sigma} q(\mathbf{x}) = \frac{1}{2} \left(-\Sigma^{-1} + \Sigma^{-1}(\mathbf{x} - \mu)(\mathbf{x} - \mu)^{\top} \Sigma^{-1} \right) q(\mathbf{x}) ,$$

²It helps to remember that $\nabla_{\Sigma} \log(q(\mathbf{x})) = (q(\mathbf{x}))^{-1} \cdot \nabla_{\Sigma} q(\mathbf{x})$.

which can be re-arranged, as we did before, in order to obtain

$$\mathbf{x}\mathbf{x}^T q(\mathbf{x}) = 2\boldsymbol{\Sigma} [\nabla_{\boldsymbol{\Sigma}} q(\mathbf{x})] \boldsymbol{\Sigma} + \left(\boldsymbol{\Sigma} + \mathbf{x}\boldsymbol{\mu}^T + \boldsymbol{\mu}\mathbf{x}^T - \boldsymbol{\mu}\boldsymbol{\mu}^T \right) q(\mathbf{x}) .$$

Multiplying both sides by $\tilde{Z}^{-1}t(\mathbf{x})$, integrating over \mathbf{x} and exploiting the linearity of the gradient operator gives

$$\begin{aligned} \left\langle \mathbf{x}\mathbf{x}^T \right\rangle_{p(\mathbf{x})} &= \boldsymbol{\Sigma} + 2\boldsymbol{\Sigma} \left(\tilde{Z}^{-1}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \nabla_{\boldsymbol{\Sigma}} \tilde{Z}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \right) \boldsymbol{\Sigma} + \langle \mathbf{x} \rangle_{p(\mathbf{x})} \boldsymbol{\mu}^T + \boldsymbol{\mu} \langle \mathbf{x} \rangle_{p(\mathbf{x})}^T - \boldsymbol{\mu}\boldsymbol{\mu}^T \\ &= \boldsymbol{\Sigma} + 2\boldsymbol{\Sigma} \left(\nabla_{\boldsymbol{\Sigma}} \log \left(\tilde{Z}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \right) \right) \boldsymbol{\Sigma} + \langle \mathbf{x} \rangle_{p(\mathbf{x})} \boldsymbol{\mu}^T + \boldsymbol{\mu} \langle \mathbf{x} \rangle_{p(\mathbf{x})}^T - \boldsymbol{\mu}\boldsymbol{\mu}^T \\ &= \boldsymbol{\Sigma} + 2\boldsymbol{\Sigma} \mathbf{G} \boldsymbol{\Sigma} + \langle \mathbf{x} \rangle_{p(\mathbf{x})} \boldsymbol{\mu}^T + \boldsymbol{\mu} \langle \mathbf{x} \rangle_{p(\mathbf{x})}^T - \boldsymbol{\mu}\boldsymbol{\mu}^T , \end{aligned}$$

where we have defined $\mathbf{G} := \nabla_{\boldsymbol{\Sigma}} \log(\tilde{Z}(\boldsymbol{\mu}, \boldsymbol{\Sigma}))$.

Matching the Covariance The update (2.4) for the covariance requires to compute

$$\left\langle \mathbf{x}\mathbf{x}^T \right\rangle_{p(\mathbf{x})} - \langle \mathbf{x} \rangle_{p(\mathbf{x})} \langle \mathbf{x} \rangle_{p(\mathbf{x})}^T = \boldsymbol{\Sigma} - \boldsymbol{\Sigma} \left(\mathbf{g}\mathbf{g}^T - 2\mathbf{G} \right) \boldsymbol{\Sigma} , \quad (4.2)$$

where we used (4.1). Substituting (4.1) and (4.2) into (2.3) and (2.4) we obtain the required updates for the mean and covariance:

$$\begin{aligned} \boldsymbol{\mu}^* &= \boldsymbol{\mu} + \boldsymbol{\Sigma} \mathbf{g} , \\ \boldsymbol{\Sigma}^* &= \boldsymbol{\Sigma} - \boldsymbol{\Sigma} \left(\mathbf{g}\mathbf{g}^T - 2\mathbf{G} \right) \boldsymbol{\Sigma} . \end{aligned}$$