
Kullback-Leibler Proximal Variational Inference

Mohammad Emtiyaz Khan*

Ecole Polytechnique Fédérale de Lausanne
Lausanne, Switzerland
emtiyaz@gmail.com

Pierre Baqué*

Ecole Polytechnique Fédérale de Lausanne
Lausanne, Switzerland
pierre.baque@epfl.ch

François Fleuret

Idiap Research Institute
Martigny, Switzerland
francois.fleuret@idiap.ch

Pascal Fua

Ecole Polytechnique Fédérale de Lausanne
Lausanne, Switzerland
pascal.fua@epfl.ch

Abstract

We propose a new variational inference method based on a proximal framework that uses the Kullback-Leibler (KL) divergence as the proximal term. We make two contributions towards exploiting the geometry and structure of the variational bound. Firstly, we propose a KL proximal-point algorithm and show its equivalence to variational inference with natural gradients (e.g. stochastic variational inference). Secondly, we use the proximal framework to derive efficient variational algorithms for non-conjugate models. We propose a splitting procedure to separate non-conjugate terms from conjugate ones. We linearize the non-conjugate terms to obtain subproblems that admit a closed-form solution. Overall, our approach converts inference in a non-conjugate model to subproblems that involve inference in well-known conjugate models. We show that our method is applicable to a wide variety of models and can result in computationally efficient algorithms. Applications to real-world datasets show comparable performance to existing methods.

1 Introduction

Variational methods are a popular alternative to Markov Chain Monte Carlo (MCMC) for Bayesian inference. They have been used extensively due to their speed and ease of use. In particular, methods based on the Evidence Lower Bound Optimization (ELBO) are quite popular since they convert a difficult integration problem to an optimization problem. This reformulation allows the application of optimization techniques for large-scale Bayesian inference.

Recently, an approach called Stochastic Variational Inference (SVI) has gained popularity for inference in conditionally-conjugate exponential family models [1]. SVI exploits the geometry of the posterior by using natural gradients, and uses a stochastic method to improve scalability. Resulting updates are simple and easy to implement.

Several generalizations of SVI have been proposed for general latent variable models where the lower bound might be intractable [2, 3, 4]. These generalizations, although important, do not take the geometry of the posterior into account.

In addition, none of these approaches exploit the structure of the lower bound. In practice, not all *factors* of the joint distribution introduce difficulty in the optimization. It is therefore desirable to treat “difficult” terms differently from “easy” terms while optimizing.

*A note on contributions: P. Baqué proposed the use of the KL proximal term and showed that the resulting proximal steps have closed-form solutions. The rest of the work was carried out by M. E. Khan.

In this context, we propose a splitting method for variational inference that exploits both the structure and the geometry of the lower bound. Our approach is based on the proximal-gradient framework. We make two important contributions. First, we propose a proximal-point algorithm that uses the Kullback-Leibler (KL) divergence as the proximal term. We show that addition of this term incorporates the geometry of the posterior. We establish equivalence of our approach to variational methods that use natural gradients (e.g. [1, 5, 6]).

Second, following the proximal-gradient framework, we propose a splitting approach for variational inference. In this approach, we linearize difficult terms such that the resulting optimization problem is easy to solve. We apply this approach to variational inference on non-conjugate models. We show that linearizing non-conjugate terms leads to subproblems that have closed-form solutions. Our approach therefore converts inference in a non-conjugate model to subproblems that involve inference in well-known conjugate models, and for which efficient implementation exists.

2 Latent Variable Models and Evidence Lower Bound Optimization

Consider a general latent variable model with data vector \mathbf{y} of length N and the latent vector \mathbf{z} of length D , following a joint distribution $p(\mathbf{y}, \mathbf{z})$ (we drop the parameters of the distribution from the notation). ELBO approximates the posterior $p(\mathbf{z}|\mathbf{y})$ by a distribution $q(\mathbf{z}|\boldsymbol{\lambda})$ that maximizes a lower bound to the marginal likelihood. Here, $\boldsymbol{\lambda}$ is the vector of parameters of the distribution q . As shown in (1), the lower bound is obtained by first multiplying and dividing by $q(\mathbf{z}|\boldsymbol{\lambda})$, and then applying Jensen’s inequality using concavity of \log . The approximate posterior $q(\mathbf{z}|\boldsymbol{\lambda})$ is obtained by maximizing the lower bound w.r.t. $\boldsymbol{\lambda}$.

$$\log p(\mathbf{y}) = \log \int q(\mathbf{z}|\boldsymbol{\lambda}) \frac{p(\mathbf{y}, \mathbf{z})}{q(\mathbf{z}|\boldsymbol{\lambda})} d\mathbf{z} \geq \max_{\boldsymbol{\lambda}} \mathbb{E}_{q(\mathbf{z}|\boldsymbol{\lambda})} \left[\log \frac{p(\mathbf{y}, \mathbf{z})}{q(\mathbf{z}|\boldsymbol{\lambda})} \right] := \underline{\mathcal{L}}(\boldsymbol{\lambda}). \quad (1)$$

It is desirable to choose $q(\mathbf{z}|\boldsymbol{\lambda})$ such that the lower bound is easy to optimize, however in general this is not the case. This might happen for several reasons, e.g. some terms in the lower bound might be intractable or may admit a form that is not easy to optimize. In addition, the optimization can be slow when N and D are large.

3 The KL Proximal-Point Algorithm for Conjugate Models

In this section, we introduce a proximal-point method based on Kullback-Leibler (KL) proximal function and establish its relation to existing approaches that are based on natural gradients [1, 5, 6]. In particular, for conditionally-conjugate exponential-family models, we show that each iteration of our proximal-point approach is equivalent to a step along the natural gradient.

The Kullback-Leibler (KL) divergence between two distributions $q(\mathbf{z}|\boldsymbol{\lambda})$ and $q(\mathbf{z}|\boldsymbol{\lambda}')$ is defined as follows: $\mathbb{D}_{KL}[q(\mathbf{z}|\boldsymbol{\lambda}) \| q(\mathbf{z}|\boldsymbol{\lambda}')] := \mathbb{E}_{q(\mathbf{z}|\boldsymbol{\lambda})}[\log q(\mathbf{z}|\boldsymbol{\lambda}) - \log q(\mathbf{z}|\boldsymbol{\lambda}')]$. Using the KL divergence as the proximal term, we introduce a proximal-point algorithm which generates a sequence of $\boldsymbol{\lambda}_k$ by solving the following subproblems:

$$\text{KL Proximal-Point : } \boldsymbol{\lambda}_{k+1} = \arg \max_{\boldsymbol{\lambda}} \underline{\mathcal{L}}(\boldsymbol{\lambda}) - \frac{1}{\beta_k} \mathbb{D}_{KL}[q(\mathbf{z}|\boldsymbol{\lambda}) \| q(\mathbf{z}|\boldsymbol{\lambda}_k)], \quad (2)$$

given an initial value $\boldsymbol{\lambda}_0$ and a bounded sequence of step-size $\beta_k > 0$,

A benefit of using the KL term is that it takes the geometry of the posterior into account. This fact has led to their extensive use in both optimization and statistics literature, e.g. to speed up expectation-maximization algorithm [7, 8], for convex optimization [9], for message-passing in graphical models [10], and for approximate Bayesian inference [11, 12, 13].

Relationship to the methods that use natural gradients: An alternative approach to incorporate the geometry of the posterior is to use natural gradients [6, 5, 1]. We now establish the relationship of our approach to this approach. The natural gradient can be interpreted as finding a descent direction that ensures a fixed amount of change in the distribution. For variational inference, this is equivalent to the following [1, 14]:

$$\arg \max_{\Delta \boldsymbol{\lambda}} \underline{\mathcal{L}}(\boldsymbol{\lambda}_k + \Delta \boldsymbol{\lambda}), \text{ s.t. } \mathbb{D}_{KL}^{sym}[q(\mathbf{z}|\boldsymbol{\lambda}_k + \Delta \boldsymbol{\lambda}) \| q(\mathbf{z}|\boldsymbol{\lambda}_k)] \leq \epsilon, \quad (3)$$

where \mathbb{D}_{KL}^{sym} is the symmetric KL divergence. It appears that the proximal-point subproblem (2) might be related to a Lagrangian of the above optimization. In fact, as we show below, the two problems are equivalent for conditionally conjugate exponential family models.

We consider the set-up described in [15] which is a bit more general than that of [1]. Consider a Bayesian network with nodes \mathbf{z}_i and a joint-distribution $\prod_i p(\mathbf{z}_i|\mathbf{pa}_i)$ where \mathbf{pa}_i are parents of \mathbf{z}_i . We assume that each factor is an exponential family distribution defined as follows:

$$p(\mathbf{z}_i|\mathbf{pa}_i) := h_i(\mathbf{z}_i) \exp [\boldsymbol{\eta}_i^T(\mathbf{pa}_i) \mathbf{T}_i(\mathbf{z}_i) - A_i(\boldsymbol{\eta}_i)], \quad (4)$$

where $\boldsymbol{\eta}_i$ is the natural parameter, $\mathbf{T}_i(\mathbf{z}_i)$ is the sufficient statistics, $A_i(\boldsymbol{\eta}_i)$ is the partition function and $h_i(\mathbf{z}_i)$ is the base measure. We seek a factorized approximation shown in (5) where each \mathbf{z}_i belongs to the same exponential family as the joint. The parameters of this distribution are denoted by $\boldsymbol{\lambda}_i$ to differentiate them from the joint-distribution parameters $\boldsymbol{\eta}_i$. Also note that the subscript refers to the factor i not to the iteration.

$$q(\mathbf{z}|\boldsymbol{\lambda}) = \prod_i q_i(\mathbf{z}_i|\boldsymbol{\lambda}_i), \quad \text{where } q_i(\mathbf{z}_i) := h_i(\mathbf{z}_i) \exp [\boldsymbol{\lambda}_i^T \mathbf{T}_i(\mathbf{z}_i) - A_i(\boldsymbol{\lambda}_i)]. \quad (5)$$

For this model, we show the following equivalence between a gradient-descent method based on natural gradients and our proximal-point approach. The proof is given in the supplementary material.

Theorem 1. *For the model shown in (4) and the posterior approximation shown in (5), the sequence $\boldsymbol{\lambda}_k$ generated by the proximal-point algorithm of (2) is equal to the one obtained using gradient-descent along the natural gradient with step lengths $\beta_k/(1 + \beta_k)$.*

Proof of convergence : Convergence of the proximal-point algorithm shown in (2) is proved in [8]. We give a summary of the results here. We assume $\beta_k = 1$, however the proof holds for any bounded sequence of β_k . Let the space of all $\boldsymbol{\lambda}$ be denoted by \mathcal{S} . Define the set $\mathcal{S}_0 := \{\boldsymbol{\lambda} \in \mathcal{S} : \underline{\mathcal{L}}(\boldsymbol{\lambda}) \geq \underline{\mathcal{L}}(\boldsymbol{\lambda}_0)\}$. Then, $\|\boldsymbol{\lambda}_{k+1} - \boldsymbol{\lambda}_k\| \rightarrow 0$ under the following conditions:

- (A) Maximum of $\underline{\mathcal{L}}$ exist and the gradient of $\underline{\mathcal{L}}$ is continuous and defined in \mathcal{S}_0 .
- (B) The KL divergence and its gradient are continuous and defined in $\mathcal{S}_0 \times \mathcal{S}_0$.
- (C) $\mathbb{D}_{KL}[q(\mathbf{z}|\boldsymbol{\lambda}) \| q(\mathbf{z}|\boldsymbol{\lambda}')] = 0$ only when $\boldsymbol{\lambda}' = \boldsymbol{\lambda}$.

In our case, conditions (A) and (B) are either assumed or satisfied, while condition (C) can be ensured by choosing an appropriate parameterization of q .

4 The KL Proximal-Gradient Algorithm for Non-Conjugate Models

The proximal-point algorithm of (2) might be difficult to optimize for non-conjugate models, e.g., due to the non-conjugate factors. In this section, we present an algorithm based on the proximal-gradient framework where we first split the objective function into “difficult” and “easy” terms, and then linearize the difficult term to simplify the optimization. See [16] for a good review of proximal methods for machine learning.

We split the ratio $p(\mathbf{y}, \mathbf{z})/q(\mathbf{z}|\boldsymbol{\lambda}) \equiv c \tilde{p}_d(\mathbf{z}|\boldsymbol{\lambda}) \tilde{p}_e(\mathbf{z}|\boldsymbol{\lambda})$, where \tilde{p}_d contains all factors that make the optimization difficult while \tilde{p}_e contains the rest (c is a constant). This result in the following split for the lower bound:

$$\underline{\mathcal{L}}(\boldsymbol{\lambda}) = \mathbb{E}_{q(\mathbf{z}|\boldsymbol{\lambda})} \left[\log \frac{p(\mathbf{y}, \mathbf{z}|\boldsymbol{\theta})}{q(\mathbf{z}|\boldsymbol{\lambda})} \right] := \underbrace{\mathbb{E}_{q(\mathbf{z}|\boldsymbol{\lambda})} [\log \tilde{p}_d(\mathbf{z}|\boldsymbol{\lambda})]}_{f(\boldsymbol{\lambda})} + \underbrace{\mathbb{E}_{q(\mathbf{z}|\boldsymbol{\lambda})} [\log \tilde{p}_e(\mathbf{z}|\boldsymbol{\lambda})]}_{h(\boldsymbol{\lambda})} + \log c, \quad (6)$$

Note that \tilde{p}_d and \tilde{p}_e can be un-normalized factors in the distribution. In the worst case, we may set $\tilde{p}_e(\mathbf{z}|\boldsymbol{\lambda}) \equiv 1$ and take the rest as $\tilde{p}_d(\mathbf{z}|\boldsymbol{\lambda})$. We give an example of the split in the next section.

The main idea is to linearize the difficult term f such that the resulting problem admits a simple form. Specifically, we use a proximal-gradient algorithm that solves the following sequence of subproblems to maximize $\underline{\mathcal{L}}$ as shown below. Here, $\nabla f(\boldsymbol{\lambda}_k)$ is the gradient of f at $\boldsymbol{\lambda}_k$.

$$\text{KL Proximal-Gradient: } \boldsymbol{\lambda}_{k+1} = \arg \max_{\boldsymbol{\lambda}} \boldsymbol{\lambda}^T \nabla f(\boldsymbol{\lambda}_k) + h(\boldsymbol{\lambda}) + \frac{1}{\beta_k} \mathbb{D}_{KL}[q(\mathbf{z}|\boldsymbol{\lambda}) \| q(\mathbf{z}|\boldsymbol{\lambda}_k)]. \quad (7)$$

Note that our linear approximation is equivalent to the one used in gradient descent. Also, the approximation is tight at λ_k . Therefore, it does not introduce any error in the optimization, rather only acts as a surrogate to take the next step. Existing variational methods have used approximations such as ours, e.g. see [17, 18, 19]. Most of these methods first approximate the $\log \tilde{p}_d(\mathbf{z}|\lambda)$ term using a linear or quadratic approximation and then compute the expectation. As a result the approximation is not tight and may result in a bad performance [20]. In contrast, our approximation is applied directly to f and therefore is tight at λ_k .

Convergence of our approach is covered under the results shown in [21] which proves convergence of a more general algorithm than ours. We summarize the results here. As before, we assume that the maximum exists and \mathcal{L} is continuous. We make three additional assumptions. First, the gradient of f is L -Lipschitz continuous in \mathcal{S} , i.e., $\|\nabla f(\lambda) - \nabla f(\lambda')\| \leq L\|\lambda - \lambda'\|$, $\forall \lambda, \lambda' \in \mathcal{S}$. Second, the function h is concave. Third, there exist an $\alpha > 0$ such that,

$$(\lambda_{k+1} - \lambda_k)^T \nabla_1 \mathbb{D}_{KL}[q(\mathbf{z}|\lambda_{k+1}) \| q(\mathbf{z}|\lambda_k)] \geq \alpha \|\lambda_{k+1} - \lambda_k\|^2, \quad (8)$$

where ∇_1 denotes the gradient w.r.t. the first argument. Under these conditions, $\|\lambda_{k+1} - \lambda_k\| \rightarrow 0$ when $0 < \beta_k < \alpha/L$. The choice of constant α is also discussed in [21].

Note that even though h is required to be concave, f could be nonconvex. The lower bound usually contains concave terms, e.g., in the entropy term. In the worst case when there are no concave terms, we can simply choose $h \equiv 0$.

5 Examples of KL Proximal-Gradient Variational Inference

In this section, we show a few examples where the subproblem (7) has a closed form solution.

Generalized linear model : We consider the generalized linear model shown in (9). Here, \mathbf{y} is the output vector (of length N) with its n 'th entry equal to y_n , while \mathbf{X} is an $N \times D$ feature matrix that contains feature vectors \mathbf{x}_n^T as rows. The weight vector \mathbf{z} is a Gaussian with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$. The linear predictor $\mathbf{x}_n^T \mathbf{z}$ is passed through $p(y_n|\cdot)$ to obtain the probability of y_n .

$$p(\mathbf{y}, \mathbf{z}) := \prod_{n=1}^N p(y_n | \mathbf{x}_n^T \mathbf{z}) \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}, \boldsymbol{\Sigma}). \quad (9)$$

We restrict the posterior to be a Gaussian $q(\mathbf{z}|\lambda) = \mathcal{N}(\mathbf{z}|\mathbf{m}, \mathbf{V})$ with mean \mathbf{m} and covariance \mathbf{V} , therefore $\lambda := \{\mathbf{m}, \mathbf{V}\}$. For this posterior family, the non-Gaussian terms $p(y_n|\mathbf{x}_n^T \mathbf{z})$ are difficult to handle, while the Gaussian term $\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is easy since it is conjugate to q . Therefore, we set $\tilde{p}_e(\mathbf{z}|\lambda) \equiv \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma})/\mathcal{N}(\mathbf{z}|\mathbf{m}, \mathbf{V})$ and let the rest of the terms go in \tilde{p}_d . The constant c is set to 1.

Substituting in (6) and using the definition of the KL divergence, we get the lower bound shown below in (10). The first term is the function f that will be linearized, and the second term is the function h .

$$\mathcal{L}(\mathbf{m}, \mathbf{V}) := \underbrace{\sum_{n=1}^N \mathbb{E}_{q(\mathbf{z}|\lambda)} [\log p(y_n | \mathbf{x}_n^T \mathbf{z})]}_{f(\mathbf{m}, \mathbf{V})} + \underbrace{\mathbb{E}_{q(\mathbf{z}|\lambda)} \left[\log \frac{\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma})}{\mathcal{N}(\mathbf{z}|\mathbf{m}, \mathbf{V})} \right]}_{h(\mathbf{m}, \mathbf{V})}. \quad (10)$$

For linearization, we compute the gradient of f using the chain rule. Denote $f_n(\tilde{m}_n, \tilde{v}_n) := \mathbb{E}_{q(\mathbf{z}|\lambda)} [\log p(y_n | \mathbf{x}_n^T \mathbf{z})]$ where $\tilde{m}_n := \mathbf{x}_n^T \mathbf{m}$ and $\tilde{v}_n := \mathbf{x}_n^T \mathbf{V} \mathbf{x}_n$. Gradients of f w.r.t. \mathbf{m} and \mathbf{V} can then be expressed in terms of gradients of f_n w.r.t. \tilde{m}_n and \tilde{v}_n :

$$\nabla_{\mathbf{m}} f(\mathbf{m}, \mathbf{V}) = \sum_{n=1}^N \mathbf{x}_n \nabla_{\tilde{m}_n} f_n(\tilde{m}_n, \tilde{v}_n), \quad \nabla_{\mathbf{V}} f(\mathbf{m}, \mathbf{V}) = \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \nabla_{\tilde{v}_n} f_n(\tilde{m}_n, \tilde{v}_n), \quad (11)$$

For notational simplicity, we denote the gradient of f_n at $\tilde{m}_{nk} := \mathbf{x}_n^T \mathbf{m}_k$ and $\tilde{v}_{nk} := \mathbf{x}_n^T \mathbf{V}_k \mathbf{x}_n$ by,

$$\alpha_{nk} := -\nabla_{\tilde{m}_n} f_n(\tilde{m}_{nk}, \tilde{v}_{nk}), \quad \gamma_{nk} := -2 \nabla_{\tilde{v}_n} f_n(\tilde{m}_{nk}, \tilde{v}_{nk}). \quad (12)$$

Using (11) and (12), we get the following linear approximation of f :

$$f(\mathbf{m}, \mathbf{V}) \approx \mathbf{m}^T [\nabla_{\mathbf{m}} f(\mathbf{m}_k, \mathbf{V}_k)] + \text{Tr}[\mathbf{V} \{\nabla_{\mathbf{V}} f(\mathbf{m}_k, \mathbf{V}_k)\}] \quad (13)$$

$$= -\sum_{n=1}^N [\alpha_{nk} (\mathbf{x}_n^T \mathbf{m}) + \frac{1}{2} \gamma_{nk} (\mathbf{x}_n^T \mathbf{V} \mathbf{x}_n)]. \quad (14)$$

Substituting the above in (7), we get the following subproblem in the k 'th iteration:

$$(\mathbf{m}_{k+1}, \mathbf{V}_{k+1}) = \arg \max_{\mathbf{m}, \mathbf{V} \succ 0} - \sum_{n=1}^N [\alpha_{nk} (\mathbf{x}_n^T \mathbf{m}) + \frac{1}{2} \gamma_{nk} (\mathbf{x}_n^T \mathbf{V} \mathbf{x}_n)] + \mathbb{E}_{q(\mathbf{z}|\lambda)} \left[\frac{\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma})}{\mathcal{N}(\mathbf{z}|\mathbf{m}, \mathbf{V})} \right] - \frac{1}{\beta_k} \mathbb{D}_{KL} [\mathcal{N}(\mathbf{z}|\mathbf{m}, \mathbf{V}) || \mathcal{N}(\mathbf{z}|\mathbf{m}_k, \mathbf{V}_k)], \quad (15)$$

Taking the gradient w.r.t. \mathbf{m} and \mathbf{V} and setting it to zero, we get the following closed-form solutions (details are given in the supplementary material):

$$\mathbf{V}_{k+1}^{-1} = r_k \mathbf{V}_k^{-1} + (1 - r_k) \left[\boldsymbol{\Sigma}^{-1} + \mathbf{X}^T \text{diag}(\boldsymbol{\gamma}_k) \mathbf{X} \right], \quad (16)$$

$$\mathbf{m}_{k+1} = \left[(1 - r_k) \boldsymbol{\Sigma}^{-1} + r_k \mathbf{V}_k^{-1} \right]^{-1} \left[(1 - r_k) (\boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} - \mathbf{X}^T \boldsymbol{\alpha}_k) + r_k \mathbf{V}_k^{-1} \mathbf{m}_k \right], \quad (17)$$

where $r_k := 1/(1 + \beta_k)$ and $\boldsymbol{\alpha}_k$ and $\boldsymbol{\gamma}_k$ are vectors of α_{nk} and γ_{nk} respectively, for all k .

Computationally efficient updates : Even though the updates are available in closed-form, they are not efficient when dimensionality D is large. In such a case, an explicit computation of \mathbf{V} is costly since the resulting $D \times D$ matrix is extremely large. We now derive a reformulation that avoids an explicit computation of \mathbf{V} .

Our reformulation involves two key steps. The first step is to show that \mathbf{V}_{k+1} can be parameterized by $\boldsymbol{\gamma}_k$. Specifically, if we initialize $\mathbf{V}_0 = \boldsymbol{\Sigma}$, then we can show that:

$$\mathbf{V}_{k+1} = \left[\boldsymbol{\Sigma}^{-1} + \mathbf{X}^T \text{diag}(\tilde{\boldsymbol{\gamma}}_{k+1}) \mathbf{X} \right]^{-1}, \text{ where } \tilde{\boldsymbol{\gamma}}_{k+1} = r_k \tilde{\boldsymbol{\gamma}}_k + (1 - r_k) \boldsymbol{\gamma}_k. \quad (18)$$

with $\tilde{\boldsymbol{\gamma}}_0 = \boldsymbol{\gamma}_0$. A detailed derivation is given in the supplementary material.

The second key step is to express the updates in terms of \tilde{m}_n and \tilde{v}_n . For this purpose, we define some new quantities. Define $\tilde{\mathbf{m}}$ to be a vector with \tilde{m}_n as its n 'th entry. Similarly, let $\tilde{\mathbf{v}}$ be the vector of \tilde{v}_n for all n . Denote the corresponding vectors in the k 'th iteration by $\tilde{\mathbf{m}}_k$ and $\tilde{\mathbf{v}}_k$, respectively. Finally, define $\tilde{\boldsymbol{\mu}} = \mathbf{X} \boldsymbol{\mu}$ and $\tilde{\boldsymbol{\Sigma}} = \mathbf{X} \boldsymbol{\Sigma} \mathbf{X}^T$.

Now, using the fact that $\tilde{\mathbf{m}} = \mathbf{X} \mathbf{m}$ and $\tilde{\mathbf{v}} = \text{diag}(\mathbf{X} \mathbf{V} \mathbf{X}^T)$ and by applying Woodbury matrix identity, we can express the updates in terms of $\tilde{\mathbf{m}}$ and $\tilde{\mathbf{v}}$, as shown below (a detailed derivation is given in the supplementary material):

$$\begin{aligned} \tilde{\mathbf{m}}_{k+1} &= \tilde{\mathbf{m}}_k + (1 - r_k) (\mathbf{I} - \tilde{\boldsymbol{\Sigma}} \mathbf{B}_k^{-1}) (\tilde{\boldsymbol{\mu}} - \tilde{\mathbf{m}}_k - \tilde{\boldsymbol{\Sigma}} \boldsymbol{\alpha}_k), \text{ where } \mathbf{B}_k := \tilde{\boldsymbol{\Sigma}} + [\text{diag}(r_k \tilde{\boldsymbol{\gamma}}_k)]^{-1}, \\ \tilde{\mathbf{v}}_{k+1} &= \text{diag}(\tilde{\boldsymbol{\Sigma}}) - \text{diag}(\tilde{\boldsymbol{\Sigma}} \mathbf{A}_k^{-1} \tilde{\boldsymbol{\Sigma}}), \text{ where } \mathbf{A}_k := \tilde{\boldsymbol{\Sigma}} + [\text{diag}(\tilde{\boldsymbol{\gamma}}_k)]^{-1}. \end{aligned} \quad (19)$$

Note that these updates depend on $\tilde{\boldsymbol{\mu}}$, $\tilde{\boldsymbol{\Sigma}}$, $\boldsymbol{\alpha}_k$, and $\boldsymbol{\gamma}_k$, whose size only depends on N and is independent of D . Most importantly, these updates avoid an explicit computation of \mathbf{V} and only requires storing $\tilde{\mathbf{m}}_k$ and $\tilde{\mathbf{v}}_k$ both of which scale linearly with N .

Also note that the matrix \mathbf{A}_k and \mathbf{B}_k differ only slightly and we can reduce computation by using \mathbf{A}_k in place of \mathbf{B}_k . In our experiments, this does not give any convergence issues.

To assess convergence, we can use the optimality condition. By taking the norm of derivative of \mathcal{L} at \mathbf{m}_{k+1} and \mathbf{V}_{k+1} and simplifying, we get the following criteria: $\|\tilde{\boldsymbol{\mu}} - \tilde{\mathbf{m}}_{k+1} - \tilde{\boldsymbol{\Sigma}} \boldsymbol{\alpha}_{k+1}\|_2^2 + \text{Tr}[\tilde{\boldsymbol{\Sigma}} \{ \text{diag}(\tilde{\boldsymbol{\gamma}}_k - \boldsymbol{\gamma}_{k+1} - \mathbf{1}) \} \tilde{\boldsymbol{\Sigma}}] \leq \epsilon$, for some $\epsilon > 0$. (derivation in the supplementary material).

Linear Basis Function Model and Gaussian Process : The algorithm presented above can be extended to linear basis function models using the *weight-space view* presented in [22]. Consider a non-linear basis function $\phi(\mathbf{x})$ that maps a D -dimensional feature vector into an N -dimensional feature space. The generalized linear model of (9) is extended to a linear basis function model by replacing $\mathbf{x}_n^T \mathbf{z}$ with the *latent function* $g(\mathbf{x}) := \phi(\mathbf{x})^T \mathbf{z}$. The Gaussian prior on \mathbf{z} then translates to a kernel function $\kappa(\mathbf{x}, \mathbf{x}') := \phi(\mathbf{x})^T \boldsymbol{\Sigma} \phi(\mathbf{x}')$ and a mean function $\tilde{\boldsymbol{\mu}}(\mathbf{x}) := \phi(\mathbf{x})^T \boldsymbol{\mu}$ in the latent function space. Given input vectors \mathbf{x}_n , we define the kernel matrix $\tilde{\boldsymbol{\Sigma}}$ whose (i, j) 'th entry is equal to $\kappa(\mathbf{x}_i, \mathbf{x}_j)$ and the mean vector $\tilde{\boldsymbol{\mu}}$ whose i 'th entry is $\tilde{\boldsymbol{\mu}}(\mathbf{x}_i)$.

Assuming a Gaussian posterior over the latent function $g(\mathbf{x})$, we can compute its mean $\tilde{m}(\mathbf{x})$ and variance $\tilde{v}(\mathbf{x})$ using the proximal-gradient algorithm. We define $\tilde{\mathbf{m}}$ to be the vector of $\tilde{m}(\mathbf{x}_n)$ for

Algorithm 1 Proximal-gradient algorithm for linear basis function models and Gaussian process

Given: Training data (\mathbf{y}, \mathbf{X}) , test data \mathbf{x}_* , kernel mean $\tilde{\boldsymbol{\mu}}$, covariance $\tilde{\boldsymbol{\Sigma}}$, step-size sequence r_k , and threshold ϵ .

Initialize: $\tilde{\mathbf{m}}_0 \leftarrow \tilde{\boldsymbol{\mu}}$, $\tilde{\mathbf{v}}_0 \leftarrow \text{diag}(\tilde{\boldsymbol{\Sigma}})$ and $\tilde{\boldsymbol{\gamma}}_0 \leftarrow \delta_1 \mathbf{1}$.

repeat

For all n in parallel: $\alpha_{nk} \leftarrow \nabla_{\tilde{m}_n} f_n(\tilde{m}_{nk}, \tilde{v}_{nk})$ and $\gamma_{nk} \leftarrow \nabla_{\tilde{v}_n} f_n(\tilde{m}_{nk}, \tilde{v}_{nk})$.

Update $\tilde{\mathbf{m}}_k$ and $\tilde{\mathbf{v}}_k$ using (19).

$\tilde{\boldsymbol{\gamma}}_{k+1} \leftarrow r_k \tilde{\boldsymbol{\gamma}}_k + (1 - r_k) \boldsymbol{\gamma}_k$.

until $\|\tilde{\boldsymbol{\mu}} - \tilde{\mathbf{m}}_k - \tilde{\boldsymbol{\Sigma}} \boldsymbol{\alpha}_k\| + \text{Tr}[\tilde{\boldsymbol{\Sigma}} \text{diag}(\tilde{\boldsymbol{\gamma}}_k - \boldsymbol{\gamma}_{k+1} - \mathbf{1}) \tilde{\boldsymbol{\Sigma}}] > \epsilon$.

Predict test inputs \mathbf{x}_* using (20).

all n and similarly $\tilde{\mathbf{v}}$ to be the vector of all $\tilde{v}(\mathbf{x}_n)$. Following the same derivation as the previous section, we can show that the updates of (19) give us the posterior mean $\tilde{\mathbf{m}}$ and variance $\tilde{\mathbf{v}}$. These updates are the *kernalized* version of (16) and (17).

For prediction, we only need the converged value of $\boldsymbol{\alpha}_k$ and $\boldsymbol{\gamma}_k$, denoted by $\boldsymbol{\alpha}^*$ and $\boldsymbol{\gamma}^*$, respectively. Given a new input \mathbf{x}_* , define $\kappa_{**} := \kappa(\mathbf{x}_*, \mathbf{x}_*)$ and $\boldsymbol{\kappa}_*$ to be a vector with n 'th entry equal to $\kappa(\mathbf{x}_n, \mathbf{x}_*)$. The predictive mean and variance can be computed as shown below:

$$\tilde{v}(\mathbf{x}_*) = \kappa_{**} - \boldsymbol{\kappa}_*^T [\tilde{\boldsymbol{\Sigma}} + (\text{diag}(\tilde{\boldsymbol{\gamma}}^*))^{-1}]^{-1} \boldsymbol{\kappa}_* \quad , \quad \tilde{m}(\mathbf{x}_*) = \tilde{\boldsymbol{\mu}} - \boldsymbol{\kappa}_*^T \boldsymbol{\alpha}^* \quad (20)$$

The final algorithm is shown in Algorithm 1. Here, we initialize $\tilde{\boldsymbol{\gamma}}$ to a small constant δ_1 , since otherwise solving the first equation may not be well-conditioned.

It is straightforward to see that these updates also work for Gaussian process (GP) with a generic kernel $k(\mathbf{x}, \mathbf{x}')$ and mean function $\tilde{\boldsymbol{\mu}}(\mathbf{x})$ and many other latent Gaussian models.

6 Experiments and Results

We now present some results on the real data. Our goal is to show that our approach gives comparable results to existing methods, while being easy to implement. We also show that, in some cases, our method is significantly faster than the alternatives due to the kernel trick.

We show results on three models: Bayesian logistic regression, GP classification with logistic likelihood, and GP regression with Laplace likelihood. For these likelihoods, expectations can be computed (almost) exactly. Specifically, we used the methods described in [23, 24]. We use a fixed step-size of $\beta_k = 0.25$ and 1 for logistic and Laplace likelihoods respectively.

We consider three datasets for each model. A summary is given in Table 1. These datasets can be found at data repository¹ of LIBSVM and UCI.

Bayesian Logistic Regression: Results for Bayesian logistic regression are shown in Table 2. We consider three datasets of various sizes. For ‘a1a’, $N > D$, for ‘Colon’, $N < D$ and for ‘gisette’ $N \approx D$. We compare our ‘proximal’ method to 3 other existing methods: ‘MAP’ method which finds the mode of the penalized log-likelihood, ‘Mean-Field’ method where the posterior is factorized across dimensions, and ‘Cholesky’ method of [25]. We implemented these methods using ‘minFunc’ software by Mark Schmidt². We used L-BFGS for optimization. All algorithms are stopped when optimality condition is below 10^{-4} . We set the Gaussian prior to $\boldsymbol{\Sigma} = \delta \mathbf{I}$ and $\boldsymbol{\mu} = \mathbf{0}$. To set the hyperparameter δ , we use cross-validation for MAP, and maximum marginal-likelihood estimate for the rest of the methods. Since we compare running time as well, we use a common set of hyperparameter values for a fair comparison. The values are shown in Table 1.

For Bayesian methods, we report the negative of the marginal likelihood approximation (‘Neg-Log-Lik’). This is (the negative of) the value of the lower bound at the maximum. We also report the log-loss computed as follows: $-\sum_n \log \hat{p}_n / N$ where \hat{p}_n are the predictive probabilities of the test data and N is the total number of test-pairs. A lower value is better and a value of 1 is equivalent to random coin-flipping. In addition, we report the total time taken for hyperparameter selection.

¹<https://archive.ics.uci.edu/ml/datasets.html> and <http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>

²Available at <https://www.cs.ubc.ca/~schmidtm/Software/minFunc.html>

| Model | Dataset | N | D | %Train | #Splits | Hyperparameter range |
|----------|------------|--------|-------|--------|---------|---|
| LogReg | ala | 32,561 | 123 | 5% | 1 | $\delta = \text{logspace}(-3,1,30)$ |
| | Colon | 62 | 2000 | 50% | 10 | $\delta = \text{logspace}(0,6,30)$ |
| | Gisette | 7,000 | 5,000 | 50% | 1 | $\delta = \text{logspace}(0,6,30)$ |
| GP class | Ionosphere | 351 | 34 | 50% | 10 | for all datasets |
| | Sonar | 208 | 60 | 50% | 10 | $\log(l) = \text{linspace}(-1,6,15)$ |
| | USPS-3vs5 | 1,540 | 256 | 50% | 5 | $\log(\sigma) = \text{linspace}(-1,6,15)$ |
| GP reg | Housing | 506 | 13 | 50% | 10 | $\log(l) = \text{linspace}(-1,6,15)$ |
| | Triazines | 186 | 60 | 50% | 10 | $\log(\sigma) = \text{linspace}(-1,6,15)$ |
| | Space_ga | 3,106 | 6 | 50% | 1 | $\log(b) = \text{linspace}(-5,1,2)$ |

Table 1: A list of models and datasets. %Train is the % of training data. The last column shows the hyperparameters values (‘linspace’ and ‘logspace’ refer to Matlab commands).

| Dataset | Methods | Neg-Log-Lik | Log Loss | Time |
|---------|------------|--------------|-------------|------------|
| ala | MAP | — | 0.499 | 27s |
| | Mean-Field | 792.8 | 0.505 | 21s |
| | Cholesky | 590.1 | 0.488 | 12m |
| | Proximal | 590.1 | 0.488 | 7m |
| Colon | MAP | — | 0.78 (0.01) | 7s (0.00) |
| | Mean-Field | 18.35 (0.11) | 0.78 (0.01) | 15m (0.04) |
| | Proximal | 15.82 (0.13) | 0.70 (0.01) | 18m (0.14) |
| Gisette | MAP | — | 0.112 | 5s |
| | Mean-Field | 1275.7 | 0.258 | 22m |
| | Proximal | 608.5 | 0.140 | 13h |

Table 2: A summary of the results obtained on Bayesian logistic regression. In all columns, a lower values implies better performance.

For MAP, this is the total cross-validation time, while for Bayesian methods it is the time taken to compute ‘Neg-Log-Lik’ for all hyperparameters values.

We summarize these results in Table 2. For all columns, a lower value is better. We see that for ‘ala’, fully Bayesian methods perform slightly better than MAP. More importantly, Proximal method is faster than Cholesky method while obtaining the same error and marginal likelihood estimate. For Proximal method, we use updates of (17) and (16) since $D \ll N$, but even in this scenario, Cholesky method is slow due to expensive line-search for large number of parameters (of the order $O(D^2)$).

For ‘Colon’ and ‘gisette’ datasets, we use the update (19) for Proximal method. Since Cholesky method is too slow for these large datasets, we do not compare to it. In Table 2, we see that for ‘Colon’ dataset our implementation is as fast as Mean-Field while performing significantly better. For ‘gisette’ dataset, however, our method is slow since both N and D are big.

Overall, we see that with the proximal approach we achieve same results as Cholesky method, while taking much less time. In some cases, we can match the running time of Mean-Field method. Note that Mean-Field does not give bad predictions overall, and the minimum value of log-loss are comparable to our approach. However, since Neg-Log-Lik values for Mean-Field are inaccurate, it ends up choosing a bad hyperparameter value. This is expected since Mean-Field makes an extreme approximation. Therefore, cross-validation is more appropriate for Mean-Field.

Gaussian process classification and regression: We compare Proximal method to expectation propagation (EP) and Laplace approximation. We use the GPML toolbox for this comparison. We used a Squared-Exponential Kernel for Gaussian process with two scale parameters σ and l (as defined in GPML toolbox). We do a grid search over these hyperparameters. The grid values are given in Table 1. We report the log-loss and running time for each method.

The left plot in Figure 1 shows the log-loss for GP classification on USPS_3vs5 dataset, where Proximal method shows very similar behaviour to EP. Results are summarized in Table 3. We see that our method performs similar to EP, sometimes a bit better. The running times of EP and Proximal are

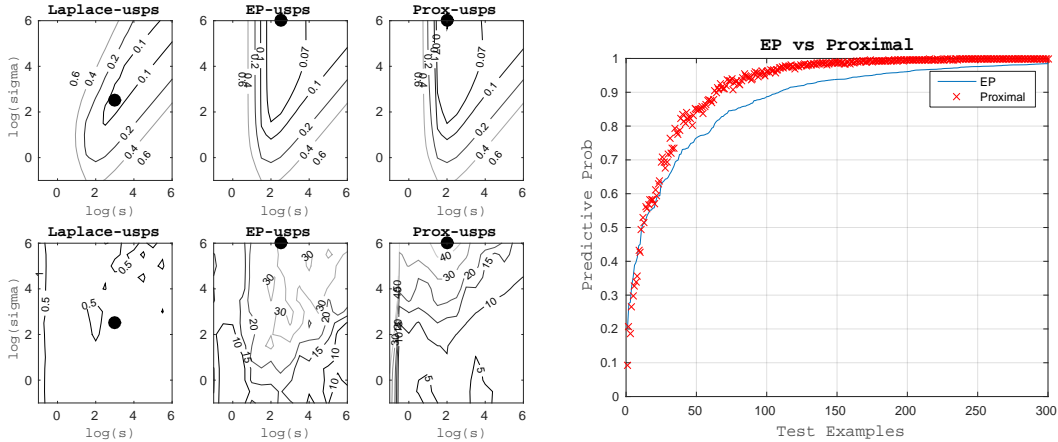


Figure 1: In the left figure, the top row shows the log-loss and the bottom row shows the running time in seconds for ‘USPS_3vs5’ dataset. In each plot, the minimum value of the log-loss is shown with a black circle. The right figure shows the predictive probabilities obtained with EP and Proximal method (a higher value implies better performance).

| Data | Log Loss | | | Time (s is sec, m is min, h is hr) | | |
|------------|-------------|-------------|-------------|------------------------------------|------------|------------|
| | Laplace | EP | Proximal | Laplace | EP | Proximal |
| Ionosphere | .285 (.002) | .234 (.002) | .230 (.002) | 10s (.3) | 3.8m (.10) | 3.6m (.10) |
| Sonar | .410 (.002) | .341 (.003) | .317 (.004) | 4s (.01) | 45s (.01) | 63s (.13) |
| USPS-3vs5 | .101 (.002) | .065 (.002) | .055 (.003) | 1m (.06) | 1h (.06) | 1h (.02) |
| Housing | 1.03 (.004) | .300 (.006) | .310 (.009) | .36m (.00) | 25m (.65) | 61m (1.8) |
| Triazines | 1.35 (.006) | 1.36 (.006) | 1.35 (.006) | 10s (.10) | 8m (.04) | 14m (.30) |
| Space_ga | 1.01 (—) | .767 (—) | .742 (—) | 2m (—) | 5h (—) | 11h (—) |

Table 3: Results for GP classification using logistic likelihood and GP regression using Laplace likelihood. For all rows, a lower value is better.

also comparable. The advantage of our approach is that it is easier to implement compared to EP and is numerically robust. The predictive probabilities obtained with EP and Proximal for ‘USPS_3vs5’ dataset are shown in the right plot of Figure 1. We see that Proximal method gives better estimates than EP in this case (higher is better). The improvement in the performance is due to the numerical error in the likelihood implementation. For Proximal method, we use the method of [23] which is quite accurate. Designing such accurate likelihood approximations for EP is challenging.

7 Discussion and Future Work

In this paper, we proposed a proximal framework that uses the KL proximal term to take the geometry of the posterior into account. We established equivalence between our proximal-point algorithm and natural-gradient methods. We proposed a proximal-gradient algorithm that exploits the structure of the bound to simplify the optimization.

An important future direction is to apply stochastic approximations to approximate gradients. This extension is discussed in [21]. It is also important to design a line-search method to set the step sizes. In addition, our proximal framework can also be used for distributed optimization in variational inference, e.g. using Alternating Direction Method of Multiplier (ADMM) [26, 11].

Acknowledgments

Emtiyaz Khan would like to thank Masashi Sugiyama and Akiko Takeda from University of Tokyo, Matthias Grossglauser and Vincent Etter from EPFL, and Hannes Nickisch from Philips Research (Hamburg) for useful discussions and feedback. Pierre Baqué was supported in part by the Swiss National Science Foundation, under the grant CRSII2-147693 “Tracking in the Wild”.

References

- [1] Matthew D Hoffman, David M Blei, Chong Wang, and John Paisley. Stochastic variational inference. *The Journal of Machine Learning Research*, 14(1):1303–1347, 2013.
- [2] Tim Salimans, David A Knowles, et al. Fixed-form variational posterior approximation through stochastic linear regression. *Bayesian Analysis*, 8(4):837–882, 2013.
- [3] Rajesh Ranganath, Sean Gerrish, and David M Blei. Black box variational inference. *arXiv preprint arXiv:1401.0118*, 2013.
- [4] Michalis Titsias and Miguel Lázaro-Gredilla. Doubly Stochastic Variational Bayes for Non-Conjugate Inference. In *International Conference on Machine Learning*, 2014.
- [5] Masa-Aki Sato. Online model selection based on the variational Bayes. *Neural Computation*, 13(7):1649–1681, 2001.
- [6] A. Honkela, T. Raiko, M. Kuusela, M. Tornio, and J. Karhunen. Approximate Riemannian conjugate gradient learning for fixed-form variational Bayes. *The Journal of Machine Learning Research*, 11:3235–3268, 2011.
- [7] Stéphane Chrétien and Alfred OIII Hero. Kullback proximal algorithms for maximum-likelihood estimation. *Information Theory, IEEE Transactions on*, 46(5):1800–1810, 2000.
- [8] Paul Tseng. An analysis of the EM algorithm and entropy-like proximal point methods. *Mathematics of Operations Research*, 29(1):27–44, 2004.
- [9] M. Teboulle. Convergence of proximal-like algorithms. *SIAM J on Optimization*, 7(4):1069–1083, 1997.
- [10] Pradeep Ravikumar, Alekh Agarwal, and Martin J Wainwright. Message-passing for graph-structured linear programs: Proximal projections, convergence and rounding schemes. In *International Conference on Machine Learning*, 2008.
- [11] Behnam Babagholami-Mohamadabadi, Sejong Yoon, and Vladimir Pavlovic. D-MFVI: Distributed mean field variational inference using Bregman ADMM. *arXiv preprint arXiv:1507.00824*, 2015.
- [12] Bo Dai, Niao He, Hanjun Dai, and Le Song. Scalable Bayesian inference via particle mirror descent. *Computing Research Repository*, abs/1506.03101, 2015.
- [13] Lucas Theis and Matthew D Hoffman. A trust-region method for stochastic variational inference with applications to streaming data. *International Conference on Machine Learning*, 2015.
- [14] Razvan Pascanu and Yoshua Bengio. Revisiting natural gradient for deep networks. *arXiv preprint arXiv:1301.3584*, 2013.
- [15] Ulrich Paquet. On the convergence of stochastic variational inference in bayesian networks. *NIPS Workshop on variational inference*, 2014.
- [16] Nicholas G Polson, James G Scott, and Brandon T Willard. Proximal algorithms in statistics and machine learning. *arXiv preprint arXiv:1502.03175*, 2015.
- [17] Harri Lappalainen and Antti Honkela. Bayesian non-linear independent component analysis by multi-layer perceptrons. In *Advances in independent component analysis*, pages 93–121. Springer, 2000.
- [18] Chong Wang and David M. Blei. Variational inference in nonconjugate models. *J. Mach. Learn. Res.*, 14(1):1005–1031, April 2013.
- [19] M. Seeger and H. Nickisch. Large scale Bayesian inference and experimental design for sparse linear models. *SIAM Journal of Imaging Sciences*, 4(1):166–199, 2011.
- [20] Antti Honkela and Harri Valpola. Unsupervised variational Bayesian learning of nonlinear models. In *Advances in neural information processing systems*, pages 593–600, 2004.
- [21] Mohammad Emtiyaz Khan, Reza Babanezhad, Wu Lin, Mark Schmidt, and Masashi Sugiyama. Convergence of Proximal-Gradient Stochastic Variational Inference under Non-Decreasing Step-Size Sequence. *arXiv preprint*, 2015.
- [22] Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian Processes for Machine Learning*. MIT Press, 2006.
- [23] B. Marlin, M. Khan, and K. Murphy. Piecewise bounds for estimating Bernoulli-logistic latent Gaussian models. In *International Conference on Machine Learning*, 2011.
- [24] Mohammad Emtiyaz Khan. Decoupled Variational Inference. In *Advances in Neural Information Processing Systems*, 2014.
- [25] E. Challis and D. Barber. Concave Gaussian variational approximations for inference in large-scale Bayesian linear models. In *International conference on Artificial Intelligence and Statistics*, 2011.
- [26] Huahua Wang and Arindam Banerjee. Bregman alternating direction method of multipliers. In *Advances in Neural Information Processing Systems*, 2014.