

# *revrand*: Technical Report

**Daniel Steinberg**

**Louis Tiao**

*DATA61, CSIRO*

*Sydney, Australia*

DANIEL.STEINBERG@DATA61.CSIRO.AU

LOUIS.TIAO@DATA61.CSIRO.AU

## Abstract

This is a technical report on the *revrand* software library. This library implements various Bayesian linear models (Bayesian linear regression), approximate Gaussian processes and generalised linear models. These algorithms have been implemented such that they can be used for large-scale inference by using stochastic gradients. All of the algorithms in *revrand* use a unified feature composition framework that allows for easy concatenation and selective application of regression basis functions.

## Contents

<b>1</b>	<b>Core Algorithms</b>	<b>1</b>
1.1	Stochastic Gradients and Variational Objective Functions . . . . .	1
1.2	Bayesian Linear Regression . . . . .	3
1.3	Bayesian Generalised Linear Models . . . . .	4
1.4	Large Scale Gaussian Process Approximation . . . . .	6
<b>2</b>	<b>Feature Composition Framework</b>	<b>7</b>
<b>3</b>	<b>Experiments</b>	<b>7</b>

## 1. Core Algorithms

### 1.1 Stochastic Gradients and Variational Objective Functions

Stochastic Gradients is now a ubiquitous method for optimisation when a whole dataset does not fit in memory, or when optimisation has to be distributed amongst many computational nodes.

When an objective function factorises over data,

$$f(\mathbf{X}, \theta) = \sum_{\mathbf{x} \in \mathbf{X}} f(\mathbf{x}, \theta), \quad (1)$$

a regular gradient descent would perform the following iterations to minimise the function w.r.t.  $\theta$ ,

$$\theta_k := \theta_{k-1} - \eta_k \sum_{\mathbf{x} \in \mathbf{X}} \nabla_{\theta} f(\mathbf{x}_n, \theta)|_{\theta=\theta_{k-1}}, \quad (2)$$

where  $\eta_k$  is the learning rate (step size) at iteration  $k$ . Stochastic gradients proposes the following update,

$$\theta_k := \theta_{k-1} - \eta_k \sum_{\mathbf{x} \in \mathbf{B}} \nabla_{\theta} f(\mathbf{x}, \theta)|_{\theta=\theta_{k-1}}, \quad (3)$$

where  $\mathbf{B} \subset \mathbf{X}$  is a mini-batch of the original dataset, where  $|\mathbf{B}| \ll |\mathbf{X}|$ .

Unfortunately some objective functions do not entirely decompose over the data, i.e.

$$f(\mathbf{X}, \theta) = \sum_{\mathbf{x} \in \mathbf{X}} f(\mathbf{x}, \theta) + g(\theta). \quad (4)$$

Let  $M = |\mathbf{B}|$  and  $N = |\mathbf{X}|$ , then we divide the contribution of the constant term amongst the mini-batches in stochastic gradients,

$$\theta_k := \theta_{k-1} - \eta_k \sum_{\mathbf{x} \in \mathbf{B}} \nabla_{\theta} f(\mathbf{x}, \theta)|_{\theta=\theta_{k-1}} - \eta_k \frac{M}{N} \nabla_{\theta} g(\theta)|_{\theta=\theta_{k-1}}. \quad (5)$$

This is particularly relevant for variational inference where the evidence lower bound objective has a component independent of the data. For example, let's consider the model,

$$\text{Likelihood: } \prod_{n=1}^N p(y_n | \theta), \quad (6)$$

$$\text{prior: } p(\theta | \alpha), \quad (7)$$

where we want to learn the values of the hyper-parameters,  $\alpha$ . Minimising negative log-marginal likelihood is a good objective in this instance, since we don't care about the value(s) of  $\theta$ ,

$$\underset{\alpha}{\operatorname{argmin}} - \log \int \prod_{n=1}^N p(y_n | \theta) p(\theta | \alpha) d\theta. \quad (8)$$

There are two problems with this objective however, (1) it may not factor over data and (2) the integral may be intractable, for instance, if the prior and likelihood are not conjugate. In variational inference we use Jensen's inequality to lower-bound log-marginal likelihood with a tractable objective function called the evidence lower bound (ELBO),

$$\begin{aligned} \log p(\mathbf{y} | \alpha) &= \log \int \prod_{n=1}^N p(y_n | \theta) p(\theta | \alpha) d\theta \\ &= \log \int \frac{\prod_n p(y_n | \theta) p(\theta | \alpha)}{q(\theta)} q(\theta) d\theta \\ &\geq \int q(\theta) \log \left[ \frac{\prod_n p(y_n | \theta) p(\theta | \alpha)}{q(\theta)} \right] d\theta \end{aligned} \quad (9)$$

where  $q(\theta)$  is an approximation of  $p(\theta | \alpha)$  that makes inference easier. This can be re-written as,

$$\mathcal{L} = \sum_{n=1}^N \langle \log p(y_n | \theta) \rangle_q - \text{KL}[q(\theta) \| p(\theta | \alpha)], \quad (10)$$

which takes the form of Equation (4), and so we can weight the Kullback-Leibler term like the constant term,  $g(\cdot)$ , from Equation (5) if we use stochastic gradients optimisation. Furthermore, if  $q(\theta) = p(\theta | \alpha)$  then the lower bound is tight, and this will be equivalent to optimising log-marginal likelihood.

## 1.2 Bayesian Linear Regression

The first machine learning algorithm in revrand is a simple Bayesian linear regressor of the following form,

$$\text{Likelihood: } \prod_{n=1}^N \mathcal{N}(y_n | \phi_n^\top \mathbf{w}, \sigma^2), \quad (11)$$

$$\text{prior: } \mathcal{N}(\mathbf{w} | \mathbf{0}, \lambda \mathbf{I}_D), \quad (12)$$

where  $\phi_n := \phi(\mathbf{x}_n, \theta)$  is a feature, or basis, function that  $\phi \mathbb{R}^d \rightarrow \mathbb{R}^D$ . This is the same algorithm described in (Rasmussen and Williams, 2006, Chapter 2). We then:

- Optimise  $\sigma^2, \lambda$  and  $\theta$  w.r.t. log-marginal likelihood,

$$\log p(\mathbf{y} | \sigma^2, \lambda, \theta) = \log \mathcal{N}(\mathbf{y} | \mathbf{0}, \sigma^2 \mathbf{I}_N + \lambda \Phi^\top \Phi), \quad (13)$$

where  $\Phi \in \mathbb{R}^{N \times D}$  is the concatenation of all the features,  $\phi_n$ . Note this results in the covariance of the log-marginal likelihood being  $N \times N$ , though we can use the Woodbury identity to simplify the corresponding matrix inversion.

- Solve analytically for the posterior over weights,  $\mathbf{w} | \mathbf{y} \sim \mathcal{N}(\mathbf{m}, \mathbf{C})$  given the above hyperparameters, where,

$$\mathbf{C} = \left[ \lambda \mathbf{I}_D + \frac{1}{\sigma^2} \Phi^\top \Phi \right]^{-1},$$

$$\mathbf{m} = \frac{1}{\sigma^2} \mathbf{C} \Phi^\top \mathbf{y}.$$

- Use the predictive distribution

$$\begin{aligned} p(y^* | \mathbf{y}, \mathbf{X}, \mathbf{x}^*) &= \int \mathcal{N}(y^* | \phi^{*\top} \mathbf{w}, \sigma^2) \mathcal{N}(\mathbf{w} | \mathbf{m}, \mathbf{C}) d\mathbf{w}, \\ &= \mathcal{N}(y^* | \phi^{*\top} \mathbf{m}, \sigma^2 + \phi^{*\top} \mathbf{C} \phi^*) \end{aligned} \quad (14)$$

for query inputs,  $\mathbf{x}^*$ . This gives us the useful expectations,

$$\mathbb{E}[y^*] = \phi^{*\top} \mathbf{m}, \quad (15)$$

$$\mathbb{V}[y^*] = \sigma^2 + \phi^{*\top} \mathbf{C} \phi^{*}. \quad (16)$$

It is actually easier to use the ELBO form with stochastic gradients for learning the parameters of this algorithm, rather than log-marginal likelihood recast using the Woodbury identity. This is because it is plainly in the same form as Equation (4), though it would give the same result as log-marginal likelihood, the “approximate” posterior is the same form as the true posterior, i.e.  $q(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}, \mathbf{C})$ . The ELBO for this model is,

$$\mathcal{L} = \sum_{n=1}^N \left\langle \log \mathcal{N}(y_n | \phi_n^\top \mathbf{w}, \sigma^2) \right\rangle_q - \text{KL}[\mathcal{N}(\mathbf{w} | \mathbf{m}, \mathbf{C}) \| \mathcal{N}(\mathbf{w} | \mathbf{0}, \lambda \mathbf{I}_D)]. \quad (17)$$

More specifically,

$$\begin{aligned} \left\langle \log \mathcal{N}(y_n | \phi_n^\top \mathbf{w}, \sigma^2) \right\rangle_q &= \log \mathcal{N}(y_n | \phi_n^\top \mathbf{m}, \sigma^2) - \frac{1}{2\sigma^2} \text{tr}(\phi_n^\top \phi_n \mathbf{C}), \\ \text{KL}[\mathcal{N}(\mathbf{w} | \mathbf{m}, \mathbf{C}) || \mathcal{N}(\mathbf{w} | \mathbf{0}, \lambda \mathbf{I}_D)] &= \frac{1}{2\lambda} \left[ \text{tr}(\mathbf{C}) + \mathbf{m}^\top \mathbf{m} \right] - \frac{1}{2} \log |\mathbf{C}| \\ &\quad + \frac{D}{2} (\log \lambda - 1). \end{aligned}$$

We have not implemented a stochastic gradient version of this algorithm since it still requires a determinant involving a  $D \times D$  matrix, and so is  $\mathcal{O}(D^3)$  in complexity, per iteration. This is true even if we optimise the posterior covariance directly (or a triangular parameterisation). The GLM presented in the next section circumvents this issue, and is more suited to really large  $N$  and  $D$  problems.

### 1.3 Bayesian Generalised Linear Models

The algorithm of primary interest in *revrand* is the Bayesian generalised linear model. The general form of the model implemented by this algorithm is,

$$\text{Likelihood: } \prod_{n=1}^N p(y_n | g(\phi_n^\top \mathbf{w}), \gamma), \quad (18)$$

$$\text{prior: } \mathcal{N}(\mathbf{w} | \mathbf{0}, \lambda \mathbf{I}_D), \quad (19)$$

for an arbitrary univariate likelihood,  $p(\cdot)$ , with an appropriate transformation (inverse link) function,  $g(\cdot)$ , and parameter(s),  $\gamma$ .

Naturally, both calculating the exact posterior over the weights,  $p(\mathbf{w} | \mathbf{y}, \mathbf{X})$ , and the log-marginal likelihood,  $p(\mathbf{y})$ , for hyperparameter learning are intractable since we may have a non-conjugate relationship between the likelihood and prior. Therefore we must resort to approximating the true posterior and the log-marginal likelihood. We use a modification of the nonparametric variational algorithm presented in [Gershman et al. \(2012\)](#) for approximating the posterior and log marginal likelihood. The posterior takes the form

$$\begin{aligned} p(\mathbf{w} | \mathbf{y}, \mathbf{X}) &\approx q(\mathbf{w}), \\ &= \frac{1}{K} \prod_{k=1}^K \mathcal{N}(\mathbf{w} | \mathbf{m}_k, \Psi_k), \end{aligned} \quad (20)$$

i.e. a mixture of  $K$  diagonal Gaussians,  $\Psi_k = \text{diag}([\Psi_{k,1}, \dots, \Psi_{k,D}]^\top)$ . This is a very flexible form for the approximate posterior, and still admits a closed-form lower bound to the log-marginal likelihood. Furthermore this has the nice property that our algorithm no longer has a  $\mathcal{O}(D^3)$  cost associated with the number of features. [Gershman et al. \(2012\)](#) make one more approximation based on local second order Taylor expansions of the joint,

$$\begin{aligned} \langle \log p(\mathbf{y}, \mathbf{w}) \rangle_q &\approx \frac{1}{K} \sum_{k=1}^K \sum_{n=1}^N \log p(y_n | g(\phi_n^\top \mathbf{m}_k), \gamma) \\ &\quad + \log \mathcal{N}(\mathbf{m}_k | \mathbf{0}, \lambda \mathbf{I}_D) + \frac{1}{2} \text{tr}(\Psi_k \mathbf{H}_k), \end{aligned} \quad (21)$$

to deal with non-conjugacy within their variational framework. Here  $\mathbf{H}_k = \nabla_{\mathbf{w}}^2 p(y, \mathbf{w})|_{\mathbf{w}=\mathbf{m}_k}$ <sup>1</sup>. Unfortunately this approximation to the ELBO actually breaks the lower bound, but works well in practice. For more details on this inference scheme we direct the reader to [Gershman et al. \(2012\)](#). The final objective function we optimise is,

$$\begin{aligned} \mathcal{L} \approx & \frac{1}{K} \sum_{k=1}^K \sum_{n=1}^N \log p(y_n | g(\phi_n^\top \mathbf{m}_k), \gamma) + \log \mathcal{N}(\mathbf{m}_k | \mathbf{0}, \lambda \mathbf{I}_D) \\ & + \frac{1}{2} \text{tr}(\mathbf{\Psi}_k \mathbf{H}_k) - \log \left( \sum_{j=1}^K \mathcal{N}(\mathbf{m}_k | \mathbf{m}_j, \mathbf{\Psi}_k + \mathbf{\Psi}_j) \right) - \log K. \end{aligned} \quad (22)$$

While this is an approximation to Equation (10), it retains the property that it easily factorises over the data ( $n$ ), and so can be optimised using stochastic gradients as described in §1.1.

In *revrand* we actually optimise Equation (22) with respect to the hyperparameters,  $\theta, \gamma$  and the parameters  $\mathbf{m}_k, \mathbf{\Psi}_k \forall k$  simultaneously. This is unlike the optimisation scheme in [Gershman et al. \(2012\)](#), since it requires 3<sup>rd</sup> derivatives of the likelihood in Equation (18) w.r.t.  $\mathbf{w}$ . However, we find that providing these derivatives is simple (for our model) compared to the more complex optimisation scheme presented in [Gershman et al. \(2012\)](#), which is not as amenable to large datasets as stochastic gradients when hyperparameter optimisation is involved.

The most simple and accurate method for approximating the predictive distribution,  $p(y^* | \mathbf{y}, \mathbf{X}, \mathbf{x}^*)$  is to Monte-Carlo sample the integral,

$$p(y^* | \mathbf{y}, \mathbf{X}, \mathbf{x}^*) \approx \int p(y | g(\phi^{*\top} \mathbf{w}), \gamma) \frac{1}{K} \sum_{k=1}^K \mathcal{N}(\mathbf{w} | \mathbf{m}_k, \mathbf{\Psi}_k) d\mathbf{w}. \quad (23)$$

However, this integral is not particularly useful unless we wish to evaluate known  $\mathbf{y}^*$  under the model. For prediction, it is more useful to compute (using Monte-Carlo integration) the predictive expectation,

$$\begin{aligned} \mathbb{E}[y^*] & \approx \int \frac{1}{K} \sum_{k=1}^K \mathcal{N}(\mathbf{w} | \mathbf{m}_k, \mathbf{\Psi}_k) \int y^* p(y^* | g(\phi^{*\top} \mathbf{w}), \gamma) dy^* d\mathbf{w} \\ & = \int \mathbb{E}[p(y^* | g(\phi^{*\top} \mathbf{w}), \gamma)] \frac{1}{K} \sum_{k=1}^K \mathcal{N}(\mathbf{w} | \mathbf{m}_k, \mathbf{\Psi}_k) d\mathbf{w}. \end{aligned} \quad (24)$$

Often we find  $\mathbb{E}[p(y^* | g(\phi^{*\top} \mathbf{w}), \gamma)] = g(\phi^{*\top} \mathbf{w})$ , however this is only true with with right choice and usage of the activation function. Furthermore, it is useful to compute quantiles of the predictive density in order to ascertain the predictive uncertainty. We start by sampling

---

1. We only need the diagonals of this Hessian matrix.

the predictive cumulative density function,  $P(\cdot)$ ,

$$\begin{aligned}
P(y^* \leq \alpha | \mathbf{y}, \mathbf{X}, \mathbf{x}^*) & \approx \int \frac{1}{K} \sum_{k=1}^K \mathcal{N}(\mathbf{w} | \mathbf{m}_k, \Psi_k) \int_{-\infty}^{\alpha} p(y^* | g(\phi^{*\top} \mathbf{w}), \gamma) dy^* d\mathbf{w} \\
& = \int P(y^* \leq \alpha | g(\phi^{*\top} \mathbf{w}), \gamma) \frac{1}{K} \sum_{k=1}^K \mathcal{N}(\mathbf{w} | \mathbf{m}_k, \Psi_k) d\mathbf{w}.
\end{aligned} \tag{25}$$

Once we have obtained sufficient samples from the (mixture) posterior we can obtain quantiles,  $\alpha$ , for some chosen level of probability,  $p$ , using root finding techniques. Specifically, we use root finding techniques to solve the following for  $\alpha$ ,

$$P(y^* \leq \alpha | \mathbf{y}, \mathbf{X}, \mathbf{x}^*) - p = 0. \tag{26}$$

#### 1.4 Large Scale Gaussian Process Approximation

TODO: re-write the following:

In *revrand* we approximate Gaussian Processes with our standard and generalised linear models by using random feature functions such as those of [Rahimi and Recht \(2007; 2008\)](#). They use Bochner’s theorem regarding the relationship between a kernel and the Fourier transform of a non-negative measure that (via Wiener-Khintchine’s theorem) establishes the duality of the covariance function of a stationary process and its spectral density,

$$k(\boldsymbol{\tau}) = \int S(\mathbf{s}) e^{i\mathbf{s}^\top \boldsymbol{\tau}} d\mathbf{s}, \tag{27}$$

$$S(\mathbf{s}) = \int k(\boldsymbol{\tau}) e^{-i\mathbf{s}^\top \boldsymbol{\tau}} d\boldsymbol{\tau}. \tag{28}$$

[Rahimi and Recht](#)’s main insight (2007) is that we can approximate the kernel by constructing ‘suitable’ random features and Monte Carlo averaging over samples from  $S(\mathbf{s})$ ,

$$k(\mathbf{x} - \mathbf{x}') = k(\boldsymbol{\tau}) \approx \frac{1}{D} \sum_{i=1}^D \phi_i(\mathbf{x})^\top \phi_i(\mathbf{x}'), \tag{29}$$

$\phi_i(\mathbf{x})$  corresponds to the  $i$ th sample from the feature map. An example of a feature vector construction in the above approximation is,

$$\begin{aligned}
[\phi_i(\mathbf{x}), \phi_{D+i}(\mathbf{x})] &= \frac{1}{\sqrt{D}} [\cos(\mathbf{s}_i^\top \mathbf{x}), \sin(\mathbf{s}_i^\top \mathbf{x})], \\
\text{with } \mathbf{s}_i &\sim \mathcal{N}(\mathbf{s}_i | \mathbf{0}, \sigma_\phi^2 \mathbf{I}_d),
\end{aligned} \tag{30}$$

for  $i = 1, \dots, D$ , which in fact is a mapping into a  $2D$ -dimensional feature space. [Rahimi and Recht \(2007\)](#) used the above feature map to approximate the commonly used (isotropic) squared exponential kernel, and showed that such an approximation converges in expectation to the true kernel.

TODO: table of kernels and sampling distributions currently in *revrand* (Laplace, Cauchy, RBF, Matern5/2, Matern3/2 etc)

TODO: FastFood

TODO: A la Carte spectral mixtures

## 2. Feature Composition Framework

TODO: We have implemented kernel addition (basis concatenation), but we haven't done kernel products.

TODO: mention advantages of using kernel and constant bases (offsets)

## 3. Experiments

### References

Samuel Gershman, Matt Hoffman, and David Blei. Nonparametric variational inference. In John Langford and Joelle Pineau, editors, *Proceedings of the 29th International Conference on Machine Learning (ICML-12)*, ICML '12, pages 663–670, New York, NY, USA, July 2012. Omnipress. ISBN 978-1-4503-1285-1.

Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *NIPS*. 2007.

Ali Rahimi and Benjamin Recht. Weighted sums of random kitchen sinks: Replacing minimization with randomization in learning. In *NIPS*. 2008.

Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian processes for machine learning*. The MIT Press, Cambridge, Massachusetts, 2006.