

The coresets variational Bayes (CVB) algorithm for mixture analysis

Qianying Liu^a, Clare A. McGrory^a and Peter W. J. Baxter^{a,b}

^a*University of Queensland*

^b*Queensland University of Technology*

Abstract. The pressing need for improved methods for analysing and coping with big data has opened up a new area of research for statisticians. Image analysis is an area where there is typically a very large number of data points to be processed per image, and often multiple images are captured over time. These issues make it challenging to design methodology that is reliable and yet still efficient enough to be of practical use. One promising emerging approach for this problem is to reduce the amount of data that actually has to be processed by extracting what we call coresets from the full dataset; analysis is then based on the coreset rather than the whole dataset. Coresets are representative subsamples of data that are carefully selected via an adaptive sampling approach. We propose a new approach called coreset variational Bayes (CVB) for mixture modelling; this is an algorithm which can perform a variational Bayes analysis of a dataset based on just an extracted coreset of the data. We apply our algorithm to weed image analysis.

1 Introduction

Finite mixture models are applied in diverse areas of science and provide a straightforward, but flexible extension of classical parametric models (Fruhwirth-Schnatter, 2006). They are used in situations where it is thought that there is more than one sub-population giving rise to points in the dataset. The sub-populations are each represented by one of the components that comprise the mixture. An obvious scenario where this type of model would be suitable is image analysis: we can associate ranges of intensity level present in the image with components of the mixture.

We take a Bayesian approach to our inference. In Bayesian analysis, the posterior distribution of the unknown parameters can be very difficult to estimate. The most commonly used Bayesian method for estimating the parameters is Markov chain Monte Carlo (MCMC), for example, a Gibbs sampling approach (Alston, Mengersen and Pettitt, 2012). However, while MCMC-based approaches are very accurate, the computational requirements associated with this approach can be prohibitive for very large datasets. Less computationally demanding algorithms for fitting the mixture models are alternative approximate Bayesian inference techniques

Key words and phrases. Mixture modelling, coresets, variational Bayes, image analysis, Bayesian statistics.

Received February 2017; accepted November 2017.

such as variational Bayes (VB) (McGrory and Titterington, 2007; Ormerod and Wand, 2010) and approximate Bayesian computation (ABC) (Marin et al., 2012).

Even though the standard VB method is highly computationally efficient in comparison to MCMC, in some cases it still might be too time-consuming for very large data problems if analysis is required within short time-frames. Consider for instance the weed-crop imaging application that we will explore in this article. Images have an extremely large number of pixels, there will likely be multiple images to analyse, and estimates are needed reasonably quickly. In order to achieve this, time-efficient techniques are required. Another avenue is to reduce the volume of data that actually has to be processed in the first place. Removing part of the dataset before running the analysis is a less drastic thing to do than we might think when we consider that much of the dataset gives us the same or very similar information. For example, consider trying to fit a mixture model to an image, we do not necessarily have to analyse all of the observed intensity levels present in the full dataset to obtain a good estimate of the mean for that particular cluster. The coresets approach (Feldman, Faulkner and Krause, 2011) is a data-reduction algorithm. It involves extracting a representative subsample of the dataset which can then be analysed in order to make inference about the whole dataset. In Feldman, Faulkner and Krause (2011) the coresets approach was used within a classical mixture modelling framework with good results. In McGrory et al. (2014) the Gibbs sampler was modified for use with coresets; the resulting algorithm was called the weighted Gibbs sampler. This was shown to give very good results when applied to analysing satellite image data and it drastically reduced the computation time required for the analysis. In a similar spirit, we wish to modify the VB algorithm to make it suitable for use with coresets of data. Since VB is more time-efficient than the Gibbs sampler, combining the concept of coresets with VB will result in even greater reductions in computational burden. We refer to this new algorithm as coresets variational Bayes (CVB).

It could be argued that a spatial mixture model is more appropriate than a finite mixture model for modelling a dataset such as an image which contains spatial information (see, e.g., McGrory et al. (2012)) as this might slightly increase the clustering accuracy. The reason we do not incorporate a spatial component into our modelling is that we cannot afford the huge extra computational burden this would incur since we require a very time efficient approach for big data settings. In this way, there is a trade-off between these two aspects.

Weeds are defined as being plants which have originated in and continue to evolve in a natural environment; but they are problematic because they do so in a manner which interferes with the growth of crops or other agriculture related activities (Zimdahl, 2009). Weeds are generally better able to compete than crops for resources like minerals, light and water. This leads to a lot of waste of agricultural investment in cases where weed plants are present, because most of these valuable inputs would be used up by them instead of by the valuable crops. There are many other harmful problems associated with weeds, such as the harbouring of extra

pests in the area which can cause plant diseases that may infect the crops in the region. Hence, it is important to be able to effectively manage weeds in farming regions if a nation's agriculture industry is to remain competitive in international markets (Sindenab et al., 2004).

A growing area of research that has the potential to be very useful in weed management is the use of statistical methodology to analyse images of weeds used in agriculture trials. This might be of use in projects where the aim is to assess and compare the effectiveness of different chemical weed killers, for instance. If calculation of the proportion of a plot of land that is weed, soil or plant after chemical applications can be done by analysing an image, this will save time by removing the need for researchers to go out into the fields and count the live and dead plants by hand. Due to the large number of pixels present in a typical image, and the fact that there will most likely be multiple images to analyse in a given trial, it is important to find a very time efficient algorithm for processing this type of data (see, e.g., Kargar and Shirzadifar (2013)).

2 Finite mixture model

Mixture models provide an excellent and flexible way to represent complex distributions. The mixture model we fit to our data comprises a linear combination of standard mixture models, these are called the components of the mode. Each component has a corresponding mixture weight which reflects the expected proportion of the data that might be captured by that particular component. In our finite mixture model for a set of continuous observations $\mathbf{y}_1, \dots, \mathbf{y}_n$, assume that the observations are all generated *i.i.d.* (independent and identically distributed) from a random variable \mathbf{Y} , which follows a mixture of K independent Gaussian distributions. In the missing-data interpretation of the mixture, we introduce an unobserved indicator variable z_{ij} for each observation; this identifies the component allocations of our observations by taking value 1 if observation i is from component j , and 0 otherwise. Since these indicators are unknown to us, the model is therefore a missing-data model. For each component, j , the value of ρ_j is the relevant weight for that particular j th component. The model density is given by:

$$p(\mathbf{y}, \mathbf{z} | \boldsymbol{\theta}) = \prod_{i=1}^n \prod_{j=1}^K \{\rho_j N_d(\mathbf{y}_i; \boldsymbol{\mu}_j, \mathbf{T}_j^{-1})\}^{z_{ij}}.$$

Here, $N_d(\cdot, \cdot)$ represents the d -dimensional multivariate normal density, where $\boldsymbol{\mu} = (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K)$ and $\mathbf{T} = (\mathbf{T}_1, \dots, \mathbf{T}_K)$, and \mathbf{T}_j denotes the j th precision matrix, which is the inverse of the j th covariance matrix.

3 Bayesian priors

We follow the standard Bayesian conjugate prior setting (Alston, Mengersen and Pettitt, 2012) for this model, and choose hyper-parameters such that they corre-

spond to non-informative prior settings, thus allowing information contained in the dataset to have more influence over the fit. The weight coefficients are assigned Dirichlet prior distributions:

$$p(\boldsymbol{\rho}) = \text{Dir}(\boldsymbol{\rho}; \alpha_1^{(0)}, \dots, \alpha_K^{(0)}).$$

The prior distributions of the means conditioned on the covariance matrices are independent multivariate normal distributions:

$$p(\boldsymbol{\mu}|\mathbf{T}) = \prod_{j=1}^K N_d(\boldsymbol{\mu}_j; \mathbf{m}_j^{(0)}, (\beta_j^{(0)} \mathbf{T}_j)^{-1}).$$

The prior of the precision matrices are given by Wishart distributions:

$$p(\mathbf{T}) = \prod_{j=1}^K W(\mathbf{T}_j; \mathbf{v}_j^{(0)}, \boldsymbol{\Sigma}_j^{(0)}).$$

Therefore, the joint distribution would finally be:

$$p(\mathbf{y}, \mathbf{z}, \boldsymbol{\theta}) = p(\mathbf{y}, \mathbf{z}|\boldsymbol{\theta})p(\boldsymbol{\rho})p(\boldsymbol{\mu}|\mathbf{T})p(\mathbf{T}).$$

The quantities of $\{\alpha_j^{(0)}\}$, $\{\mathbf{m}_j^{(0)}\}$, $\{\beta_j^{(0)}\}$ and $\{\boldsymbol{\Sigma}_j^{(0)}\}$ are all hyper-parameters.

4 Bayesian posterior distributions

4.1 The variational approach

The variational Bayesian method is a time-efficient approach for estimating Bayesian mixture models (Faes, Ormerod and Wand, 2011; McGrory and Titterton, 2007; Wand et al., 2012) and can be thought of as an alternative to MCMC. The main difference between the two approaches is that instead of estimating the parameter directly by sampling from the posterior distribution, variational Bayesian methods approximate it. This is done by artificially “introducing” a more amenable distribution $q(\boldsymbol{\theta}, \mathbf{z})$ which is often referred to as the variational approximating function. This function will end up becoming an approximate the joint conditional distribution of $\boldsymbol{\theta}$ and \mathbf{z} given the observations \mathbf{y} after the variational method is applied to it. In the following, we explain how $q(\boldsymbol{\theta}, \mathbf{z})$ should be chosen and integrated into the variational framework in order to achieve this outcome.

The distribution $q(\boldsymbol{\theta}, \mathbf{z})$ is chosen to minimise the Kullback–Leibler (KL) divergence between the approximating density $q(\boldsymbol{\theta}, \mathbf{z})$ and the true joint density $p(\boldsymbol{\theta}, \mathbf{z}|\mathbf{y})$. In doing so, we are trying to obtain a relatively tight lower bound on the marginal density, $p(\mathbf{y})$. Essentially we manipulate and re-express the joint density to allow us to introduce the variational approximating function in such a way that we can then use a maximisation approach to estimate parameters of that target ap-

proximating function (see also [McGrory and Titterington \(2007\)](#)). We begin then by showing that the joint density is lower bounded as:

$$\begin{aligned} \log p(\mathbf{y}) &= \log \int \sum_{\{z\}} p(\mathbf{y}, z, \boldsymbol{\theta}) d\boldsymbol{\theta} \\ &= \log \int \sum_{\{z\}} q(\boldsymbol{\theta}, z) \frac{p(\mathbf{y}, z, \boldsymbol{\theta})}{q(\boldsymbol{\theta}, z)} d\boldsymbol{\theta} \\ &\geq \int \sum_{\{z\}} \log \frac{p(\mathbf{y}, z, \boldsymbol{\theta})}{q(\boldsymbol{\theta}, z)} d\boldsymbol{\theta} \quad \text{by Jensen's inequality.} \end{aligned} \tag{1}$$

Another way of viewing this is that finding the tightest lower bound is the same as minimising the Kullback–Leibler divergence between the variational distribution and the true target posterior. It is exactly minimised when we take $q(\boldsymbol{\theta}, z) = p(\boldsymbol{\theta}, z|\mathbf{y})$. However, as we are trying to simplify the problem, $q(\boldsymbol{\theta}, z)$ should be a close enough approximation to the true density, yet have a simple form for computational purposes. Normally, to achieve this $q(\boldsymbol{\theta}, z)$ is restricted to have the factorised form $q(\boldsymbol{\theta}, z) = q_{\boldsymbol{\theta}}(\boldsymbol{\theta})q_z(z)$.

Unlike the MCMC approach, the variational method approximates the parameters in the finite mixture model. This difference may cause a slight decrease in accuracy of the variational method when compared with MCMC. It has been demonstrated that in many contexts, including mixture modelling (e.g., [McGrory and Titterington \(2007\)](#); [Wand et al. \(2012\)](#); [Faes, Ormerod and Wand \(2011\)](#)) that the variational method can largely reduce operating time and yet the loss in accuracy that arises from the approximation is not terribly great. Indeed the approximate result is typically adequate for practical purposes. When computational efficiency is an important consideration, it is worthwhile for us to pursue this approximate approach rather than MCMC.

4.2 Variational posterior

After we maximise the lower bound (equation (1)), the posteriors are

$$\begin{aligned} q_{\boldsymbol{\rho}}(\boldsymbol{\rho}) &= \text{Dir}(\boldsymbol{\rho}; \alpha_1, \dots, \alpha_k), \\ q_{\boldsymbol{\mu}|T}(\boldsymbol{\mu}|T) &= \prod_{j=1}^K N_d(\boldsymbol{\mu}_j; \mathbf{m}_j, (\beta_j T_j)^{-1}), \end{aligned}$$

and

$$q_T(T) = \prod_{j=1}^K W(T_j; \mathbf{v}_j, \boldsymbol{\Sigma}_j),$$

the hyperparameters will be updated as:

$$\alpha_j = \alpha_j^{(0)} + \sum_{i=1}^n q_{ij}, \tag{2}$$

$$\beta_j = \beta_j^{(0)} + \sum_{i=1}^n q_{ij}, \tag{3}$$

$$\mathbf{m}_j = \frac{\beta_j^{(0)} \mathbf{m}_j^{(0)} + \sum_{i=1}^n q_{ij} \mathbf{y}_i}{\beta_j}, \tag{4}$$

$$\boldsymbol{\Sigma}_j = \boldsymbol{\Sigma}_j^{(0)} + \sum_{i=1}^n q_{ij} \mathbf{y}_i \mathbf{y}_i^T + \beta_j^{(0)} \mathbf{m}_j^{(0)} \mathbf{m}_j^{(0)T} - \beta_j \mathbf{m}_j \mathbf{m}_j^T, \tag{5}$$

$$v_j = v_j^{(0)} + \sum_{i=1}^n q_{ij}, \tag{6}$$

with q_{ij} being the variational posterior expected probability that the indicator variable $z_{ij} = 1$. The form of q_{ij} is:

$$q_{ij} = \frac{\exp\{\langle \log \rho_j \rangle + \frac{1}{2} \{\langle \log |\mathbf{T}_j| \rangle\} - \frac{1}{2} \text{tr}(\langle \mathbf{T}_j \rangle (\mathbf{y}_i - \mathbf{m}_j)(\mathbf{y}_i - \mathbf{m}_j)^T + 1/\beta_j) \mathbf{I}_d\}}{s_i} \tag{7}$$

$$= \frac{\phi_{ij}}{s_{ij}},$$

where s_{ij} is the normalization constant and $\langle \cdot \rangle$ denotes the expected values required to evaluate the expressions in Equation 7. These are given by:

$$\langle \boldsymbol{\mu}_j \rangle = \mathbf{m}_j,$$

$$\langle \mathbf{T}_j \rangle = v_j \boldsymbol{\Sigma}_j^{-1},$$

$$\langle |\log \mathbf{T}_j| \rangle = \sum_{s=1}^d \Psi \left(\frac{v_j + 1 - s}{2} \right) + d \log(2) + \log |\boldsymbol{\Sigma}_j|,$$

$$\langle |\log(\rho_j)| \rangle = \Psi(\hat{\alpha}_j) + \Psi(\hat{\alpha}_j),$$

where $\Psi(\cdot)$ is the digamma function and $\hat{\alpha}_j = \sum_i \hat{\alpha}_{ij}$.

4.3 The standard variational Bayes (VB) algorithm

In Algorithm 1, we outline the pseudo-code for the VB algorithm for obtaining the posterior estimates. We draw the reader’s attention to what we call the component elimination property of VB. By this we mean that the algorithm can automatically determine the complexity of the model since:

1. the initial value for the number of components is set to be greater than what the user would reasonably expect in the final fit,
2. components which converge to have similar estimated parameter values will be dominated by only one of them, then components with small mixing weights can be removed leading to automatic complexity assessment.

Algorithm 1 The standard VB algorithm

Set initial number of components K .

Set initial values for hyperparameters $\alpha^{(0)}$, $\beta^{(0)}$, $\Sigma^{(0)}$, $v^{(0)}$, $m^{(0)}$.

Specify initial allocation of observations to components and get q_{ij} .

while Not Converged **do**:

 Update variational posterior expressions for model parameters: Equations 2–6

 Update variational posterior for q_{ij} : Equation 7

if any component has a mixing weight $\leq \varepsilon$

 remove the component from model

end if

if the algorithm has converged

 exit loop.

end if

end while

This feature is very useful in applications involving large amounts of data as it greatly reduces computing costs: no need to perform different runs for various numbers of mixture components and compare them, which is what we would have to do if we used an MCMC approach in order to estimate the dimension. The other alternative would be to use the computationally burdensome reversible jump MCMC. This feature is also particularly useful for practical applications where the operator would like the analysis to run in an unsupervised manor because there is no need for the user to manually control searches over dimensionality.

Of course initial settings for the hyperparameters and epsilon (the threshold for component removal) can have some influence on estimation some cases, although previous research has shown the method to be generally fairly robust to initial parameter settings. We chose values for the hyperparameters $\alpha^{(0)}$, $\beta^{(0)}$, $\Sigma^{(0)}$, $v^{(0)}$, $m^{(0)}$ to correspond to vague non-informative priors. Note that these are the hyperparameters of the Gaussian mixture model, and they need to be chosen in the initial settings of the variational algorithm. In this, we follow the standard guidelines for prior settings used in any Bayesian analysis. Naturally, if a user has specific prior knowledge in their particular application, they might decide an informative prior is suitable in their case, this is a user driven choice.

Note that the choice of epsilon determines how small a component's allocated weighting has to be at a given iteration of the algorithm in order for it to be selected for removal, and again is the user's choice, within reason of course. The simplest, and perhaps "safest" choice is to set that equal to 1. That is a "safe" option because clearly once number of allocations is less than 1, the user can be sure that they are not excluding a component that has any real significant contribution within the model. However, that might not be the most time efficient choice, if you imagine a dataset with thousands of observations, it might be that any component with less than say 50–100 observations allocated is unlikely to be useful and on the

way to being eventually removed with a gradual reduction in allocations at each subsequent iteration. However, choosing epsilon in the range 50–100 might be less “safe” because it does incur higher risk that a component that is in fact potentially useful in representing some features of the data will get removed in error. Some users might decide on a proportional value of the size of the observation set which can be used as the cut off, for example, 1% of the dataset size. Users must select a suitable value according to their particular application.

5 Adapting the variational Bayes approach for use with coresets of data

In McGrory et al. (2014), it was shown how the Gibbs sampler could be adapted for use with coresets of data. In a similar spirit, we will adapt the VB method for use with coresets. We first describe the basic procedure for finding coresets, as outlined in Feldman, Faulkner and Krause (2011).

5.1 Finding coresets

The coreset method described in Feldman, Faulkner and Krause (2011) can be used to find an appropriate weighted subset to represent the information in the complete dataset. The starting point is to first sample uniformly a small number of points, then remove half of the data points which are closest to the sampled points. Next, sample again from the rest of the points and remove half of the points lying closest until all of the data points are labeled as removed or sampled.

By doing this, we construct a hierarchy of data points and the importance-weight of the sampled points is associated with the log-likelihood. The weights are set to be optimal if the estimated log-likelihood is of the least variance. This construction of a sampled set gives a higher probability to observations that are further away from the initial cluster center, and the sampling bias would be fixed by adapting the weight which is to be associated with the sampling probabilities. We can then finally build an appropriate coreset from the whole dataset based on the weights.

This algorithm for coreset construction is more formally summarised in pseudo-code form in Algorithm 2.

5.2 VB inference using coreset sampling

In this section, we propose a new algorithm in which we adapt the variational Bayes method in order that it can be used in conjunction with a coreset sampling approach. We will use the standard prior settings as described in Section 3 and the posterior of the model will be adjusted using the coreset samples and coreset weights. This novel modification of the algorithm makes it suitable for use in analysing a coreset of the image. The update equations that have to be adapted for

Algorithm 2 Algorithm for finding a coreset (see Feldman, Faulkner and Krause (2011))

input: Whole dataset D , ϵ , δ , K .

set: $D' \leftarrow D$; $B \leftarrow \emptyset$

Specify initial allocation of observations to components.

while $|D'| \geq 10dK \ln(1/\delta)\epsilon$ **do** :

Sample set S of $\beta = 10dK \ln(1/\delta)$ points uniformly at random from D' ;

Remove $\lceil |D'|/2 \rceil$ points $\mathbf{x} \in D'$ closest to S (i.e., minimising $\text{dist}(\mathbf{x}, S)$) from D' ;

Set $B \leftarrow B \cup S$;

Set $B \leftarrow B \cup D'$

for every $b \in B$ **do**

$D_b \leftarrow$ the points in D whose closest point B is b .

for every $b \in B$ and every $\mathbf{x} \in D_b$ **do**

$$m(\mathbf{x}) \leftarrow \lceil \frac{5}{|D_b|} + \frac{\text{dist}(\mathbf{x}, B)^2}{\sum_{\mathbf{x}' \in D} \text{dist}(\mathbf{x}', B)^2} \rceil;$$

Pick a non-uniform random sample C of $10 \lceil dk|B|^2 \ln(1/\delta)/\epsilon^2 \rceil$ points from D , where for every $\mathbf{x}' \in C$ and $\mathbf{x}' \in D$, we have $\mathbf{x}' = \mathbf{x}$ with probability

$$m(\mathbf{x}) / \sum_{\mathbf{x}' \in D} m(\mathbf{x}');$$

for each $\mathbf{x}' \in C$ **do** $\gamma(\mathbf{x}') \leftarrow \frac{\sum_{\mathbf{x} \in D} m(\mathbf{x})}{|C| \cdot m(\mathbf{x}'')}$

output: Coreset $C = \{(\gamma(\mathbf{x}_1), \mathbf{x}_1), (\gamma(\mathbf{x}_2), \mathbf{x}_2), \dots, (\gamma(\mathbf{x}_{|C|}), \mathbf{x}_{|C|})\}$.

use with the coreset data are Equations 2–7. In those equations, \mathbf{y}_i is replaced by $\hat{\mathbf{y}}_i$, where $\hat{\mathbf{y}}_i$ corresponds to the weighted observations and is defined as follows:

$$\hat{\mathbf{y}}_i = \frac{\gamma_i \times \mathbf{y}_i}{\frac{1}{n} \sum_{i=1}^n \gamma_i}.$$

The expected values required to update the expressions remain unaltered from the form they take in the standard VB algorithm. Pseudo-code for the weighted VB algorithm which we call coreset variational Bayes (CVB) is outlined in Algorithm 3.

6 Application to weed-crop image segmentation using coreset variational Bayes (CVB)

6.1 Weed image

We illustrate the effectiveness of the CVB algorithm for analysing an image of a weed plant amongst soil and dead leaves, see Figure 1. The size of the coreset is 2599 pixels, which is 1/100 of the size of the original dataset. The idea is to use the CVB algorithm to segment the image and classify the pixels as representing

Algorithm 3 Coreset variational Bayes (CVB) algorithm

input: $C = \{(\gamma(\mathbf{x}_1), \mathbf{x}_1), (\gamma(\mathbf{x}_2), \mathbf{x}_2), \dots, (\gamma(\mathbf{x}_N), \mathbf{x}_N)\}$ from Algorithm 2.
 Set initial values for hyper-parameters, $\varepsilon, \alpha^{(0)}, \beta^{(0)}, \Sigma^{(0)}, v^{(0)}, m^{(0)}$.
 Specify initial allocation of observations to components via initial q_{ij} .

while Not Converged, **do**:
 Update modified variational posterior expressions for model parameters
 Update modified variational posterior for q_{ij}
if any component has a mixing weight $\leq \varepsilon$
 remove the component from model
end if
if the algorithm has converged
 exit loop.
end if

end while

either living plant, background soil, or dead leaves. This is the type of classification that would be required for research purposes in agricultural trials. In running the analysis, the hyper-parameters are chosen to correspond to vague prior settings.

Figure 2 shows the segmentation of each component into the three different types of matter. As we can see, the area of the living part of the weed is clearly defined, and the algorithm can also distinguish between soil and dead leaves. The numerical segmentation results are in Table 1 and for comparison we also show



Figure 1 *The original image of a weed against background soil to be analysed.*

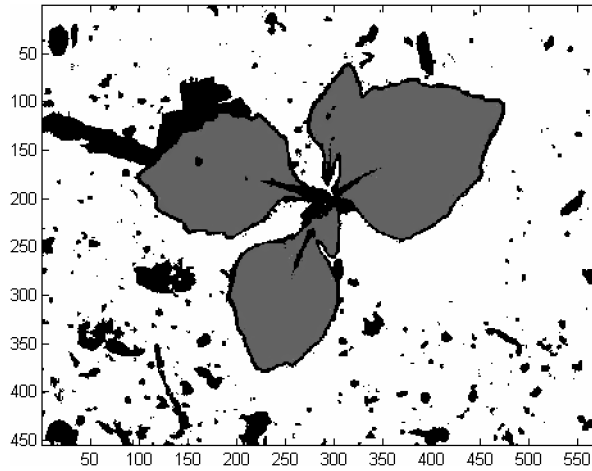


Figure 2 Result of the CVB analysis of the weed image. The pixels in the image have been segmented into 3 different components representing the types background soil, plant and dead leaf.

the results that we would have obtained from an analysis of the full dataset. There is close agreement between these results showing that the coreset modification of the VB algorithm is reliable and useful. Significantly, the VB coreset algorithm is greatly more time-efficient than the standard VB as it runs around 18 times faster (around 40 minutes for CVB compared to around 12 hours for VB on the full dataset). This is very impressive when we consider how similar the final results were. Hence, CVB is a practical and useful alternative to standard VB for the image segmentation based on finite Gaussian mixture models.

7 Discussion

We have presented a new algorithm (CVB) for analysing data using variational Bayes based on a representative coreset of the data. This allows us to perform re-

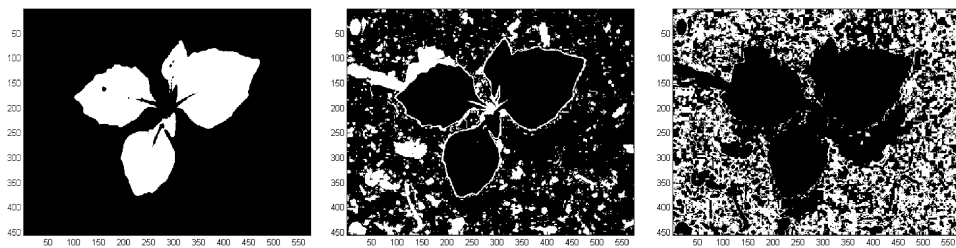


Figure 3 The three different components identified via the CVB algorithm shown in three separate plots. The first plot shows, in white, the area of living plant (weed), the second shows dead leaves. The third plot shows regions of background soil.

Table 1 *Result comparison*

Component type	Standard VB		VB coresets	
	Mixing weight	Mean	Mixing weight	Mean
Background soil	0.7005	$\begin{bmatrix} 0.1075 \\ 0.0986 \\ 0.0941 \end{bmatrix}$	0.6989	$\begin{bmatrix} 0.1093 \\ 0.0999 \\ 0.0952 \end{bmatrix}$
Plant	0.1765	$\begin{bmatrix} 0.3200 \\ 0.4396 \\ 0.1753 \end{bmatrix}$	0.1894	$\begin{bmatrix} 0.3233 \\ 0.4422 \\ 0.1810 \end{bmatrix}$
Dead leaves	0.1231	$\begin{bmatrix} 0.2660 \\ 0.2447 \\ 0.1986 \end{bmatrix}$	0.1117	$\begin{bmatrix} 0.2768 \\ 0.2561 \\ 0.2059 \end{bmatrix}$

liable inference in a highly time-efficient way. However, the running time after algorithm modification is still around 40 minutes for a medium size (67 KB) image. While this is good in comparison to other existing approaches, there is still much scope for further research if these ideas are to be put to routine use. Consider that agricultural activities may need weed detection means applied on a large area of land, and thus the image can be more than 1 GB, and contain even more detail than the presented example. One option to explore might be looking at improving computational speed through the use of more efficient programming algorithms, more sophisticated computers or GPU programming for increased efficiency Suchard et al. (2010). This is a topic for future research.

References

- Alston, C., Mengersen, K. and Pettitt, A. N. (2012). *Case Studies in Bayesian Statistical Modelling and Analysis*, 1st ed. New York: John Wiley & Sons.
- Faes, C., Ormerod, J. and Wand, M. (2011). Variational Bayesian inference for parametric and non-parametric regression with missing data. *Journal of the American Statistical Association* **106**, 959–971. [MR2894756](#)
- Feldman, D., Faulkner, M. and Krause, A. (2011). Scalable training of mixture models via coresets. In *Advances in Neural Information Processing Systems 24: 25th Annual Conference on Neural Information Processing Systems 2011* (J. Shawe-Taylor, R. S. Zemel, P. Bartlett, F. Pereira and K. Q. Weinberger, eds.) 2142–2150. NY: Curran Associates, Inc.
- Fruhwirth-Schnatter, S. (2006). *Finite Mixture and Markov Switching Models*. New York: Springer.
- Kargar, A. H. B. and Shirzadifar, A. M. (2013). Automatic weed detection system and smart herbicide sprayer robot for corn fields. In *Proceeding of the 2013 RSI/ISM International Conference on Robotics and Mechatronics. February 13–15, Tehran, Iran*.
- Marin, J. M., Pudlo, P., Robert, C. P. and Ryder, R. (2012). Approximate Bayesian computation methods. *Statistics and Computing* **22**, 1167–1180.
- McGrory, C. A., Ahfock, D., Horsley, J. and Alston, C. L. Weighted Gibbs sampling for mixture modelling of massive datasets via coresets. *Stat* **3**, 291–299.

- McGrory, C. A., Pettitt, A. N., Reeves, R., Griffin, M. and Dwyer, M. (2012). Variational Bayes and the reduced dependence approximation for the autologistic model on an irregular grid with applications. *Journal of Computational and Graphical Statistics* **21**, 781–796. [MR2970919](#)
- McGrory, C. A. and Titterington, D. M. (2007). Variational approximations in Bayesian model selection for finite mixture distributions. *Computational Statistics & Data Analysis* **51**, 5352–5367. [MR2370876](#)
- Ormerod, J. T. and Wand, M. P. (2010). Explaining variational approximations. *American Statistician* **64**, 140–153. [MR2757005](#)
- Sindenab, J., Jonesbc, R., Hesterba, S., Odomba, D., Kalischda, C., Jamese, R. and Cacho, O. (2004). *The Economic Impact of Weeds in Australia. CRC for Australian Weed Management Technical Series No. 8.*
- Suchard, M. A., Wang, Q., Chan, C., Frelinger, J., Cron, A. J. and West, M. (2010). Understanding GPU programming for statistical computation: Studies in massively parallel massive mixtures. *Journal of Computational and Graphical Statistics* **19**, 419–438.
- Wand, M., Ormerod, J., Padoan, S. and Fruhwirth, R. (2012). Mean field variational Bayes for elaborate distributions. *Bayesian Analysis* **6**, 847–900.
- Zimdahl, R. (2009). *Fundamentals of Weed Science*. San Diego: Academic Press.

Q. Liu
C. A. McGrory
Centre for Applications in Natural
Resource Mathematics
School of Mathematics
University of Queensland
Brisbane, Qld 4072
Australia
E-mail: c.mcgrory@uq.edu.au

P. W. J. Baxter
Centre for Applications in Natural
Resource Mathematics
School of Mathematics
University of Queensland
Brisbane, Qld 4072
Australia
and
Queensland University of Technology
Brisbane, Qld 4001
Australia