

Particle Swarm Optimization Assisted Metropolis Hastings Algorithms

Abstract

Fitting dependent data models is often a challenging endeavor and often requires some form of dimension reduction or customized estimation algorithm. Particle swarm optimization (PSO) refers to a class of heuristic optimization algorithms that exploit analogies with animal flocking behavior in order to obtain optima without strong conditions on the objective function. We introduce two new classes of PSO algorithms, termed adaptively tuned PSO (AT-PSO) and adaptively tuned bare bones PSO (AT-BBPSO). In both algorithms we add a dynamically tuned parameter to previously existing PSO algorithms. We propose using these PSO algorithms to approximate Bayesian posterior distributions in order to construct efficient proposals for independent Metropolis-Hastings and independent Metropolis-Hastings within Gibbs algorithms. For the latter, we propose using a global approximation to the posterior in order to construct an approximation to the conditional distribution which requires a Metropolis step, thereby requiring the optimization algorithm only once rather than every iteration of the Markov chain Monte Carlo (MCMC) algorithm. In order to illustrate our method and compare it to alternatives, we provide a simulation study and apply it to constructing MCMC algorithms to estimate reduced rank spatial models of American Community Survey (ACS) 5-year estimates of county populations in the United States.

KEY WORDS: Bayesian estimation; Markov chain Monte Carlo; Official statistics; Optimization; Spatial data

1 Introduction

[REWRITE THIS PARAGRAPH - IS TALL DATA WORTH MENTIONING?] One common type of big data problem may be called “tall” data — that is many observations of a relatively small number of variables. Bayesian estimation of tall data models can be particularly challenging in the dependent data and non-Gaussian settings due to the need to estimate or integrate out a high dimensional latent process, e.g., with dimension equal to the size of the dataset. Sometimes the latent process can be written as a function of a relatively small number of latent random variables, but, even in this case, standard Markov chain Monte Carlo (MCMC) techniques to approximate the posterior are often still slow. A common alternative is to leverage normal approximations to the posterior to obtain approximations to the marginal distributions of each parameter, e.g., INLA (Rue et al., 2009). This approximation is often very good and much faster than MCMC, though neither advantage is universally true, particularly when features of the joint posterior distribution rather than just the marginals are desired or when the parameter space is too large (e.g., see Taylor and Diggle (2014)). [CHECK ON THE JOINT VS MARGINAL THING AND OTHER INLA LIMITATIONS — CITATIONS!!!] Another common strategy is to use Hamiltonian Monte Carlo (HMC) (Neal et al., 2011), particularly the No-U-Turn sampler (NUTS) (Homan and Gelman, 2014) which is implemented in the Stan software (Carpenter et al., 2015). However HMC requires many log posterior evaluations per iteration and in the tall data setting these are typically expensive.

We propose using old MCMC technology improved by new optimization techniques. It is well known that a Bayesian posterior distributions for a fixed set of parameters tends asymptotically to a normal distribution centered at the posterior mode as the sample size increases; see Schervish (1997, Chapter 7.4). This normal approximation, often called a Laplace approximation, is used by INLA but is also frequently used to construct independent

Metropolis-Hastings (IMH) samplers (Metropolis et al., 1953; Hastings, 1970) — typically using a t proposal instead of a normal. IMH samplers based on a good approximation to the posterior distribution typically have high acceptance rates along with fast mixing and convergence, but it is often impractical to find an adequate approximation. Even in cases where the normal approximation is appropriate the posterior may be too high dimensional for this approach to be practical. In larger models numerical methods are usually required to find the posterior mode and sometimes also to integrate out a latent process, but standard methods for doing so are usually prohibitively slow or fail outright. However, new heuristic optimization algorithms tend to work reasonably well over a wider class of objective functions and in larger parameter spaces, for example, genetic algorithms (Goldberg and Holland, 1988) and particle swarm optimization (PSO) (Clerc and Kennedy, 2002; Blum and Li, 2008; Clerc, 2010). We propose using PSO to obtain the posterior mode so that effective independent Metropolis samplers can be constructed in a wider range of models. Other heuristic optimization algorithms may be useful for the same task, but we do not explore them here. We also develop several novel PSO algorithms that help obtain good estimates of the posterior mode and that may also have utility in other optimization contexts. We illustrate our method on two different modeling applications. The first application considers a group of reduced rank spatial models of American Community Survey (ACS) 5-year period estimates of county population estimates for the U.S. in 2014. In the second application, we adapt a hierarchical model for predicting the outcome of the 1988 presidential election from Gelman and Hill (2006). We fit these models using our proposed independent Metropolis-Hastings algorithm, an independent Metropolis within Gibbs algorithm based on the same Laplace approximation, and a variety of other MCMC algorithms. We compare the computational cost for each of the algorithms and discuss their ease of use.

[REWRITE THIS PARAGRAPH - SECTIONS HAVE CHANGED]

2 Particle swarm optimization

We briefly describe PSO here; refer to Blum and Li (2008) for an excellent introduction and Clerc (2010) for more detail. Suppose that we wish to optimize some objective function $Q(\boldsymbol{\theta}) : \mathbb{R}^D \rightarrow \mathbb{R}$ — without loss of generality we will assume the goal is maximization. Let $i = 1, 2, \dots, n$ index a set of particles over time, $t = 1, 2, \dots, T$, where in every period each particle consists of a location $\boldsymbol{\theta}_i(t) \in \mathbb{R}^D$, a velocity $\mathbf{v}_i(t) \in \mathbb{R}^D$, a personal best location $\mathbf{p}_i(t) \in \mathbb{R}^D$, and a group best location $\mathbf{g}_i(t) \in \mathbb{R}^D$. Here we mean “best” in the sense of maximizing Q , so $Q(\mathbf{p}_i(t)) \geq Q(\boldsymbol{\theta}_i(s))$ for any $s \leq t$. The group best location is defined with respect to some neighborhood \mathcal{N}_i of particle i ; that is, $\mathbf{g}_i(t) = \arg \max_{\{\mathbf{p}_j | j \in \mathcal{N}_i\}} Q(\mathbf{p}_j(t))$. In the simplest case where the entire swarm is the neighborhood of each particle, $\mathbf{g}_i(t) \equiv \mathbf{g}(t) = \arg \max_{j \in 1:n} Q(\mathbf{p}_j(t))$. The generic PSO algorithm updates as follows:

$$\begin{aligned}
 \mathbf{v}_i(t+1) &= \omega \mathbf{v}_i(t) + \phi_1 \mathbf{r}_{1i}(t) \circ \{\mathbf{p}_i(t) - \boldsymbol{\theta}_i(t)\} + \phi_2 \mathbf{r}_{2i}(t) \circ \{\mathbf{g}_i(t) - \boldsymbol{\theta}_i(t)\}, \\
 \boldsymbol{\theta}_i(t+1) &= \boldsymbol{\theta}_i(t) + \mathbf{v}_i(t+1), \\
 \mathbf{p}_i(t+1) &= \begin{cases} \mathbf{p}_i(t) & \text{if } Q(\mathbf{p}_i(t)) \geq Q(\boldsymbol{\theta}_i(t+1)) \\ \boldsymbol{\theta}_i(t+1) & \text{otherwise,} \end{cases} \\
 \mathbf{g}_i(t+1) &= \arg \max_{\{\mathbf{p}_j(t+1) | j \in \mathcal{N}_i\}} Q(\mathbf{p}_j(t+1)), \tag{1}
 \end{aligned}$$

where \circ denotes the Hadamard product (element-wise product), $\mathbf{r}_{1i}(t)$ and $\mathbf{r}_{2i}(t)$ are each vectors of D random variates independently generated from the $U(0, 1)$ distribution, and $\omega > 0$, $\phi_1 > 0$, and $\phi_2 > 0$ are user-defined parameters. The term $\omega \mathbf{v}_i(t)$ controls the particle’s tendency to keep moving in the direction it is already going, so ω is called the inertia parameter. For $\omega < 1$ velocities tend to decrease over time, while for $\omega > 1$ they tend to increase over time. Similarly $\phi_1 \mathbf{r}_{1i}(t) \circ (\mathbf{p}_i(t) - \boldsymbol{\theta}_i(t))$ controls the particle’s tendency to move towards its personal best location while $\phi_2 \mathbf{r}_{2i}(t) \circ (\mathbf{g}_i(t) - \boldsymbol{\theta}_i(t))$ controls its tendency to move toward the group’s best location, so ϕ_1 and ϕ_2 are called the cognitive correction

factor and social correction factor, respectively (Blum and Li, 2008). This version of PSO is equivalent to Clerc and Kennedy (2002)’s constriction type I particle swarm. There are many variants of the PSO algorithm, often obtained through choosing (ω, ϕ_1, ϕ_2) in special ways or sometimes even dynamically. A default version of the algorithm sets $\omega = 0.7298$ and $\phi_1 = \phi_2 = 1.496$; see Clerc and Kennedy (2002) and Blum and Li (2008) for justification of these choices. Even when $\omega < 1$, if ϕ_1 and ϕ_2 are set high enough the velocities of the particles can continually increase and cause the swarm to make jumps that are much too large. A heavy handed way to solve this problem is by setting an upper bound on the velocity of any particle in any given direction, called velocity clamping. However, the default parameter values suggested by Clerc and Kennedy (2002) are also designed to prevent exactly this sort of velocity explosion.

Any PSO variant can also be combined with various neighborhood topologies that control how the particles communicate to each other. The default global topology allows each particle to see each other particle’s previous best location for the social components of their respective velocity updates, but this can cause inadequate exploration and premature convergence. Alternative neighborhood topologies limit how many other particles each particle can communicate with. For example, particle 5 may only look at itself and particles 4 and 6 when determining what its group best location is. This allows information about high value locations in the domain of the objective function to eventually reach every particle in the swarm, but much more slowly so that each particle has an opportunity to explore the space more fully first. Appendix A contains a short description of the ring topologies, but there are many alternatives in the literature.[CITATION]

A variant of PSO, the bare bones PSO algorithm (BBPSO) is a PSO algorithm introduced by Kennedy (2003) that strips away the velocity term and removes the need for the user to choose parameters outside of the swarm size. Let $\theta_{ij}(t)$ denote the j th coordinate of the position for the i th particle in period t , and similarly for $p_{ij}(t)$ and $g_{ij}(t)$. Then the BBPSO

algorithm obtains a new position coordinate θ_{ij} via

$$\theta_{ij}(t+1) \sim N\left(\frac{p_{ij}(t) + g_{ij}(t)}{2}, |p_{ij}(t) - g_{ij}(t)|^2\right) \quad (2)$$

where $N(\mu, \sigma^2)$ is the normal distribution with mean μ and standard deviation σ^2 . The updates of $\mathbf{p}_i(t)$ and $\mathbf{g}_i(t)$ are the same as in (1). There are several variants of this algorithm proposed including using distributions from different location-scale families — e.g., using the t -distribution and modifying the location or the scale parameters, for example Krohling et al. (2009), Hsieh and Lee (2010), Richer and Blackwell (2006), and Campos et al. (2014). Appendix A contains more detail on the BBPSO variants in the literature which we employ.

2.1 Adaptively tuned BBPSO

BBPSO adapts the size effective search space of the swarm over time through the variance term: $|p_{ij}(t) - g_{ij}(t)|^2$. As the personal best locations of the swarm move closer together, these variances decrease and the swarm tries locations which are closer to known high value areas in the space. This behavior is desirable, but the adaptation is forced to occur through one channel: the personal and group best locations. If the personal best locations of the swarm are arranged in a rough ring around the global optimum, smaller variances are desirable so that the new particle locations have a tendency to be in the center of the ring. On the other hand, if the personal best locations of the swarm are all to one side of the optimum and fairly far away, larger variances are desirable so that the particles can quickly approach the neighborhood of the optimum. BBPSO cannot distinguish between these two cases. [CAN ADD A GRAPHIC ILLUSTRATING THIS. WORTH IT?]

In order to allow it to adapt in a more flexible manner, we modify the BBPSO variance to $\sigma^2(t)|p_{ij}(t) - g_{ij}(t)|^2$ and tune $\sigma^2(t)$ in a manner similar to adaptive random walk Metropolis MCMC algorithms (Andrieu and Thoms, 2008). Define the improvement rate of the swarm in period t as $R(t) = \#\{i : Q(\mathbf{p}_i(t)) > Q(\mathbf{p}_i(t-1))\}/n$ where if A is a set then $\#A$ is

the number of members of that set, and let R^* denote the target improvement rate. Then what we term adaptively tuned BBPSO (AT-BBPSO) updates personal best and group best locations as in (1), then updates particle locations as follows:

$$\begin{aligned}\theta_{ij}(t+1) &\sim t_{df} \left(\frac{p_{ij}(t) + g_{ij}(t)}{2}, \sigma^2(t) |p_{ij}(t) - g_{ij}(t)|^2 \right), \\ \log \sigma^2(t+1) &= \log \sigma^2(t) + c \times \{R(t+1) - R^*\},\end{aligned}\tag{3}$$

where df is a user chosen degrees of freedom parameter and c is another user chosen parameter that controls the speed of adaptation. We use a t kernel instead of a Gaussian in order to allow for more flexibility, and smaller values of df appear to combine well with adaptively tuning $\sigma^2(t)$ — in particular $df = 1$. Since the target rate in AT-BBPSO is similar to the target rate in an adaptive random walk Metropolis algorithm, a priori values around 0.25 and in particular below 0.5 seem reasonable (Gelman et al., 1996). In practice we find that values from $R^* = 0.3$ to $R^* = 0.5$ tend to yield good AT-BBPSO algorithms. The parameter c controls the speed of adaptation so that larger values of c mean the algorithm adapts $\sigma^2(t)$ faster. We find that $c = 0.1$ to be a good value, though anything within an order of magnitude often yields similar results. The initial value $\sigma^2(0)$ also needs to be chosen, though this does not have much impact on the algorithm as long as c is not too small. We use $\sigma^2(0) = 1$ to initialize the algorithm at the standard BBPSO algorithm.

Both using a t kernel and adding a fixed scale parameter have been discussed in the BBPSO literature, but as Kennedy (2003) notes, something about setting $\sigma = 1$ is special that causes the algorithm to work well. Another similar BBPSO algorithm in the literature comes from Hsieh and Lee (2010). They propose a modified version of BBPSO with

$$\theta_{ij}(t+1) \sim N \left(\omega \frac{p_{ij}(t) + g_{ij}(t)}{2}, \sigma^2 |p_{ij}(t) - g_{ij}(t)|^2 \right),$$

where $\omega \leq 1$ and $\sigma^2 \leq 1$ are constriction parameters that are eventually both set to one after enough iterations of the algorithm. The authors suggest dynamically adjusting the

constriction parameters in the early stage of the algorithm before they are set to one, but give no suggestion for how to do this. The default BBPSO algorithm essentially sets $\sigma^2(t) = 1$ for all t and Hsieh and Lee (2010)’s algorithm deterministically adjusts $\sigma^2(t)$, but our AT-BBPSO algorithm is able to adapt its value on the fly based on local knowledge about the objective function. If too much of the swarm is failing to find new personal best locations, AT-BBPSO proposes new locations closer to known high value areas. If too much of the swarm is improving, AT-BBPSO proposes bolder locations in an effort to make larger improvements. This ability to adapt to local information about the objective function allows AT-BBPSO to more quickly traverse the search space towards the global optimum, though by using local information AT-BBPSO does risk premature convergence to a local optimum. The adaptively tuned component of AT-BBPSO can also be combined with most BBPSO variants, some of which are outlined in Appendix A. In Appendix B we conduct a simulation study on several test functions that compares AT-BBPSO variants to other PSO and BBPSO variants in order to justify the parameter settings discussed above and demonstrate that making a BBPSO variant into an AT-BBPSO variant will tend to improve it.

2.2 Adaptively tuned PSO

In AT-BBPSO variants, the parameter $\sigma(t)$ partially controls the effective size of the swarm’s search area, and we increase or decrease $\sigma(t)$ and consequently the search area depending on how much of the swarm is finding new personal best locations. In standard PSO there is no direct analogue to $\sigma(t)$, though the inertia parameter, ω in (1), is related. It controls the effective size of the swarm’s search area by controlling how the magnitude of the velocities evolve over time — larger values of ω allow for larger magnitude velocities in future periods. In AT-BBPSO we use an analogy with tuning a random walk Metropolis-Hastings MCMC algorithm in order to build intuition about how to tune $\sigma(t)$. The analogy is much weaker

in this case; nonetheless, we allow $\omega(t)$ to be time-varying and use the same mechanism in order to tune it as we did for $\sigma(t)$.

The idea of time-varying $\omega(t)$ has been in the PSO literature for some time. An early suggestion was to set $\omega(0) = 0.9$ and deterministically decrease it until it reaches $\omega(T) = 0.4$ after the maximum number of iterations allowed (Eberhart and Shi, 2000). In particular, Tuppadung and Kurutach (2011) suggest defining $\omega(t)$ via the parameterized inertia weight function

$$\omega(t) = \frac{1}{1 + \left(\frac{t}{\alpha}\right)^\beta} \quad (4)$$

where α and β are user-defined parameters. Roughly, α controls how low $\omega(t)$ can go and β controls how fast it gets there, so α and β can be thought of as intercept and slope parameters respectively. The suggestion in Tuppadung and Kurutach (2011) is to set α to a small fraction of the total amount of iterations in which the algorithm is allowed to run (e.g., 10% or 20%), and set β between one and four.

This approach tends to improve on standard PSO if $\omega(t)$'s progression is set appropriately, but it invariably makes using PSO more difficult for the average user. Additionally, depending on the problem, it may be more useful to let the swarm explore the space for more or less iterations, necessitating different progressions of $\omega(t)$. A priori it may not be clear exactly which approach is best for any given problem, so an automatic method is desirable. Adaptively tuned PSO (AT-PSO) is just that — it provides an automatic method to adjust the value of $\omega(t)$ depending on local information obtained by the particle swarm. Formally, AT-PSO updates personal and group best locations as in (1) and updates $\omega(t)$ and particle

locations as follows:

$$\begin{aligned}
\mathbf{v}_i(t+1) &= \omega(t+1)\mathbf{v}_i(t) + \phi_1\mathbf{r}_{1i}(t) \circ \{\mathbf{p}_i(t) - \boldsymbol{\theta}_i(t)\} + \phi_2\mathbf{r}_{2i}(t) \circ \{\mathbf{g}_i(t) - \boldsymbol{\theta}_i(t)\}, \\
\boldsymbol{\theta}_i(t+1) &= \boldsymbol{\theta}_i(t) + \mathbf{v}_i(t+1), \\
\log \omega(t+1) &= \log \omega(t) + c \times \{R(t+1) - R^*\},
\end{aligned} \tag{5}$$

where $R(t)$ is the improvement rate of the swarm in iteration t , R^* is the target improvement rate, and c controls how much $R(t)$ changes on a per iteration basis. For AT-BBPSO we used an analogy with random walk Metropolis-Hastings algorithms to suggest that a good value for the target improvement rate is smaller than 0.5 but not too small and this turned out to be correct. The analogy does not apply as well here, though we still find in Appendix B that $R^* = 0.3$ or 0.5 still seems to work well for AT-PSO. The value of c controls the speed of adaptation, and in particular if c is small and $\omega(0)$ is large, AT-PSO can mimic DI-PSO to some extent. This turns out to produce poor AT-PSO algorithms, however. We use $c = 0.1$ as a default value and in simulations not reported here, we find that the gains from optimizing c appear to be small. However, very small values like those suggested by an attempt to mimic DI-PSO turn out to cause the algorithm to perform very poorly.

A major strength of AT-PSO relative to DI-PSO and standard PSO is that AT-PSO can increase $\omega(t)$ when information from the swarm suggests there is an unexplored high value region of the space — when too much of the swarm is improving on their personal best locations AT-PSO increases $\omega(t)$ until velocities start increasing, the swarm starts exploring a larger amount of the nearby space, and more of the particles fail to find improvements on their personal best. Just like in AT-BBPSO, this mechanism provides a way for the swarm to adapt its behavior on the fly based on local conditions but it can also cause premature convergence to a local optimum. While DI-PSO monotonically decreases $\omega(t)$ toward some minimum value, AT-PSO typically oscillates $\omega(t)$ so that the algorithm alternates between exploring and exploiting more, relative to standard PSO. [CAN ADD A PLOT OF INERTIA

OVER TIME HERE. WORTH IT?] Appendix B contains an extended simulation study comparing a variety of these PSO and BBPSO algorithms on a suite of test functions that motivates some of the recommendations detailed above and demonstrates that AT-PSO tends to be consistently one of the best performing PSO or BBPSO algorithms and often is the best.

3 Laplace Approximations and Independent Metropolis-Hastings

Ultimately we will use PSO algorithms to help construct MCMC algorithms for posterior sampling, in particular independent Metropolis-Hastings (IMH) algorithms based on Laplace approximations. To do this, first we introduce the Laplace approximation. Let $\boldsymbol{\theta}$ be the model parameter and $[\boldsymbol{\theta}|\mathbf{Y}_{1:n}]$ be its posterior distribution, available up to a normalizing constant. We use the generic notation $[X|Y]$ to denote the density or mass function of X given Y . Further let $\boldsymbol{\theta}_n^*$ denote the posterior mode of $[\boldsymbol{\theta}|\mathbf{Y}_{1:n}]$ and let $\mathbf{H}_n(\boldsymbol{\theta}_n^*)$ denote the Hessian matrix of $\log[\boldsymbol{\theta}|\mathbf{Y}_{1:n}]$ evaluated at $\boldsymbol{\theta}_n^*$. Then under suitable regularity conditions (Schervish, 1997, Sections 7.4.2 and 7.4.3) $\boldsymbol{\theta}$'s posterior distribution is asymptotically normal so that for a fixed but large value of n , $\boldsymbol{\theta}|\mathbf{Y}_{1:n} \stackrel{a}{\sim} N(\boldsymbol{\theta}_n^*, -\mathbf{H}_n^{-1}(\boldsymbol{\theta}_n^*))$ where the notation $\stackrel{a}{\sim}$ means “approximately distributed as.” We do not discuss the precise technical conditions necessary for this result here, but intuitively for the approximation to be good each element of $\boldsymbol{\theta}$ must have enough “observations” for asymptotics to kick in. We put “observations” in scare quotes because often rather than data, it is other elements of $\boldsymbol{\theta}$ are directly informing on another element. In order to see this more clearly we add structure to the model by considering a class of generalized linear mixed models (GLMMs) (Stroup, 2012), but by no means is this the only class of models where the Laplace approximation works well.

3.1 Laplace approximations for GLMMs

Generalized linear mixed models with latent Gaussian processes (LGP) provide a plethora of examples where PSO assisted IMH algorithms are attractive for MCMC. That latent parameters are Gaussian distributed often increases the quality of the normal approximation to the posterior, but it is not a necessary feature of the model for the approximation to work well. We will conceptualize our model class using the strategy of Berliner (1996) and Wikle et al. (2003), that is hierarchically with a data model conditional on parameters and a latent process, a process model conditional on parameters, and finally a parameter model. Let $\mathbf{Z} = (Z_1, Z_2, \dots, Z_n)$ denote a vector of response variables. We assume a conditionally independent data model given a vector of location parameters $\boldsymbol{\mu} = (\mu_1, \mu_2, \dots, \mu_n)$ and a common dispersion parameter ϕ , i.e. $Z_i \stackrel{ind}{\sim} [Z|\mu_i, \phi]$ where $[z|\mu, \phi]$ is some known family of density functions indexed by (μ, ϕ) , often an exponential family. The location parameters are a known function g of a LGP, so $g^{-1}(\boldsymbol{\mu}) = \mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{S}\boldsymbol{\delta}$. Here $\mathbf{X}'\boldsymbol{\beta}$ represents fixed effects and $\mathbf{S}'\boldsymbol{\delta}$ represents random effects, where \mathbf{X} is a known $n \times p$ matrix, \mathbf{S} is a known $n \times r$ matrix, $\boldsymbol{\beta}$ is an unknown p -dimensional vector and $\boldsymbol{\delta}$ is an unknown r -dimensional vector. Often $r \ll n$, but in some cases $r = n$. $\boldsymbol{\delta}$ is further modeled as Gaussian with mean zero and a covariance matrix $\boldsymbol{\Sigma}$ with a structure appropriate to the specific problem. Finally a prior on $(\phi, \boldsymbol{\beta}, \boldsymbol{\Sigma})$ serves as the parameter model where typically $\boldsymbol{\beta}$ is assumed normally distributed a priori so that $g^{-1}(\boldsymbol{\mu})$ is an LGP. The model can be generalized to multivariate Z_i and multivariate Y_i , but we omit this complication here.

Three key factors affect the quality of the Laplace approximation for GLMMs: the sample size, how close the data model is to Gaussian at reasonable parameter values for the given dataset, and the structure of $\boldsymbol{\Sigma}$. Due to the asymptotic nature of the Laplace approximation, the larger the sample size the better the approximation is for $\boldsymbol{\beta}$ and ϕ . To the extent that the data model is already approximately normal, smaller sample sizes will still

yield reasonably good approximations. [I CAN ILLUSTRATE THIS WITH SOME SIMULATIONS IN A GLMM WITH A BETA DATA MODEL. NECESSARY?] The last factor can best be illustrated by considering two extremes: iid random effects where $\Sigma = \sigma^2 \mathbf{I}_r$, and fully correlated random effects where Σ is fully parameterized and unknown. In the iid case when r is moderately large, e.g. $r = 10$, σ^2 or more likely $\log \sigma^2$ is often approximately normal. The “observations” directly informing on σ^2 are the estimated random effects in δ rather than data, but the Laplace approximation works all the same. In the fully correlated case, the only “observation” relevant to estimating Σ is δ , and there is only one observation to estimate $r(r + 1)/2$ unique parameters. In this case it is unrealistic to expect Σ or any commonly used transformation of Σ to be approximately normal in the posterior unless there is strong information in the prior suggesting otherwise.

When $r = n$ it is more important that the data model is already approximately normal and that the random effects are a priori modeled as Gaussian because there are only n data points informing on n latent random effects. Since the dimension of δ is growing at the same rate as the sample size in that case, asymptotic arguments about δ ’s posterior distribution do not apply. Additionally, when n is large and $r = n$ it may be difficult to find the posterior mode in a reasonable amount of time using any optimization algorithm. Many GLMMs feature e.g. group specific random effects, so r is often much smaller than n and these concerns do not apply. The same principles apply to other models outside of the GLMM framework — each model parameter and random effect needs enough effective observations informing it so that the asymptotics can take hold or else other features of the model need to ensure that the parameters are approximately Gaussian in the posterior.

3.2 Laplace approximations for IMH algorithms

Given any reasonable approximation to the posterior distribution, we can construct an IMH MCMC sampler using that approximation (Robert and Casella, 2013, Chapter 7.4), though typically a t_{df} proposal is used instead with df set small enough so that its tails dominate the posterior's in order to ensure that the Markov chain is uniformly ergodic — see Robert and Casella (2013, Theorem 7.8). As long as the conditions mentioned in Section 3.1 are satisfied, the Laplace approximation is likely to be good and the IMH algorithm is likely to yield efficient MCMC. In particular, when the approximation is good, IMH has high acceptance rates, converges quickly, and mixes well. This is well known, but the problem lies in finding good approximations. The Laplace approximation is often poor even for the class of GLMMs we discussed in Section 3.1 because of the random effect covariance matrix. Often there are a small number of random effects or they are allowed to have less structured covariance matrices, so the asymptotics never kick in. Yet the random effects themselves and the fixed effect parameter are often still approximately normal in their conditional posterior. Further, the dispersion parameter and the random effect covariance matrix often have easy conditional posteriors which are easy to sample from.

This suggests an obvious alternative: construct a Metropolis within Gibbs MCMC algorithm (Robert and Casella, 2013, Chapter 10.3) with an independent Metropolis step for $(\boldsymbol{\beta}, \boldsymbol{\delta})$, i.e an independent Metropolis-Hastings within Gibbs (IMHwG) algorithm. This algorithm consists of three steps: 1) sample $\boldsymbol{\Sigma}$ directly from its full conditional distribution, 2) sample ϕ directly from its full conditional distribution, and 3) sample $(\boldsymbol{\beta}, \boldsymbol{\delta})$ from its full conditional distribution using an independent Metropolis-Hastings step with the Laplace approximation as a proposal. Typically steps 1) and 2) are independent of each other, so the algorithm only consists of two Gibbs steps.

The naive way to implement this algorithm is to construct the Laplace approximation

to $[\boldsymbol{\beta}, \boldsymbol{\delta} | \boldsymbol{\Sigma}, \phi, \mathbf{Y}, \mathbf{X}]$ every iteration of the MCMC algorithm to use it as the IMH step proposal. This cost get multiplied by the number of MCMC iterations and can therefore be prohibitive when the optimization is costly. We call this approximation strategy the local conditional Laplace approximation. Alternatively, the joint Laplace approximation to the full joint posterior implies a conditional distribution for $\boldsymbol{\beta}, \boldsymbol{\delta} | \boldsymbol{\Sigma}, \phi$, which we call the global conditional Laplace approximation. Let $\boldsymbol{\theta}_1 = (\boldsymbol{\beta}, \boldsymbol{\delta})$ and $\boldsymbol{\theta}_2 = (\phi, \boldsymbol{\Sigma})$, or rather that $\boldsymbol{\theta}_2$ is a vector containing ϕ and the unique elements of $\boldsymbol{\Sigma}$. Then suppose we have the joint Laplace approximation to the posterior of $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)$ given by $\boldsymbol{\theta} | \mathbf{Y}, \mathbf{X} \stackrel{a}{\sim} N(\boldsymbol{\theta}^*, (-\mathbf{H}^*)^{-1})$ where $\boldsymbol{\theta}^* = (\boldsymbol{\theta}_1^*, \boldsymbol{\theta}_2^*)$ and

$$\boldsymbol{\Omega}^* = (-\mathbf{H}^*)^{-1} = \begin{bmatrix} \boldsymbol{\Omega}_{11}^* & \boldsymbol{\Omega}_{12}^* \\ \boldsymbol{\Omega}_{21}^* & \boldsymbol{\Omega}_{22}^* \end{bmatrix}.$$

Then the global conditional Laplace approximation to the conditional posterior of $\boldsymbol{\theta}_1$ is $\boldsymbol{\theta}_1 | \boldsymbol{\theta}_2, \mathbf{Y}, \mathbf{X} \stackrel{a}{\sim} N(\tilde{\boldsymbol{\theta}}_1, \tilde{\boldsymbol{\Omega}}_{11})$, where $\tilde{\boldsymbol{\theta}}_1 = \boldsymbol{\theta}_1^* + \boldsymbol{\Omega}_{12}^* (\boldsymbol{\Omega}_{22}^*)^{-1} (\boldsymbol{\theta}_2^{(t+1)} - \boldsymbol{\theta}_2^*)$ and $\tilde{\boldsymbol{\Omega}}_{11} = \boldsymbol{\Omega}_{11}^* - \boldsymbol{\Omega}_{12}^* (\boldsymbol{\Omega}_{22}^*)^{-1} \boldsymbol{\Omega}_{21}^*$.

The global approximation may be worse than the local approximation since the local approximation directly approximates the target conditional posterior, but the computation savings realized from only having to run the optimization algorithm once are typically significant and make it worthwhile. A probably with using either approximation is that by moving from an IMH algorithm to an IMHwG algorithm the mixing and convergence properties of the chain do deteriorate to the extent that $\boldsymbol{\theta}_1$ and $\boldsymbol{\theta}_2$ are dependent in the posterior. In practice this is often not an issue, or if it is it can typically be solved by reparameterizing the model, e.g. using a non-centered parameterization (Gelfand et al., 1995; Roberts and Sahu, 1997; Van Dyk and Meng, 2001; Bernardo et al., 2003).

4 Spatially Modeling County Population Estimates

The American Community Survey (ACS) provides 5-year period estimates of county populations as recently as 2014. In 2014 there were 3,142 counties in the United States, including the District of Columbia and counties in Alaska and Hawaii. We use two separate data models in order to illustrate when both the joint and global conditional Laplace approximations work well. The first is a Poisson data model, i.e., $Z_i \sim \text{Poisson}(\lambda_i)$ where $\lambda_i = \exp(Y_i)$. Through visual inspection of a histogram, it was determined that log county populations look approximately normally distributed. So our second data model is $\log Z_i \sim N(\mu_i, \phi^2)$, where $\mu_i = Y_i$.

The process model in both cases is a reduced rank spatial model $Y = \mathbf{X}\boldsymbol{\beta} + \mathbf{S}\boldsymbol{\delta}$, where $\mathbf{X}\boldsymbol{\beta}$ represents the process mean at each county and $\mathbf{S}\boldsymbol{\delta}$ implies the spatial correlation across counties. The spatial correlation term consists of a set of r basis functions evaluated at each of the $n = 3,142$ counties, denoted by the $n \times r$ matrix \mathbf{S} , and a common random effect $\boldsymbol{\delta}$. We assume that $\boldsymbol{\delta}$ is r -dimensional with $r \ll n$ so that the model is reduced rank. Any set of spatial basis functions could be used for \mathbf{S} but we use the Moran's I (MI) basis set, described below, but see Hughes and Haran (2013), Porter et al. (2015), Bradley et al. (2015) and references therein for additional discussion. Another possibility is to define a reduced rank model for a point-level spatial process using a basis function expansion and compute the implied set of basis functions for each of the census tracts by integrating the point level basis functions appropriately. See Sections 2.1, 3.1, and 4 of Bradley et al. (2016) for details.

The MI basis functions are defined through the orthogonal projection matrix $\mathbf{P}_\mathbf{X} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$. Let \mathbf{A} denote the binary adjacency matrix with $a_{ij} = 1$ if counties i and j are neighbors, $a_{ij} = 0$ otherwise, and $a_{ii} = 0$ along the diagonal, and define the MI operator \mathbf{G} as

$$\mathbf{G} = (\mathbf{I}_n - \mathbf{P}_\mathbf{X})\mathbf{A}(\mathbf{I}_n - \mathbf{P}_\mathbf{X})$$

where \mathbf{I}_n is the $n \times n$ identity matrix. The spectral decomposition of \mathbf{G} is $\mathbf{G} = \mathbf{\Phi}\mathbf{\Lambda}\mathbf{\Phi}'$. To use a reduced rank version of the MI basis functions we truncate the basis function expansion and take \mathbf{S} to be the $n \times r$ matrix formed by the r columns of $\mathbf{\Phi}$ corresponding to the r largest in magnitude eigenvalues of \mathbf{G} . The random effect $\boldsymbol{\delta}$ is further modeled as $\boldsymbol{\delta} \sim N(\mathbf{0}_r, \boldsymbol{\Sigma}(\boldsymbol{\theta}))$ where $\mathbf{0}_r$ denotes an r -dimensional vector of zeroes and $\boldsymbol{\Sigma}(\boldsymbol{\theta})$ is an unknown covariance matrix. The covariance matrix of $\mathbf{S}\boldsymbol{\delta}$ is then $\mathbf{S}\boldsymbol{\Sigma}(\boldsymbol{\theta})\mathbf{S}'$.

The process model depends on choices for \mathbf{X} and r . In practice r can be chosen using a sensitivity analysis. Since our goal is to illustrate computational methods, we will elide choosing r in a principled way and instead use several values for r in order to illustrate when the parameter space becomes too high dimensional for our method to be advantageous. For simplicity we choose an intercept only model, but all derivations for the model will assume that \mathbf{X} is $n \times p$.

Finally, we consider the two extreme parameterizations of $\boldsymbol{\Sigma}(\boldsymbol{\theta})$ from Section 3.2 — the iid parameterization where $\boldsymbol{\Sigma} = \sigma^2\mathbf{I}_r$ and the full parameterization where $\boldsymbol{\Sigma}$ is fully parameterized. In the iid parameterization we assume that $\sigma^2 \sim IG(a_\sigma, b_\sigma)$ in the prior, while in the full parameterization we assume that $\boldsymbol{\Sigma} \sim IW(d, \mathbf{E})$. The prior for $\boldsymbol{\beta}$ in all models is $\boldsymbol{\beta} \sim N(\mathbf{b}, v^2\mathbf{I}_p)$, and in the lognormal models the prior for ϕ^2 is $\phi^2 \sim IG(a_\phi, b_\phi)$. We assume that the parameters $\boldsymbol{\beta}$, $\boldsymbol{\Sigma}$, and when applicable ϕ^2 are mutually independent in the prior. For our examples we assume that $a_\sigma = a_\phi = b_\sigma = b_\phi = 1$, $\mathbf{b} = \mathbf{0}_p$, $v = 10$, $d = r + 1$ and $\mathbf{E} = \mathbf{I}_p$. Often a more complicated prior is appropriate on variance or covariance matrix parameters so that the marginal posterior is not sensitive to arbitrary choices in the prior. We use these conditionally conjugate priors because they allow for a fair comparison between MCMC algorithms — most alternatives will complicate alternative Gibbs samplers with extra steps or necessitate Metropolis steps making IMH or IMHwG relatively more attractive.

Between the two possible parameterizations of $\boldsymbol{\Sigma}$ and the two choices for the data model

we consider four possible classes of models. Both for implementing PSO algorithms to find the Laplace approximations and for running the IMH and IMHwG algorithms, we reparameterize the variance and covariance matrix parameters so that they have support on an unconstrained space. For σ^2 and ϕ^2 we use the log transformation, and for Σ we use the Cholesky factorization of Σ^{-1} allowing the diagonal elements to be negative. See Appendix C for details about the posterior distributions, including this transformation and relevant full conditionals.

5 Predicting the 1988 presidential election

Gelman and Hill (2006, Chapter 14) describes a model used to predict state-level opinions about the 1988 presidential candidates from national polls in order to predict the outcome of the election. They model the responses to a series of seven polls conducted by CBS News during the week before the 1988 presidential election. The variable of interest is binary: $Z_i = 1$ if the i th respondent said they supported the Republican candidate and $Z_i = 0$ if they said they supported the Democratic candidate, with undecideds being excluded. Focusing on the last poll, they ultimately estimate a logistic regression model with fixed effects for race (whether the respondent was African American or not), sex, and race \times sex, and random effects for four age categories, four education categories, and 16 age \times education categories, as well as for the respondent’s state of residence (including District of Columbia). The mean of the state random effect distribution is one of five region random effects plus the proportion of the state that voted republican in the last election times a slope coefficient.

The model is somewhat overparameterized since it has age, education, and age \times education random effects, so we reduce its size somewhat by omitting the age and education random effects, but keeping the age \times education interaction terms. In our model the age \times education random effects now represent the random effects for each age \times education category rather

than an interaction term. The single poll data model is

$$P(Z_i = 1) = \theta_i, \quad \theta_i = \exp(Y_i) / \{1 + \exp(Y_i)\},$$

$$Y_i = \beta_0 + f_i\beta_f + b_i\beta_b + f_ib_i\beta_{fb} + \alpha_{ae}[ae_i] + \alpha_s[s_i] \quad (\text{single poll data model}), \quad (6)$$

where f_i indicates whether respondent i identified as female, b_i indicates whether respondent i identified as African American, ae_i indicates respondent i 's age \times education category, and s_i indicates respondent i 's state of residence. Here we use $\alpha_s[k]$ to denote the k th element of the vector α_s , so α_{ae} contains 16 elements, and α_s contains 51 elements (50 states plus the District of Columbia). The single poll process model is

$$\alpha_s[k] \stackrel{ind}{\sim} N(\alpha_r[r_k] + prev_k\beta_{prev}, \sigma_s^2) \text{ for } k = 1, 2, \dots, 51,$$

$$\alpha_{ae}[k] \stackrel{iid}{\sim} N(0, \sigma_{ae}^2) \text{ for } k = 1, 2, \dots, 16,$$

$$\alpha_r[k] \stackrel{iid}{\sim} N(0, \sigma_r^2) \text{ for } k = 1, 2, \dots, 5 \quad (\text{single poll process model}), \quad (7)$$

where $\alpha_r[r_k]$ denotes the region containing state k and $prev_k$ denotes the average vote share for the Republicans in the previous three presidential elections. This model expands the class of models discussed in Section 3.1 by allowing the mean of δ to depend on random effects that are further modeled, but adding a level to the hierarchy does not fundamentally change the applicability of the Laplace approximations we use.

The last poll had 2,015 respondents, but together all seven polls have 11,566 respondents. Using each poll with a minimal number of additional parameters to account for poll to poll variability should increase the quality of the model and result in a posterior with better Laplace approximations since there is so much more data. We analyze a model for all of the polls using the following data model

$$P(Z_i = 1) = \theta_i, \quad \theta_i = \exp(Y_i) / \{1 + \exp(Y_i)\},$$

$$Y_i = \beta_0 + f_i\beta_f + b_i\beta_b + f_ib_i\beta_{fb} + \alpha_{ae}[ae_i] + \alpha_s[s_i] + \alpha_p[p_i] \quad (\text{all polls data model}), \quad (8)$$

where p_i denotes which poll respondent i was surveyed in. The process model is given by

$$\begin{aligned}
\alpha_s[k] &\stackrel{iid}{\sim} N(\alpha_r[r_k] + prev_k \beta_{prev}, \sigma_s^2) \text{ for } k = 1, 2, \dots, 51, \\
\alpha_{ae}[k] &\stackrel{iid}{\sim} N(0, \sigma_{ae}^2) \text{ for } k = 1, 2, \dots, 16, \\
\alpha_r[k] &\stackrel{iid}{\sim} N(0, \sigma_r^2) \text{ for } k = 1, 2, \dots, 5, \\
\alpha_p[k] &\stackrel{iid}{\sim} N(0, \sigma_p^2) \text{ for } k = 1, 2, \dots, 7 \quad (\text{all polls process model}).
\end{aligned} \tag{9}$$

In both models we assume each of the β s have independent $N(0, 1000)$ priors, and each of the σ^2 s have $IG(1, 1)$ priors. Including the random effects the single poll model contains 80 parameters while the all polls model contains 89 parameters. Writing down the log posteriors and deriving the Hessians is straightforward but tedious for these models, so we omit these steps, though note that in both the PSO and IMH algorithms the variances should be transformed to the log scale. [SHOULD WE ADD THIS INTO AN APPENDIX]

6 PSO results for finding posterior modes

[THIS STUFF ALL GETS PUT INTO THE LATTER TWO SECTIONS]

[I HAVE RESULTS USING 10 RANDOM EFFECTS iid MODELS AND 5 IN full; WORTH FINDING SPACE FOR?]

We conduct a simulation study using R (R Development Core Team, 2008) to compare various PSO algorithms at finding the posterior mode in each of the example models from Sections 4 and 5. Based on the results of Appendix B, we limit the study to 7 PSO algorithms: standard PSO algorithm using parameter values suggested by Blum and Li (2008) and Clerc and Kennedy (2002) (PSO in Figures 1, 2, and [REFERENCE TO POLL BOX PLOT FIGURE]), the standard BBPSOxp-MC algorithm (BBPSO), DI-PSO with $\alpha = 0.2 \times n_{iter} = 200$ and $\beta = 1$ (DI-PSO), AT-PSO with $c = 0.1$ and either $R^* = 0.3$ or 0.5 (AT-PSO-0.3 and AT-PSO-0.5), and AT-BBPSOxp-MC with $df = 1$, $c = 0.1$, and either $R^* = 0.3$ or 0.5 (AT-

BBPSO-0.3 and AT-BBPSO-0.5). Each algorithm was tried with one of two initializations. In the “BFGS” initialization, we first ran the BFGS algorithm using R’s `optim` function (R Development Core Team, 2008) until convergence using default settings to obtain an initial guess of the argmax, $\hat{\theta}$, then initialized the swarm with one particle at this initial guess and the rest uniformly in a length 2 hypercube centered on $\hat{\theta}$, i.e. $\theta_1 = \hat{\theta}$ and $\theta_{ij} = \hat{\theta}_j + U(-1, 1)$ for $i = 2, 3, \dots, 50$ and $j = 1, 2, \dots, n_{par}$ where n_{par} is the number of parameters in the model and the $U(-1, 1)$ random variates are drawn independently. In the “no BFGS” initialization, each particle was initialized uniformly in a length 200 hypercube centered at zero, i.e. $\theta_{ij} \stackrel{iid}{\sim} U(-100, 100)$ for $i = 1, 2, \dots, 50$ and $j = 1, 2, \dots, n_{par}$. In addition, each algorithm was run using both the ring-1 and ring-3 neighborhood topologies, for a total of four combinations of initializations and neighborhoods, each for 20 replications of 1,000 iterations using 50 particles.

Figure 1 contains boxplots of the results of the simulations for the Poisson county population models of Section 4. Models with iid random effects had 30 random effects while models with fully correlated random effects had 15. Each box plot was created using the 10, 25, 50, 75, and 90 percentiles of the maximum value of the log posterior found in all 20 replications after 1,000 iterations. Figure 2 is similar, except for the lognormal models. Across all models we see that the ring-3 topology and BFGS initialization both improve each of the algorithms, especially in combination. In Appendix B the ring-3 neighborhood performed the best on our suite of test functions, and the same seems to be true here. The BFGS initialization is cheap, taking essentially no time to compute, yet seems to drastically improve the quality of every PSO algorithms. In Appendix B we speculated that the AT-PSO algorithms would benefit significantly from some sort of stage one optimization, and our results confirm this for all of the algorithms. Another lesson from these boxplots is that the AT-PSO and AT-BBPSO algorithms tend to do the best, especially when $R^* = 0.5$ in both cases.

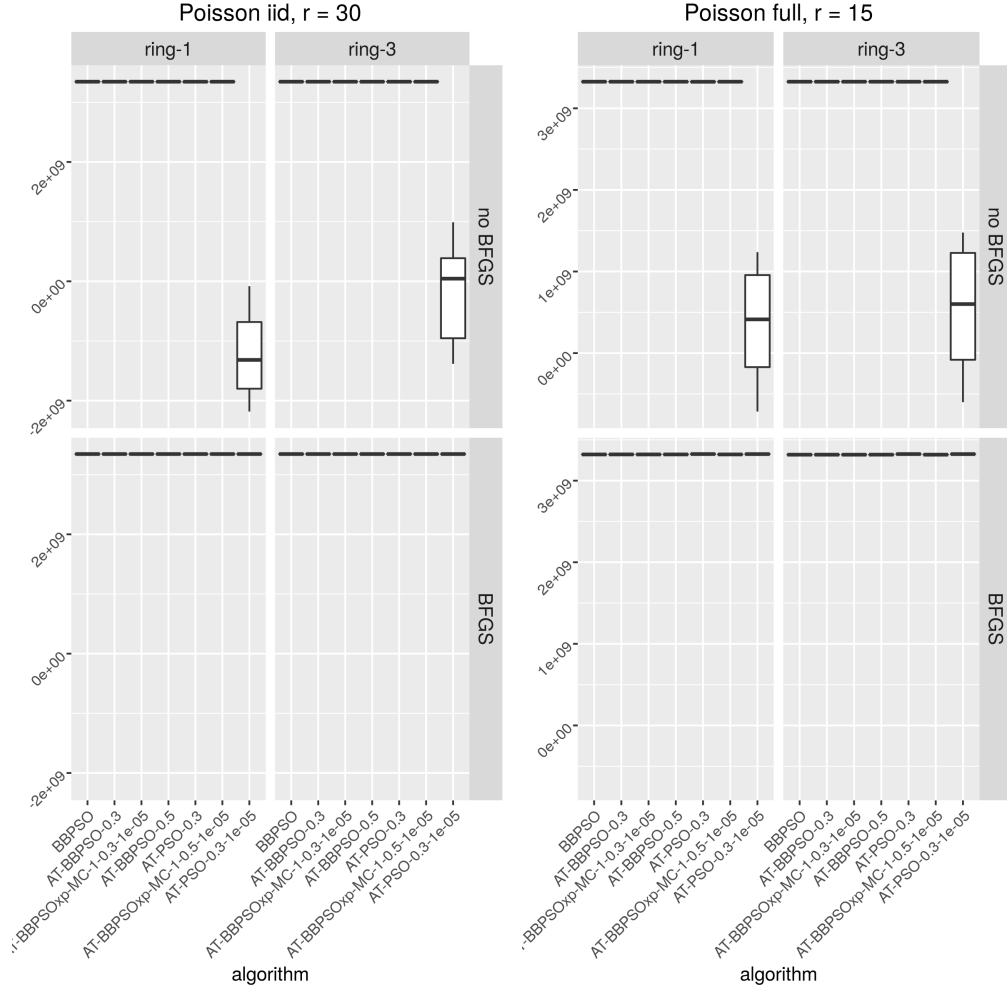


Figure 1: Boxplots of the maximum value of the log posterior found for Poisson models by random effect type, neighborhood topology, optimization initialization, and algorithm. Each box plot was created using the 10, 25, 50, 75, and 90 percentiles of the maximum found from 20 replications of 1,000 iterations with 50 particles for each factor combination.

[PARAGRAPH OR TWO AND PLOTS FOR THE ELECTION MODELS]

For our purposes, the PSO algorithms are useful only insofar as they allow us to construct IMH and IMHwG algorithms. So we conduct another simulation study to see when each algorithm gets close enough to the posterior mode for the IMH and IMHwG algorithms to have high acceptance rates. To this end we conduct another simulations study using the

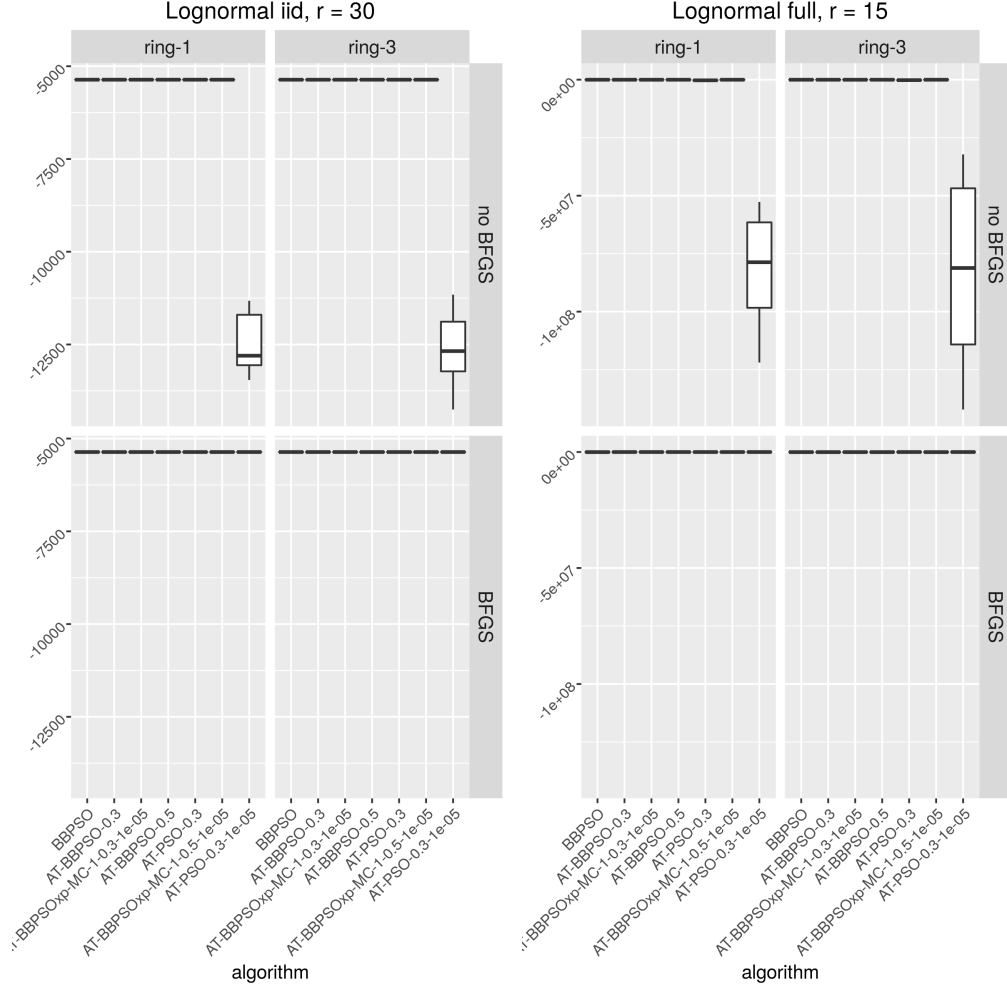


Figure 2: Boxplots of the maximum value of the log posterior found for lognormal models by random effect type, number of random effects, neighborhood topology, optimization initialization, and algorithm. Each box plot was created using the 10, 25, 50, 75, and 90 percentiles of the maximum found from 20 replications of 1,000 iterations with 50 particles for each factor combination.

same 7 PSO algorithms from the previous study, except only using the ring-3 neighborhood topology and BFGS initialization. Next, we run the IMH and IMHwG MCMC algorithms based on 0, 100, 500, 1,000, 1,500, and 2,000 iterations of the PSO algorithm and compute the acceptance rate after 1,000 iterations of the MCMC algorithm. Each MCMC algorithm

is initialized at the PSO estimate of the posterior mode used to construct the Laplace approximation. When we run the IMH and IMHwG algorithms after 0 iterations of a PSO algorithm we still run the BFGS initialization, so this serves as a control to see if running the PSO algorithm is necessary to get reasonable acceptance rates. Additionally, we run each IMH and IMHwG algorithm with different choices for the degrees of freedom parameter in the Laplace approximation.

In single poll model of Section 5, we draw $(\beta_0, \beta_f, \beta_b, \boldsymbol{\alpha}_s, \boldsymbol{\alpha}_{ae})$ in the Metropolis step, while in the all polls model we draw all of those parameters and additionally $\boldsymbol{\alpha}_p$ in the Metropolis step. Then for both models there are two additional Gibbs steps — one where each random effect variance is drawn, and one where $(\beta_{prev}, \boldsymbol{\alpha}_r)$ is drawn. These full conditionals are straightforward to derive, so we do not reproduce them here. Likewise, the Hessian is easy though tedious to derive, so we omit it as well.

[SIMULATION STUDY TO BE COMPLETED AND PUT HERE - WON'T TAKE VERY LONG, 1 DAY]

7 Comparison of MCMC techniques

Next we move to comparing various MCMC algorithms for both classes of models. The other MCMC algorithms we consider in this section are single move random walk Metropolis within Gibbs (RWwG), block random walk Metropolis within Gibbs (B-RWwG), Hamiltonian Monte Carlo as implemented by the Stan software (Carpenter et al., 2015) using the No U-turn Sampler (Homan and Gelman, 2014) and for the lognormal model for county populations, a two step Gibbs sampler (Gibbs). The last algorithm is only considered for the lognormal model because only in that case are all of the full conditionals tractable. We use standard settings to fit each of the models in Stan, though minimal tweaking was required in some cases. In all of the Gibbs algorithms, the blocks are the same as in the IMHwG

algorithms detailed in the previous section. Both the RWwG and B-RWwG algorithms are tuned during the burn-in using an adaptive method. Appendix D explains the details of the two random walk algorithms we use including how the adaptation was performed.

We compare these algorithms in Table [TO BE CONSTRUCTED] in terms of two measures: the minimum estimated effective sample size for all parameters in the model (n_{eff} , see Robert and Casella (2013, Section 12.3.5)), and time in seconds per n_{eff} . The effective sample size is the size of an iid sample which yields the same standard error for estimating the mean of some function of the parameters in our MCMC simulations, so we take the minimum estimated n_{eff} among all elements of the parameter vector and latent process.

[SIMULATIONS TO BE CONDUCTED - SHOULD TAKE A COUPLE DAYS ONCE THEY'RE RUNNING. FOR THE IMH AND IMHwG ALGORITHMS, THESE WILL ONLY USE ONE PSO ALGORITHM]

8 Discussion

[WILL NEED TO BE REWRITTEN IN LIGHT OF FORTHCOMING RESULTS]

In practice PSO-assisted Metropolis-Hastings algorithms are useful to the extent that the problem falls in something of a “sweet spot”. The model must be complex enough that the posterior mode cannot be found analytically and that a fully conjugate Gibbs sampler is not available. The combined parameter space and latent space must be large enough that traditional numerical methods fail to find the posterior mode, but small enough that the heuristic PSO algorithms have some chance without being too costly. Finally, both model parameters and latent random variables must be approximately Gaussian in the posterior, though there is some leeway for model parameters using the IMHwG algorithm. Under these conditions, using PSO to obtain the Laplace approximation for use as a proposal for IMH or IMH within Gibbs improves on other more commonly used approaches for MCMC. The

comparison to Stan is also relevant — even though our PSO assisted IMH algorithms beat Stan in terms of total computational time for a desired level of precision for some of the models we considered, the difference will be small enough in many settings for most users to stick with the more general algorithm. [SOMETHING ABOUT THE ADVANTAGE BEING IN BIG / TALL DATA SITUATIONS - STAN REQUIRES A LOT OF LOG POSTERIOR EVALUATIONS]

Our approach is fairly hands off, so long as the conditions listed above appear to be satisfied. In that case the standard PSO algorithm with default parameter values and the ring-1 neighborhood usually does a great job of finding the posterior mode quickly relative to the most competitive alternatives. Alternatives such as the adaptively tuned BBPSO variants we introduced only seem to do better when neither algorithm has converged on the mode yet or when the Laplace approximation is poor. So in practice we recommend using standard PSO first.

The AT-BBPSOxp-MC algorithms we introduced still require a larger swarm size or more iterations or both to do as well as standard PSO when standard PSO works well. With the same swarm size and number of iterations AT-BBPSOxp-MC does only slightly worse than standard PSO in terms of maximizing the objective function, but that slight difference often amounts to a large difference in the Metropolis acceptance rate for the resulting MCMC algorithms. When AT-BBPSOxp-MC does do better than standard PSO, the resulting MCMC algorithms do not always have high acceptance rates, at least not without running AT-BBPSOxp-MC with a high swarm size for many iterations. So in practice we recommend using standard PSO first. If the resulting MCMC algorithms have poor acceptance rates or if it seems likely a priori that PSO might have trouble, then it can be useful to explore the options provided by AT-BBPSOxp-MC with the caveat that a large swarm size and number of iterations may be required.

References

- Andrieu, C. and Thoms, J. (2008). “A tutorial on adaptive MCMC.” *Statistics and Computing*, 18, 4, 343–373.
- Berliner, L. M. (1996). “Hierarchical Bayesian time series models.” In *Maximum entropy and Bayesian methods*, 15–22. Springer.
- Bernardo, J., Bayarri, M., Berger, J., Dawid, A., Heckerman, D., Smith, A., and West, M. (2003). “Non-centered Parameterisations for Hierarchical Models and Data Augmentation.” In *Bayesian Statistics 7: Proceedings of the Seventh Valencia International Meeting*, 307–326. London: Oxford University Press.
- Blum, C. and Li, X. (2008). “Swarm Intelligence in Optimization.” In *Swarm Intelligence: Introduction and Applications*, eds. C. Blum and D. Merkle. Springer.
- Bradley, J. R., Holan, S. H., and Wikle, C. K. (2015). “Multivariate spatio-temporal models for high-dimensional areal data with application to longitudinal employer-household dynamics.” *Annals of Applied Statistics*, 9, 4, 1761–1791.
- Bradley, J. R., Wikle, C. K., and Holan, S. H. (2016). “Regionalization of multiscale spatial processes by using a criterion for spatial aggregation error.” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*.
- Campos, M., Krohling, R. A., and Enriquez, I. (2014). “Bare bones particle swarm optimization with scale matrix adaptation.” *Cybernetics, IEEE Transactions on*, 44, 9, 1567–1578.
- Carpenter, B., Gelman, A., Hoffman, M., Lee, D., Goodrich, B., Betancourt, M., Brubaker, M. A., Guo, J., Li, P., and Riddell, A. (2015). “Stan: a probabilistic programming language.” *Journal of Statistical Software*.

- Clerc, M. (2010). *Particle swarm optimization*. John Wiley & Sons.
- Clerc, M. and Kennedy, J. (2002). “The particle swarm-explosion, stability, and convergence in a multidimensional complex space.” *Evolutionary Computation, IEEE Transactions on*, 6, 1, 58–73.
- Eberhart, R. C. and Shi, Y. (2000). “Comparing inertia weights and constriction factors in particle swarm optimization.” In *Evolutionary Computation, 2000. Proceedings of the 2000 Congress on*, vol. 1, 84–88. IEEE.
- Gelfand, A. E., Sahu, S. K., and Carlin, B. P. (1995). “Efficient Parametrisations for Normal Linear Mixed Models.” *Biometrika*, 82, 3, 479–488.
- Gelman, A. and Hill, J. (2006). *Data analysis using regression and multilevel/hierarchical models*. Cambridge University Press.
- Gelman, A., Roberts, G., and Gilks, W. (1996). “Efficient Metropolis jumping rules.” *Bayesian statistics*, 5, 599–608, 42.
- Goldberg, D. E. and Holland, J. H. (1988). “Genetic algorithms and machine learning.” *Machine learning*, 3, 2, 95–99.
- Hastings, W. K. (1970). “Monte Carlo sampling methods using Markov chains and their applications.” *Biometrika*, 57, 1, 97–109.
- Homan, M. D. and Gelman, A. (2014). “The no-U-turn sampler: Adaptively setting path lengths in Hamiltonian Monte Carlo.” *The Journal of Machine Learning Research*, 15, 1, 1593–1623.
- Hsieh, H.-I. and Lee, T.-S. (2010). “A modified algorithm of bare bones particle swarm optimization.” *International Journal of Computer Science Issues*, 7, 11.

- Hughes, J. and Haran, M. (2013). “Dimension reduction and alleviation of confounding for spatial generalized linear mixed models.” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 75, 1, 139–159.
- Kennedy, J. (2003). “Bare bones particle swarms.” In *Swarm Intelligence Symposium, 2003. SIS’03. Proceedings of the 2003 IEEE*, 80–87. IEEE.
- Krohling, R., Mendel, E., et al. (2009). “Bare bones particle swarm optimization with Gaussian or Cauchy jumps.” In *Evolutionary Computation, 2009. CEC’09. IEEE Congress on*, 3285–3291. IEEE.
- Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H., and Teller, E. (1953). “Equation of state calculations by fast computing machines.” *The journal of chemical physics*, 21, 6, 1087–1092.
- Neal, R. M. et al. (2011). “MCMC using Hamiltonian dynamics.” *Handbook of Markov Chain Monte Carlo*, 2, 113–162.
- Porter, A. T., Holan, S. H., and Wikle, C. K. (2015). “Bayesian semiparametric hierarchical empirical likelihood spatial models.” *Journal of Statistical Planning and Inference*, 165, 78–90.
- R Development Core Team (2008). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0.
- Richer, T. J. and Blackwell, T. M. (2006). “The Lévy particle swarm.” In *Evolutionary Computation, 2006. CEC 2006. IEEE Congress on*, 808–815. IEEE.
- Robert, C. and Casella, G. (2013). *Monte Carlo statistical methods*. 2nd ed. Springer Science & Business Media.

- Roberts, G. O. and Sahu, S. K. (1997). “Updating Schemes, Correlation Structure, Blocking and Parameterization for the Gibbs Sampler.” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 59, 2, 291–317.
- Rue, H., Martino, S., and Chopin, N. (2009). “Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximations.” *Journal of the royal statistical society: Series B (statistical methodology)*, 71, 2, 319–392.
- Schervish, M. J. (1997). *Theory of statistics*. Springer Science & Business Media.
- Stroup, W. W. (2012). *Generalized linear mixed models: modern concepts, methods and applications*. CRC press.
- Taylor, B. M. and Diggle, P. J. (2014). “INLA or MCMC? A tutorial and comparative evaluation for spatial prediction in log-Gaussian Cox processes.” *Journal of Statistical Computation and Simulation*, 84, 10, 2266–2284.
- Tuppadung, Y. and Kurutach, W. (2011). “Comparing nonlinear inertia weights and constriction factors in particle swarm optimization.” *International Journal of Knowledge-based and Intelligent Engineering Systems*, 15, 2, 65–70.
- Van Dyk, D. and Meng, X.-L. (2001). “The Art of Data Augmentation.” *Journal of Computational and Graphical Statistics*, 10, 1, 1–50.
- Wikle, C. K. et al. (2003). “Hierarchical Models in Environmental Science.” *International Statistical Review*, 71, 2, 181–199.

Particle Swarm Optimization Assisted Metropolis Hastings Algorithms

A PSO and BBPSO details

[THIS IS JUST PASTED TEXT THAT CAME FROM THE MAIN BODY OF THE DOCUMENT - SECTION NEEDS TO BE WRITTEN IN EARNEST. INCLUDE RING TOPOLOGY STUFF HERE.]

A downside of both versions of BBPSO above is that any particle currently at its group best location does not move due to the definition of the standard deviation term. Several methods have been proposed to overcome this; e.g., see Hsieh and Lee (2010) and Zhang et al. (2011). Zhang et al. (2011) propose using mutation and crossover operations for the group best particle. To do this, the group best particle randomly selects three other distinct particles i_1 , i_2 , and i_3 , and updates:

$$\theta_{ij}(t+1) = p_{i_1j}(t) + 0.5(p_{i_2j}(t) - p_{i_3j}(t)). \quad (1)$$

When combined with BBPSOxp, this becomes

$$\theta_{ij}(t+1) = \begin{cases} p_{i_1j}(t) + 0.5(p_{i_2j}(t) - p_{i_3j}(t)) & \text{with probability 0.5} \\ g_{ij}(t) & \text{otherwise,} \end{cases} \quad (2)$$

but it also introduces another problem — a particle which moves to its group best location on a single coordinate will have a standard deviation of zero for that coordinate during the next iteration. This is easily overcome by manually setting the standard deviation to a small value (e.g., 0.001), which allows us to define the two base BBPSO algorithms that we will consider: BBPSO-MC and BBPSOxp-MC. BBPSO-MC is standard BBPSO where each particle evolves according to (2) except any current group best particle evolves according to (1). BBPSOxp-MC evolves every particle according to

$$\theta_{ij}(t+1) = \begin{cases} N\left(\frac{p_{ij}(t) + g_{ij}(t)}{2}, \sigma_{ij}^2(t)\right) & \text{with probability 0.5} \\ g_{ij}(t) & \text{otherwise,} \end{cases} \quad (3)$$

where $\sigma_{ij}(t) = |p_{ij}(t) - g_{ij}(t)|$ if $|p_{ij}(t) - g_{ij}(t)| > 0$ and $\sigma_{ij}(t) = 0.001$ otherwise, except current group best particles evolve according to (2).

B Comparing AT-BBPSO, BBPSO, and PSO algorithms

In order to compare AT-BBPSO to other PSO variants, we employ a subset of test functions used in Hsieh and Lee (2010). Each function is listed in Table 1 along with the global maximum and argmax, and the initialization range for the simulations. Further description of many of these functions can be found in Clerc (2010). For each function, we set $D = 20$ so the domain of each function is \Re^{20} . For each function, the PSO algorithms are initialized in a range that does not contain the true maximum.

Equation	ArgMax	Maximum	Initialization
$Q_1(\boldsymbol{\theta}) = -\sum_{i=1}^D \theta_i^2$	$\boldsymbol{\theta}^* = \mathbf{0}$	$Q_1(\boldsymbol{\theta}^*) = 0$	$(50, 100)^D$
$Q_2(\boldsymbol{\theta}) = -\sum_{i=1}^D \left(\sum_{j=1}^i \theta_j\right)^2$	$\boldsymbol{\theta}^* = \mathbf{0}$	$Q_2(\boldsymbol{\theta}^*) = 0$	$(50, 100)^D$
$Q_3(\boldsymbol{\theta}) = -\sum_{i=1}^{D-1} [100\{\theta_{i+1} + 1 - (\theta_i + 1)^2\} + \theta_i^2]$	$\boldsymbol{\theta}^* = \mathbf{0}$	$Q_3(\boldsymbol{\theta}^*) = 0$	$(15, 30)^D$
$Q_4(\boldsymbol{\theta}) = 9D - \sum_{i=1}^D \{\theta_i^2 - \cos(2\pi\theta_i) + 10\}$	$\boldsymbol{\theta}^* = \mathbf{0}$	$Q_4(\boldsymbol{\theta}^*) = 0$	$(2.56, 5.12)^D$
$Q_5(\boldsymbol{\theta}) = -\frac{1}{4000}\ \boldsymbol{\theta}\ ^2 + \prod_{i=1}^D \cos\left(\frac{\theta_i}{\sqrt{i}}\right) - 1$	$\boldsymbol{\theta}^* = \mathbf{0}$	$Q_5(\boldsymbol{\theta}^*) = 0$	$(300, 600)^D$
$Q_6(\boldsymbol{\theta}) = 20 \exp\left(-0.2\sqrt{\frac{1}{D}\ \boldsymbol{\theta}\ }\right) + \exp\left\{\frac{1}{D}\sum_{i=1}^D \cos(2\pi\theta_i)\right\} - 20 - \exp(1)$	$\boldsymbol{\theta}^* = \mathbf{0}$	$Q_6(\boldsymbol{\theta}^*) = 0$	$(16, 32)^D$

Table 1: Test functions for evaluating PSO algorithms. The dimension of $\boldsymbol{\theta}$ is D and $\|\cdot\|$ is the Euclidean norm: $\|\boldsymbol{\theta}\| = \sqrt{\sum_{i=1}^D \theta_i^2}$.

We use several PSO algorithms in the simulation study. The standard PSO algorithm uses the parameter values suggested by Blum and Li (2008) and Clerc and Kennedy (2002). The AT-BBPSO variants are implemented a wide variety of parameter values, but all have

the scale parameter initialized at $\sigma(0) = 1$, and both set $c = 0.1$. The AT-PSO variants are initialized at $\omega(0) = 1$ and also set $c = 0.1$. In addition, each algorithm is implemented using each of three neighborhood structures — the global, ring-3, and ring-1 neighborhoods. Each algorithm was used to optimize each objective function for 500 iterations over 50 replications using 20 particles. Initializations were changed across replications but held constant across algorithms. The standard PSO, DI-PSO, and AT-PSO algorithms initialized their velocity terms using the same method as a function of the initial locations of the particles, which we will denote by \mathbf{x}_i for particle i . Let $x_{max,j}$ be the maximum initial value of coordinate j of \mathbf{x}_i for each particle i , and let $x_{min,j}$ be the corresponding minimum. Then let $d_{max} = \max_j x_{max,j} - x_{min,j}$. Then we initialize the velocities with $v_{ij}(0) \stackrel{iid}{\sim} U(-d_{max}/2, d_{max}/2)$. Tables 2-7 contain the simulation results for objective functions 1-6 respectively (OF1, OF2, etc.). We use several measures to quantify how well each algorithm finds the global maximum. First, each table includes the mean and standard deviation of the absolute difference between the true global maximum and the algorithm's estimated global maximum across all 50 replications, denoted by Mean and SD. Second, each table includes a convergence criterion — the proportion of the replications that came within 0.01 of the true global maximum, denoted by \hat{p} . Finally, \hat{t} denotes the median number of iterations until the algorithm reaches the convergence criterion. When $\hat{p} < 0.5$ then $\hat{t} = \infty$ since greater than 50% of the replications did not converge in the maximum number of iterations allowed. Mean, \hat{p} , and \hat{t} can be thought of how close the algorithm gets to the global maximum, what proportion of the time it converges, and long it takes to converge respectively.

We highlight only some of the features of these tables. First and foremost, PSO, BBPSO-MC, and BBPSOxp-MC almost always do worse than their AT cousins. Adaptively tuning either the scale or inertia parameter leads to gains in all three of our measures, sometimes large. Second, for most algorithms the more restrictive neighborhood appears to be result

in algorithms which do a better job of finding the global max. This is not true for the AT algorithms. For the AT-PSO algorithms, it appears that the target improvement rate (R^*) and the neighborhood interact. When the rate is high a more restrictive neighborhood is preferable, while when the rate is low a less restrictive neighborhood is preferable. On the other hand, for the AT-BBPSO algorithms, the opposite appears to be true — when the target improvement rate is high, a less restrictive neighborhood is desirable. In addition, the best AT-BBPSO algorithms often use the global neighborhood.

For the DI-PSO algorithms often there is a parameter-neighborhood combination that does well, typically from setting $\alpha = 200$ (20% of the 500 iterations) and $\beta = 1$ and using either the ring-1 or ring-3 neighborhood. One of these combinations typically does the best of all the DI-PSO algorithms but they can sometimes still do much worse than the best alternatives (e.g., for OF2). In the AT-BBPSO algorithms, it is not always clear what the best parameter settings are, but good default values appear to be $df = 3$ or 5 and $R^* = 0.3$ or 0.5 . When these values are good, they often lead to the best performing PSO algorithms we consider. However, it is not always clear whether to use AT-BBPSO-MC or AT-BBPSOxp-MC. For some objective functions, e.g. OF4, the xp version is consistently better than the non-xp alternative, but the opposite is true for others, e.g. OF2.. AT-PSO is also often very competitive with $R^* = 0.3$ or $R^* = 0.5$, though again sometimes different parameter settings also appear to work well.

The DI-PSO and AT-PSO algorithms are similar conceptually, but often yield very different results. DI-PSO deterministically reduces the inertia parameter over time in the same manner given a fixed set of parameter values (α and β), while AT-PSO dynamically adjusts the inertia parameter to hit a target improvement rate. Figure 1 plots the inertia over time for the DI-PSO algorithm with $\alpha = 200$, i.e. 20% of the total number of PSO iterations, and $\beta = 1$, and observed inertia over time for one replication of the *AT – PSO* algorithm with target rate $R^* = 0.5$ and ring-1 neighborhood for OF1 and one replication for OF6.

All three algorithms have an initial inertia of $\omega(0) = 1$. While DI-PSO smoothly decreases its inertia with a slowly decreasing rate, for OF1 AT-PSO very quickly drops its inertia to about 0.5 then bounces around around near that point. It also jumps up above 1 initially, imploring the particles to cast a wider net in search of higher value areas of the search space. This is pretty typical behavior for the inertia parameter of AT-PSO — it tends to bounce around a level which is approximately the average over time of the DI-PSO’s inertia, though lower values of R^* will result in higher inertias. In this way, AT-PSO alternates periods of exploration (relatively high inertia) and periods of exploitation (relatively low inertia). The main exception to this pattern is when AT-PSO converges around a local maximum. In this case, inertia plummets to zero as the particles settle down. This is precisely what happens for OF6 in Figure 1, though in this case the maximum is not global — Table 7 indicates that ring-1 AT-PSO with $R^* = 0.5$ never converged to the global max. In optimization problems with multiple local optima, both AT-PSO and AT-BBPSO variants can exhibit this behavior prematurely converge to a local optima, so they may not be advantageous for those problems.

Based on these simulations, our default recommendation is to use AT-BBPSO-MC or AT-BBPSOxp-MC with $R^* = 0.3$ or 0.5 and $df = 3$ or $df = 5$ along with the global neighborhood. These algorithms will not always be the best of the PSO algorithms, but they will often be very good. AT-PSO with $R^* = 0.3$ or $R^* = 0.5$ with a restrictive neighborhood topology such as ring-3 also tends to be a very good choice, though perhaps less consistent than the AT-BBPSO variants. Default PSO also performs rather well and is a good baseline algorithm to use for comparisons.

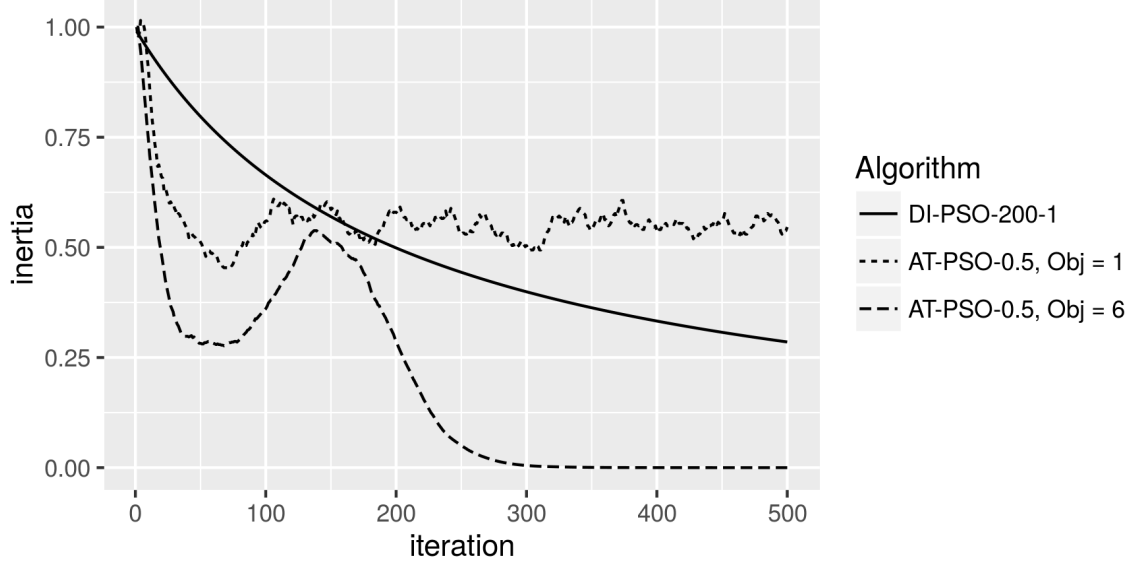


Figure 1: Inertia over time for the DI-PSO algorithm with $\alpha = 200$ and $\beta = 1$, and for one replication of the AT-PSO-0.5 algorithm for each of OFs 1 and 6.

C County Population Model Details

We consider two possible data models, Poisson and lognormal, and two possible random effect distributions, iid normal and normal with a fully correlated covariance matrix. The posterior distributions for the models with iid random effects can be written as

$$\begin{aligned}
p(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\delta}, \phi^2 | \mathbf{z}, \mathbf{X}, \mathbf{S}) &\propto (\phi^2)^{-n/2 - a_\phi - 1} \exp \left[-\frac{1}{\phi^2} \left\{ \frac{(\log \mathbf{z} - \mathbf{y})'(\log \mathbf{z} - \mathbf{y})}{2} + b_\phi \right\} \right] \\
&\times (\sigma^2)^{-\frac{r}{2} - a_\sigma - 1} \exp \left\{ -\frac{1}{\sigma^2} \left(\frac{\boldsymbol{\delta}'\boldsymbol{\delta}}{2} + b_\sigma \right) \right\} \exp \left\{ -\frac{(\boldsymbol{\beta} - \mathbf{b})'(\boldsymbol{\beta} - \mathbf{b})}{2v} \right\} \quad (\text{iid lognormal}), \\
p(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\delta} | \mathbf{z}, \mathbf{X}, \mathbf{S}) &\propto \prod_{i=1}^n \frac{\exp\{-\exp(y_i)\} \exp(y_i z_i)}{z_i!} \\
&\times (\sigma^2)^{-\frac{r}{2} - a_\sigma - 1} \exp \left\{ -\frac{1}{\sigma^2} \left(\frac{\boldsymbol{\delta}'\boldsymbol{\delta}}{2} + b_\sigma \right) \right\} \exp \left\{ -\frac{(\boldsymbol{\beta} - \mathbf{b})'(\boldsymbol{\beta} - \mathbf{b})}{2v} \right\} \quad (\text{iid Poisson}).
\end{aligned}$$

For fully parameterized Σ we write the posteriors in terms of the precision matrix $\Omega = \Sigma^{-1}$, yielding

$$\begin{aligned}
p(\beta, \Omega, \delta, \phi^2 | \mathbf{z}, \mathbf{X}, \mathbf{S}) &\propto (\phi^2)^{-n/2 - a_\phi - 1} \exp \left\{ -\frac{1}{\phi^2} \left(\frac{(\log \mathbf{z} - \mathbf{y})'(\log \mathbf{z} - \mathbf{y})}{2} + b_\phi \right) \right\} \\
&\times |\Omega|^{(d-r)/2} \exp \left\{ -\frac{1}{2} \left(\delta' \Omega \delta + \frac{(\beta - \mathbf{b})'(\beta - \mathbf{b})}{v} + \text{tr}(\mathbf{E} \Omega) \right) \right\} \quad (\text{full lognormal}), \\
p(\beta, \Omega, \delta | \mathbf{z}, \mathbf{X}, \mathbf{S}) &\propto \prod_{i=1}^n \frac{\exp\{-\exp(y_i)\} \exp(y_i z_i)}{z_i!} \\
&\times |\Omega|^{(d-r)/2} \exp \left\{ -\frac{1}{2} \left(\delta' \Omega \delta + \frac{(\beta - \mathbf{b})'(\beta - \mathbf{b})}{v} + \text{tr}(\mathbf{E} \Omega) \right) \right\} \quad (\text{full Poisson})
\end{aligned}$$

where $\text{tr}(\cdot)$ denotes the trace operator. For the PSO algorithm and some MCMC algorithms it will be easier to work with the Cholesky decomposition of Ω given by $\mathbf{L}\mathbf{L}' = \Omega$ where \mathbf{L} is lower triangular. It is often more convenient to put the prior distribution directly on \mathbf{L} rather than on Ω or Ω^{-1} and solving for the Jacobian, but this depends in part on which MCMC algorithm is used to fit the model. So while we use the prior distribution on \mathbf{L} implied by a Wishart prior on Ω , in practice it is advantageous to use one of the priors suggested by Chen and Dunson (2003) or Frühwirth-Schnatter and Tüchler (2008). We allow the diagonal entries of \mathbf{L} to be negative in the PSO and IMH algorithms, so \mathbf{L} is not strictly speaking a Cholesky decomposition. The determinant of the Jacobian is the same in both cases up to a proportionality constant. The signs of the elements of \mathbf{L} are not identified, therefore care needs to be taken when interpreting the results of MCMC, but transforming back to the precision matrix in a post processing step is sufficient. Let ℓ_{ij} denote the (i, j) th element of \mathbf{L} . Then the Jacobian of $\Omega \rightarrow \mathbf{L}$ is given by

$$|J(\Omega \rightarrow \mathbf{L})| \propto \prod_{k=1}^r |\ell_{kk}|^{r+1-k}$$

where $\mathbf{\Omega}$ is $r \times r$. Under this parameterization the full posteriors can be written as

$$\begin{aligned}
p(\boldsymbol{\beta}, \mathbf{L}, \boldsymbol{\delta}, \phi^2 | \mathbf{z}, \mathbf{X}, \mathbf{S}) &\propto (\phi^2)^{-n/2 - a_\phi - 1} \exp \left[-\frac{1}{\phi^2} \left\{ \frac{(\log \mathbf{z} - \mathbf{y})'(\log \mathbf{z} - \mathbf{y})}{2} + b_\phi \right\} \right] \\
&\times \prod_{k=1}^r (\ell_{kk}^2)^{(d-k+1)/2} \\
&\times \exp \left[-\frac{1}{2} \left\{ \boldsymbol{\delta}' \mathbf{L} \mathbf{L}' \boldsymbol{\delta} + \frac{(\boldsymbol{\beta} - \mathbf{b})'(\boldsymbol{\beta} - \mathbf{b})}{v} + \text{tr}(\mathbf{E} \mathbf{L} \mathbf{L}') \right\} \right] \quad (\text{full lognormal}), \\
p(\boldsymbol{\beta}, \mathbf{L}, \boldsymbol{\delta} | \mathbf{z}, \mathbf{X}, \mathbf{S}) &\propto \prod_{i=1}^n \frac{\exp\{-\exp(y_i)\} \exp(y_i z_i)}{z_i!} \times \prod_{k=1}^r (\ell_{kk}^2)^{(d-k+1)/2} \\
&\times \exp \left[-\frac{1}{2} \left\{ \boldsymbol{\delta}' \mathbf{L} \mathbf{L}' \boldsymbol{\delta} + \frac{(\boldsymbol{\beta} - \mathbf{b})'(\boldsymbol{\beta} - \mathbf{b})}{v} + \text{tr}(\mathbf{E} \mathbf{L} \mathbf{L}') \right\} \right] \quad (\text{full Poisson}).
\end{aligned}$$

In Appendix C.2 we derive the Hessian for the fully parameterized Poisson model. The other models are analogous, though the variances in the iid models and in the lognormal models should be transformed to the log scale first.

C.1 Full conditionals for MCMC

We consider alternative several Gibbs sampling algorithms, detailed in Appendix D, which draw the covariance matrix from its full conditional distribution. When the covariance matrix is fully parameterized the full conditional distribution of the precision matrix $\mathbf{\Omega}$ is Wishart, that is $\mathbf{\Omega} \sim W(\tilde{d}, \tilde{\mathbf{E}})$. This draw is usually accomplished via the Bartlett decomposition (Smith and Hocking, 1972). Let $\tilde{\mathbf{C}}$ be the lower triangular Cholesky decomposition of $\tilde{\mathbf{E}}$. Then let \mathbf{A} be an $r \times r$ random lower triangular matrix with independent elements $\{a_{ij} : 0 < i \leq j \leq r\}$ where $a_{ii} \sim \sqrt{\chi_{\tilde{d}-i+1}^2}$ for $i = 1, 2, \dots, r$ and $a_{ij} \sim N(0, 1)$ for $0 < i < j \leq r$. Then $\mathbf{\Omega} = \mathbf{L} \mathbf{L}' \sim W(\tilde{d}, \tilde{\mathbf{E}})$ where $\mathbf{L} = \tilde{\mathbf{C}} \mathbf{A}$. In the process of drawing $\mathbf{\Omega}$ we must first draw \mathbf{L} , so we construct our Gibbs samplers in terms of \mathbf{L} instead of $\mathbf{\Omega}$.

IMH applies straightforwardly to each of these models, though see Appendix C.2 for a detailed derivation of the Hessian for the county population models with a fully parameterized covariance matrix associated. To apply IMHwG in the county population

models, we draw $(\boldsymbol{\beta}, \boldsymbol{\delta})$ in the Metropolis step and either σ^2 or \mathbf{L} in the conditionally conjugate step, depending on the model. The full conditional distribution of σ^2 in both iid models is $IG(a_\sigma + r/2, b_\sigma + \boldsymbol{\delta}'\boldsymbol{\delta}/2)$. The full conditional distribution of $\boldsymbol{\Omega}$ in both fully correlated models is $W(r + 1, (\mathbf{E} + \boldsymbol{\delta}\boldsymbol{\delta}')^{-1})$. Then a draw from the full conditional distribution of \mathbf{L} can be obtained via the Bartlett decomposition as described in the previous paragraph. In the lognormal models the full conditional distribution of ϕ^2 is $IG(a_\phi + n/2, b_\phi + (\log \mathbf{z} - \mathbf{X}\boldsymbol{\beta} - \mathbf{S}\boldsymbol{\delta})'(\log \mathbf{z} - \mathbf{X}\boldsymbol{\beta} - \mathbf{S}\boldsymbol{\delta})/2)$.

C.2 Deriving the Hessian

The log posterior in the fully parameterized Poisson case can be written as

$$\begin{aligned} \log p(\boldsymbol{\beta}, \mathbf{L}, \boldsymbol{\delta} | \mathbf{z}, \mathbf{X}, \mathbf{S}) &= \text{constant} + \sum_{i=1}^n (y_i z_i - \exp(y_i)) + \sum_{k=1}^r \frac{d - k + 1}{2} \log \ell_{kk}^2 \\ &\quad - \frac{1}{2} \left\{ \boldsymbol{\delta}' \mathbf{L} \mathbf{L}' \boldsymbol{\delta} + \frac{(\boldsymbol{\beta} - \mathbf{b})'(\boldsymbol{\beta} - \mathbf{b})}{v^2} + \text{tr}(\mathbf{E} \mathbf{L} \mathbf{L}') \right\} \\ &= \text{constant} + \sum_{i=1}^n \{z_i y_i - \exp(y_i)\} + \frac{1}{2} \mathbf{R}_r \text{vech}(\log \mathbf{L}^2) \\ &\quad - \frac{1}{2} \left[\text{vech}(\mathbf{L})' \mathbf{M}_r \{ \mathbf{K}_r' (\boldsymbol{\delta} \boldsymbol{\delta}' \otimes \mathbf{I}_r) \mathbf{K}_r + (\mathbf{I}_r \otimes \mathbf{E}) \} \mathbf{M}_r' \text{vech}(\mathbf{L}) + \frac{(\boldsymbol{\beta} - \mathbf{b})'(\boldsymbol{\beta} - \mathbf{b})}{v^2} \right] \end{aligned}$$

where $y_i = \mathbf{x}_i' \boldsymbol{\beta} + \mathbf{s}_i' \boldsymbol{\delta}$, \otimes is the Kronecker product, \mathbf{I}_r is the $r \times r$ identity matrix, and \mathbf{R}_r is an $r(r+1)/2 \times r(r+1)/2$ matrix of zeroes except the $(r+1)k - k(k+1)/2 + k - \text{rst}$ diagonal element is equal to $d - k + 1$ for $k = 1, 2, \dots, r$, corresponding to locations in $\text{vech}(\mathbf{L})$ that store the diagonal elements of \mathbf{L} . In $\text{vech}(\log \mathbf{L}^2)$ the log and power are applied element-wise. The expansions of $\boldsymbol{\delta}' \mathbf{L} \mathbf{L}' \boldsymbol{\delta}$ and $\text{tr}(\mathbf{E} \mathbf{L} \mathbf{L}')$ can be derived from the properties of the trace operator, Kronecker products, and the following identities.

$$\boldsymbol{\delta}' \mathbf{L} \mathbf{L}' \boldsymbol{\delta} = \text{vec}(\mathbf{L}' \boldsymbol{\delta})' \text{vec}(\mathbf{L}' \boldsymbol{\delta})$$

$$\text{vec}(\mathbf{L}' \boldsymbol{\delta}) = (\boldsymbol{\delta}' \otimes \mathbf{I}_r) \text{vec}(\mathbf{L}')$$

$$\text{vec}(\mathbf{L}') = \mathbf{K}_r \text{vec}(\mathbf{L}) = \mathbf{K}_r \mathbf{M}_r' \text{vech}(\mathbf{L})$$

and assuming $\mathbf{E} = \mathbf{C}\mathbf{C}'$, then $\text{tr}(\mathbf{L}'\mathbf{C}\mathbf{C}'\mathbf{L}) = \text{vec}(\mathbf{C}'\mathbf{L})' \text{vec}(\mathbf{C}'\mathbf{L})$ where $\text{vec}(\cdot)$ is the vectorization operation which stacks each column of its argument on top of each other into a single column vector, $\text{vech}(\cdot)$ is the half-vectorization operator which is similar but omits elements above the diagonal, tr is the trace operator, \mathbf{K}_r is the $r^2 \times r^2$ commutation matrix such that for any $r \times r$ matrix \mathbf{L} , $\text{vec}(\mathbf{L}') = \mathbf{K}_r \text{vec}(\mathbf{L})$, \mathbf{M}_r is the $r^2 \times r(r+1)/2$ elimination matrix such that for $\text{vech}(\mathbf{L}) = \mathbf{M}_r \text{vec}(\mathbf{L})$ and if additionally \mathbf{L} is lower triangular, $\text{vec}(\mathbf{L}) = \mathbf{M}_r' \text{vech}(\mathbf{L})$ (Magnus and Neudecker, 1980; Magnus, 1988). Let \mathbf{E}_{ij} be an $r \times r$ matrix of zeroes with a one in only its (i, j) th position, $\mathbf{u}_{ij} = \text{vech}(\mathbf{E}_{ij})$, and \mathbf{e}_i be a r -dimensional column vector of zeroes with a one in the i th row. Then \mathbf{K}_r and \mathbf{M}_r can be written as

$$\mathbf{K}_r = \sum_{i=1}^r \sum_{j=1}^r \mathbf{E}_{ij} \otimes \mathbf{E}_{ij}' \quad \text{and} \quad \mathbf{M}_r = \sum_{i \geq j}^r \mathbf{u}_{ij} \otimes \mathbf{e}_j' \otimes \mathbf{e}_i'.$$

Now the first derivatives of the log posterior are given by

$$\begin{aligned} \frac{\partial \log p}{\partial \boldsymbol{\beta}} &= \sum_{i=1}^n (z_i - e^{y_i}) \mathbf{x}_i' - \frac{(\boldsymbol{\beta} - \mathbf{b})'}{v^2}, \\ \frac{\partial \log p}{\partial \boldsymbol{\delta}} &= \sum_{i=1}^n (z_i - e^{y_i}) \mathbf{s}_i' - \boldsymbol{\delta}' \mathbf{L} \mathbf{L}', \\ \frac{\partial \log p}{\partial \text{vech}(\mathbf{L})} &= -\text{vech}(\mathbf{L})' \mathbf{M}_r [\mathbf{K}_r' (\boldsymbol{\delta} \boldsymbol{\delta}' \otimes \mathbf{I}_r) \mathbf{K}_r + (\mathbf{I}_r \otimes \mathbf{E})] \mathbf{M}_r' + \mathbf{R}_r \text{vech}(1/\mathbf{L}), \end{aligned}$$

where in $\text{vech}(1/\mathbf{L})$ the division is applied element-wise. Next the second derivatives are

$$\begin{aligned}
\frac{\partial^2 \log p}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}'} &= - \sum_{i=1}^n e^{y_i} \mathbf{x}_i \mathbf{x}_i' - \frac{1}{v^2} \mathbf{I}_p \\
\frac{\partial^2 \log p}{\partial \boldsymbol{\delta} \partial \boldsymbol{\delta}'} &= - \sum_{i=1}^n e^{y_i} \mathbf{s}_i \mathbf{s}_i' - \mathbf{L} \mathbf{L}' \\
\frac{\partial^2 \log p}{\partial \text{vech}(\mathbf{L}) \partial \text{vech}(\mathbf{L})'} &= - \mathbf{M}_r [\mathbf{K}_r' (\boldsymbol{\delta} \boldsymbol{\delta}' \otimes \mathbf{I}_r) \mathbf{K}_r + (\mathbf{I}_r \otimes \mathbf{E})] \mathbf{M}_r' - \mathbf{R}_r \text{diag}(\text{vech}(\mathbf{L})^{-2}) \\
\frac{\partial^2 \log p}{\partial \boldsymbol{\beta} \partial \boldsymbol{\delta}'} &= - \sum_{i=1}^n e^{y_i} \mathbf{s}_i \mathbf{x}_i' \\
\frac{\partial^2 \log p}{\partial \boldsymbol{\beta} \partial \text{vech}(\mathbf{L})'} &= \mathbf{0}_{p \times r(r+1)/2} \\
\frac{\partial^2 \log p}{\partial \boldsymbol{\delta} \partial \text{vech}(\mathbf{L})'} &= - \frac{\partial \boldsymbol{\delta}' \mathbf{L} \mathbf{L}'}{\partial \text{vech}(\mathbf{L})'} = - (\boldsymbol{\delta}' \otimes \mathbf{I}_r) (\mathbf{I}_{r^2} + \mathbf{K}_r) (\mathbf{L} \otimes \mathbf{I}_r) \mathbf{M}_r',
\end{aligned}$$

where $\text{diag}(\text{vech}(\mathbf{L})^{-2})$ is a diagonal matrix with the elements of $\text{vech}(\mathbf{L})$ raised to the power -2 along the diagonal. The last second derivative matrix can be derived using repeated application of the chain rule and from the following facts for any $r \times r$ matrix \mathbf{A} (Magnus and Neudecker, 2005):

$$\begin{aligned}
\frac{\partial \text{vec}(\mathbf{A} \mathbf{A}')}{\partial \text{vec}(\mathbf{A})} &= (\mathbf{I}_{r^2} + \mathbf{K}_r) (\mathbf{A} \otimes \mathbf{I}_r) \\
\frac{\partial \mathbf{A} \boldsymbol{\delta}}{\partial \text{vec}(\mathbf{A})} &= \boldsymbol{\delta}' \otimes \mathbf{I}_r.
\end{aligned}$$

Then, using the chain rule for lower triangular \mathbf{L} we have

$$\begin{aligned}
\frac{\partial \boldsymbol{\delta}' \mathbf{L} \mathbf{L}'}{\partial \text{vech}(\mathbf{L})'} &= \frac{\partial \boldsymbol{\delta}' \mathbf{L} \mathbf{L}'}{\partial \mathbf{L} \mathbf{L}' \boldsymbol{\delta}} \frac{\partial \mathbf{L} \mathbf{L}' \boldsymbol{\delta}}{\partial \text{vec}(\mathbf{L} \mathbf{L}')} \frac{\partial \text{vec}(\mathbf{L} \mathbf{L}')}{\partial \text{vec}(\mathbf{L})} \frac{\partial \text{vec}(\mathbf{L})}{\partial \text{vech}(\mathbf{L})} \frac{\partial \text{vech}(\mathbf{L})}{\partial \text{vech}(\mathbf{L})'} \\
&= \mathbf{I}_r \frac{\partial \mathbf{L} \mathbf{L}' \boldsymbol{\delta}}{\partial \text{vec}(\mathbf{L} \mathbf{L}')} \frac{\partial \text{vec}(\mathbf{L} \mathbf{L}')}{\partial \text{vec}(\mathbf{L})} \mathbf{M}_r' \mathbf{I}_{r(r+1)/2} \\
&= (\boldsymbol{\delta}' \otimes \mathbf{I}_r) (\mathbf{I}_{r^2} + \mathbf{K}_r) (\mathbf{L} \otimes \mathbf{I}_r) \mathbf{M}_r'.
\end{aligned}$$

Finally, the Hessian is

$$\mathbf{H} = \begin{bmatrix} \frac{\partial^2 \log p}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}'} & \frac{\partial^2 \log p}{\partial \boldsymbol{\beta} \partial \boldsymbol{\delta}'} & \frac{\partial^2 \log p}{\partial \boldsymbol{\beta} \partial \text{vech}(\mathbf{L})'} \\ \frac{\partial^2 \log p}{\partial \boldsymbol{\delta} \partial \boldsymbol{\beta}'} & \frac{\partial^2 \log p}{\partial \boldsymbol{\delta} \partial \boldsymbol{\delta}'} & \frac{\partial^2 \log p}{\partial \boldsymbol{\delta} \partial \text{vech}(\mathbf{L})'} \\ \frac{\partial^2 \log p}{\partial \text{vech}(\mathbf{L}) \partial \boldsymbol{\beta}'} & \frac{\partial^2 \log p}{\partial \text{vech}(\mathbf{L}) \partial \boldsymbol{\delta}'} & \frac{\partial^2 \log p}{\partial \text{vech}(\mathbf{L}) \partial \text{vech}(\mathbf{L})'} \end{bmatrix}.$$

D Alternative MCMC algorithms

[CURRENTLY THIS SECTION ONLY HAS GENERIC VERSIONS OF THE ALGORITHMS. DETAILS FOR OUR CASES?] This section describes the alternative MCMC algorithms used for comparisons with PSO assisted Metropolis-Hastings algorithms. In particular, we use two types of adaptive random walk Metropolis within Gibbs algorithms. In the generic problem, suppose we wish to sample from the posterior distribution $p(\boldsymbol{\theta}|\mathbf{y})$ using MCMC methods. Suppose that $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)$ and that the full conditional $p(\boldsymbol{\theta}_1|\boldsymbol{\theta}_2, \mathbf{y})$ can be sampled from easily while the full conditional for $\boldsymbol{\theta}_2$ may be intractable. Then a single move random walk Metropolis within Gibbs (RWwG) algorithm for this problem is as follows.

Algorithm 1 (Single move random walk Metropolis within Gibbs) *Given target posterior $p(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2|\mathbf{y})$ where $\boldsymbol{\theta}_2 = (\theta_{21}, \theta_{22}, \dots, \theta_{2n})$, it is easy to sample from $p(\boldsymbol{\theta}_1|\boldsymbol{\theta}_2, \mathbf{y})$, and given that the support of $p(\boldsymbol{\theta}_2|\boldsymbol{\theta}_1, \mathbf{y})$ is \mathbb{R}^n , iteration $t + 1$ is obtained from iteration t via*

1. Draw $\boldsymbol{\theta}_1^{(t+1)} \sim p(\boldsymbol{\theta}_1|\boldsymbol{\theta}_2^{(t)}, \mathbf{y})$

2. For $k = 1, 2, \dots, n$

Draw $\theta_{2k}^{(prop)} \sim N(\theta_{2k}^{(t)}, \eta_k)$ and form the Metropolis acceptance ratio

$$a_k^{(t)} = \frac{p(\theta_{2k}^{(prop)}|\boldsymbol{\theta}_1^{(t+1)}, \theta_{21}^{(t+1)}, \dots, \theta_{2k-1}^{(t+1)}, \theta_{2k+1}^{(t)}, \dots, \theta_{2n}^{(t)}, \mathbf{y})}{p(\theta_{2k}^{(t)}|\boldsymbol{\theta}_1^{(t+1)}, \theta_{21}^{(t+1)}, \dots, \theta_{2k-1}^{(t+1)}, \theta_{2k+1}^{(t)}, \dots, \theta_{2n}^{(t)}, \mathbf{y})}.$$

Then set $\theta_{2k}^{(t+1)} = \theta_{2k}^{(prop)}$ with probability $r_k^{(t)} = \min(a_k, 1)$ and otherwise set

$$\theta_{2k}^{(t+1)} = \theta_{2k}^{(t)}$$

A major problem with this algorithm is that the η_k s must be selected so that the algorithm Metropolis steps accept a reasonable amount of time. According to Gelman et al. (1996) the optimal acceptance rate in a narrow set of problems is 0.44 for a single dimensional random walk, though this is often used as a guideline for more complex problems. We

adaptively tune η_k during the burn-in period of the chain in a manner discussed in Andrieu and Thoms (2008). The computed Metropolis acceptance probability for θ_{2k} is denoted by $r_k^{(t)} = \min(a_k^{(t)}, 1)$; let r^* denote the target acceptance probability. Then we evolve η_k over time via

$$\log \eta_k^{(t+1)} = \log \eta_k^{(t)} + \gamma^{(t+1)}(r_k^{(t)} - r^*)$$

where $\gamma^{(t)}$ is a fixed sequence decreasing to zero fast enough to enough that the algorithm converges. The intuition here is that if the acceptance rate is too high, we can increase the effective search space by increasing the random walk standard deviation, while if it is too low, we can decrease the effective search space by decreasing the random walk standard deviation. We use $r^* = 0.44$. A tuning method such as this implies that $(\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}, \dots)$ is no longer a Markov chain, and additional conditions are required to ensure convergence to the target posterior distribution. So in practice our RWwG algorithms run Algorithm 1 with the tuning method described above until convergence and until the η_k s settle into a rough equilibrium, then we use Algorithm 1 without tuning but using the last known values of the η_k s. In practice this is equivalent setting $\gamma^{(t)} = 0$ after the burn-in period. During the burn-in we set $\gamma^{(t)} = 1$ and we initialize at $\eta_k^{(0)} = 1$ for all k .

An alternative to RWwG algorithms are block random walk Metropolis within Gibbs algorithms (B-RWwG). Suppose we have an estimate of the covariance structure between the elements of $\boldsymbol{\theta}_2$ in the posterior. Then

Algorithm 2 (Block random walk Metropolis within Gibbs) *Given target posterior $p(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2 | \mathbf{y})$ where it is easy to sample from $p(\boldsymbol{\theta}_1 | \boldsymbol{\theta}_2, \mathbf{y})$, the support of $p(\boldsymbol{\theta}_2 | \boldsymbol{\theta}_1, \mathbf{y})$ is \mathbb{R}^n , and given an estimate $\boldsymbol{\Sigma}$ of $\text{cov}(\boldsymbol{\theta}_2 | \mathbf{y})$, iteration $t + 1$ is obtained from iteration t via*

1. Draw $\boldsymbol{\theta}_1^{(t+1)} \sim p(\boldsymbol{\theta}_1 | \boldsymbol{\theta}_2^{(t)}, \mathbf{y})$

2. Draw $\boldsymbol{\theta}_2^{(prop)} \sim N(\boldsymbol{\theta}_2^{(t)}, \eta \boldsymbol{\Sigma})$ and form the Metropolis acceptance ratio

$$a = \frac{p(\boldsymbol{\theta}_2^{(prop)} | \boldsymbol{\theta}_1^{(t+1)}, \mathbf{y})}{p(\boldsymbol{\theta}_2^{(t)} | \boldsymbol{\theta}_1^{(t+1)}, \mathbf{y})}.$$

Then set $\boldsymbol{\theta}_2^{(t+1)} = \boldsymbol{\theta}_2^{(prop)}$ with probability $\min(a, 1)$ and otherwise set $\boldsymbol{\theta}_2^{(t+1)} = \boldsymbol{\theta}_2^{(t)}$.

We make this algorithm adaptive as well by tuning η in the same fashion as the η_k s above. The main difference is that this algorithm must be initialized with a crude estimate of $\boldsymbol{\Sigma}$ by running, e.g., Algorithm 1.

OF1	Global nbhd				Ring-3 nbhd				Ring-1 nbhd			
Algorithm	Mean	SD	\hat{p}	\hat{t}	Mean	SD	\hat{p}	\hat{t}	Mean	SD	\hat{p}	\hat{t}
PSO	3147.42	2726.79	0.00	∞	0.00	0.00	1.00	360.50	0.01	0.01	0.58	482.50
BBPSO-MC	88237.79	4013.36	0.00	∞	89130.57	4684.99	0.00	∞	88274.59	4528.45	0.00	∞
BBPSOxp-MC	81476.62	3295.56	0.00	∞	75471.96	4929.58	0.00	∞	59817.42	5244.39	0.00	∞
AT-BBPSO-MC												
$df = 1, R^* = 0.1$	1.71	0.55	0.00	∞	2.01	0.57	0.00	∞	2.04	0.66	0.00	∞
$df = 1, R^* = 0.3$	0.00	0.00	1.00	449.00	0.00	0.00	1.00	463.00	0.01	0.00	0.46	∞
$df = 1, R^* = 0.5$	0.00	0.00	1.00	286.00	0.00	0.00	1.00	308.00	0.00	0.00	1.00	353.50
$df = 1, R^* = 0.7$	0.00	0.00	1.00	223.00	0.00	0.00	1.00	254.50	0.00	0.00	1.00	328.00
$df = 3, R^* = 0.1$	7.17	2.25	0.00	∞	6.91	1.65	0.00	∞	5.11	1.45	0.00	∞
$df = 3, R^* = 0.3$	0.02	0.01	0.10	∞	0.02	0.01	0.02	∞	0.03	0.01	0.00	∞
$df = 3, R^* = 0.5$	0.00	0.00	1.00	330.00	0.00	0.00	1.00	348.00	0.00	0.00	1.00	390.50
$df = 3, R^* = 0.7$	0.00	0.00	1.00	263.00	0.00	0.00	1.00	288.50	0.00	0.00	1.00	388.00
$df = 5, R^* = 0.1$	12.49	3.78	0.00	∞	12.42	3.39	0.00	∞	7.43	3.03	0.00	∞
$df = 5, R^* = 0.3$	0.03	0.01	0.00	∞	0.03	0.01	0.00	∞	0.07	0.03	0.00	∞
$df = 5, R^* = 0.5$	0.00	0.00	1.00	357.00	0.00	0.00	1.00	373.00	0.00	0.00	1.00	420.50
$df = 5, R^* = 0.7$	0.00	0.00	1.00	290.00	0.00	0.00	1.00	325.50	195.94	336.65	0.06	∞
$df = \infty, R^* = 0.1$	43.64	10.08	0.00	∞	37.75	8.20	0.00	∞	19.58	6.50	0.00	∞
$df = \infty, R^* = 0.3$	0.14	0.03	0.00	∞	0.17	0.05	0.00	∞	0.23	0.06	0.00	∞
$df = \infty, R^* = 0.5$	0.00	0.00	1.00	413.00	0.00	0.00	1.00	430.50	0.01	0.00	0.82	491.00
$df = \infty, R^* = 0.7$	0.00	0.00	1.00	357.00	0.00	0.00	1.00	443.00	16625.74	3842.72	0.00	∞
AT-BBPSOxp-MC												
$df = 1, R^* = 0.1$	3.77	0.91	0.00	∞	2.22	0.76	0.00	∞	0.43	0.20	0.00	∞
$df = 1, R^* = 0.3$	0.02	0.01	0.00	∞	0.01	0.01	0.22	∞	0.00	0.00	1.00	448.50
$df = 1, R^* = 0.5$	0.00	0.00	1.00	349.50	0.00	0.00	1.00	349.50	0.00	0.00	1.00	323.00
$df = 1, R^* = 0.7$	0.00	0.00	1.00	283.50	0.00	0.00	1.00	295.00	0.00	0.00	1.00	285.00
$df = 3, R^* = 0.1$	11.64	2.95	0.00	∞	6.57	2.30	0.00	∞	1.52	0.91	0.00	∞
$df = 3, R^* = 0.3$	0.08	0.03	0.00	∞	0.05	0.02	0.00	∞	0.01	0.01	0.20	∞
$df = 3, R^* = 0.5$	0.00	0.00	1.00	399.50	0.00	0.00	1.00	407.50	0.00	0.00	1.00	374.00
$df = 3, R^* = 0.7$	0.00	0.00	1.00	340.50	0.00	0.00	1.00	365.50	0.00	0.00	1.00	365.00
$df = 5, R^* = 0.1$	18.53	4.90	0.00	∞	9.29	2.96	0.00	∞	1.86	0.75	0.00	∞
$df = 5, R^* = 0.3$	0.14	0.04	0.00	∞	0.10	0.03	0.00	∞	0.02	0.01	0.04	∞
$df = 5, R^* = 0.5$	0.00	0.00	1.00	430.50	0.00	0.00	1.00	436.00	0.00	0.00	1.00	404.00
$df = 5, R^* = 0.7$	0.00	0.00	1.00	389.50	0.01	0.02	0.92	451.50	0.04	0.13	0.74	454.50
$df = \infty, R^* = 0.1$	43.01	9.65	0.00	∞	23.58	7.63	0.00	∞	4.82	1.83	0.00	∞
$df = \infty, R^* = 0.3$	0.42	0.12	0.00	∞	0.32	0.09	0.00	∞	0.07	0.04	0.00	∞
$df = \infty, R^* = 0.5$	0.01	0.00	0.96	482.50	0.01	0.00	0.84	489.00	0.00	0.00	1.00	454.00
$df = \infty, R^* = 0.7$	0.06	0.15	0.24	∞	1766.73	1214.77	0.00	∞	2237.55	1691.20	0.00	∞
DI-PSO												
$\alpha = 50, \beta = 1$	5800.59	3155.97	0.00	∞	698.71	1183.33	0.00	∞	35.24	79.82	0.00	∞
$\alpha = 50, \beta = 2$	8130.93	5431.61	0.00	∞	2254.74	1775.35	0.00	∞	707.41	570.80	0.00	∞
$\alpha = 50, \beta = 4$	14881.90	8467.41	0.00	∞	6133.60	3281.73	0.00	∞	3051.71	2314.76	0.00	∞
$\alpha = 100, \beta = 1$	3574.49	3170.07	0.00	∞	146.17	331.79	0.00	∞	0.34	0.91	0.22	∞
$\alpha = 100, \beta = 2$	7903.23	4250.18	0.00	∞	1248.39	1326.94	0.00	∞	83.35	140.17	0.00	∞
$\alpha = 100, \beta = 4$	13455.03	6484.96	0.00	∞	4378.86	3004.31	0.00	∞	1591.49	1194.36	0.00	∞
$\alpha = 200, \beta = 1$	2746.62	2890.48	0.00	∞	27.67	126.87	0.00	∞	0.00	0.00	0.96	384.50
$\alpha = 200, \beta = 2$	8261.84	4537.27	0.00	∞	180.57	334.02	0.00	∞	3.38	7.17	0.02	∞
$\alpha = 200, \beta = 4$	12100.60	6066.27	0.00	∞	2947.83	2171.07	0.00	∞	581.14	1413.97	0.00	∞
AT-PSO												
$R^* = 0.1$	118.60	228.55	0.00	∞	188.52	718.43	0.00	∞	10169.59	7405.93	0.00	∞
$R^* = 0.3$	70.10	326.06	0.04	∞	0.00	0.00	1.00	339.00	0.01	0.01	0.86	464.50
$R^* = 0.5$	373.93	894.14	0.00	∞	0.00	0.00	1.00	275.00	0.00	0.00	1.00	302.50
$R^* = 0.7$	1932.49	1906.13	0.00	∞	0.15	0.99	0.90	395.50	0.00	0.00	0.98	343.00

Table 2: Simulation results for OF1. See text for description.

OF2		Global nbhd				Ring-3 nbhd				Ring-1 nbhd			
Algorithm		Mean	SD	\hat{p}	\hat{t}	Mean	SD	\hat{p}	\hat{t}	Mean	SD	\hat{p}	\hat{t}
PSO		9058.41	9685.19	0.00	∞	1264.42	799.32	0.00	∞	5339.50	3143.93	0.00	∞
BBPSO-MC		11883975.00	628634.30	0.00	∞	12276069.13	699375.46	0.00	∞	11493578.31	1172502.58	0.00	∞
BBPSOxp-MC		11010472.51	544947.26	0.00	∞	10140651.29	619325.29	0.00	∞	7004850.88	1225368.46	0.00	∞
AT-BBPSO-MC													
$df = 1, R^* = 0.1$		945.06	802.75	0.00	∞	710.80	469.14	0.00	∞	774.49	747.78	0.00	∞
$df = 1, R^* = 0.3$		327.20	240.51	0.00	∞	926.29	1132.84	0.00	∞	1008.87	1034.00	0.00	∞
$df = 1, R^* = 0.5$		469.09	503.21	0.00	∞	1468.95	1561.26	0.00	∞	3253.60	2878.65	0.00	∞
$df = 1, R^* = 0.7$		1406.45	1261.92	0.00	∞	3977.07	3875.19	0.00	∞	10836.34	7853.42	0.00	∞
$df = 3, R^* = 0.1$		219.80	111.62	0.00	∞	202.99	142.50	0.00	∞	255.86	157.95	0.00	∞
$df = 3, R^* = 0.3$		19.11	13.28	0.00	∞	63.30	60.89	0.00	∞	145.93	153.16	0.00	∞
$df = 3, R^* = 0.5$		29.80	35.95	0.00	∞	214.82	214.48	0.00	∞	846.76	697.18	0.00	∞
$df = 3, R^* = 0.7$		144.09	103.71	0.00	∞	1033.32	796.45	0.00	∞	8079.02	20596.37	0.00	∞
$df = 5, R^* = 0.1$		185.82	78.05	0.00	∞	187.84	85.91	0.00	∞	237.84	122.56	0.00	∞
$df = 5, R^* = 0.3$		10.85	8.18	0.00	∞	24.21	18.38	0.00	∞	55.58	39.08	0.00	∞
$df = 5, R^* = 0.5$		11.60	10.52	0.00	∞	54.75	40.96	0.00	∞	468.53	369.98	0.00	∞
$df = 5, R^* = 0.7$		72.38	51.38	0.00	∞	422.19	277.83	0.00	∞	8889.40	7833.07	0.00	∞
$df = \infty, R^* = 0.1$		282.38	102.62	0.00	∞	262.49	79.80	0.00	∞	281.59	115.66	0.00	∞
$df = \infty, R^* = 0.3$		7.83	6.04	0.00	∞	16.43	10.41	0.00	∞	37.22	23.14	0.00	∞
$df = \infty, R^* = 0.5$		6.84	5.16	0.00	∞	35.82	32.97	0.00	∞	261.12	152.46	0.00	∞
$df = \infty, R^* = 0.7$		50.86	36.80	0.00	∞	484.76	270.15	0.00	∞	1167787.27	401166.29	0.00	∞
AT-BBPSOxp-MC													
$df = 1, R^* = 0.1$		1872.52	1038.69	0.00	∞	1799.53	862.91	0.00	∞	942.70	590.06	0.00	∞
$df = 1, R^* = 0.3$		1204.75	814.82	0.00	∞	1864.65	1185.93	0.00	∞	961.45	928.89	0.00	∞
$df = 1, R^* = 0.5$		1881.72	1041.49	0.00	∞	3108.80	1482.50	0.00	∞	2519.47	1307.36	0.00	∞
$df = 1, R^* = 0.7$		5161.81	3212.03	0.00	∞	7854.74	4468.87	0.00	∞	7775.57	4525.87	0.00	∞
$df = 3, R^* = 0.1$		565.47	256.64	0.00	∞	665.17	307.78	0.00	∞	508.77	268.54	0.00	∞
$df = 3, R^* = 0.3$		191.75	130.61	0.00	∞	363.46	260.03	0.00	∞	267.88	190.70	0.00	∞
$df = 3, R^* = 0.5$		438.28	273.27	0.00	∞	654.47	404.48	0.00	∞	1110.91	1424.90	0.00	∞
$df = 3, R^* = 0.7$		1384.56	909.83	0.00	∞	3503.63	3161.16	0.00	∞	4256.67	3547.43	0.00	∞
$df = 5, R^* = 0.1$		480.63	187.88	0.00	∞	504.39	220.77	0.00	∞	400.66	214.64	0.00	∞
$df = 5, R^* = 0.3$		115.80	68.47	0.00	∞	197.49	138.10	0.00	∞	225.22	142.07	0.00	∞
$df = 5, R^* = 0.5$		181.62	128.13	0.00	∞	333.87	159.91	0.00	∞	472.13	276.68	0.00	∞
$df = 5, R^* = 0.7$		660.28	313.88	0.00	∞	1817.45	834.00	0.00	∞	4632.54	4302.74	0.00	∞
$df = \infty, R^* = 0.1$		586.91	176.77	0.00	∞	595.51	179.25	0.00	∞	541.43	255.32	0.00	∞
$df = \infty, R^* = 0.3$		81.20	47.02	0.00	∞	128.97	94.06	0.00	∞	176.75	109.79	0.00	∞
$df = \infty, R^* = 0.5$		119.58	66.07	0.00	∞	232.02	163.01	0.00	∞	559.93	315.23	0.00	∞
$df = \infty, R^* = 0.7$		971.24	600.68	0.00	∞	118742.43	156281.96	0.00	∞	34066.02	28609.67	0.00	∞
DI-PSO													
$\alpha = 50, \beta = 1$		40389.25	24692.95	0.00	∞	18720.02	11239.63	0.00	∞	8622.85	4980.17	0.00	∞
$\alpha = 50, \beta = 2$		44518.25	22261.16	0.00	∞	21028.28	10707.64	0.00	∞	12726.73	5322.96	0.00	∞
$\alpha = 50, \beta = 4$		51851.57	25120.63	0.00	∞	29532.84	13000.01	0.00	∞	16540.60	8846.75	0.00	∞
$\alpha = 100, \beta = 1$		24843.60	15674.35	0.00	∞	10917.91	7294.55	0.00	∞	5011.53	2384.94	0.00	∞
$\alpha = 100, \beta = 2$		37890.78	24040.40	0.00	∞	17226.82	8980.99	0.00	∞	10305.47	5737.74	0.00	∞
$\alpha = 100, \beta = 4$		52019.62	25824.86	0.00	∞	30086.99	20743.14	0.00	∞	12073.62	5451.03	0.00	∞
$\alpha = 200, \beta = 1$		22007.15	13078.38	0.00	∞	5785.75	4635.98	0.00	∞	3694.82	2434.96	0.00	∞
$\alpha = 200, \beta = 2$		32862.70	16618.65	0.00	∞	13173.95	6095.13	0.00	∞	8033.58	4532.79	0.00	∞
$\alpha = 200, \beta = 4$		50850.73	21907.04	0.00	∞	21202.31	11345.72	0.00	∞	14321.18	7476.13	0.00	∞
AT-PSO													
$R^* = 0.1$		7359.75	6803.21	0.00	∞	8340.43	7908.49	0.00	∞	56057.72	26008.24	0.00	∞
$R^* = 0.3$		7582.23	12259.21	0.00	∞	593.85	754.83	0.00	∞	5252.44	4221.37	0.00	∞
$R^* = 0.5$		17909.32	16070.66	0.00	∞	1397.43	1598.40	0.00	∞	1846.47	1312.13	0.00	∞
$R^* = 0.7$		43178.25	24887.81	0.00	∞	6301.86	5349.85	0.00	∞	5760.60	3355.07	0.00	∞

Table 3: Simulation results for OF2. See text for description.

OF3		Global nbhd				Ring-3 nbhd				Ring-1 nbhd			
Algorithm		Mean	SD	\hat{p}	\hat{t}	Mean	SD	\hat{p}	\hat{t}	Mean	SD	\hat{p}	\hat{t}
PSO		970637.16	1719993.18	0.00	∞	144.21	197.18	0.00	∞	207.96	306.95	0.00	∞
BBPSO-MC		121023363.91	13812838.21	0.00	∞	154011414.01	16232288.40	0.00	∞	174968862.42	23613990.88	0.00	∞
BBPSOxp-MC		138645962.86	15657816.44	0.00	∞	132653861.04	13853873.26	0.00	∞	89975682.91	15709701.24	0.00	∞
AT-BBPSO-MC													
$df = 1, R^* = 0.1$		502.81	549.75	0.00	∞	489.99	659.01	0.00	∞	417.60	453.85	0.00	∞
$df = 1, R^* = 0.3$		238.76	403.40	0.00	∞	241.03	308.07	0.00	∞	221.48	324.13	0.00	∞
$df = 1, R^* = 0.5$		195.84	235.99	0.00	∞	122.20	115.62	0.00	∞	230.69	508.80	0.00	∞
$df = 1, R^* = 0.7$		226.04	526.49	0.00	∞	274.61	575.37	0.00	∞	225.18	415.80	0.00	∞
$df = 3, R^* = 0.1$		872.47	954.91	0.00	∞	603.26	714.65	0.00	∞	657.78	863.02	0.00	∞
$df = 3, R^* = 0.3$		187.54	155.35	0.00	∞	140.43	115.93	0.00	∞	267.09	432.14	0.00	∞
$df = 3, R^* = 0.5$		179.88	235.11	0.00	∞	156.27	214.31	0.00	∞	172.14	335.49	0.00	∞
$df = 3, R^* = 0.7$		217.20	455.21	0.00	∞	123.77	165.76	0.00	∞	207.58	550.05	0.00	∞
$df = 5, R^* = 0.1$		781.19	711.43	0.00	∞	679.51	634.98	0.00	∞	541.45	530.74	0.00	∞
$df = 5, R^* = 0.3$		221.15	284.28	0.00	∞	214.66	361.21	0.00	∞	221.80	250.50	0.00	∞
$df = 5, R^* = 0.5$		206.31	228.63	0.00	∞	144.48	112.27	0.00	∞	165.02	319.30	0.00	∞
$df = 5, R^* = 0.7$		182.94	257.84	0.00	∞	48.60	73.43	0.00	∞	108.13	360.85	0.00	∞
$df = \infty, R^* = 0.1$		1382.42	1071.14	0.00	∞	1089.55	708.37	0.00	∞	704.85	580.05	0.00	∞
$df = \infty, R^* = 0.3$		269.15	320.43	0.00	∞	141.57	144.43	0.00	∞	160.88	290.82	0.00	∞
$df = \infty, R^* = 0.5$		191.11	147.95	0.00	∞	194.51	297.34	0.00	∞	165.00	212.70	0.00	∞
$df = \infty, R^* = 0.7$		146.08	127.14	0.00	∞	61.55	73.27	0.00	∞	317.42	999.55	0.00	∞
AT-BBPSOxp-MC													
$df = 1, R^* = 0.1$		611.07	434.02	0.00	∞	402.89	264.11	0.00	∞	216.75	152.59	0.00	∞
$df = 1, R^* = 0.3$		221.23	231.63	0.00	∞	129.90	92.90	0.00	∞	92.61	60.58	0.00	∞
$df = 1, R^* = 0.5$		170.81	394.57	0.00	∞	107.35	136.60	0.00	∞	104.84	110.07	0.00	∞
$df = 1, R^* = 0.7$		136.87	248.84	0.00	∞	131.50	432.59	0.00	∞	98.33	135.63	0.00	∞
$df = 3, R^* = 0.1$		720.92	215.53	0.00	∞	560.46	249.45	0.00	∞	275.19	192.65	0.00	∞
$df = 3, R^* = 0.3$		174.45	244.90	0.00	∞	127.44	58.88	0.00	∞	150.36	127.50	0.00	∞
$df = 3, R^* = 0.5$		151.95	244.27	0.00	∞	92.78	78.27	0.00	∞	109.43	148.82	0.00	∞
$df = 3, R^* = 0.7$		104.03	181.86	0.00	∞	168.15	436.88	0.00	∞	36.17	52.63	0.00	∞
$df = 5, R^* = 0.1$		946.03	410.19	0.00	∞	555.51	250.59	0.00	∞	278.16	196.26	0.00	∞
$df = 5, R^* = 0.3$		121.66	80.43	0.00	∞	120.35	74.89	0.00	∞	109.47	91.42	0.00	∞
$df = 5, R^* = 0.5$		106.68	49.70	0.00	∞	129.63	124.28	0.00	∞	76.17	45.36	0.00	∞
$df = 5, R^* = 0.7$		54.18	64.54	0.00	∞	39.95	84.35	0.00	∞	37.37	59.62	0.00	∞
$df = \infty, R^* = 0.1$		1429.56	774.39	0.00	∞	896.51	333.33	0.00	∞	390.78	225.22	0.00	∞
$df = \infty, R^* = 0.3$		109.28	76.19	0.00	∞	112.02	95.78	0.00	∞	115.01	89.85	0.00	∞
$df = \infty, R^* = 0.5$		103.74	35.01	0.00	∞	118.78	171.39	0.00	∞	90.54	74.14	0.00	∞
$df = \infty, R^* = 0.7$		43.35	75.92	0.00	∞	18.13	12.93	0.00	∞	19.99	20.18	0.00	∞
DI-PSO													
$\alpha = 50, \beta = 1$		3854178.52	3563712.76	0.00	∞	476166.48	1068595.33	0.00	∞	49808.17	90745.61	0.00	∞
$\alpha = 50, \beta = 2$		6562468.86	8285814.85	0.00	∞	1595754.13	1624066.71	0.00	∞	1327935.46	1951123.16	0.00	∞
$\alpha = 50, \beta = 4$		19701402.98	12961045.13	0.00	∞	10884269.41	8299038.15	0.00	∞	5985429.81	3917709.91	0.00	∞
$\alpha = 100, \beta = 1$		1870977.66	2324247.39	0.00	∞	57070.98	164749.98	0.00	∞	1760.17	3897.70	0.00	∞
$\alpha = 100, \beta = 2$		7530967.66	8856235.46	0.00	∞	866151.25	1270221.03	0.00	∞	207921.17	463007.58	0.00	∞
$\alpha = 100, \beta = 4$		22584490.35	22402688.80	0.00	∞	7448049.14	6263068.66	0.00	∞	4916024.35	4602692.81	0.00	∞
$\alpha = 200, \beta = 1$		1217548.01	1884033.63	0.00	∞	3621.58	8954.06	0.00	∞	483.36	1129.34	0.00	∞
$\alpha = 200, \beta = 2$		4461093.42	4947722.54	0.00	∞	120433.51	307696.11	0.00	∞	3403.48	4242.88	0.00	∞
$\alpha = 200, \beta = 4$		17327523.00	13868599.90	0.00	∞	3061892.74	3273826.33	0.00	∞	1189001.99	1583478.70	0.00	∞
AT-PSO													
$R^* = 0.1$		5135.53	14617.98	0.00	∞	446469.66	2260240.37	0.00	∞	34766227.37	32913538.16	0.00	∞
$R^* = 0.3$		1577.07	4053.03	0.00	∞	176.34	324.07	0.00	∞	168.65	155.33	0.00	∞
$R^* = 0.5$		19154.96	47380.82	0.00	∞	315.71	711.01	0.00	∞	110.71	98.15	0.00	∞
$R^* = 0.7$		1297015.56	3307403.64	0.00	∞	395.90	933.41	0.00	∞	239.90	394.68	0.00	∞

Table 4: Simulation results for OF3. See text for description.

OF4		Global nbhd				Ring-3 nbhd				Ring-1 nbhd			
Algorithm		Mean	SD	\hat{p}	\hat{t}	Mean	SD	\hat{p}	\hat{t}	Mean	SD	\hat{p}	\hat{t}
PSO		26.31	11.97	0.00	∞	4.45	2.63	0.02	∞	3.17	1.81	0.00	∞
BBPSO-MC		151.12	10.61	0.00	∞	144.75	15.12	0.00	∞	124.80	14.36	0.00	∞
BBPSOxp-MC		165.73	5.44	0.00	∞	159.62	8.67	0.00	∞	132.89	13.23	0.00	∞
AT-BBPSO-MC													
$df = 1, R^* = 0.1$		5.11	2.09	0.00	∞	4.13	1.43	0.00	∞	3.82	1.43	0.00	∞
$df = 1, R^* = 0.3$		4.87	1.98	0.02	∞	3.05	1.54	0.02	∞	2.88	1.66	0.06	∞
$df = 1, R^* = 0.5$		6.74	1.77	0.00	∞	4.62	2.09	0.00	∞	3.96	1.84	0.00	∞
$df = 1, R^* = 0.7$		8.20	2.12	0.00	∞	7.04	2.17	0.00	∞	5.75	1.94	0.00	∞
$df = 3, R^* = 0.1$		5.58	1.31	0.00	∞	5.60	1.33	0.00	∞	5.77	1.80	0.00	∞
$df = 3, R^* = 0.3$		3.24	1.93	0.08	∞	1.41	0.97	0.16	∞	1.15	0.94	0.24	∞
$df = 3, R^* = 0.5$		4.30	1.80	0.00	∞	2.49	1.56	0.06	∞	2.38	1.23	0.06	∞
$df = 3, R^* = 0.7$		5.18	2.00	0.00	∞	3.41	1.44	0.00	∞	3.84	1.82	0.02	∞
$df = 5, R^* = 0.1$		6.30	1.42	0.00	∞	6.44	1.59	0.00	∞	6.60	1.63	0.00	∞
$df = 5, R^* = 0.3$		2.63	1.46	0.04	∞	1.51	1.14	0.18	∞	1.07	1.05	0.20	∞
$df = 5, R^* = 0.5$		3.60	2.01	0.04	∞	2.13	1.28	0.08	∞	1.77	1.22	0.14	∞
$df = 5, R^* = 0.7$		4.36	1.40	0.00	∞	3.01	1.58	0.02	∞	3.10	1.47	0.02	∞
$df = \infty, R^* = 0.1$		9.32	1.28	0.00	∞	9.44	1.44	0.00	∞	9.11	1.73	0.00	∞
$df = \infty, R^* = 0.3$		1.54	1.02	0.00	∞	0.86	1.08	0.00	∞	0.68	0.86	0.02	∞
$df = \infty, R^* = 0.5$		2.68	1.60	0.10	∞	1.41	1.07	0.16	∞	1.39	1.11	0.20	∞
$df = \infty, R^* = 0.7$		3.23	1.87	0.04	∞	1.88	1.51	0.18	∞	1.69	1.31	0.24	∞
AT-BBPSOxp-MC													
$df = 1, R^* = 0.1$		3.25	1.09	0.00	∞	2.05	0.79	0.00	∞	0.57	0.41	0.00	∞
$df = 1, R^* = 0.3$		0.66	0.76	0.26	∞	0.26	0.46	0.62	490.50	0.59	0.72	0.54	445.00
$df = 1, R^* = 0.5$		1.43	0.98	0.20	∞	1.20	1.10	0.32	∞	1.47	1.25	0.22	∞
$df = 1, R^* = 0.7$		2.78	1.75	0.08	∞	2.45	1.45	0.06	∞	2.58	1.56	0.06	∞
$df = 3, R^* = 0.1$		4.33	1.11	0.00	∞	3.33	1.04	0.00	∞	0.80	0.53	0.00	∞
$df = 3, R^* = 0.3$		0.22	0.44	0.02	∞	0.12	0.31	0.38	∞	0.31	0.49	0.70	458.50
$df = 3, R^* = 0.5$		0.51	0.72	0.60	335.50	0.51	0.93	0.64	327.50	0.84	0.99	0.48	∞
$df = 3, R^* = 0.7$		1.24	0.95	0.22	∞	1.18	0.99	0.24	∞	1.58	1.21	0.20	∞
$df = 5, R^* = 0.1$		5.12	1.04	0.00	∞	4.14	1.26	0.00	∞	0.91	0.62	0.00	∞
$df = 5, R^* = 0.3$		0.17	0.35	0.00	∞	0.05	0.19	0.24	∞	0.27	0.43	0.72	458.00
$df = 5, R^* = 0.5$		0.34	0.57	0.70	334.00	0.30	0.49	0.70	327.50	0.74	0.82	0.46	∞
$df = 5, R^* = 0.7$		1.10	0.99	0.34	∞	0.89	0.95	0.40	∞	1.94	1.57	0.14	∞
$df = \infty, R^* = 0.1$		6.79	1.32	0.00	∞	5.75	1.36	0.00	∞	1.29	0.60	0.00	∞
$df = \infty, R^* = 0.3$		0.03	0.01	0.00	∞	0.08	0.23	0.00	∞	0.20	0.38	0.78	470.50
$df = \infty, R^* = 0.5$		0.25	0.46	0.76	340.50	0.21	0.40	0.78	334.00	0.63	0.76	0.52	323.00
$df = \infty, R^* = 0.7$		0.71	0.83	0.46	∞	0.95	1.02	0.42	∞	1.50	1.11	0.22	∞
DI-PSO													
$\alpha = 50, \beta = 1$		32.00	14.44	0.00	∞	12.58	5.43	0.00	∞	6.46	3.78	0.00	∞
$\alpha = 50, \beta = 2$		35.82	12.86	0.00	∞	18.69	7.41	0.00	∞	9.51	4.88	0.00	∞
$\alpha = 50, \beta = 4$		54.12	18.93	0.00	∞	28.25	10.10	0.00	∞	15.24	5.62	0.00	∞
$\alpha = 100, \beta = 1$		23.66	12.32	0.00	∞	11.01	6.22	0.00	∞	4.88	2.75	0.00	∞
$\alpha = 100, \beta = 2$		32.56	13.79	0.00	∞	13.56	5.63	0.00	∞	7.34	4.10	0.00	∞
$\alpha = 100, \beta = 4$		52.63	21.66	0.00	∞	24.35	9.50	0.00	∞	13.91	5.97	0.00	∞
$\alpha = 200, \beta = 1$		19.10	6.98	0.00	∞	8.20	4.44	0.00	∞	4.53	2.38	0.00	∞
$\alpha = 200, \beta = 2$		36.31	12.85	0.00	∞	11.45	5.21	0.00	∞	5.46	2.87	0.00	∞
$\alpha = 200, \beta = 4$		51.61	20.50	0.00	∞	20.01	7.50	0.00	∞	10.15	3.55	0.00	∞
AT-PSO													
$R^* = 0.1$		7.78	3.44	0.00	∞	5.20	2.00	0.00	∞	23.42	13.89	0.00	∞
$R^* = 0.3$		15.00	6.19	0.00	∞	7.59	3.13	0.00	∞	4.11	2.31	0.00	∞
$R^* = 0.5$		22.59	8.87	0.00	∞	11.60	5.14	0.00	∞	6.01	2.99	0.00	∞
$R^* = 0.7$		35.35	14.56	0.00	∞	15.32	4.51	0.00	∞	9.60	4.15	0.00	∞

Table 5: Simulation results for OF4. See text for description.

OF5		Global nbhd				Ring-3 nbhd				Ring-1 nbhd			
Algorithm		Mean	SD	\hat{p}	\hat{t}	Mean	SD	\hat{p}	\hat{t}	Mean	SD	\hat{p}	\hat{t}
PSO		28.03	26.04	0.00	∞	0.04	0.05	0.24	∞	0.10	0.09	0.00	∞
BBPSO-MC		874.83	38.75	0.00	∞	871.43	40.28	0.00	∞	841.29	53.32	0.00	∞
BBPSOxp-MC		777.25	29.42	0.00	∞	723.32	36.24	0.00	∞	580.79	51.18	0.00	∞
AT-BBPSO-MC													
$df = 1, R^* = 0.1$		0.86	0.09	0.00	∞	0.90	0.07	0.00	∞	0.91	0.06	0.00	∞
$df = 1, R^* = 0.3$		0.02	0.01	0.34	∞	0.02	0.01	0.38	∞	0.03	0.02	0.00	∞
$df = 1, R^* = 0.5$		0.01	0.01	0.54	432.50	0.01	0.01	0.56	457.00	0.01	0.01	0.54	464.50
$df = 1, R^* = 0.7$		0.02	0.02	0.44	∞	0.01	0.02	0.60	413.50	0.04	0.05	0.10	∞
$df = 3, R^* = 0.1$		1.06	0.02	0.00	∞	1.06	0.02	0.00	∞	1.04	0.02	0.00	∞
$df = 3, R^* = 0.3$		0.08	0.03	0.00	∞	0.12	0.04	0.00	∞	0.31	0.13	0.00	∞
$df = 3, R^* = 0.5$		0.01	0.01	0.52	461.00	0.01	0.01	0.72	467.50	0.10	0.09	0.02	∞
$df = 3, R^* = 0.7$		0.02	0.01	0.44	∞	0.09	0.12	0.08	∞	225.75	51.35	0.00	∞
$df = 5, R^* = 0.1$		1.11	0.03	0.00	∞	1.11	0.03	0.00	∞	1.08	0.03	0.00	∞
$df = 5, R^* = 0.3$		0.17	0.05	0.00	∞	0.30	0.11	0.00	∞	0.64	0.14	0.00	∞
$df = 5, R^* = 0.5$		0.02	0.01	0.34	∞	0.02	0.01	0.22	∞	0.73	0.20	0.00	∞
$df = 5, R^* = 0.7$		0.02	0.01	0.30	∞	31.72	19.26	0.00	∞	476.26	46.12	0.00	∞
$df = \infty, R^* = 0.1$		1.46	0.10	0.00	∞	1.47	0.13	0.00	∞	1.22	0.07	0.00	∞
$df = \infty, R^* = 0.3$		0.68	0.13	0.00	∞	0.84	0.09	0.00	∞	0.97	0.05	0.00	∞
$df = \infty, R^* = 0.5$		0.04	0.01	0.00	∞	0.25	0.11	0.00	∞	6.49	6.73	0.00	∞
$df = \infty, R^* = 0.7$		0.88	0.18	0.00	∞	377.73	43.53	0.00	∞	655.88	49.74	0.00	∞
AT-BBPSOxp-MC													
$df = 1, R^* = 0.1$		0.98	0.04	0.00	∞	0.89	0.10	0.00	∞	0.48	0.18	0.00	∞
$df = 1, R^* = 0.3$		0.05	0.02	0.00	∞	0.06	0.03	0.00	∞	0.02	0.02	0.22	∞
$df = 1, R^* = 0.5$		0.01	0.01	0.74	419.50	0.01	0.01	0.66	438.50	0.01	0.01	0.52	482.50
$df = 1, R^* = 0.7$		0.02	0.03	0.58	439.50	0.02	0.02	0.60	437.50	0.02	0.02	0.48	∞
$df = 3, R^* = 0.1$		1.10	0.03	0.00	∞	1.06	0.02	0.00	∞	0.93	0.09	0.00	∞
$df = 3, R^* = 0.3$		0.47	0.10	0.00	∞	0.47	0.12	0.00	∞	0.17	0.07	0.00	∞
$df = 3, R^* = 0.5$		0.04	0.03	0.00	∞	0.07	0.05	0.02	∞	0.04	0.03	0.02	∞
$df = 3, R^* = 0.7$		1.57	1.96	0.00	∞	74.80	39.86	0.00	∞	95.06	41.38	0.00	∞
$df = 5, R^* = 0.1$		1.19	0.06	0.00	∞	1.10	0.03	0.00	∞	1.00	0.05	0.00	∞
$df = 5, R^* = 0.3$		0.70	0.10	0.00	∞	0.71	0.11	0.00	∞	0.28	0.12	0.00	∞
$df = 5, R^* = 0.5$		0.14	0.07	0.00	∞	0.34	0.15	0.00	∞	0.22	0.13	0.00	∞
$df = 5, R^* = 0.7$		116.60	42.47	0.00	∞	306.40	53.78	0.00	∞	261.33	43.65	0.00	∞
$df = \infty, R^* = 0.1$		1.49	0.12	0.00	∞	1.26	0.09	0.00	∞	1.05	0.03	0.00	∞
$df = \infty, R^* = 0.3$		0.96	0.06	0.00	∞	0.93	0.07	0.00	∞	0.64	0.11	0.00	∞
$df = \infty, R^* = 0.5$		0.73	0.12	0.00	∞	0.95	0.07	0.00	∞	0.89	0.15	0.00	∞
$df = \infty, R^* = 0.7$		395.68	33.71	0.00	∞	504.84	39.71	0.00	∞	389.26	40.73	0.00	∞
DI-PSO													
$\alpha = 50, \beta = 1$		60.30	38.12	0.00	∞	7.26	10.66	0.00	∞	1.14	0.61	0.00	∞
$\alpha = 50, \beta = 2$		66.12	33.41	0.00	∞	21.06	15.97	0.00	∞	6.18	4.13	0.00	∞
$\alpha = 50, \beta = 4$		127.46	57.77	0.00	∞	56.24	25.74	0.00	∞	24.65	13.98	0.00	∞
$\alpha = 100, \beta = 1$		28.60	21.21	0.00	∞	1.62	1.71	0.00	∞	0.27	0.24	0.00	∞
$\alpha = 100, \beta = 2$		73.41	67.24	0.00	∞	7.98	8.39	0.00	∞	1.90	1.34	0.00	∞
$\alpha = 100, \beta = 4$		129.76	51.11	0.00	∞	36.80	25.83	0.00	∞	14.49	11.96	0.00	∞
$\alpha = 200, \beta = 1$		25.75	23.08	0.00	∞	0.91	2.29	0.00	∞	0.06	0.07	0.28	∞
$\alpha = 200, \beta = 2$		72.48	45.02	0.00	∞	4.38	6.49	0.00	∞	0.66	0.34	0.00	∞
$\alpha = 200, \beta = 4$		126.67	76.77	0.00	∞	23.59	19.07	0.00	∞	4.98	4.04	0.00	∞
AT-PSO													
$R^* = 0.1$		2.25	1.98	0.00	∞	1.91	1.98	0.00	∞	98.02	57.26	0.00	∞
$R^* = 0.3$		1.09	1.89	0.00	∞	0.02	0.03	0.44	∞	0.05	0.05	0.16	∞
$R^* = 0.5$		5.65	11.07	0.00	∞	0.06	0.08	0.28	∞	0.02	0.02	0.46	∞
$R^* = 0.7$		32.51	39.20	0.00	∞	0.13	0.18	0.14	∞	0.07	0.17	0.26	∞

Table 6: Simulation results for OF5. See text for description.

OF6	Global nbhd				Ring-3 nbhd				Ring-1 nbhd			
Algorithm	Mean	SD	\hat{p}	\hat{t}	Mean	SD	\hat{p}	\hat{t}	Mean	SD	\hat{p}	\hat{t}
PSO	19.59	1.20	0.00	∞	15.41	8.13	0.04	∞	17.89	5.35	0.00	∞
BBPSO-MC	20.16	0.47	0.00	∞	20.24	0.50	0.00	∞	20.49	0.45	0.00	∞
BBPSOxp-MC	20.11	0.19	0.00	∞	20.23	0.16	0.00	∞	20.06	0.11	0.00	∞
AT-BBPSO-MC												
$df = 1, R^* = 0.1$	17.10	6.89	0.00	∞	9.21	9.18	0.00	∞	6.23	8.00	0.00	∞
$df = 1, R^* = 0.3$	18.66	4.77	0.00	∞	15.62	8.27	0.00	∞	13.05	9.54	0.00	∞
$df = 1, R^* = 0.5$	19.05	3.82	0.02	∞	17.29	6.59	0.06	∞	16.41	7.57	0.04	∞
$df = 1, R^* = 0.7$	19.31	3.54	0.00	∞	18.57	4.73	0.00	∞	19.05	3.24	0.00	∞
$df = 3, R^* = 0.1$	19.61	2.52	0.00	∞	18.65	4.81	0.00	∞	17.12	7.06	0.00	∞
$df = 3, R^* = 0.3$	19.82	0.05	0.00	∞	19.87	0.28	0.00	∞	19.08	3.97	0.00	∞
$df = 3, R^* = 0.5$	19.85	0.17	0.00	∞	19.60	2.85	0.02	∞	20.08	0.55	0.00	∞
$df = 3, R^* = 0.7$	19.89	0.22	0.00	∞	20.13	0.40	0.00	∞	20.35	0.42	0.00	∞
$df = 5, R^* = 0.1$	19.99	0.06	0.00	∞	19.75	2.52	0.00	∞	17.90	6.25	0.00	∞
$df = 5, R^* = 0.3$	19.83	0.03	0.00	∞	19.87	0.24	0.00	∞	19.66	2.85	0.00	∞
$df = 5, R^* = 0.5$	19.84	0.09	0.00	∞	19.90	0.25	0.00	∞	20.12	0.62	0.00	∞
$df = 5, R^* = 0.7$	19.83	0.05	0.00	∞	20.00	0.33	0.00	∞	20.33	0.39	0.00	∞
$df = \infty, R^* = 0.1$	20.10	0.07	0.00	∞	20.10	0.09	0.00	∞	20.47	0.29	0.00	∞
$df = \infty, R^* = 0.3$	19.83	0.02	0.00	∞	19.84	0.15	0.00	∞	19.68	2.83	0.00	∞
$df = \infty, R^* = 0.5$	19.83	0.03	0.00	∞	19.92	0.26	0.00	∞	20.10	0.39	0.00	∞
$df = \infty, R^* = 0.7$	19.83	0.03	0.00	∞	19.99	0.33	0.00	∞	20.36	0.38	0.00	∞
AT-BBPSOxp-MC												
$df = 1, R^* = 0.1$	14.13	8.56	0.00	∞	11.18	9.07	0.00	∞	12.71	9.05	0.00	∞
$df = 1, R^* = 0.3$	18.04	6.04	0.00	∞	16.06	7.84	0.00	∞	16.02	7.64	0.00	∞
$df = 1, R^* = 0.5$	18.96	4.84	0.06	∞	18.16	5.26	0.00	∞	17.59	5.83	0.06	∞
$df = 1, R^* = 0.7$	20.17	0.17	0.00	∞	20.04	0.87	0.00	∞	19.60	2.01	0.00	∞
$df = 3, R^* = 0.1$	20.22	0.11	0.00	∞	20.17	0.11	0.00	∞	19.70	2.21	0.00	∞
$df = 3, R^* = 0.3$	19.70	2.82	0.00	∞	20.12	0.15	0.00	∞	19.99	0.16	0.00	∞
$df = 3, R^* = 0.5$	20.17	0.13	0.00	∞	20.20	0.12	0.00	∞	20.02	0.12	0.00	∞
$df = 3, R^* = 0.7$	20.20	0.12	0.00	∞	20.21	0.10	0.00	∞	20.06	0.12	0.00	∞
$df = 5, R^* = 0.1$	20.22	0.10	0.00	∞	20.19	0.11	0.00	∞	20.03	0.12	0.00	∞
$df = 5, R^* = 0.3$	20.08	0.16	0.00	∞	20.15	0.13	0.00	∞	20.03	0.12	0.00	∞
$df = 5, R^* = 0.5$	20.12	0.13	0.00	∞	20.21	0.12	0.00	∞	20.03	0.12	0.00	∞
$df = 5, R^* = 0.7$	20.19	0.12	0.00	∞	20.22	0.10	0.00	∞	20.05	0.11	0.00	∞
$df = \infty, R^* = 0.1$	20.18	0.09	0.00	∞	20.19	0.09	0.00	∞	20.03	0.11	0.00	∞
$df = \infty, R^* = 0.3$	20.03	0.14	0.00	∞	20.09	0.16	0.00	∞	20.02	0.13	0.00	∞
$df = \infty, R^* = 0.5$	20.08	0.14	0.00	∞	20.15	0.13	0.00	∞	20.02	0.09	0.00	∞
$df = \infty, R^* = 0.7$	20.12	0.15	0.00	∞	20.19	0.10	0.00	∞	20.04	0.10	0.00	∞
DI-PSO												
$\alpha = 50, \beta = 1$	19.81	0.52	0.00	∞	18.11	3.34	0.00	∞	18.80	3.52	0.00	∞
$\alpha = 50, \beta = 2$	19.72	1.24	0.00	∞	18.36	2.48	0.00	∞	19.59	1.65	0.00	∞
$\alpha = 50, \beta = 4$	19.85	0.60	0.00	∞	19.54	0.81	0.00	∞	20.03	0.42	0.00	∞
$\alpha = 100, \beta = 1$	19.09	2.07	0.00	∞	15.60	6.01	0.00	∞	18.20	4.53	0.00	∞
$\alpha = 100, \beta = 2$	19.39	1.51	0.00	∞	18.33	3.97	0.00	∞	19.48	2.42	0.00	∞
$\alpha = 100, \beta = 4$	20.18	0.62	0.00	∞	19.71	1.32	0.00	∞	20.22	0.42	0.00	∞
$\alpha = 200, \beta = 1$	18.96	2.05	0.00	∞	14.66	6.67	0.00	∞	17.95	4.74	0.00	∞
$\alpha = 200, \beta = 2$	20.15	0.64	0.00	∞	19.78	1.46	0.00	∞	20.29	0.35	0.00	∞
$\alpha = 200, \beta = 4$	20.32	0.69	0.00	∞	20.17	0.81	0.00	∞	20.42	0.41	0.00	∞
AT-PSO												
$R^* = 0.1$	19.31	2.85	0.00	∞	16.16	6.93	0.00	∞	20.25	0.32	0.00	∞
$R^* = 0.3$	19.58	1.24	0.00	∞	14.64	7.42	0.04	∞	16.28	7.41	0.04	∞
$R^* = 0.5$	19.84	0.67	0.00	∞	19.00	2.47	0.00	∞	19.84	2.12	0.00	∞
$R^* = 0.7$	20.00	0.23	0.00	∞	19.61	1.19	0.00	∞	20.11	0.29	0.00	∞

Table 7: Simulation results for OF6. See text for description.

References

- Andrieu, C. and Thoms, J. (2008). “A tutorial on adaptive MCMC.” *Statistics and Computing*, 18, 4, 343–373.
- Blum, C. and Li, X. (2008). “Swarm Intelligence in Optimization.” In *Swarm Intelligence: Introduction and Applications*, eds. C. Blum and D. Merkle. Springer.
- Chen, Z. and Dunson, D. B. (2003). “Random effects selection in linear mixed models.” *Biometrics*, 59, 4, 762–769.
- Clerc, M. (2010). *Particle swarm optimization*. John Wiley & Sons.
- Clerc, M. and Kennedy, J. (2002). “The particle swarm-explosion, stability, and convergence in a multidimensional complex space.” *Evolutionary Computation, IEEE Transactions on*, 6, 1, 58–73.
- Frühwirth-Schnatter, S. and Tüchler, R. (2008). “Bayesian parsimonious covariance estimation for hierarchical linear mixed models.” *Statistics and Computing*, 18, 1, 1–13.
- Gelman, A., Roberts, G., and Gilks, W. (1996). “Efficient Metropolis jumping rules.” *Bayesian statistics*, 5, 599-608, 42.
- Hsieh, H.-I. and Lee, T.-S. (2010). “A modified algorithm of bare bones particle swarm optimization.” *International Journal of Computer Science Issues*, 7, 11.
- Magnus, J. R. (1988). *Linear structures*. No. 42 in Griffin’s Statistical Monographs and Courses. New York: Oxford University Press.
- Magnus, J. R. and Neudecker, H. (1980). “The elimination matrix: some lemmas and applications.” *SIAM Journal on Algebraic Discrete Methods*, 1, 4, 422–449.

- (2005). *Matrix differential calculus with applications in statistics and econometrics*. 3rd ed. New York: John Wiley & Sons.
- Smith, W. and Hocking, R. (1972). “Algorithm AS 53: Wishart variate generator.” *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 21, 3, 341–345.
- Zhang, H., Kennedy, D. D., Rangaiah, G. P., and Bonilla-Petriciolet, A. (2011). “Novel bare-bones particle swarm optimization and its performance for modeling vapor–liquid equilibrium data.” *Fluid Phase Equilibria*, 301, 1, 33–45.