

1 A practical guide to non-parametric approximate Bayesian computation
2 with improved implementation and error characterization

3 William R. Holmes

4 *Center for Complex Biological Systems and Department of Mathematics, University of California Irvine,*
5 *Department of Mathematics and Statistics, The University of Melbourne, Melbourne, Victoria 3010 Australia -*
6 *(613) 8344 5550.*

7 **Abstract**

8 A critical task in modelling is to determine how well the theoretical assumptions encoded in
9 a model account for observations. Bayesian methods are an ideal framework for doing just this.
10 Existing approximate Bayesian computation (ABC) methods however rely on often insufficient
11 “summary statistics”. Here, I present and analyze a highly efficient extension of a recently pro-
12 posed (Turner et al., 2014) non-parametric approximate Bayesian computation (npABC) algorithm,
13 which circumvents this insufficiency. This method combines Markov Chain Monte Carlo simulation
14 with tools from non-parametric statistics to improve upon existing ABC methods. The primary
15 contributions of this article: 1) A more efficient implementation of this method is described, that
16 substantially improves computational performance and chain mixing. 2) Theoretical results describ-
17 ing the influence of methodological approximation errors on posterior estimation are discussed. In
18 particular, while this method is highly accurate, even small errors have a strong influence on model
19 comparisons when using standard statistical approaches (such as deviance information criterion).
20 Thus care must be taken when using this (or any other ABC) method for model comparison. 3) An
21 augmentation of the standard MCMC procedure, termed “Resampled MCMC”, that reduces the
22 negative influence of approximation errors on performance and accuracy, is presented. 4) In order
23 to make this method accessible to a broader audience, a number of examples of varying complexity
24 are presented along with supplementary code for their implementation.

25 *Keywords:* Non-parametric Approximate Bayesian Computation, Approximate Likelihood,
26 Kernel Density Estimate, Markov Chain Monte Carlo, Linear Ballistic Accumulator Model

Email address: wrholmes@uci.edu (William R. Holmes)

27 **1. Introduction**

28 Numerous models of cognitive processes have been developed over the decades, varying in detail
29 and complexity from neuro-physiological models (Brunel and Wang, 2001; Albantakis and Deco,
30 2009) to high level models (Ratcliff, 1978; Usher and McClelland, 2001; Brown and Heathcote,
31 2008). A central benefit of modelling over experimentation alone is the ability to translate formal
32 theories into testable predictions. Subsequent testing can be performed in a number of ways. In
33 more detailed models, qualitative comparisons between model and data are often made to determine
34 if the theory can account for the qualitative trends observed in data. While such comparisons can
35 quickly refute a theory, they can rarely provide more than weak support. Higher level models on the
36 other hand are typically tractable enough for quantitative comparisons against data. This typically
37 involves “fitting” a model to data and estimating the values of model parameters. Quantitative
38 methods provide two distinct benefits. First, quantitative measures provide a more fine grained
39 account of how well a theory accounts for observations. Second, access to parameter values and
40 their uncertainty provides further inferences about underlying behaviour.

41 Historically, parameters have been estimated using linear programming methods (Dantzig and
42 Thapa, 1997) designed to minimize some statistic such as sum squared error. This frequentist
43 approach however provides no information about estimation uncertainty. For this reason, these
44 methods are often supplemented with some form of sensitivity analysis (Saltelli et al., 2008). More
45 recently however, following algorithmic improvements and computing advances, Bayesian methods
46 have moved to the forefront. These have a number of powerful benefits (Abelson, 2008; Gallistel,
47 2009; Lee, 2008; Lee and Wagenmakers, 2013) which are beyond the scope of this article, but in short
48 they provide a more principled way to account for uncertainty and incorporate prior knowledge.

49 Numerous methods for performing Bayesian analysis have been devised. Standard Markov
50 Chain Monte Carlo (MCMC) techniques (Gelman et al., 2003; Robert and Casella, 2004; Cappé
51 et al., 2004; Del Moral et al., 2006) have proven very powerful, but can only be applied to the
52 simplest problems where a model can be analytically described by a closed form probability den-
53 sity function. Other approximate Bayesian computation (ABC) methods have been developed to
54 circumvent this requirement (Csilléry et al., 2010; Turner and Van Zandt, 2012). Unfortunately
55 they have a number of downsides, particularly that in most cases they do not actually estimate pa-
56 rameters of the desired model (see below for further discussion). Recently, a new method, referred
57 to here as a non-parameteric ABC method (or npABC), was developed (Turner and Sederberg,
58 2014) that bridges the gap between these exact and approximate methods, alleviating some of
59 these issues.

60 The goals of this article are three fold: 1) describe this new methodology in detail and determine
61 its strengths, weaknesses, and limitations, 2) improve upon it, and 3) present it in an accessible
62 way so it can be utilized by a broader audience of end users. In particular, toward this third
63 goal, a number of examples of this method are presented along with documented MATLAB code
64 for two of the examples. This is not intended as a “plug in” software package, but rather to aid
65 implementation of this method by others. Also toward this goal, a detailed procedural overview of
66 this method (with improvements presented here) is provided in Section 4.

67 *1.1. Brief introduction to Bayesian methods*

The canonical Bayesian model estimation problem is to determine the posterior probability distribution of a set of model parameters conditioned on observed data $\pi(\theta|X)$. This typically involves three critical steps. First, a *prior* distribution on the parameters $\pi(\theta)$, which encodes pre-existing knowledge of the system, must be supplied. This is rarely a problem as some form of subjective prior belief is typically available from which a prior can be constructed. Second, a probability density (aka. likelihood) function $L(X|\theta)$, which describes the likelihood that the model will give rise to the observed data given the parameters θ , must be computed. When observations are identical and independently distributed, this reduces to

$$L(X|\theta) = \prod_{i=1}^I L_i(\theta),$$

68 where $L_i(\theta) := L(x_i|\theta)$, so only the likelihood of each individual observation is required. The
 69 following exposition will be confined to this simplified setting. Third, with these model components
 70 provided, the models posterior distribution is computed via Bayes' theorem

$$\pi(\theta|X) = \frac{L(\theta|X)\pi(\theta)}{\int L(\theta|X)\pi(\theta)}. \quad (1.1)$$

71 In almost all practical situations, this third step is impossible to perform analytically since the
 72 required integral is usually not solvable. For this reason, numerous MCMC methods have been
 73 developed to circumvent this third step.

74 In many cases however, the second step is problematic as well since the model may either 1) not
 75 emit an analytic density function or 2) emit one that is too cumbersome to compute. Approximate
 76 Bayesian computation methods have been developed to deal with this difficulty, see (Csilléry et al.,
 77 2010; Turner and Van Zandt, 2012) for existing reviews. Generally speaking, ABC deals with the
 78 absence of a likelihood by prescribing a surrogate measure for how likely or plausible a particular
 79 parameter set (θ) is. To accomplish this, a large number of simulated data observations (\tilde{X}) are
 80 drawn from the model. The observed (X) and simulated (\tilde{X}) data are then compared in some
 81 way to determine how likely that parameter set is. Typically this comparison is accomplished by
 82 compressing both data sets into a set of summary statistics $S(X)$ and then defining a “distance”
 83 between them $\rho(S(X), S(\tilde{X}))$.

84 This method raises two distinct issues. First, a reasonable distance function $\rho(\cdot, \cdot)$ must be
 85 prescribed. A more serious issue however is that the summary statistics must adequately represent
 86 the models output, often referred to as a sufficiency condition. ABC methods do not approximate
 87 the models posterior $\pi(\theta|X)$, but rather a posterior augmented by the choice of S , $\pi(\theta|S(X))$. So
 88 ABC only estimates the posterior distribution of the intended model if $\pi(\theta|S(X)) = \pi(\theta|X)$, which
 89 is often not possible to verify. The essential problem here is that the use of summary statistics
 90 force assumptions on the structure of the underlying likelihood function, which if inaccurate lead to
 91 potentially serious errors (Robert et al., 2011) that no amount of computational effort will correct.
 92 As an extreme example, using mean and variance as summary statistics to describe a distribution
 93 implies a normality assumption, which could be very poor if the underlying model is multimodal or

94 heavily skewed. In a more cognitive context, choice response time distributions are often described
 95 by quantile summary statistics (Heathcote et al., 2002; Ratcliff and Tuerlinckx, 2002; Heathcote
 96 and Brown, 2004). This was however recently shown to be an insufficient summary of the data
 97 (Turner and Sederberg, 2014), leading to substantial posterior inaccuracies.

98 In the broader statistics field, such issues have been overcome through the development of “non-
 99 parametric statistical” methods, which free the user from having to make potentially erroneous
 100 assumptions (e.g. summary statistics) on the structure of their data / model. Recently, non-
 101 parametric methods have been incorporated into the ABC context (Turner and Sederberg, 2014)
 102 to improve Bayesian estimation methods (e.g. npABC). These methods begin the same way as
 103 canonical ABC by first simulating a large number of samples from the underlying distribution.
 104 Next however, they construct an approximation of the underlying likelihood $\hat{L}(X|\theta)$. In this case,
 105 no summary statistics are prescribed and much weaker assumptions on the structure of $L(X|\theta)$ are
 106 made. The approximate likelihood is then substituted ($L \rightarrow \hat{L}$) into the chosen MCMC framework
 107 and an approximate posterior is determined.

108 In the following sections, the implementation details required to apply this method will be
 109 discussed. First, the “kernel density estimate” (or KDE), which is the core of the likelihood
 110 estimation, is described. An improvement of this method that yields substantial efficiency gains
 111 over that in (Turner and Sederberg, 2014) is also provided. Second, the influence of likelihood
 112 estimation errors on posterior estimation and MCMC efficiency will be discussed from a theoretical
 113 perspective. Third, the full npABC will be demonstrated through three examples of increasing
 114 complexity. Using these examples, the strengths and weaknesses of the method will be described.
 115 Finally, a stand alone section (4) describing the core implementation steps of this method is provided
 116 for the user interested primarily in implementing this method.

117 2. Methods

118 2.1. The kernel density estimate

119 The critical step in npABC is the construction of the approximate likelihood $\hat{L}(X|\theta)$ that will
 120 replace L in canonical MCMC estimation. The kernel density estimate (KDE) is a powerful tool
 121 for doing just this (Silverman, 1982, 1986; Epanechnikov, 1969). The first step in this process is
 122 to simulate N_s draws (\tilde{X}) from $Model(\theta)$. This step is of course dependent on the model under
 123 consideration, which will dictate how these samples are produced. The second and final step is to
 124 use kernel density estimation (KDE) to extract likelihood estimates for the observed data X from
 125 the simulated \tilde{X} . For purposes of generality, the KDE process will first be discussed independent
 126 of posterior estimation.

127 The basic problem is to estimate the probability density $f(x_i)$ (a placeholder for $L_i(\theta)$) of each
 128 individual observation x_i from the samples $\tilde{X} = \{\tilde{x}_j\}$, where $j = 1 \dots N_s$. The KDE of this quantity
 129 is given by

$$f(x_i) \approx \hat{f}(x_i) := \frac{1}{N_s} \sum_{j=1}^{N_s} K_h(x_i - \tilde{x}_j). \quad (2.1)$$

From here on the $\hat{\cdot}$ will reference a kernel density estimate of the underlying likelihood, or a quantity

derived from it. Here K_h is a “smoothing kernel” defined by

$$K_h(z) = \frac{1}{h} K\left(\frac{z}{h}\right),$$

where K is a continuous function that is symmetric about $z = 0$ and integrates to 1. The parameter h , commonly referred to as a “bandwidth” size, determines the smoothing properties of the kernel: large h heavily smoothes the sampled data while small h provides less smoothing. To illustrate this, consider the uniform kernel $K(z) = \chi_{[-0.5, 0.5]}(z)$ where χ is the standard indicator function that is one on the prescribed interval and zero elsewhere. This kernel produces a standard histogram estimator with h corresponding to the size of the histogram bins. Histograms with small bins (i.e. small h) of course produce noisy plots while those with large bins produce smoother but less refined plots. See (Silverman, 1986; Epanechnikov, 1969) for a full review of KDE theory.

There are a few critical properties of the KDE estimator relevant to this discussion. \hat{f} is an approximation to f and as such can be thought of as an estimator with some underlying distribution. For the standard class of first order kernel functions (e.g. biweight, Gaussian, Epanechnikov, etc.), this distribution is approximately normal with intrinsic bias and variance

$$\text{Bias}(\hat{f}(x)) \approx \frac{h^2}{2} f''(x) M_2(K), \quad \text{Var}(\hat{f}(x)) \approx \frac{1}{N_s h} f(x) \|K\|_2, \quad (2.2)$$

where $M_2(K)$ and $\|K\|_2$ denote the second moment and Euclidean (or L^2) norm of K respectively (Silverman, 1986). From these estimates we see that the bandwidth, number of samples, and specific choice of K all affect accuracy of this approximation. In practice, the specific choice of K has only marginal effects on accuracy, though the Epanechnikov kernel is known to minimize mean integrated square error (Epanechnikov, 1969). The bandwidth (h) and number of samples (N_s) however are of critical importance and will each have different effects on the posterior estimation process. For now, simply note that N_s plays a role in variance control while h modulates a classic bias-variance tradeoff.

2.1.1. An improved, more efficient KDE implementation

Before discussing the influences of this approximation procedure on posterior estimation, I will discuss a technical improvement on the classic implementation of KDE that substantially speeds computational implementation of npABC. Before continuing with this section however, note that it is technical in nature. This improved procedure has the same accuracy, the same bias / variance issues as the standard KDE procedure, and will lead to the same results as the standard KDE when embedded into the npABC procedure. This improvement does however substantially improve efficiency for reasons that will be discussed.

Direct computation of \hat{f} from Equ. (2.1), while simple, is inefficient. Given a set of N_d observations, the kernel function must be evaluated $N_s \cdot N_d$ times. While a single evaluation of this size is reasonable, this becomes a computational bottle neck in MCMC applications where the likelihood must be evaluated at many chain iterations. The improvement presented here, first proposed in (Silverman, 1982), takes advantage of the observation that the KDE formula in Equ. (2.1) resembles

163 a convolution. The discrete model samples \tilde{X} can be represented by the following function

$$d(x) = \frac{1}{N_s} \sum_{j=1}^{N_s} \delta_{\tilde{x}_j}(x), \quad (2.3)$$

164 where δ is the Dirac delta function. It is then direct to show that

$$d \star K_h(x) = \frac{1}{N_s} \sum_{j=1}^{N_s} K_h(x - \tilde{x}_j), \quad (2.4)$$

165 where \star denotes the standard convolution. This is precisely the KDE formula in Equ. (2.1). The
 166 KDE thus resembles a canonical smoothing operation (with K as the smoother), proposed as early
 167 as 1944 in partial differential equations literature (Friedrichs, 1944).

168 While convolutions are well know to be intensive to compute directly, this burden can be
 169 greatly reduced by making use of techniques from signal processing theory, where convolutions are
 170 common. The ‘‘convolution theorem’’ states that $\mathcal{F}(f \star g) = \mathcal{F}(f) \cdot \mathcal{F}(g)$, where \mathcal{F} is the continuous
 171 Fourier transform. Since multiplication is much more efficient than convolution, the basic idea
 172 of this method is to transform both d and K_h into the spectral domain, multiply in the spectral
 173 domain (which effectively is the convolution), then transform back. Given the high efficiency of Fast
 174 Fourier Transform methods (FFT), transferring to and from the spectral domain is fast relative to
 175 the convolution. This was originally proposed as an efficient method for generating a high resolution
 176 PDF on a regular grid, particularly for plotting purposes. Likelihood values of observations can
 177 however readily be interpolated from this regular grid.

178 There is a technical point that must be addressed before applying this method; the FFT is only
 179 efficient if the data being transformed is on an regular grid, which is not the case for the samples
 180 $\{\tilde{x}_j\}$. To circumvent this, the samples should first be binned to a very fine grid with 2^n points (a
 181 power of 2 is used for technical reasons related to FFT efficiency). This grid should be much more
 182 finely spaced ($n > 8$ typically) than a typical histogram grid for reasons discussed in a moment.
 183 The improved FFT based KDE procedure is then as follows:

- 184 1. Bin the simulated samples to a very fine grid, $d \rightarrow \tilde{d}$.
- 185 2. Transform the resulting data to the spectral domain ($\tilde{d}(x) \rightarrow \mathcal{F}[\tilde{d}](s)$) using a FFT (where
 186 $\mathcal{F}[\tilde{d}](s)$ is the contribution of wave number s , i.e. the frequency spectrum of \tilde{d}).
- 187 3. Carry out the convolution operation in the spectral domain

$$\mathcal{F}[\tilde{d} \star K_h](s) = \mathcal{F}[\tilde{d}](s) \cdot \mathcal{F}[K_h](s). \quad (2.5)$$

- 188 4. Using an inverse FFT, transform the resulting expression back to obtain the likelihood esti-
 189 mate on the same 2^n grid

$$\hat{f} = \mathcal{F}^{-1} \left(\mathcal{F}[\tilde{d}] \cdot \mathcal{F}[K_h] \right) \quad (2.6)$$

- 190 5. Interpolate the density from this grid to the observed data points to obtain $\hat{f}(x_i)$. Linear
 191 interpolation should be used here since higher order methods (such as cubic splines) can
 192 generate negative likelihood values in the tail of a distribution.

193 Since FFT, multiplication, and interpolation are each highly efficient and usually optimized within
194 programming languages, this procedure is vastly more efficient than direct computation of the
195 convolution.

196 A few notes about this procedure are in order. First, the FFT itself introduces some errors
197 into the approximation, but these are orders of magnitude smaller than the primary error sources.
198 Second, the binning and interpolation steps will introduce errors as well. However these will
199 again be very small provided $n > 8$ ($n = 10$ is used in all following applications). Third, in
200 principal, any kernel (K) can still be used in this process. However the canonical Gaussian kernel
201 is particularly useful in this case since its Fourier transform is another Gaussian, $\mathcal{F}[K_h](s) \propto$
202 $\exp(-0.5h^2s^2)$. Fourth, one must be careful when applying FFT's since various operations (scalings,
203 shifts, etc.) must be performed to correctly prepare the data. Such details are not mentioned here
204 as they are specific to programming language and the FFT implementation being called. Instead,
205 supplementary files that demonstrate implementation in MATLAB are provided.

206 It also interesting to note that this view of the kernel density estimate is somewhat of a math-
207 ematical departure from the original view. In its original form, the KDE was essentially designed
208 as an extension of the histogram to be a method of pooling information from nearby samples in a
209 weighted manner to make a more accurate density estimate. This procedure however more closely
210 resembles filtering of a noisy signal to determine the underlying "true" trend. In this way, the
211 transformed function $\mathcal{F}[K_h]$ acts as a low pass filter in the spectral domain that attenuates high
212 frequency noise. To illustrate this, and more generally how this procedure works, consider the
213 following simple example.

214 *2.1.2. Example 1: Reconstructing a Gaussian distribution*

215 In this section, the KDE procedure is demonstrated through a simple example. We begin with a
216 known normal distribution with known mean ($\mu = 5$) and variance ($\sigma = 1$). Next, this FFT based
217 KDE procedure is used to reconstruct the underlying distribution from $N_s = 10,000$ simulated
218 draws. The binning of these N_s observations to 2^{10} grid points yields a very noisy distribution
219 (Figure 1a, grey). Application of the spectral filter (e.g. the smoothing step) attenuates the
220 high frequency noise, revealing a smoothed normal distribution that agrees well with the exact
221 distribution (Figure 1a, black).

222 To test the accuracy of log likelihood estimation, which is critical in MCMC applications, a fixed
223 set of $N_d = 1,000$ "observations" from the known normal are drawn as a synthetic data set. Then,
224 100 independent reconstructions of this normal are performed, with the resulting approximate log
225 likelihood computed. Results show the mean error in this example is 0.3% with a maximum error
226 of 0.8%.

227 This example also demonstrates the critical effect of the bandwidth parameter (h) on the like-
228 lihood construction (Figure 1b). Recall that small bandwidths produce a less biased but higher
229 variance estimate while larger bandwidths produce lower variance but higher bias. These results
230 confirm this tradeoff and indicate the source of the increased bias at higher bandwidths. Specifically,
231 with larger bandwidths, the peak of the distribution is attenuated and the tails are overestimated,
232 due to over-smoothing. Essentially mass from the peak of the PDF is transferred to the tail in
233 the smoothing process. For this reason, the bandwidth must be chosen carefully so that it is small
234 enough to account for the most refined feature of the model / data but still large enough to pro-

235 duce a reliable estimate. When likelihood distributions are nearly normal, automated bandwidth
 236 determination methods can choose nearly optimal values (Silverman, 1986), however these auto-
 237 mated methods can lead to poor results when distributions are more complicated (multi-modal for
 238 example, see the next example).

239 2.2. npABC: Incorporating the KDE into ABC

240 The above sections describe a manner of approximating the likelihood $L(X|\theta)$. This however is
 241 only an approximation and it is important to ask the question, how do inevitable errors influence
 242 the MCMC procedure and posterior estimation. In this section, I will discuss critical points that
 243 must be considered when embedding KDE in a Bayesian MCMC framework. The primary results
 244 of this section are as follows. 1) Small likelihood estimation errors will inevitably propagate into
 245 small posterior errors. While this will in many cases have little effect on parameter estimation,
 246 these small errors can have dramatic effects on hypothesis testing and model comparison through
 247 the use of AIC, BIC, or DIC. 2) Variance in the KDE approximation influences Metropolis-Hastings
 248 acceptance probabilities in a manner that substantially degrades MCMC chain mixing. While a
 249 rigorous characterization of these points in a general setting is likely very difficult and beyond the
 250 scope of this article, I will first outline the theoretical reasoning behind each and subsequently
 251 demonstrate them through simple examples.

252 2.2.1. The influence of KDE on likelihood estimation

253 The likelihood $L(X|\theta)$ and choice of priors fully determine the posterior distribution for a model.
 254 In the context here however, we only have access to the approximate likelihood obtained as

$$\hat{L}(X|\theta) = \prod_{i=1}^{N_d} \hat{L}_i(\theta), \quad (2.7)$$

255 which is itself a stochastic quantity. Define the estimation error for the likelihood of observation i
 256 as

$$\epsilon_i = \hat{L}_i - L_i, \quad (2.8)$$

257 where for brevity, the dependence on θ has been omitted. The following relation then connects the
 258 approximate and true likelihood

$$\hat{L} := \prod_{i=1}^{N_d} \hat{L}_i = L \prod_{i=1}^{N_d} \left(1 + \frac{\epsilon_i}{L_i}\right). \quad (2.9)$$

259 From Equ. (2.2), we know that $1 + \epsilon_i/L_i \sim N(1 + \mu_i, \sigma_i)$ where

$$\mu_i = \frac{h^2}{2} M_2(K) \frac{L_i''}{L_i}, \quad \sigma_i^2 = \frac{\|K\|_2}{N_s h}, \quad (2.10)$$

260 and $L_i'' = L''(x_i|\theta)$. Using basic facts about normal distributions, we know that the product of
 261 independent normals is again normal so that

$$\frac{\hat{L}}{\bar{L}} \sim N(\mu_{1\dots N_d}, \sigma_{1\dots N_d}), \quad (2.11a)$$

262 where

$$\mu_{1\dots N_d} = 1 + \langle \mu_i \rangle, \quad \sigma_{1\dots N_d}^2 = \frac{\|K\|_2}{N_d N_s h}, \quad (2.11b)$$

263 and $\langle \mu_i \rangle$ is the mean of the set $\{\mu_i\}$.

264 Unfortunately these quantities are not rigorous quantitative estimates of the mean and variance
 265 of the likelihood ratio and cannot be used to make post hoc error estimates. A critical assumption in
 266 this derivation was that $\{\epsilon_i\}$ are uncorrelated, which is not the case here since the density estimation
 267 of two nearby points will pool information from common samples. Nonetheless, these expressions
 268 provide valuable insights into the scaling behavior of the bulk mean and variance. In particular,
 269 the same bias variance tradeoff we saw in the estimate of an individual likelihood appears in the
 270 estimate of the full likelihood. So again, a tradeoff must be made between accuracy and precision
 271 of the estimate. Similarly, increasing the number of samples N_s drawn for estimation does not
 272 improve accuracy, but instead decreases the estimation variance and improves precision.

273 2.2.2. Influence of the KDE on model comparison statistics

274 Most common measures used for model comparison (except Bayes Factors) utilize values of the
 275 log likelihood $LL(X|\theta)$ (AIC, BIC, DIC). Thus, likelihood estimation errors will influence these
 276 model comparison statistics. To see this, note that

$$\hat{LL}(X|\theta) - LL(X|\theta) = \sum_{i=1}^{N_d} \hat{L}(x_i|\theta) - LL(x_i|\theta) \approx \sum_{i=1}^{N_d} \frac{\epsilon_i}{L_i}, \quad (2.12)$$

277 where ϵ_i, L_i are as above and the approximation results from a Taylor expansion of the log function
 278 near 1 (assuming ϵ_i/L_i is small). We thus see that small relative errors in the approximate likelihood
 279 of each individual observation translate directly into small relative errors in the log likelihood.

280 While these errors are small in a relative sense, they raise a substantial problem from a model
 281 comparison standpoint. Model comparison tests will often use differences in these statistics of as
 282 little as $\Delta DIC = 10$ to conclude evidence for or against a model. However, log likelihoods and in
 283 turn DIC values are often on the order of 100 – 1000. So these small relative errors in LL can easily
 284 be larger than the differential commonly taken as “significant”. For this reason, care should be
 285 taken when using these measures for model comparison in this or any context where approximations
 286 are used. Examples in subsequent sections will further elucidate this issue.

287 Also note that, in contrast to estimates in the previous section, this estimate provides a quanti-
 288 tative approximation of the estimation bias being made. In particular, from Equ. (2.12) it is direct
 289 to show that

$$E \left(\hat{LL}(X|\theta) - LL(X|\theta) \right) \approx \sum_{i=1}^{N_d} \frac{\mu_i}{L_i} = \frac{h^2}{2} M_2(K) \sum_{i=1}^{N_d} \frac{L_i''}{L_i}, \quad (2.13)$$

290 which can be used for post hoc estimates of the expected error. Unfortunately an estimate of the
 291 underlying variance cannot be obtained, again because $\{\epsilon_i\}$ are not independent. However, if such
 292 an estimate is required, both the bias and variance of this estimate for any particular parameter
 293 set (the mean of the posterior for example) can easily be assessed through repeated simulation.

294 2.2.3. Influence of the KDE on rejection rates

295 A critical quantity in any MCMC procedure is the Metropolis-Hastings probability (α^n) of
 296 accepting a parameter set θ^n at the n^{th} chain iteration, defined as

$$\alpha^n = \frac{L(X|\theta^n)\pi(\theta^n)}{L(X|\theta^{n-1})\pi(\theta^{n-1})}. \quad (2.14)$$

297 Given this procedure provides only an estimate of $L(X|\theta^n)$, this quantity will be stochastic as well.
 298 While a simple description of the density for this quantity is not available, some intuition into the
 299 influence of the KDE on it is possible. Define $r^n = (\hat{L}^n - L^n)/L^n$ to be the relative sample error
 300 between the exact and approximated likelihood values. Note this quantity is distinct from ϵ_i , which
 301 is the absolute error in the approximate likelihood of observation i , where this is the relative error
 302 in the full approximate likelihood. Then it is direct to show that

$$\frac{\hat{\alpha}^n}{\alpha^n} = \frac{1 + r^n}{1 + r^{n-1}}, \quad (2.15)$$

303 where $\hat{\alpha}$ is obtained by substituting the likelihood estimate for the true value.

304 This raises the practical problem for MCMC efficiency. Suppose \hat{L}^n is an overly optimistic
 305 estimate of the likelihood of L^n so that $r^n > 0$. This over estimate will increase the probability
 306 of accepting this particular chain iteration. However, since this over estimate enters into the
 307 acceptance probability of subsequent chain iterations, it will reduce the acceptance probability of
 308 every subsequent iteration of that chain until a new parameter set is accepted. If r^n is fairly near
 309 0, this will likely only have a marginal effect on the time spent by that chain at the current state.
 310 However, if it departs significantly from 0, a significant number of chain iterations will be required
 311 to displace it, causing chains to stagnate and impairing exploration of parameter space.

312 This issue is exacerbated by the fact that the next acceptance in that chain is likely to result
 313 from yet another over estimation. Generally speaking, this process will lead to a net increase
 314 in r^n as n increases due to acceptance being influenced by the variance. Underestimates will
 315 rarely be accepted and quickly discarded, while overestimates are more likely to be accepted and
 316 rarely discarded. Performance will thus be further degraded and posterior estimates might become
 317 skewed. This issue will be demonstrated through practical examples in subsequent sections and an
 318 augmentation of the standard MCMC to correct this deficiency will be discussed.

319 3. Results

320 The capabilities of this method will next be demonstrated through three examples of increasing
 321 complexity. In the first, a bimodal posterior will be fit to demonstrate the substantial effects the
 322 bandwidth parameter h can have on accuracy. Second, a canonical response time model in decision

323 making will be fit using this method. Finally, a third example will be presented where standard
 324 ABC methods based on summary statistics are insufficient for posterior estimation.

325 3.1. Example 2: Fitting a mixture of Gaussians distribution

326 In this first example, the npABC algorithm is used to estimate the posterior of a mixture model

$$X \sim (1 - p)N(\mu_1, \sigma) + pN(\mu_2, \sigma), \quad (3.1)$$

327 where N indicates the normal distribution and p is a weighting parameter indicating the probability
 328 that an observation is derived from the normal centred at μ_2 . Mixture models are common in a
 329 number of applications where transitioning between discrete strategies might occur. For example
 330 participants might transition between actively participating in a task and simply guessing or may
 331 follow instructions in some instances while disregarding them in others (Cassey et al., 2014; Ollman,
 332 1966; Yellott Jr, 1967, 1971; Yantis et al., 1991; Vandekerckhove and Tuerlinckx, 2007).

333 We begin by creating a data set with $N_d = 1,000$ simulated observations drawn from this
 334 distribution with $p = 0.6, \mu_1 = -6, \mu_2 = 4, \sigma = 1$. Next, the following prior distributions for the
 335 parameters are prescribed

$$p \sim U(0, 1), \quad \mu_1 \sim U(-10, 0), \quad \mu_2, \sigma \sim U(0, 10), \quad (3.2)$$

336 where $U(a, b)$ indicates the uniform distribution on the interval $[a, b]$. For a simple model such as
 337 this, any standard MCMC procedure should be sufficient. Since subsequent examples require
 338 more sophisticated techniques though, a differential evolution MCMC procedure (DE-MCMC)
 339 (Ter Braak, 2006; Storn and Price, 1997; Turner et al., 2013) is used for consistency. For all
 340 simulations of this model, 15 chains are propagated for 500 burn in iterations followed by 2000
 341 recorded iterations. All that is left now is to specify the density estimation parameters N_s and h .
 342 Rather than specify a single set of KDE parameters, different combinations of N_s and h are used
 343 to determine the influence of these parameters on results.

344 3.1.1. Example 2: Results

345 As discussed previously, h mediates a bias variance tradeoff. To determine how this parameter
 346 effects posteriors, the model is fit to data for different values of h . Figure 2a shows the posterior
 347 distribution for (μ_1, μ_2) for two values of h , where *Silv* indicates the value derived by “Silverman’s
 348 rule of thumb” (Silverman, 1986)

$$h = 1.06\bar{\sigma}N_s^{-0.2}. \quad (3.3)$$

349 Here, $\bar{\sigma}$ is the sample variance of the data, and for $N_s = 10,000$ is $h = 0.67$. These results show
 350 that posterior estimates are visually identical for the two values. This however is misleading. To
 351 investigate the influence of h further, 1) the quality of posterior model fit to the data and 2) the
 352 computed log likelihood were determined by comparing to the analytic solution, Figure 2b.

353 In this figure, the mean parameter values from the posterior for each value of h were used
 354 to construct the likelihood. These values are effectively identical, differing by $< 0.1\%$ between
 355 the two simulations. Yet, when the model’s PDF is constructed from these parameter sets, we
 356 see significant deviation of the $h = \textit{Silv}$ case from the analytic solution while $h = 0.2$ faithfully

357 captures the analytic solution. This results from over-smoothing of the density in the $h = Silv$
358 case, since the underlying distribution is bimodal. Essentially, the smoothing kernel is too broad, so
359 when it is applied to smooth the simulated PDF, the peak of the distribution is attenuated leading
360 to fatter tails.

361 These results also show the choice of h influences log likelihood estimates. While the difference
362 between the computed and actual log likelihood is small in a relative sense ($\sim 2\%$ for $h = Silv$), it
363 is still relatively large in an absolute sense (~ 38). This raises a substantial problem for hypothesis
364 testing and model comparison based on AIC, BIC, and DIC, all of which rely on log likelihoods. It
365 is important to remember however that this is a problem with ABC in general since the posterior
366 being estimated is always only an approximation of the actual posterior.

367 To further determine the extent to which small posterior errors influence these statistics, DIC
368 was computed for different pairings of (N_s, h) . For each, 100 independent posterior estimation
369 simulations were performed (all on the same synthetic data set), and the mean and standard
370 deviation of the DIC computed from those simulations is shown, Figure 2c. While even the worst
371 model fit ($N_s = 5,000, h = 0.8$) leads to a relative DIC error of $< 2\%$, the resulting absolute DIC
372 error is > 200 . Given that DIC differences between different models of as little as $\Delta DIC = 10$
373 is often taken as strong evidence for a particular model, clearly these small relative errors can
374 overwhelm standard hypothesis tests.

375 3.1.2. *The influence of the bias-variance tradeoff*

376 These results (Figure 2c) also further illustrate the influence of h (and the bias variance tradeoff
377 it mediates) on results. Broadly speaking, as h increases, error in the DIC increases as well, largely
378 independent of N_s which has no influence on estimation bias. Further, as h decreases, the DIC
379 error decreases while the DIC variance increases. This is consistent with the fact that increasing
380 h reduces estimation variance but increases estimation bias. Thus the MCMC procedure does not
381 abrogate this tradeoff and h should be chosen carefully. Unfortunately there is no universal way
382 of choosing this value and for example, the Silverman value is usually too large for multimodal
383 distributions and too small for heavy tailed distributions. Thus some trial and error is required for
384 choosing this bandwidth.

385 There is one last point to consider here. Recall from the previous section that variability in the
386 likelihood estimation is hypothesized to impair MCMC performance by causing chains to get stuck
387 when a likelihood is significantly overestimated. To determine the extent of this problem, proposal
388 acceptance rates as a function of N_s and h are computed, Figure 2d. Again, 100 independent
389 simulations of the posterior are used for each KDE parameter set. Results show a clear decrease in
390 the acceptance rate as the number of samples N_s decreases, consistent with the supposition that
391 likelihood variability leads to poor MCMC performance. One way to ameliorate this issue is to
392 simply increase the number of samples used for estimation. In many cases however this will not be
393 possible for performance reasons. In the next example, an alternative correction that ameliorates
394 this performance issue is presented.

395 3.2. *Example 3: Fitting the Linear Ballistic Accumulator (LBA)*

396 In this example, the canonical Linear Ballistic Accumulator (LBA) model (Brown and Heath-
397 cote, 2008) is considered as an example of a large of class of evidence accumulation models in

398 decision making literature. A number of accumulator models, including Ratcliff’s drift diffusion
399 model (Ratcliff, 1978; Ratcliff and Rouder, 1998), the leaky competing accumulator (Usher and
400 McClelland, 2001), the ballistic accumulator (Brown and Heathcote, 2005), and decision field the-
401 ory (Busemeyer and Townsend, 1993), have been developed over the years to account for different
402 aspects of decision making. What differentiates the LBA from the other models, is that evidence
403 accumulation for different choice alternatives are independent, linear, and deterministic. This sim-
404 plicity allows for a closed form solution for the simplest settings. However as will be discussed in
405 the next example, even simple variations of this model make it impossible to obtain a tractable
406 likelihood function. This example will thus be used to demonstrate the potential power of these
407 methods in response time modelling. A brief description of this model will be provided and the
408 interested reader can find further details in (Brown and Heathcote, 2008).

409 The basic assumptions of the LBA are that following the presentation of information, evidence
410 for each of a set of choice alternatives accumulates linearly and deterministically until an evidence
411 threshold b is reached. The rate of evidence accumulation for choice alternative i , given by v_i , is
412 assumed to be fixed within a trial (this is the deterministic assumption) but to vary among trials.
413 This rate is sampled from an underlying normal distribution $v_i \sim N(\mu_i, \sigma)$, while the start point
414 $x_{0,i}$ for the i^{th} accumulator, which is also assumed to vary across trials, is uniformly distributed
415 $x_{0,i} \sim U(0, A)$. Additionally, a non-decision time τ_{er} is included to account for encoding and motor
416 response delays. This simplest LBA variant is thus fully parameterized by the parameters b, A, σ, τ_{er}
417 and the collection of mean drift rates $\{\mu_i\}$. The likelihood $L(c_i, \tau_i | \theta)$ of a option c_i being chosen at
418 time τ_i can then be described by an analytic function (Brown and Heathcote, 2005).

419 This method will be used to perform parameter recovery, as was done previously in (Turner and
420 Sederberg, 2014), and assess the properties of this method in a cognitive modelling context. To
421 begin, a synthetic data set for a two choice experiment, consisting of $N_d = 1,000$ observations, is
422 created by simulating N_d trials with $A = 1.6, b = 2.7, \mu_1 = 3.4, \mu_2 = 2.1, \tau_{er} = 0.1$. The canonical
423 assumption $\sigma = 1$ is further made to identify the model. To place the model in a Bayesian
424 framework, the following priors on the parameters are further prescribed

$$b, A \sim U(0, 10), \quad \mu_1, \mu_2 \sim U(-10, 10), \quad \tau_{er} \sim U(0, 1). \quad (3.4)$$

425 The same differential evolution MCMC procedure used previously is used here as well, again with
426 15 chains, a 500 iteration burn in, and 2000 recorded chain iterations.

427 3.2.1. Example 3: Results

428 The posterior of this model is fit both analytically and using this FFT based npABC procedure.
429 Again, for each combination of (N_s, h) , 100 independent fits are performed to determine how
430 estimation variability influences various quantities, Figure 3. For all but the largest value of h ,
431 the posterior estimated by the two methods was visually indistinguishable, and so they are not
432 shown. Panel *a* shows the quality of fit for the two methods, analytic MCMC and npABC (using
433 $N_s = 10,000$ and $h = Silv$ where *Silv* again indicates $h = 0.028$ was chosen according to Silverman’s
434 rule of thumb). In each case, the mean value of the parameters from the associated posterior were
435 determined and the PDF was constructed from those values. The resulting PDF’s are virtually
436 indistinguishable and in this case, the log likelihoods are very close.

437 Again however, we see that small errors propagate into the DIC measure, Figure 3b. For all
438 but the worst case fits ($h = 0.07$), the relative DIC error is $\sim 1 - 2\%$. This translates into absolute
439 errors of $\Delta DIC \sim 10 - 20$, which will again have a strong influence on model comparison. We
440 again see that N_s has effectively no influence on the DIC error. This along with results from the
441 previous example confirms that errors cannot be reduced by increasing the number of samples used
442 in the likelihood reconstruction. Only reductions in h can improve posterior estimates.

443 We also again see a problem with acceptance rates, Figure 3c, which generally decrease when
444 either N_s or h decreases (black dots). Further recall that both of these parameter changes lead to
445 increased likelihood estimation variance. These results are thus again consistent with the fact that
446 increased estimation variance reduces acceptance rates and hence MCMC efficiency. In particular,
447 the acceptance rate for this procedure with $N_s = 10,000, h = Silv$ (which are the same estimation
448 parameters used in (Turner and Sederberg, 2014)) is only $\sim 6\%$. Fortunately, this can be ameliorated
449 to a significant extent with a minor augmentation of the MCMC procedure and a little more
450 computation.

451 3.2.2. Resampled MCMC for npABC

452 The central problem that leads to chain stagnation and poor performance is that the likelihood
453 of a particular parameter set θ^n can, on a rare occasion, be significantly overestimated. While this
454 will be rare, it will substantially degrade performance. This overestimation will increase the chance
455 of that parameter set being accepted. Subsequently, acceptance probabilities will significantly favor
456 keeping that state on further MCMC chain iterations. A simple way to “unstick” chains that become
457 stagnant for this methodological reason is to simply resample that likelihood value frequently. This
458 will of course increase computational cost, but it will ensure that no chain becomes stuck due to
459 mis-estimation of the likelihood. This will have the additional benefit of reducing contamination
460 of the posterior by oversampling less likely parameters. From here on, this augmented MCMC will
461 be referred to as a “resampled MCMC”.

462 This adjustment was added so that the likelihood of every chain is resampled every three chain
463 iterations, Figure 3c (grey dots). That is for each of the N_c chains, the likelihood of the current state
464 of that chain is resampled every third MCMC iteration, independent of the history of the chain.
465 Why was this frequency chosen? It is well established that the theoretical acceptance rate for this
466 form of MCMC is $\sim 25\%$ for five or more parameters. The resampling rate was chosen to be faster
467 than the theoretical frequency of chain movement. Results show this augmentation substantially
468 improves acceptance rates, increasing them to $\sim 17 - 18\%$. Furthermore, the resulting acceptance
469 rate is only weakly dependent on N_s and h , suggesting the effects of variance on performance have
470 been removed. The exception to this is that for large $h = 0.07$, there is a substantial drop in
471 performance. It is unclear what is causing this, but this value is well above any reasonable choice
472 for h .

473 3.2.3. A note on performance

474 A brief note is in order regarding practical performance of this algorithm. To assess performance,
475 both the standard MCMC with the analytic LBA likelihood and the npABC algorithm with the
476 resampled MCMC were timed for a single posterior fit. To obtain consistent timings, all but one
477 active computational core on a Mac Pro computer were turned off. Results show it takes ~ 137

478 seconds for the npABC algorithm to complete while it takes ~ 315 seconds for the standard MCMC
479 to complete. Thus the npABC with resampled MCMC and FFT based density estimation is more
480 than twice as fast as the standard MCMC using the analytic likelihood.

481 Profiling of the codes shows the primary reason for this is that computation of the cumulative
482 density function (CDF) of the normal distribution, which is required to evaluate the LBA likelihood,
483 is cumbersome. While we have a tendency consider functions such as the normal CDF to be
484 “analytic”, the term analytic has little meaning in computational settings. In fact, a complex
485 numerical procedure is required to approximate the normal CDF, which in this application is
486 slower than Monte Carlo sampling of the full likelihood. To be clear, this is not meant to advocate
487 for abandoning the standard methods, since accuracy should always be favored over efficiency when
488 reasonable. Rather, this is intended to demonstrate that for many types of models, with proper
489 coding techniques, computational cost can be very reasonable.

490 3.3. Example 4: Fitting the piecewise Linear Ballistic Accumulator (pLBA)

491 In this final example, a piecewise LBA type model will be considered. The canonical LBA
492 describes decision process that might be described as stationary in the sense that the information
493 available to the decision maker remains the same over time. In many cases however, information
494 may change during the course of the decision process. In (Huk and Shadlen, 2005; Kiani et al.,
495 2008; Thura et al., 2012; Tsetsos et al., 2012; Winkel et al., 2014) for example, a random dot motion
496 paradigm where the direction motion of dots change at discrete times during the course of individual
497 trials was utilized. In these cases, the information itself is non-stationary and one would expect
498 the decision process to change in response to the new information. To account for this, a piecewise
499 variant of the standard LBA was first presented in (Holmes et al., 2014). This non-stationary
500 model, which has no tractable closed form likelihood function, will be used to demonstrate this
501 method in a context where existing methods are insufficient.

502 Briefly, to account for the changes in information, this model makes two assumptions on top
503 of those of the standard LBA. First, that changes in information influence the rate of evidence
504 accumulation so that the rates prior to the change are $v_i \sim N(\mu_{vi}, \sigma)$ while those after the change are
505 $w_i \sim N(\mu_{wi}, \sigma)$. This model is referred to as “piecewise LBA” since evidence accumulation is linear
506 and deterministic on each of two segments corresponding to the two separate pieces of information.
507 Second, there is some delay (t_{delay}) between onset of new information and its incorporation into the
508 decision process, which is assumed fixed across trials. In the context of a two choice decision, after
509 setting $\sigma = 1$, the model is fully described by the eight parameters $A, b, \mu_{v1}, \mu_{v2}, \mu_{w1}, \mu_{w2}, t_{er}$, and
510 t_{delay} . See Holmes et al. (2014) for further details. To begin, a data set consisting of $N_d = 1,000$
511 observations is created, assuming $A = 1.6, b = 2.7, \mu_{v1} = 3.4, \mu_{v2} = 2.5, \mu_{w1} = 1.5, \mu_{w2} = 3.6, t_{er} =$
512 $0.1, t_{delay} = 0.3$.

513 While this model is simple to describe, it does not have an analytic description. Nonetheless,
514 the methods described here can be applied to this model without too much augmentation. In fact,
515 the density estimation procedure itself is identical to that used in the previous examples. The only
516 changes that are required for this application are entirely in the resampled MCMC procedure itself.
517 In this example, KDE parameters $h = 0.02, N_s = 10,000$ will be used with no other changes to the
518 likelihood approximation. A slightly more complex MCMC procedure must however be used. The
519 DE-MCMC procedure will again be used, this time with 24 chains. However, since the size of the

parameter space has increased (8 parameters), a blocked variant must be used to improve sampling performance. Here, as in Holmes et al. (2014), the parameters will be grouped into two blocks: $(A, b, \mu_{v1}, \mu_{v2}, t_{er})$ which describe the accumulation process prior to the change of information, and $(\mu_{w1}, \mu_{w2}, t_{delay})$ which describe the process after the change.

3.3.1. Example 4: Results

Estimated posteriors for this piecewise LBA model are shown in Figure 4a for all model parameters. At first glance, it may appear the method has performed poorly since the posteriors are quite broad. This is not however the case. First, the mean parameter set from these posteriors provides a good fit to the data, Figure 4c. Second, it is well known that the LBA model exhibits significant parameter correlations, which commonly lead to poorly localized posteriors. In biological and physics literature, this is commonly referred to as a “sloppy model” (Gutenkunst et al., 2007; Apgar et al., 2010) since the likelihood is nearly unchanged over a wide range of parameters, Figure 5. To confirm parameter correlations are the source of this posterior spread, principal component analysis (PCA) was performed on the saved MCMC chain data. This reveals that the first and second principal components account for $\sim 92\%$ and 5% of the variability in the posterior respectively. Furthermore, the eigenvector of the principal component shows this correlated direction involves only the pre-switch model parameter A, b, μ_{v1}, μ_{v2} .

To determine how the log likelihood varies along this principal component, the mean parameter set $\vec{\mu}$ and eigenvector for the principal component \vec{v}_1 were extracted and the log likelihood was computed at values \vec{p} along the affine linear subspace

$$\vec{p} = \mu + k\vec{v}_1, \tag{3.5}$$

see Figure 5 for a schematic depiction. For even a relatively large displacement ($k = 4$, Figure 4d), the log likelihood and quality of fit change only marginally. Furthermore, it is simple to check that the exact parameter set used to construct the data lies nearly on this subspace, and that it provides only a marginally better fit, Figure 4e. This supports the supposition that there is a single, strong correlation within the model and that the poor localization of the posterior is intrinsic to the standard LBA model.

Since the model degeneracy (i.e. correlation) involves only a one dimensional subspace of the 8 dimensional parameter space, fixing a single parameter in that subspace should in theory fully localize the posterior. To test this, the threshold parameter was fixed to the value $b = 2.7$, which was originally chosen to produce the data. The same procedure was carried out to sample the posterior (Figure 4c), and results indeed show it becomes substantially more localized. Furthermore, the mean parameter set from the constrained model is almost identical to the exact parameters used to construct the data. This confirms the spread in the posterior is a result of the strong correlation.

3.3.2. npABC and sloppy models

These observations do however raise an important issue. Recall from the previous two examples that this approximation procedure yields small errors in the log likelihood of on the order of $\sim 1 - 2\%$. While this might not seem too large, it can have a substantial effect on estimation of sloppy models such as this. The issue is that along this correlated parameter dimension, the

558 variation of the log likelihood is the same size or slightly larger than the log likelihood estimation
 559 variance. The npABC will thus explore the corridor along this correlated dimension more so than
 560 an MCMC with an analytic likelihood would. This variance can of course be reduced with extra
 561 computational power, but in practice this will not be practical.

562 If the goal is to understand the behavior of the model and its capacity to account for observa-
 563 tions, this may not be an issue. However, if the context being considered requires one to extract
 564 a single parameter set, more must be done. Strategies for dealing with sloppiness in models have
 565 been discussed extensively in other literature (Gutenkunst et al., 2007; Apgar et al., 2010), but
 566 such exposition is beyond the scope of this article. It is important to reiterate though that the
 567 sloppiness of posterior estimates here is more a reflection of an underlying model property that
 568 prevents accurate estimation.

569 4. A procedural overview of npABC for the practitioner

570 The previous sections outline many points that must be considered when using this method.
 571 Here, a procedural overview building on these results is provided for the interested practitioner.
 572 Familiarity with MCMC methods is assumed and only the details that relate to non-parametric
 573 component of this method is provided. It is impossible to list all details, however I again note that
 574 MATLAB codes have been supplied to aid the interested practitioner fill in the technical details.

575 0) Choose the number of samples to be used in the estimation process (N_s) and the kernel
 576 bandwidth (h). A minimum of $N_s = 10,000$ should generally be used. Choosing h will
 577 require trial and error, but Silverman’s rule of thumb (Silverman, 1986) provides a good
 578 starting point. As a general rule however, err on the side of smaller h since this will reduce
 579 estimation bias (at the expense of performance).

580 1) Loop over chains.

581 a) Generate a proposal θ^n . In the applications here this was done using DE-MCMC
 582 (Ter Braak, 2006; Turner et al., 2013), but any MCMC procedure can be used.

583 b) Compute $\hat{L}L(X|\theta^n)$ using the KDE.

584 i) Generate N_s samples from the model.

585 ii) Create a discrete representation of the likelihood by binning those samples into 2^n
 586 ($n > 8$) equally spaced bins with centers z_0, \dots, z_l . Set these bin centers so that
 587 $z_0 < \min(X) - 3h$ and $z_l > \max(X) + 3h$. This pads both sides of the histogram
 588 with zeros so the FFT is more accurate.

589 iii) Apply a FFT to map the data into the spectral domain.

590 iv) Apply the Gaussian smoothing filter. This is essentially the convolution step in the
 591 spectral domain.

592 v) Map the filtered signal back to the data space, producing a likelihood function
 593 $\hat{L}(z_i|\theta^n)$ on the regularly spaced grid.

594 vi) Interpolate this likelihood on the grid to the observation values, $\hat{L}(z_i|\theta^n) \rightarrow \hat{L}(x_i|\theta^n)$,
 595 using linear interpolation. Do not use cubic splines or anything higher order than
 596 linear as they can induce negative values in the tail of the distribution.

- 597 vii) Replace any zero values of $\hat{L}(x_i|\theta^n)$ with a minimum value, say $L_{min} = 1/(10 * N_s)$.
 598 viii) Compute the approximate log likelihood as

$$\hat{LL}(X|\theta^n) = \sum_{i=1}^{N_d} \log \left(\hat{L}(x_i|\theta^n) \right). \quad (4.1)$$

- 599 c) Compute the acceptance probability $\hat{\alpha}^n$ and accept or reject the proposal.
- 600 2) Resample the log likelihood of any previous chain. The algorithms here resample each chain
 601 every third MCMC iteration, though more efficient schemes are certainly possible. For ex-
 602 ample, the length of time a chain remains stuck can be recorded and used to determine when
 603 to resample.

604 Steps 1,2 define a single update of every chain in the re-sampled MCMC procedure. Simply iterate
 605 these steps the desired number of times and apply burn-in rules. Note that steps such as compu-
 606 tation of the prior, its incorporation into the acceptance probability, and specifically how to call
 607 the FFT have been neglected for brevity. Details can be found in the codes associated with these
 608 examples.

609 5. Discussion

610 This article presents a non-parametric approximate Bayesian computation (npABC) algorithm.
 611 This method, which combines non-parametric statistical methods with Bayesian inference tech-
 612 niques, is an extensible methodology for performing Bayesian posterior estimation. The purpose
 613 of this article is to elaborate this methodology in detail, discuss its pitfalls, improve its efficiency,
 614 and make it accessible a broader audience of end users.

615 A great many algorithms, ranging from Markov Chain Monte Carlo (MCMC) (Gelman et al.,
 616 2003; Robert and Casella, 2004) to particle filtering methods (Cappé et al., 2004; Del Moral et al.,
 617 2006), have been developed for the purpose of posterior estimation in contexts where Bayes' formula
 618 cannot be computed directly. These methods however typically require a closed form description
 619 of the model's likelihood. More recently, numerous approximate Bayesian computation (ABC)
 620 methods have extended these to likelihood free contexts. These methods however require the user to
 621 prescribe a set of summary statistics that describe the model / data. Unfortunately, these summary
 622 statistics are rarely sufficient to describe the model, and so the model that is fit is different from
 623 the one intended, by a substantial margin in some cases. More recently, non-parametric methods
 624 have been incorporated into ABC to circumvent this requirement (Turner and Sederberg, 2014).

625 Both npABC and ABC are similar in that they seek to determine the likelihood or plausibility
 626 of a given parameter set by first simulating a large number of model realizations, and second
 627 comparing those model realizations to the data. The central feature that differentiates npABC from
 628 other ABC methods however is that for each parameter set under consideration, a non-parametric
 629 approximation of the underlying likelihood is constructed, as opposed to some surrogate based on
 630 summary statistics. This provides two distinct benefits. First, the user does not have to make a
 631 possibly erroneous assumption about the form of the underlying model distribution. Second, this

632 method more fully utilizes the data since it does not compress it into a small number of summary
633 statistics.

634 The key step in this method of course is to construct an approximation of the underlying models
635 density function $L(x|\theta)$. This is accomplished using a kernel density estimation (KDE) technique
636 (Silverman, 1986), which is a method of directly computing an approximate of the likelihood of any
637 particular observation $L(x_i|\theta)$ from a collection of simulated model observations. The KDE can
638 thus be used to directly compute an approximation (\hat{L}) to the models log likelihood $LL(X|\theta)$,
639 which is the key piece of information needed for MCMC sampling. Thus a third practical benefit of
640 this method, in addition to theoretical benefits mentioned above, is that this KDE can be directly
641 integrated into standard MCMC techniques, since the likelihood itself is being assessed rather than
642 some surrogate.

643 Results here and elsewhere (Holmes et al., 2014; Turner and Sederberg, 2014) show this method-
644 ology is highly efficient and performs well. There are however a number of implementation details
645 that must be considered. First, the standard KDE procedure is highly inefficient and can itself
646 become a computational bottleneck. For this reason, a highly efficient implementation of KDE,
647 which utilizes only standard and highly optimized fast Fourier transform and linear interpolation
648 subroutines, is presented. In the applications discussed here, this implementation improved com-
649 putation times by a factor of 10 or more. Second, while this method can be directly plugged into
650 standard MCMC procedures, doing so can lead to inefficiencies. This stems from the fact that
651 the KDE is a statistical estimator of the underlying likelihood and as a result has an intrinsic
652 variance. To overcome this issue, a “resampled MCMC” procedure is proposed, which accounts for
653 the variability in this estimator and substantially improves performance.

654 While these investigations demonstrate the efficacy and efficiency of this methodology, like any
655 approximate method, it does come with drawbacks that must be kept in mind. First, the KDE
656 likelihood estimator is inherently biased. In applications discussed here, this bias is quite small,
657 being on the order of 1% or less. Unfortunately, these very small errors can have a profound effect on
658 model comparison and hypothesis testing. The essential problem is that standard quantities such as
659 AIC, BIC, or DIC are inherently flawed as they are absolute measures of model comparison rather
660 than relative measures. It is commonly accepted that Δ DIC of 10 is interpreted as “significant”
661 evidence for or against a model. However, if DIC measures are on the order of 1,000 (which
662 is common), a difference between two models of 1% would be considered “significant”, which is
663 rarely sensible. Since KDE approximation errors are on the order of 1 – 2%, those errors will
664 often overwhelm these model comparison statistics since the distinguishing difference is within the
665 methods margin of error.

666 A second issue is that this method can have difficulties with models containing very strong param-
667 eter correlations, which in other fields are commonly referred to as “sloppy models” (Gutenkunst
668 et al., 2007). The essential issue here is that the models with strong correlations are under-
669 determined in the sense that large parameter variations along the correlated dimension can lead
670 to very small changes in log likelihoods. In the final example presented here, varying parameters
671 by a factor of 10 along the correlated dimension leads to a $\sim 1\%$ variation in log likelihood and
672 nearly indistinguishable fits to the models data. Given these model fit variations are within the
673 small margin of error of the KDE approximation, posteriors become broadened. Thus, care must
674 be taken in interpreting the results of this method when such under-determined, highly correlated

675 models are being considered.

676 Despite these issues, this method has distinct benefits over existing ABC methods. With stan-
677 dard methods, it is rarely possible to know how good or bad the summary statistics being used
678 are. Using npABC however, the quality of an approximation can be controlled in a predictable
679 way by varying kernel density estimation parameters. Furthermore, since the types of errors being
680 made with npABC are somewhat predictable and quantifiable, their influence on results is also
681 reasonably predictable. Additionally, the efficiency of this methodology is comparable to existing
682 methods, especially with the more efficient KDE implementation presented here. Thus it can be
683 applied in almost any context where ABC methods are currently being or might be used. For these
684 reasons, this method should be added to the toolbox of any researcher performing Bayesian anal-
685 ysis of complex models beyond the reach of existing toolboxes such as JAGS (Plummer, 2003) or
686 WinBUGS (Lunn et al., 2000). The hope is that this article (along with the supporting MATLAB
687 codes) will make this method more accessible to those who could benefit from its use.

688 **Acknowledgements**

689 I would like to thank Joachim Vandekerckhove and Jennifer S. Trueblood for a critical reading of
690 this manuscript.

691 **References**

- 692 Abelson, R. P., 2008. Statistics as principled argument. Psychology Press.
- 693 Albantakis, L., Deco, G., 2009. The encoding of alternatives in multiple-choice decision-making.
694 BMC Neuroscience 10 (Suppl 1), P166.
- 695 Apgar, J. F., Witmer, D. K., White, F. M., Tidor, B., 2010. Sloppy models, parameter uncertainty,
696 and the role of experimental design. Mol. BioSyst. 6, 1890–1900.
- 697 Brown, S., Heathcote, A., 2005. A ballistic model of choice response time. Psychological review
698 112 (1), 117.
- 699 Brown, S. D., Heathcote, A., 2008. The simplest complete model of choice response time: Linear
700 ballistic accumulation. Cognitive psychology 57 (3), 153–178.
- 701 Brunel, N., Wang, X.-J., 2001. Effects of neuromodulation in a cortical network model of ob-
702 ject working memory dominated by recurrent inhibition. Journal of computational neuroscience
703 11 (1), 63–85.
- 704 Busemeyer, J. R., Townsend, J. T., 1993. Decision field theory: a dynamic-cognitive approach to
705 decision making in an uncertain environment. Psychological review 100 (3), 432–459.
- 706 Cappé, O., Guillin, A., Marin, J.-M., Robert, C. P., 2004. Population monte carlo. Journal of
707 Computational and Graphical Statistics 13 (4).

- 708 Cassey, P., Heathcote, A., Brown, S. D., 07 2014. Brain and behavior in decision-making. PLoS
709 Comput Biol 10 (7), e1003700.
- 710 Csilléry, K., Blum, M. G., Gaggiotti, O. E., François, O., 2010. Approximate bayesian computation
711 (abc) in practice. Trends in ecology & evolution 25 (7), 410–418.
- 712 Dantzig, G. B., Thapa, M. N., 1997. Linear Programming 1: 1: Introduction. Vol. 1. Springer.
- 713 Del Moral, P., Doucet, A., Jasra, A., 2006. Sequential monte carlo samplers. Journal of the Royal
714 Statistical Society: Series B (Statistical Methodology) 68 (3), 411–436.
- 715 Epanechnikov, V., 1969. Non-parametric estimation of a multivariate probability density. Theory
716 of Probability and Its Applications 14 (1), 153–158.
- 717 Friedrichs, K. O., 1944. The identity of weak and strong extensions of differential operators. Trans-
718 actions of the American Mathematical Society 55 (1), 132–151.
- 719 Gallistel, C., 2009. The importance of proving the null. Psychological review 116 (2), 439.
- 720 Gelman, A., Carlin, J. B., Stern, H. S., Rubin, D. B., 2003. Bayesian Data Analysis, 2nd Edition.
721 Chapman and Hall/CRC.
- 722 Gutenkunst, R. N., Waterfall, J. J., Casey, F. P., Brown, K. S., Myers, C. R., Sethna, J. P., 10
723 2007. Universally sloppy parameter sensitivities in systems biology models. PLoS Comput Biol
724 3 (10), e189.
- 725 Heathcote, A., Brown, S., 2004. Reply to speckman and roudner: A theoretical basis for qml.
726 Psychonomic Bulletin & Review 11 (3), 577–578.
- 727 Heathcote, A., Brown, S., Mewhort, D., 2002. Quantile maximum likelihood estimation of response
728 time distributions. Psychonomic Bulletin & Review 9 (2), 394–401.
- 729 Holmes, W. R., Trueblood, J. S., Heathcoat, A., 2014. Asymmetric updating and hysteresis in
730 perceptual decision-making with changing information. In review.
- 731 Huk, A. C., Shadlen, M. N., 2005. Neural activity in macaque parietal cortex reflects temporal inte-
732 gration of visual motion signals during perceptual decision making. The Journal of neuroscience
733 25 (45), 10420–10436.
- 734 Kiani, R., Hanks, T. D., Shadlen, M. N., 2008. Bounded integration in parietal cortex underlies
735 decisions even when viewing duration is dictated by the environment. The Journal of Neuroscience
736 28 (12), 3017–3029.
- 737 Lee, M. D., 2008. Three case studies in the bayesian analysis of cognitive models. Psychonomic
738 Bulletin & Review 15 (1), 1–15.
- 739 Lee, M. D., Wagenmakers, E.-J., 2013. Bayesian cognitive modeling: A practical course. Cambridge
740 University Press.

- 741 Lunn, D. J., Thomas, A., Best, N., Spiegelhalter, D., 2000. Winbugs a bayesian modelling frame-
742 work: concepts, structure, and extensibility. *Statistics and computing* 10 (4), 325–337.
- 743 Ollman, R., 1966. Fast guesses in choice reaction time. *Psychonomic Science* 6 (4), 155–156.
- 744 Plummer, M., 2003. Jags: A program for analysis of Bayesian graphical models using gibbs sam-
745 pling. In: *Proceedings of the 3rd International Workshop on Distributed Statistical Computing*.
- 746 Ratcliff, R., 1978. A theory of memory retrieval. *Psychological Review* 85, 59–108.
- 747 Ratcliff, R., Rouder, J. N., 1998. Modeling response times for two-choice decisions. *Psychological*
748 *Science* 9 (5), 347–356.
- 749 Ratcliff, R., Tuerlinckx, F., 2002. Estimating parameters of the diffusion model: Approaches to
750 dealing with contaminant reaction times and parameter variability. *Psychonomic bulletin & re-*
751 *view* 9 (3), 438–481.
- 752 Robert, C., Casella, G., 2004. *Monte Carlo statistical methods*. Springer, New York, NY.
- 753 Robert, C. P., Cornuet, J.-M., Marin, J.-M., Pillai, N. S., 2011. Lack of confidence in approximate
754 bayesian computation model choice. *Proceedings of the National Academy of Sciences* 108 (37),
755 15112–15117.
- 756 Saltelli, A., Ratto, M., Andres, T., Campolongo, F., Cariboni, J., Gatelli, D., Saisana, M., Taran-
757 tola, S., 2008. *Global sensitivity analysis: the primer*. John Wiley & Sons.
- 758 Silverman, B. W., 1982. Algorithm as 176: Kernel density estimation using the fast fourier trans-
759 form. *Journal of the Royal Statistical Society. Series C (Applied Statistics)* 31 (1), pp. 93–99.
- 760 Silverman, B. W., 1986. *Density estimation for statistics and data analysis*. Vol. 26. CRC press.
- 761 Storn, R., Price, K., 1997. Differential evolution—a simple and efficient heuristic for global opti-
762 mization over continuous spaces. *Journal of global optimization* 11 (4), 341–359.
- 763 Ter Braak, C. J., 2006. A markov chain monte carlo version of the genetic algorithm differential
764 evolution: easy bayesian computing for real parameter spaces. *Statistics and Computing* 16 (3),
765 239–249.
- 766 Thura, D., Beauregard-Racine, J., Fradet, C.-W., Cisek, P., 2012. Decision making by urgency
767 gating: theory and experimental support. *Journal of Neurophysiology* 108 (11), 2912–2930.
- 768 Tsetsos, K., Gao, J., McClelland, J. L., Usher, M., 2012. Using time-varying evidence to test models
769 of decision dynamics: Bounded diffusion vs. the leaky competing accumulator model. *Frontiers*
770 *in Neuroscience* 6.
- 771 Turner, B., Sederberg, P., 2014. A generalized, likelihood-free method for posterior estimation.
772 *Psychonomic Bulletin and Review* 21 (2), 227–250.

- 773 Turner, B. M., Sederberg, P. B., Brown, S. D., Steyvers, M., 2013. A method for efficiently sampling
774 from distributions with correlated dimensions. *Psychological methods* 18 (3), 368–384.
- 775 Turner, B. M., Van Zandt, T., 2012. A tutorial on approximate bayesian computation. *Journal of*
776 *Mathematical Psychology* 56 (2), 69–85.
- 777 Usher, M., McClelland, J. L., 2001. The time course of perceptual choice: the leaky, competing
778 accumulator model. *Psychological Review* 108 (3), 550–592.
- 779 Vandekerckhove, J., Tuerlinckx, F., 2007. Fitting the ratcliff diffusion model to experimental data.
780 *Psychonomic Bulletin & Review* 14 (6), 1011–1026.
- 781 Winkel, J., Keuken, M. C., Van Maanen, L., Wagenmakers, E.-J., Forstmann, B. U., 2014. Early
782 evidence affects later decisions: Why evidence accumulation is required to explain response time
783 data. *Psychon Bull Rev* 21, 777–784.
- 784 Yantis, S., Meyer, D. E., Smith, J. K., 1991. Analyses of multinomial mixture distributions: new
785 tests for stochastic models of cognition and action. *Psychological bulletin* 110 (2), 350.
- 786 Yellott Jr, J. I., 1967. Correction for guessing in choice reaction time. *Psychonomic Science* 8 (8),
787 321–322.
- 788 Yellott Jr, J. I., 1971. Correction for fast guessing and the speed-accuracy tradeoff in choice reaction
789 time. *Journal of Mathematical Psychology* 8 (2), 159–199.

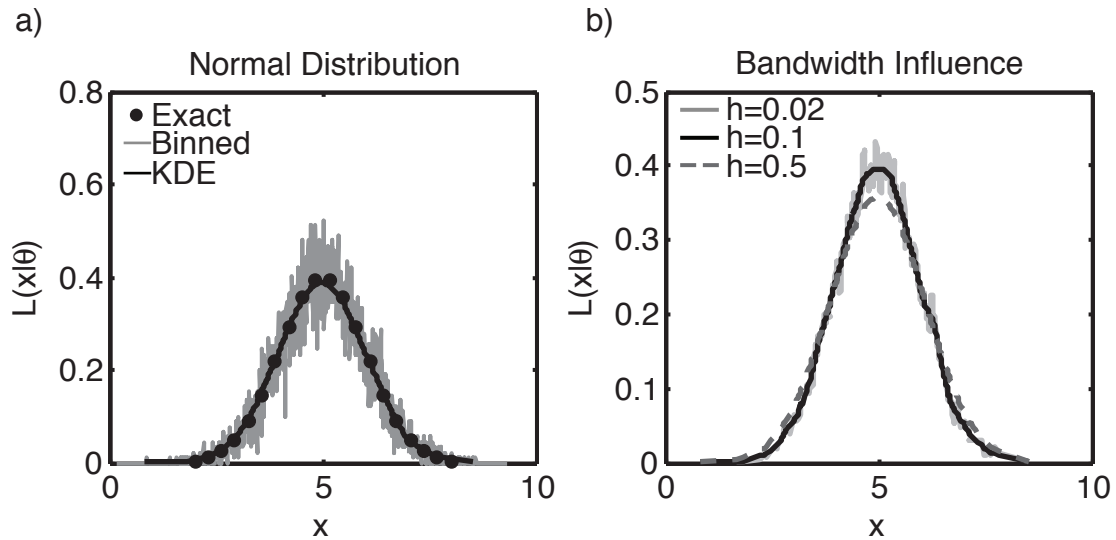


Figure 1: **Reconstructing a Gaussian: Panel a)** Reconstruction of a normal distribution using the FFT based KDE method. Gray lines indicate the noisy density estimate derived from binning $N_s = 10,000$ sampled points into 2^{10} bins and normalizing to produce a density. Black line indicates the smoothed version after convolving in the spectral domain. Circles show the exact likelihood at a few values, indicating agreement between the constructed and analytic likelihood. The bandwidth $h = 0.1$ is used here. **Panel b)** Reconstructed likelihood for three choices of the bandwidth parameter h . The mean and standard deviation parameters used are $\mu = 5, \sigma = 1$.

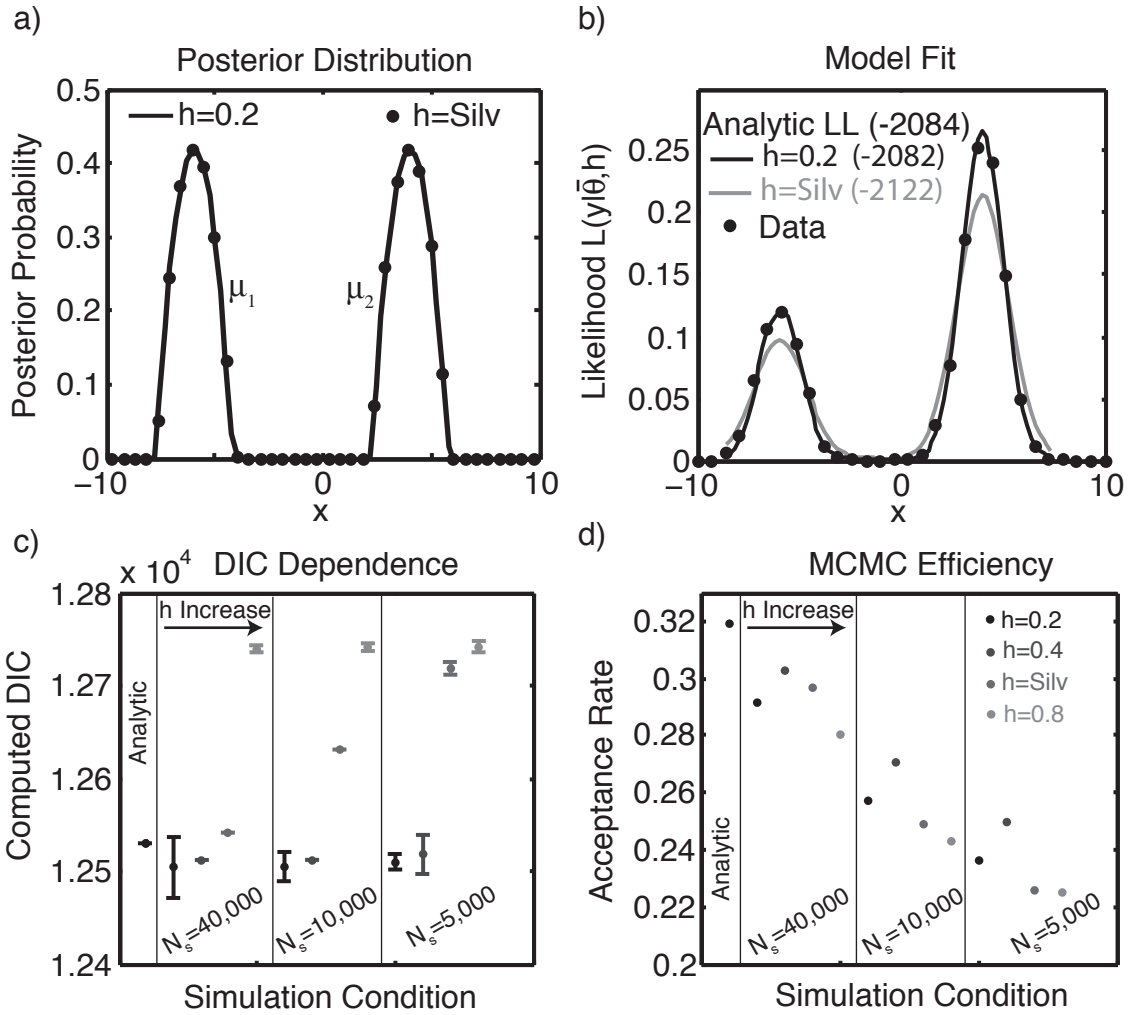


Figure 2: **Fitting a bimodal distribution:** *Panel a)* Posterior distribution obtained using $N_s = 10,000$ and two separate values of $h = 0.2, 0.68$, the latter is computed from Silverman’s rule of thumb. Posterior for the two means μ_1, μ_2 are shown. *Panel b)* Quality of model fit. For each value of h (and $N_s = 10,000$). The mean of the posterior for each parameter was computed, which was used to reconstruct the likelihood. In both cases, the approximate likelihood method with the associated h was used to construct the likelihood, but with $N_s = 1,000,000$ to reduce variance. In both cases the log likelihood along with that computed analytically are reported. *Panel c)* Dependence of DIC on N_s and h . DIC was computed by fitting the posterior using the analytic likelihood. Then, 100 fits of the posterior were obtained for each of different combinations of N_s and h . Within each value of N_s , the h values increase from left to right $h = 0.2, 0.4, Silv, 0.8$ where *Silv* indicates the bandwidth computed from Silverman’s rule of thumb. *Panel d)* Acceptance rates as a function of N_s, h . Data from the 100 posterior fits in (c) were used, though variance was so small it is not shown. The reported values of h increase from left to right, as in (c). Note the reduced efficiency with decreased N_s . A correction to the MCMC that alleviates this performance reduction is discussed in Section 3.2.2.

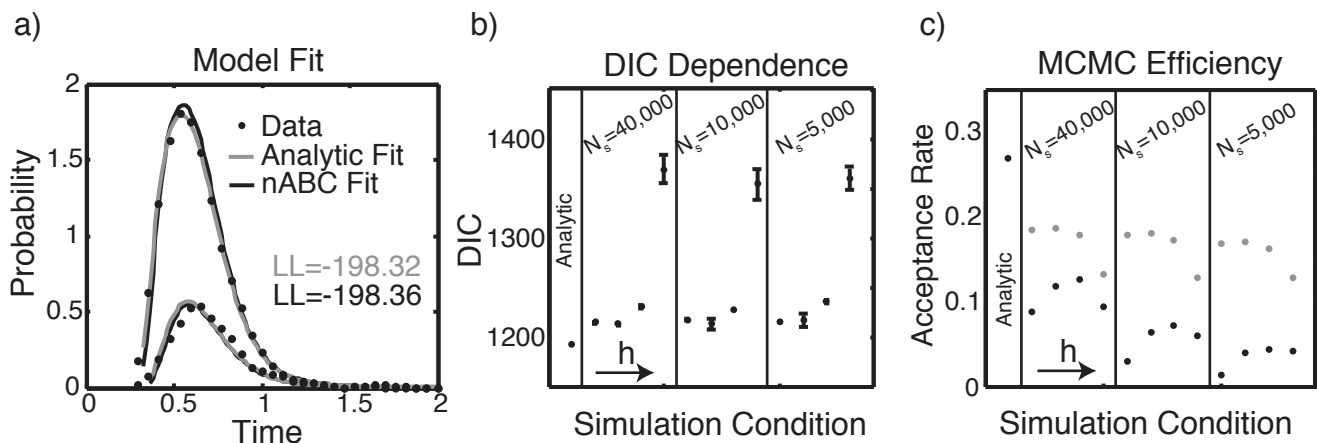


Figure 3: **Fitting the LBA:** *Panel a)* Quality of model fit using the analytic likelihood and approximate likelihood methods (with resampled MCMC). For the npABC fit, the approximate likelihood function was used to reconstruct the likelihood using the same value of h , but $N_s = 1,000,000$ to reduce variance. The analytic LBA likelihood was used for the analytic case. The high and low peaked curves correspond to correct and incorrect choice options respectively. In both cases, the log likelihood is reported. *Panel b)* Dependence of DIC on N_s and h . DIC was computed by fitting the posterior using the analytic likelihood. Then, 100 fits of the posterior were obtained for each of different combinations of N_s and h . Within each value of N_s , the h values increase from left to right $h = 0.01, Silv, 0.04, 0.07$ where *Silv* indicates the bandwidth computed from Silverman's rule of thumb. Mean and variance of the DIC samples is shown for each. The resampled MCMC is again used here. *Panel c)* Acceptance rates as a function of N_s, h for the standard (black) and resampled MCMC (gray). The same values of h as in panel *b*, again increasing from left to right. Resampling occurred every third chain iteration for the resampled MCMC.

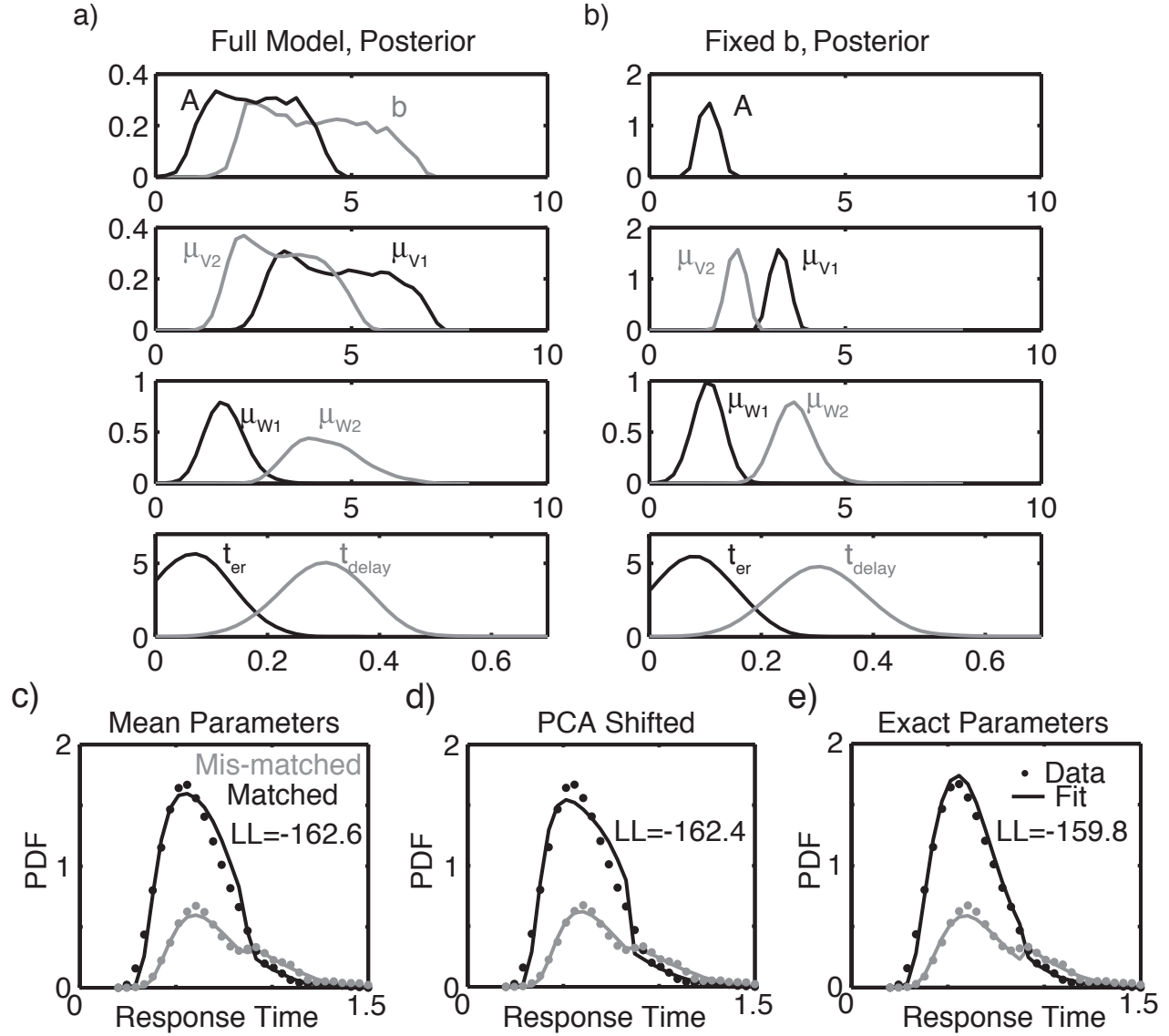


Figure 4: **Fitting the piecewise LBA:** *Panels a,b)* Posteriors for each parameter in the full piecewise LBA model and the restricted model with $b = 2.7$ fixed. *Panels c,d,e)* Fit to data for three parameter sets: c) the mean parameter set taken from (a), d) a translation from the mean parameter set along the first principal component, and e) the parameter set used to produce the data. Matched (dark) and mis-matched (gray) refer to response times for correct / incorrect prior to the information change. Note that in computing the quoted log likelihoods, h was taken very small and N_s very large to effectively remove any bias / variance in the estimates.

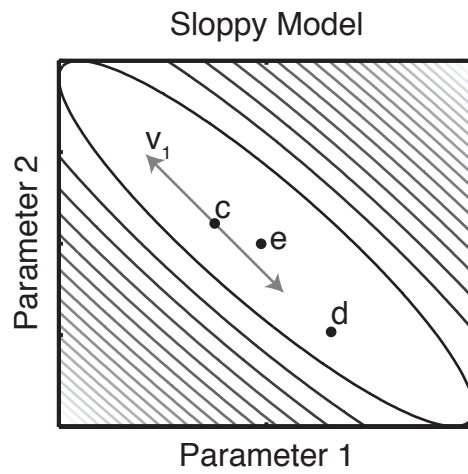


Figure 5: **Sloppy models:** Schematic depiction of a sloppy model showing a strong one dimensional correlation in the likelihood space. v_1 indicates the first principal component while the point c would be akin to the mean parameter set determined from MCMC chain samples. Point d indicates a point along the same principal component subspace while point e indicates the “best fit” parameter (e.g. maximum likelihood). These points schematically indicate the parameters used to produce Figures 4c-e respectively.