


Towards end-to-end likelihood-free inference with convolutional neural networks

Stefan T. Radev^{1*} , Ulf K. Mertens^{1*}, Andreas Voss¹ and Ullrich Köthe²

¹Institute of Psychology, Heidelberg University, Germany

²Heidelberg Collaboratory for Image Processing (HCI), Interdisciplinary Center for Scientific Computing (IWR), Heidelberg University, Germany

Complex simulator-based models with non-standard sampling distributions require sophisticated design choices for reliable approximate parameter inference. We introduce a fast, end-to-end approach for approximate Bayesian computation (ABC) based on fully convolutional neural networks. The method enables users of ABC to derive simultaneously the posterior mean and variance of multidimensional posterior distributions directly from raw simulated data. Once trained on simulated data, the convolutional neural network is able to map real data samples of variable size to the first two posterior moments of the relevant parameter's distributions. Thus, in contrast to other machine learning approaches to ABC, our approach allows us to generate reusable models that can be applied by different researchers employing the same model. We verify the utility of our method on two common statistical models (i.e., a multivariate normal distribution and a multiple regression scenario), for which the posterior parameter distributions can be derived analytically. We then apply our method to recover the parameters of the leaky competing accumulator (LCA) model and we reference our results to the current state-of-the-art technique, which is the probability density estimation (PDA). Results show that our method exhibits a lower approximation error compared with other machine learning approaches to ABC. It also performs similarly to PDA in recovering the parameters of the LCA model.

1. Introduction

A major goal in statistical modelling is to describe data-generation processes by a finite parameter vector θ . Because there are different sources of uncertainty (Goodfellow, Bengio, Courville, & Bengio, 2016, p. 52), such descriptions are inherently stochastic. Naturally, researchers are interested in quantifying the uncertainty in their parametric estimates (Kendall & Gal, 2017; Lee, 2008; Schoot *et al.*, 2014). By drawing on the mathematical framework of probability theory, one way to represent uncertainty is to place probability distributions over parameters, and to update the parameters of these

*Correspondence should be addressed to Stefan T. Radev and Ulf Mertens, Institute of Psychology, Heidelberg University, Hauptstr. 47-51, 69117 Heidelberg, Germany (emails: stefan.radev93@gmail.com and ulf.mertens@psychologie.uni-heidelberg.de).

The co-authors Stefan T. Radev and Ulf K. Mertens contributed equally to the work.

DOI:10.1111/bmsp.12159

distributions as the current state of knowledge changes. In fact, this is the quintessential idea incorporated in Bayesian statistics (Gelman *et al.*, 2013).

Suppose we have collected a dataset $\mathbf{x} = \{\mathbf{x}^{(i)}\}_{i=1}^N$ in an observational or an experimental setting. At the core of all Bayesian data analysis lies Bayes' rule:

$$p(\theta|\mathbf{x}) = \frac{p(\mathbf{x}|\theta)\pi(\theta)}{\int p(\mathbf{x}|\theta)\pi(\theta)d\theta}. \quad (1)$$

In the above expression, $p(\theta|\mathbf{x})$ denotes the posterior distribution of the model parameters, $p(\mathbf{x}|\theta)$ denotes the likelihood function, $\pi(\theta)$ denotes the prior probability distribution and the denominator $p(\mathbf{x}) = \int p(\mathbf{x}|\theta)\pi(\theta)d\theta$ denotes the marginal probability of the data. Considered as a normalization constant, the computation of the marginal probability can often be bypassed, so the posterior distribution can be expressed as being proportional to the prior times the likelihood:

$$p(\theta|\mathbf{x}) \propto p(\mathbf{x}|\theta)\pi(\theta). \quad (2)$$

If the likelihood belongs to a known family of probability distributions, and if the prior is conjugate to the likelihood, then the posterior belongs to the same distribution family as the prior. Therefore, by choosing mathematically convenient priors, we can obtain a closed-form expression for the posterior, which has the same algebraic form as the prior, merely with modified parameter values. In the case of non-conjugate priors, it is necessary to resort to Monte Carlo sampling methods, such as Markov chain Monte Carlo (MCMC) algorithms.

A different problem arises when the likelihood $p(\mathbf{x}|\theta)$ is computationally intractable, or cannot be specified algebraically. This situation occurs when models become increasingly complex, or rely on a simulation-based mechanism for generating the data (Turner & Van Zandt, 2012). In this case, approximation methods are required in order to perform any type of Bayesian inference.

1.1. Approximate Bayesian computation methods

Approximate Bayesian computation (ABC) methods are part of a larger class of techniques aimed at bypassing the intractability problem arising when parametric models require a non-standard solution based on a likelihood function.

The ABC method was first successfully implemented in genetics research (Tavaré, Balding, Griffiths, & Donnelly, 1997) and the utility of the method has been steadily increasing across research domains ever since (Pritchard, Seielstad, Perez-Lezaun, & Feldman, 1999). The main idea of ABC methods is to approximate the likelihood function by repeatedly sampling parameters from prior distributions and simulating artificial datasets conditioned on the sampled parameters. Even though exact numerical evaluation of the likelihood $p(\mathbf{x}|\theta)$ might be intractable, a generative model of the form $q(\cdot|\theta)$, usually specified as a function code in a programming language, is necessary for performing the simulations.

The starting point for all ABC algorithms is the creation of the so-called reference table (Mertens, Voss, & Radev, 2018; Raynal *et al.*, 2016). The reference table is a special data structure used to store a large number of simulated datasets or summary-statistics of such produced from an approximation of the posterior distribution given by $p(\theta|\mathbf{x}) \propto p(\mathbf{x}|\theta)\pi(\theta) \approx q(\cdot|\theta)\pi(\theta)$ as well as the corresponding samples from the prior

distributions(s) used to generate these datasets. In the following, the details of creating the reference table are discussed.

In line with the notation embraced by the literature on ABC, let θ denote a random vector of parameters and let $\pi(\theta)$ denote the joint prior distribution over the parameter vector. Let $\{\tilde{\mathbf{x}}^{(i)}\}_{i=1}^N$ represent a collection of artificial datasets simulated by a generative model $q(\cdot|\theta)\pi(\theta)$, let \mathbf{x} denote the observed dataset and let $\eta(\tilde{\mathbf{x}}^{(i)})$ denote a (vector) summary statistic obtained by some form of dimensionality reduction applied to each $\tilde{\mathbf{x}}^{(i)}$. Consequently, Algorithm 1 describes the standard procedure for generating the ABC reference table.

Algorithm 1. Generation of the reference table

for $i \leftarrow 1$ to N

 Sample $\theta^{(i)} \sim \pi(\cdot)$

 Simulate $\tilde{\mathbf{x}}^{(i)} \sim q(\cdot|\theta^{(i)})$

 Compute $\eta(\tilde{\mathbf{x}}^{(i)})$

 Store the pair $(\eta(\tilde{\mathbf{x}}^{(i)}), \theta^{(i)})$ in row i of the reference table

end for

There are multiple ways to utilize the reference table for the purpose of parameter inference. The rejection algorithm, initially proposed by Rubin (1984), and refined by Tavaré *et al.* (1997) and Pritchard *et al.* (1999), requires a distance function $d(\eta(\tilde{\mathbf{x}}^{(i)}), \eta(\mathbf{x}))$ to be specified. This quantifies the deviation of a given simulated sample from the observed data as well as a threshold (also called the tolerance level), $0 < \epsilon \leq 1$, according to which a fraction of the parameters, $1 - \epsilon$, is rejected, and the remaining are considered as samples from an approximate posterior distribution $p_\epsilon(\theta|\eta(\tilde{\mathbf{x}}))$. Even though, in theory, the rejection method is doubly asymptotically consistent (Frazier, Martin, Robert, & Rousseau, 2016), its practical limitations have been repeatedly noted by researchers (Blum, 2010; Marin, Pudlo, Robert, & Ryder, 2012; Raynal *et al.*, 2016), necessitating the use of more sophisticated methods. First, the rejection method suffers from the curse of dimensionality, which means that as the dimensions of the parameter space increase, the number of simulations required to obtain reasonable convergence grows exponentially large. Secondly, in order for $p_\epsilon(\theta|\eta(\tilde{\mathbf{x}}))$ to approximate the true posterior distribution $p(\theta|\mathbf{x})$ within a reasonable amount of computational time, the parameter ϵ needs to be carefully tuned. Thirdly, the summary statistic $\eta(\cdot)$ has to be sufficient, meaning that no loss of information from the sample should be incurred by computing it. Finally, the choice of a particular distance metric $d(\cdot, \cdot)$ is often arbitrary (e.g., Euclidian distance) and thus presents another non-trivial hyperparameter of the method requiring further attention.

1.2. Machine learning approaches to ABC

Recently, it has been proposed that techniques can be ‘borrowed’ from the machine learning literature in order to confront the above-mentioned shortcomings of traditional ABC methods. In the next two subsections, we briefly review two recent applications of the supervised machine learning approach to ABC. We address both the advantages and limitations of these approaches. In Section 1.5, we introduce our method.

1.3. Random forest methodology

The random forest (RF) algorithm, originally developed by Breiman (2001), appears to provide an especially promising approach for estimating the first two moments (see Raynal *et al.*, 2016 for details) and quantiles of a posterior parameter distribution. In particular, RF regression is a supervised learning algorithm, which, given a sample of input–output pairs, learns a highly non-linear mapping from the input to the output by training an ensemble of decision trees (a standard non-parametric algorithm for classification and regression) and aggregating their regression function outputs.

In the context of ABC, the inputs to the RF algorithm are the simulated and summarized datasets $\{\eta(\tilde{\mathbf{x}}^{(i)})\}_{i=1}^N$ from the reference table, while the outputs are the corresponding $\{\theta^{(i)}\}_{i=1}^N$ parameter values. In other words, the algorithm learns to recover the unknown parameters used to generate a given dataset. Algorithm 2 describes the ABC–RF procedure (Raynal *et al.*, 2016).

The ABC–RF methodology provides reliable and theoretically sound approximations of the posterior mean $\mathbb{E}[\theta|\eta(\tilde{\mathbf{x}})]$ and variance $\text{Var}[\theta|\eta(\tilde{\mathbf{x}})]$ of the parameter vector of interest. However, the user still needs to manually decide on and compute a (potentially large) collection of summary statistics. A further feature of the ABC–RF methodology is the fact that existing implementations involve the training of a separate RF ensemble for each individual model parameter θ_k of the parameter vector θ . Even though this approach could easily be parallelized across parameters, sequential solutions might experience a computational bottleneck, especially when dealing with high-dimensional parameter spaces, combined with a large reference table.

Algorithm 2. ABC–RF methodology (Raynal *et al.*, 2016)

Construct a reference table $\left\{\left(\eta(\tilde{\mathbf{x}}^{(i)}), \theta^{(i)}\right)\right\}_{i=1}^N$ using **Algorithm 1**.

$K \leftarrow \#$ of model parameters

for $k \leftarrow 1$ to K

Train a random forest regression $\hat{f}_{\theta_k}(\eta(\tilde{\mathbf{x}}))$ to predict parameter θ_k

end for

Output a group of random forests $\{\hat{f}_{\theta_k}\}_{k=1}^K$

Posterior uncertainty estimation using OOB samples:

for $k \leftarrow 1$ to K

Estimate $\mathbb{E}[\theta_k|\eta(\tilde{\mathbf{x}}), \Theta_k]$, $\text{Var}[\theta_k|\eta(\tilde{\mathbf{x}}), \Theta_k]$ and posterior quantiles $\mathcal{Q}_p(\theta_k|\eta(\tilde{\mathbf{x}}), \Theta_k)$

end for

1.4. Deep neural network methodology

Another promising approach to ABC inspired by machine learning utilizes deep neural networks (DNNs) for automatically learning summary statistics from the reference table (Jiang, Wu, Zheng, & Wong, 2015; Sheehan & Song, 2016). At a general level, a neural network defines a mapping $y = f_W(x)$ from input x to output y and learns the parameters W that lead to the best function approximation (Goodfellow *et al.*, 2016). As in the case of ABC–RF, the challenge of ABC is cast as a supervised learning task, with artificial datasets $\{\tilde{\mathbf{x}}^{(i)}\}_{i=1}^N$ as inputs, and dataset-generating parameters $\{\theta^{(i)}\}_{i=1}^N$ as outputs. It is important

to note that raw simulated data are used as input to the neural network, because the goal is to learn the functional form of $\eta(\cdot)$ implicitly from the data, treating the predicted parameters as the summary statistic. Algorithm 3 describes the ABC–DNN procedure (Jiang *et al.*, 2015).

Algorithm 3. ABC-DNN methodology (Jiang *et al.*, 2015)

Construct a reference table $\left\{(\tilde{\mathbf{x}}^{(i)}, \theta^{(i)})\right\}_{i=1}^N$ using **Algorithm 1**.

Train a DNN f_W with $\left\{\tilde{\mathbf{x}}^{(i)}\right\}_{i=1}^N$ as input and $\left\{\theta^{(i)}\right\}_{i=1}^N$ as output

ABC rejection algorithm with DNN outputs as summary statistics

for $t \leftarrow 1$ to M

 Sample $\theta^{(t)} \sim \pi(\cdot)$

 Simulate $\tilde{\mathbf{x}}^{(t)} \sim q(\cdot | \theta^{(t)})$

 Compute $\hat{\theta}^{(t)} = f_W(\tilde{\mathbf{x}}^{(t)})$ as a summary statistic

end for

Compute summary statistic of observed data $\theta = f_W(\mathbf{x})$

Select $\hat{\theta}^{(t)}$ such that $d(\theta, \hat{\theta}^{(t)}) < \epsilon$ where $d(\cdot, \cdot)$ is a distance function and ϵ is a pre-defined tolerance level

Use selected $\hat{\theta}^{(t)}$ to estimate features of $p(\theta | \mathbf{x})$

The most important advantage of using DNNs as supervised approximators is their huge representational power (Goodfellow *et al.*, 2016; Hornik, 1991; Jiang *et al.*, 2015). According to the universal approximation theorem (Cybenko, 1989; Hornik, 1991), neural networks are, in principle, able to approximate any continuous function to an arbitrary degree of accuracy, given the appropriate conditions and parameters. Furthermore, neural network architectures can easily be designed to incorporate multi-output regression, as is the case when θ is multidimensional, thus learning each mapping from $\left\{\tilde{\mathbf{x}}^{(i)}\right\}_{i=1}^N$ to individual parameters simultaneously.

One disadvantage of the ABC–DNN method compared to the ABC–RF approach is that it does not offer a straightforward way to quantify the uncertainty inherent in the DNN estimates, which is captured by the estimate of the posterior variance $\text{Var}[\theta | \eta(\tilde{\mathbf{x}}), \Theta]$ within the ABC–RF framework. As a consequence, Jiang *et al.* (2015) used the DNN estimates of the posterior expectation as summary statistics in the classical ABC rejection algorithm, thus inheriting the difficulties in specifying a convenient distance metric $d(\cdot, \cdot)$ and guessing an appropriate tolerance level ϵ . However, Sheehan and Song (2016) forgo rejection sampling altogether and rely solely on the DNN estimates for parameter inference. Despite requiring fewer steps, this method does not allow the quantification of uncertainty in parameter estimates. Finally, Blum and François (2010) use a powerful combination of kernel-based and DNN-based methods to estimate the posterior parameter variance along the posterior mean. These estimates are used in an adaptive algorithm aimed at adjusting local parameter estimates towards the posterior distribution $p_\epsilon(\theta | \eta(\tilde{\mathbf{x}}))$. Even though this approach ameliorates some of the issues surrounding rejection sampling by employing importance sampling, the practical implementation requires the selection of summary statistics and two separate DNNs for estimating the posterior mean and variances.

Another, more subtle, yet serious shortcoming of all previous supervised learning approaches to ABC, is the inability to deal with variable input sizes (e.g., the observed sample size). This inevitably confines the dimensionality of the simulated data to the dimensionality of the observed data sample. In addition, it is not possible to reuse the same model in a context where the observed data sample has a different size. Hence, researchers interested in applying current machine learning models need to generate a completely new reference table and train the models from scratch.

1.5. DeepInference

The approach we investigate in this paper is an attempt to fuse the advantages of both the ABC–RF and the ABC–DNN methodologies discussed so far (Jiang *et al.*, 2015; Raynal *et al.*, 2016). It is also inspired by recent advances in modelling uncertainty in deep learning predictions (Kendall & Gal, 2017).

We propose to train a fully convolutional neural network (FCN) on simulated data distributions in order to approximate the posterior mean of the relevant model parameters and, at the same time, to estimate the uncertainty in the posterior approximations directly from the available data by minimizing the heteroscedastic loss (Kendall & Gal, 2017). We argue that this approach overcomes most of the shortcomings of previous machine learning approaches to ABC.

The outline of the rest of this paper is as follows. In Section 2, we introduce the building blocks of our approach. Then, in Section 3, we confirm, using two toy examples, that the predictions and uncertainty estimates obtained by the optimization procedure employed in the convolutional neural network (CNN) closely approximate the analytically computed posterior expectations and posterior variances of the data-generating parameters. Also, we demonstrate the usefulness of the method in recovering the parameters of the leaky competitive accumulator (LCA; Usher & McClelland, 2001) model, which is used widely in cognitive science and psychology. Then, we compare our results with the latest state-of-the-art method (Miletić, Turner, Forstmann, & van Maanen, 2017). In Section 4, we discuss the advantages of the current approach.

2. Methods

2.1. Inference as optimization

In general, a regression model is trained by minimizing the mean squared error (MSE) loss across N training examples (simulations) between actual parameters $\theta^{(i)}$ and parameters predicted by the model $\hat{\theta}^{(i)} = f_W(\tilde{\mathbf{x}}^{(i)})$. The loss function for a single-parameter model is thus given by:

$$\mathcal{L}(W) = \frac{1}{N} \sum_{i=1}^N \|\hat{\theta}^{(i)} - \theta^{(i)}\|_2^2, \quad (3)$$

and optimized via stochastic gradient descent using the back propagation algorithm. Because the MSE loss is proportional to the cross-entropy between the empirical distribution of the training set and a Gaussian model, minimizing the MSE loss is equivalent to minimizing the negative log-likelihood of a conditional Gaussian model¹ defined as:

¹ Another interpretation of maximum likelihood is that it minimizes the cross-entropy between the data distribution and the model distribution (see Goodfellow *et al.*, 2016, p. 129).

$$p(\theta|\tilde{\mathbf{x}}; W) = \mathcal{N}(\theta|f_W(\tilde{\mathbf{x}}), \sigma^2), \quad (4)$$

where the prediction of the neural network $\hat{\theta} = f_W(\tilde{\mathbf{x}})$ corresponds to the mean of the Gaussian distribution. The negative log-likelihood thus becomes

$$-\log p(\theta|\tilde{\mathbf{x}}; W) = -\sum_{i=1}^N \log p(\theta^{(i)}|\tilde{\mathbf{x}}^{(i)}; W), \quad (5)$$

$$= -\sum_{i=1}^N \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(\hat{\theta}^{(i)} - \theta^{(i)})^2} \right), \quad (6)$$

$$= \sum_{i=1}^N \frac{\|\hat{\theta}^{(i)} - \theta^{(i)}\|_2^2}{2\sigma^2} + \frac{N}{2} \log(2\pi\sigma^2), \quad (7)$$

and we see that the minimizing the MSE criterion in equation (3) with respect to the DNN weights W is equivalent to minimizing equation (7), where σ^2 is considered to be a fixed constant. Because setting $\hat{\theta}$ to be equal to the posterior expectation $\mathbb{E}[\theta|\tilde{\mathbf{x}}]$ leads to the best possible prediction of θ given simulated data $\tilde{\mathbf{x}}$ and hence minimizes the expected MSE, the optimization procedure of the supervised learning approach effectively approximates the posterior mean of a given parameter θ (and similarly for a parameter vector θ in a multidimensional context). This observation also implies that any supervised learning approach that optimizes the MSE criterion can potentially be trained on an ABC reference table to approximate $\mathbb{E}[\theta|\mathbf{x}]$, thus reframing the problem of inference as a problem of optimization.

2.2. Heteroscedastic loss

As in the previous machine learning approaches to ABC, we frame the problem of parameter estimation as a supervised learning task, provided that a reference table is available. We propose to train a one-dimensional (1D) CNN on the raw simulated datasets $\{\tilde{\mathbf{x}}^{(i)}\}_{i=1}^N$, optimizing the heteroscedastic loss criterion (Kendall & Gal, 2017; Nix & Weigend, 1994):

$$\mathcal{L}(W) = \frac{1}{N} \sum_{i=1}^N \frac{1}{2\sigma^2(\tilde{\mathbf{x}}^{(i)})} \|\hat{\theta}^{(i)} - \theta^{(i)}\|_2^2 + \frac{1}{2} \log \sigma^2(\tilde{\mathbf{x}}^{(i)}). \quad (8)$$

Here, $\hat{\theta}^{(i)}$ and $\sigma^2(\tilde{\mathbf{x}}^{(i)})$ are the outputs of the network and depend on the network's weights W . It is immediately obvious that the heteroscedastic loss is proportional to the negative log-likelihood of a Gaussian model with constant variance as defined by equation (7). However, now we also train the network to explicitly predict the σ^2 term, corresponding to the variance of the target parameter distribution. The target distribution from equation (4) then becomes:

$$p(\theta|\tilde{\mathbf{x}}; W) = \mathcal{N}(\theta|f_W(\tilde{\mathbf{x}}), \sigma^2(\tilde{\mathbf{x}})), \quad (9)$$

where we explicitly denote the dependence of the σ^2 variance term on the simulated input dataset $\tilde{\mathbf{x}}$. Thus, the network is trained to predict $\sigma^2(\tilde{\mathbf{x}}^{(i)})$ along with $\hat{\theta}^{(i)}$. The

network output then becomes $(\hat{\theta}, \sigma^2(\tilde{\mathbf{x}})) = f_w(\tilde{\mathbf{x}})$. Assuming a Gaussian regression model, another interpretation of the heteroscedastic loss suggests that the network's predictions approximate both the posterior mean $\mathbb{E}[\theta|\mathbf{x}]$ and the posterior variance $\text{Var}[\theta|\mathbf{x}]$ of the parameter distribution.

Heteroscedastic regression comes with two important advantages: (1) it does not assume a fixed variance parameter across all inputs $\tilde{\mathbf{x}}^{(i)}$; (2) it enables us to learn the variance term $\sigma^2(\tilde{\mathbf{x}}^{(i)})$ during training as a function of the data. Intuitively, if the model makes a particularly bad prediction, the squared $L2$ loss term $\|\hat{\theta}^{(i)} - \theta^{(i)}\|_2^2$ is going to be large, thereby inducing an increase in the $\sigma^2(\tilde{\mathbf{x}}^{(i)})$ uncertainty term to keep the overall error low. However, the $\log \sigma^2(\tilde{\mathbf{x}}^{(i)})$ term prevents the model from 'cheating' by not allowing the uncertainty to grow to infinity, which would inevitably reduce the squared error term, but would not result in the model learning any meaningful information from the data.

2.3. Neural network architecture

For both toy examples, we used a FCN with three hidden layers (see O'Shea & Nash, 2015 for an in-depth discussion on convolutional networks) where the last hidden layer computes the average value of each kernel's output, also known as global average pooling (Lin, Chen, & Yan, 2013). If m is the number of parameters of the model from which samples are generated, the output layer of the CNN consists of m linear activation units for predicting each component of θ , and further m units predicting each corresponding $\sigma^2(\tilde{\mathbf{x}})$. The latter units apply the softplus activation function, defined by

$$\xi(z) = \log(1 + e^z). \quad (10)$$

The softplus activation is introduced in order to ensure that estimates of $\sigma^2(\tilde{\mathbf{x}})$ are always non-negative. Thus, for each application, the number of output units is double the number of parameters to be estimated. For all examples, we trained the CNN for five epochs using the Adam optimization algorithm (Kingma & Ba, 2014) with the learning rate set to 0.001. We did not perform an extensive hyperparameter search over the CNN parameter space, and we strived to keep the model as small as possible.

The core idea of our convolutional approach is to fully exploit the grid-like structure inherent in most datasets. In the simplest case of 1D independent and identically distributed (i.i.d.) datasets, each data point is statistically independent of all others, so it makes sense to use convolutional kernels (also called filters) of size 1 and stride (step size) 1 in the first layer of the network. Importantly, due to the properties of convolutional layers (see Section 4 for more details on why CNN properties are useful for extracting information from raw data samples), shuffling the data sample does not dramatically change the output of a particular kernel, which has hopefully learned to encode a given range of values, regardless of the position of the values in the input sample. In later layers, we allow for kernel sizes >1 in order to extract further informative characteristics of the data, thus implicitly learning a hierarchy of summary statistics.

In the more complicated case where each simulated dataset forms a multidimensional array, and data points along a given axis (e.g., rows) are connected in some way, the first hidden layer is built from kernels of size appropriate for encoding information contained in combinations or sequences of data points.

To make the latter description more concrete, consider a dataset with n rows and k columns containing $k - 1$ predictor variables and an outcome variable in a multiple regression scenario (see also Figure 1). Each row is composed of one single value per

predictor variable and the corresponding outcome. We exploit this structure by using 1D convolutional filters with a kernel size corresponding to the number of columns and stride of 1 as the first layers. These settings (kernel size = k , stride = 1) ensure that the network is fed all the relevant information inherent in this particular data structure. Algorithm 4 lays out the essential steps of our method.

Algorithm 4. DeepInference algorithm

Construct a reference table $\left\{ \left(\tilde{\mathbf{x}}^{(t)}, \theta^{(t)} \right) \right\}_{t=1}^N$ using **Algorithm 1**.

Train a convolutional neural network f_W with $\left\{ \tilde{\mathbf{x}}^{(t)} \right\}_{t=1}^N$ as input and $\left\{ \theta^{(t)} \right\}_{t=1}^N$ as output by optimizing the heteroscedastic loss (equation 8)

Uncertainty estimation:

Perform a forward pass on the observed data $(\theta, \sigma^2(\mathbf{x})) = f_W(\mathbf{x})$ to obtain estimates of $\mathbb{E}[\theta|\mathbf{x}]$ and $\text{Var}[\theta|\mathbf{x}]$.

All CNN models were implemented via the high-level neural networks API keras (<https://keras.io/>; see also Chollet, 2017) written in the Python programming language. Reference tables used for training the models and model implementation code are readily available at <https://github.com/mertensu/deep-inference>.

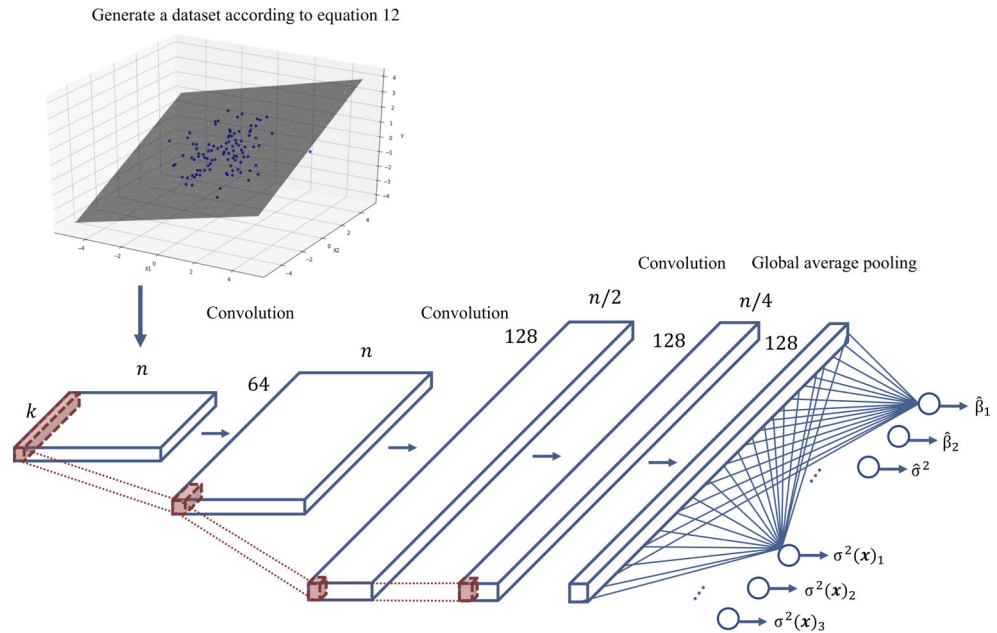


Figure 1. Graphical illustration of the CNN architecture used throughout this paper. The input represents datasets generated from the Bayesian regression model (toy example 2), organized as an $n \times k$ grid, with k equal to the number of predictors (in this case, 2) + the outcome, and n equal to the number of simulated samples. The CNN then performs a series of non-linear transformations, each layer applying convolution as a linear transformation followed by a rectified linear (ReLU) non-linearity, which is typical of modern convolutional architectures. The last layer then outputs the estimates of the posteriors mean and variance (one tuple of the form $(\hat{\theta}, \sigma^2(\tilde{\mathbf{x}}))$ for each model parameter). The loss is then evaluated via equation (8). [Colour figure can be viewed at wileyonlinelibrary.com]

In the following section, we evaluate the method proposed above on two artificial statistical models with known posterior distributions. Then, we apply the method to the LCA model from cognitive psychology.

3. Results

3.1. Toy example 1: Multivariate normal with unknown mean and variance

We demonstrate the utility of our method on the simple conjugate multivariate normal (MVN) model introduced in Gelman *et al.* (2013, pp. 72–73). The conjugate prior distributions for the mean μ and covariance matrix Σ of the MVN model are given by

$$\begin{aligned}\Sigma &\sim \text{InvWishart}(\Lambda_0, v_0) \\ \mu|\Sigma &\sim \mathcal{N}_d(\mu_0, \Sigma/k_0) \\ \mathbf{x}_i|\mu, \Sigma &\sim \mathcal{N}_d(\mu, \Sigma), \quad i = 1, \dots, m.\end{aligned}\tag{11}$$

In this model, $\text{InvWishart}(\Lambda_0, v_0)$ denotes an inverse Wishart distribution parametrized by a d -by- d scale matrix Λ_0 and degrees of freedom v_0 , and $\mathcal{N}_d(\mu, \Sigma)$ denotes a MVN distribution with a d -dimensional mean vector μ and a d -by- d covariance matrix Σ . As these priors are conjugate to the normal likelihood, it is possible to derive closed-form solutions for the joint posterior distributions (Gelman *et al.*, 2013). Thus, to obtain samples from the true joint posterior distribution of μ and Σ , we use the following procedure. First, draw $\Sigma|\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m \sim \text{InvWishart}(\Lambda_m, v_m)$, and then draw $\mu|\Sigma, \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m \sim \mathcal{N}_d(\mu_m, \Sigma/k_m)$, where $(\Lambda_m, v_m, \mu_m, k_m)$ represent the posterior model parameters, computed as functions of the data and the prior parameters (see Gelman *et al.*, 2013, pp. 72–73).

For the current example, we set $d = 2$, $m = 100$, $\Lambda_0 = I_2$, $v_0 = 5$, $k_0 = 3$, $\mu_0 = (0, 0)^T$ and we simulated 100,000 datasets from the model defined by equation (11). Thus, each row i of the reference table contains the raw sample $\{\tilde{\mathbf{x}}_k^{(i)}\}_{k=1}^{100}$, each element $\tilde{\mathbf{x}}_k^{(i)}$ being a realization of a two-dimensional random vector. Given the referenced table as input, the CNN model was trained to predict the corresponding values of $\{\mu^{(i)}, \text{diag}(\Sigma)^{(i)}\}$. Note that we are only interested in the diagonal elements of the covariance matrix, as we are only estimating the variance of the MVN model. We trained the model on a training set of 99,000 simulated datasets, holding out the remaining 1,000 datasets for online validation during training. Finally, to evaluate the performance of the model, we generated 100 independent datasets \mathbf{X}_{test} on which we computed the true values of $\mathbb{E}[\mu|\mathbf{x}]$, $\mathbb{E}[\text{diag}(\Sigma)|\mathbf{x}]$, $\text{Var}[\mu|\mathbf{x}]$ and $\text{Var}[\text{diag}(\Sigma)|\mathbf{x}]$ by sampling from the true joint posterior as described above. Figure 2 compares the estimates obtained from the true joint distribution with the estimates obtained from our CNN model.

Table 1 compares the performance of our DeepInference approach to the performance of the ABC-RF and ABC-DNN methods as specified by Algorithms 2 and 3 in terms of root mean squared error (RMSE). We used the same training and test sets for all three methods. For the ABC-RF approach, we computed the sample means, sample variances, the sample covariance and sample median absolute deviations, as well as all possible sums and products with these four metrics as summary statistics. We also used the hyperparameter settings for the RF regression algorithm recommended by Raynal *et al.* (2016, p. 16). For the ABC-DNN approach, we created a fully connected neural network with three hidden layers, each hidden layer consisting of 100 units, and we trained it to minimize the MSE loss between true and predicted parameters (Jiang *et al.*, 2015). For the

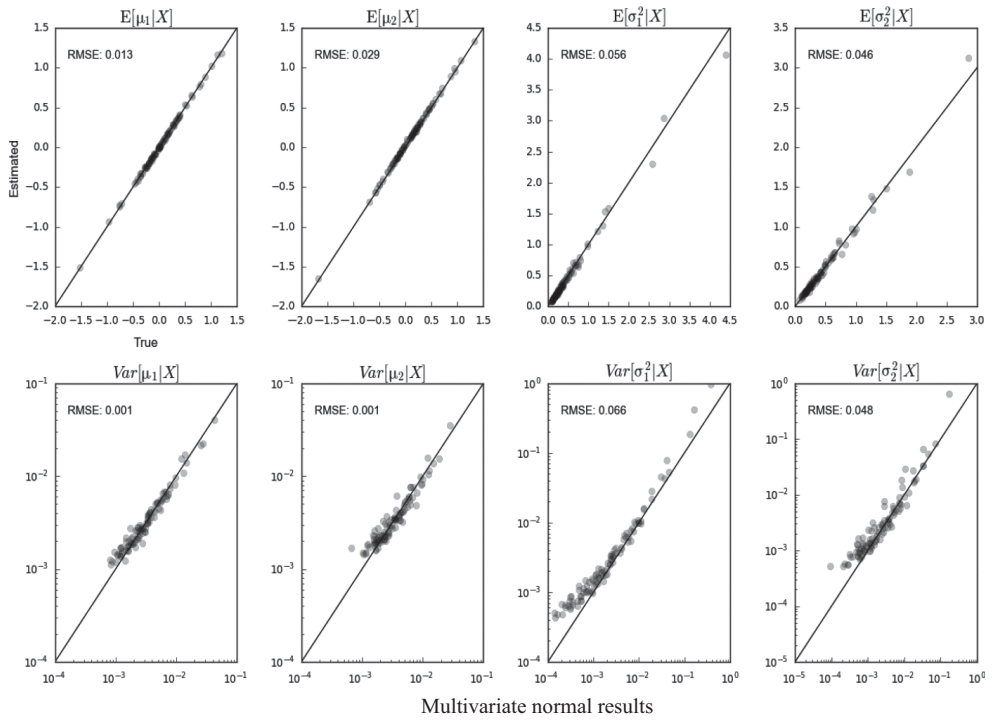


Figure 2. Comparison of the estimates obtained by the CNN model with the true parameter values on the multivariate normal toy example. True parameter values are plotted on the x -axis, and corresponding estimates on the y -axis. The first row depicts the posterior expectations of the model parameters. Posterior variances are depicted in the second row. The RMSE is also given for each estimate.

subsequent rejection sampling scheme, we used the estimates of the trained neural network as summary statistics $\eta(\tilde{\mathbf{x}})$ with tolerance level $\epsilon = .01$, and Euclidian distance as the distance function $d(\theta, \hat{\theta})$. Posterior moments were computed by considering accepted samples from the approximate posterior distribution $p_\epsilon(\theta|\eta(\tilde{\mathbf{x}}))$.

Both graphical and numerical evaluations of the performance of our CNN model suggest a decent approximation of the first two posterior moments of the MVN model. We observe a slight overestimation of the posterior variance $\text{Var}[\text{diag}(\Sigma)|\mathbf{x}]$ in the lower and upper regions of the parameter space. Moreover, our method clearly outperforms the ABC–DNN approach and yields approximations comparable to the ABC–RF approach. However, it is important to note that, in this example, it is easy to come up with summary-statistics that capture all the relevant information, thus giving the ABC–RF approach a clear advantage.

3.2. Toy example 2: Bayesian regression

We now turn to a slightly more complex example based on Zellner’s hierarchical regression model (Marin & Robert, 2014, chapter 3; Raynal *et al.*, 2016). Given a simulated $m \times d$ design matrix of covariates $\mathbf{X} = (x_1, x_2, \dots, x_d)$, the hierarchical conjugate model is defined by:

Table 1. Comparison of parameter estimation methods on the multivariate normal toy example with the RMSE metric

Method	$\mathbb{E}[\mu_1 \mathbf{x}]$	$\mathbb{E}[\mu_2 \mathbf{x}]$	$\mathbb{E}[\sigma_1^2 \mathbf{x}]$	$\mathbb{E}[\sigma_2^2 \mathbf{x}]$	$\text{Var}[\mu_1 \mathbf{x}]$	$\text{Var}[\mu_2 \mathbf{x}]$	$\text{Var}[\sigma_1^2 \mathbf{x}]$	$\text{Var}[\sigma_2^2 \mathbf{x}]$	$\mathbb{E}[\mu_1 \mathbf{x}]$
Deep Inference	0.013	0.029	0.056	0.046	0.001	0.001	0.066	0.048	0.013
ABC-DNN	0.046	0.047	0.361	0.334	0.009	0.1	0.857	0.368	0.046
ABC-RF	0.013	0.014	0.028	0.015	0.002	0.002	0.012	0.013	0.013

Note. RMSE scores given in bold highlight the best approximation for each parameter (i.e., each column).

$$\begin{aligned}
 \sigma^2 &\sim IG(a_0, b_0) \\
 \boldsymbol{\beta}|\sigma^2 &\sim \mathcal{N}_d(\mathbf{0}, n\sigma^2(\mathbf{X}^T\mathbf{X})^{-1}) \\
 y|\boldsymbol{\beta}, \sigma^2 &\sim \mathcal{N}_m(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I}),
 \end{aligned} \tag{12}$$

where $IG(a_0, b_0)$ denotes an inverse gamma distribution with shape parameter a_0 and scale parameter b_0 , the regression weights vector $\boldsymbol{\beta} = (\beta_1, \beta_2, \dots, \beta_d)^T$ follows a d -dimensional normal distribution, and the likelihood follows an m -dimensional normal distribution with an isotropic covariance matrix. As in the previous example, the conjugacy property of the model allows for an analytical computation of the marginal posterior distributions of $\boldsymbol{\beta}$ and σ^2 (see Marin & Robert, 2014, chapter 3).

For this example, we set $d = 2$, $a_0 = 4$, $b_0 = 3$, $m = 100$ and generate a reference table containing 100,000 rows of raw i.i.d. samples $\{\tilde{\mathbf{x}}_k^{(i)}\}_{k=1}^{100}$ from the model defined by equation (12), where each element $\tilde{\mathbf{x}}_k^{(i)}$ is an ordered triple containing the two covariates and the outcome generated using the covariates $(x_1^{(i)}, x_2^{(i)}, y^{(i)})_k$. The CNN is then used to predict the corresponding data-generating distribution parameters $\{\beta_1^{(i)}, \beta_2^{(i)}, \sigma^2^{(i)}\}$. In contrast to the previous example, each kernel in the input layer should learn to encode a particular combination of covariates and outcome (tripe) leading to a set of parameters instead of a single data point.

As in the previous toy example, we trained the model on a training set of 99,000 simulated datasets, holding out the remaining 1,000 datasets for validation. Again, for the purpose of evaluation, we generated 100 independent datasets \mathbf{X}_{test} , on which we computed the true values of $\mathbb{E}[\boldsymbol{\beta}|\mathbf{x}]$, $\mathbb{E}[\sigma^2|\mathbf{x}]$, $\text{Var}[\boldsymbol{\beta}|\mathbf{x}]$, and $\text{Var}[\sigma^2|\mathbf{x}]$. Figure 3 compares the analytically computed posterior moments with the estimates obtained from our CNN model.

Once again, we observe a very low approximation error (RMSE) of the posterior expectations of the regression weights $\boldsymbol{\beta}$. The approximation error of the posterior expectation of the residual variance σ^2 is slightly higher. The approximation of the posterior variance of $\boldsymbol{\beta}$ is also good, whereas we make the observation that the posterior variance of σ^2 is overestimated.

Table 2 compares the performance of our approach to the performance of the ABC-RF and ABC-DNN methods as specified by Algorithms 2 and 3 in terms of RMSE. As in the previous example, we used the same training and test sets for all three methods. The hyperparameter settings of the algorithms were identical to those used in the previous example. The summary statistics we computed for the input to the ABC-RF were the same as those employed by Raynal *et al.* (2016): the maximum likelihood estimates of $\boldsymbol{\beta}$, the

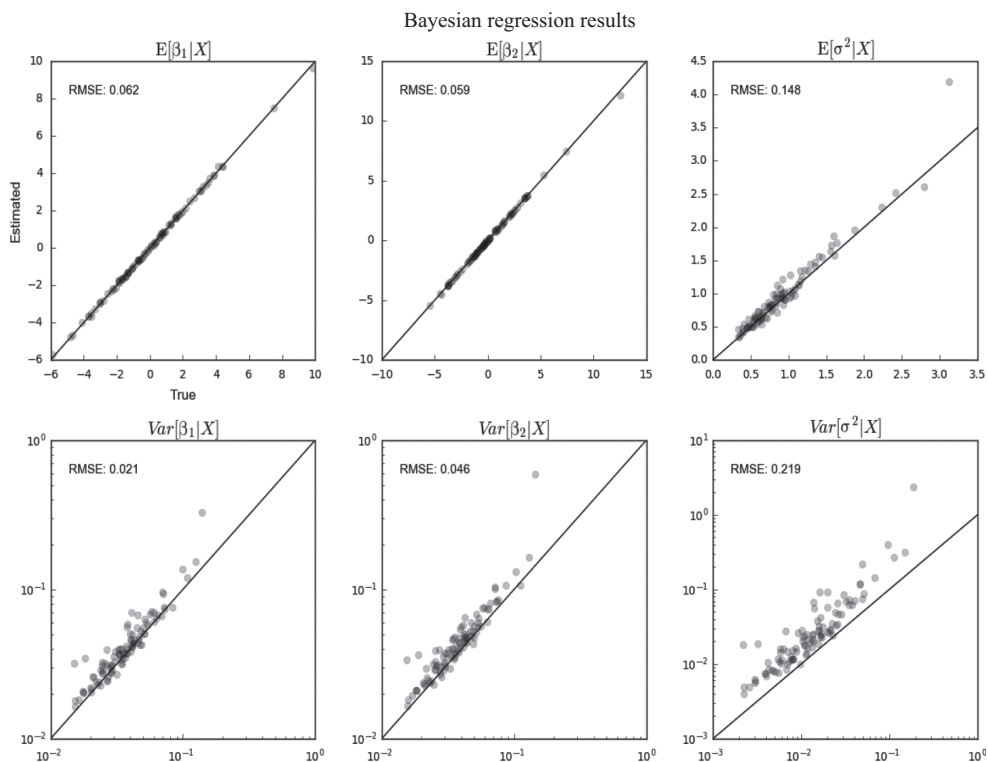


Figure 3. Comparison of the estimates obtained by the CNN model with the true posterior parameter values on the Bayesian linear regression example. True parameter values are plotted on the x -axis, and corresponding estimates on the y -axis. The first row depicts the posterior expectations of the model parameters. Posterior variances are depicted in the second row. The RMSE is also given for each estimate.

Table 2. Comparison of parameter estimation methods on the Bayesian regression toy example with the RMSE metric

Method	$E[\beta_1 \mathbf{x}]$	$E[\beta_2 \mathbf{x}]$	$E[\sigma^2 \mathbf{x}]$	$Var[\beta_1 \mathbf{x}]$	$Var[\beta_2 \mathbf{x}]$	$Var[\sigma^2 \mathbf{x}]$
DeepInference	0.062	0.059	0.148	0.021	0.046	0.219
ABC-DNN	0.456	0.674	0.369	0.232	0.203	0.076
ABC-RF	0.093	0.087	0.062	0.055	0.099	0.037

Note. RMSE scores given in bold highlight the best approximation for each parameter (i.e., each column).

residual sum of squares, the empirical covariance and correlation between y and \mathbf{X} , the sample mean, variance, and median of y , resulting in a total of 10 metrics. Network architecture, tolerance level and distance function for the ABC-DNN did not differ from the previous example.

Compared to the ABC-RF approach, our method achieves a slightly better approximation of the posterior expectation and variance of β , whereas the posterior expectation and variance of σ^2 are better approximated by the ABC-RF approach. Compared to the

ABC–DNN approach, our model yields better approximations for all parameters except for the variance of σ^2 .

3.3. Example 3: Leaky competing accumulator model

As a final example, we apply our approach to a realistic model instance from the cognitive modelling literature: the LCA model (Miletić *et al.*, 2017; Usher & McClelland, 2001). The LCA is rooted in the theoretical framework of sequential sampling models (SSMs), which attempt to describe perceptual decision-making via the mechanism of noisy evidence accumulation. SSMs assume that, for each decision path, there is a latent process at work, termed an accumulator, which samples evidence from sensory input through time. A decision process terminates when an accumulator exceeds a certain threshold of activation, meaning that evidence for a given option has culminated in a decision in favour of the given option. The inherent stochasticity of evidence accumulation ensures variation in predicted response times given identical stimuli and also allows for an explanation of the occurrence of erroneous responses.

Despite the fact that different researchers have proposed different variants of SSMs (Ratcliff & McKoon, 2008; Usher & McClelland, 2001; Voss, Nagler, & Lerche, 2013), most flavours of SSMs are specified in terms of a finite set of parameters, each interpretable in light of some assumed neurocognitive process or the property of such a process. Importantly, all SSMs provide a way to decompose empirical response times into a decision time component and a non-decision time component, the latter accounting for pre-decisional perceptual (encoding time) and post-decisional motor processes (response execution). Parameters of SSMs are typically estimated from empirical response time distributions via an assumption- and application-dependent estimation procedure (i.e., ML estimation).

In particular, the LCA defines evidence accumulation via the following stochastic differential equation (Miletić *et al.*, 2017):

$$\begin{aligned} dx_i &= \left[I_i - kx_i - \beta \sum_{i' \neq i} x_{i'} \right] \frac{dt}{\tau} + \xi_i \sqrt{\frac{dt}{\tau}} \\ x_i &= \max(0, x_i). \end{aligned} \quad (13)$$

In this model equation, dx_i denotes the change in activation of accumulator i , I_i is the input accumulator i receives, k is leakage, β is inhibition, $\frac{dt}{\tau}$ denotes the time-step size and ξ represents Gaussian noise.

The LCA model assumes one accumulator process for each option in a decision-making task. Each accumulator competes with all others for reaching a common threshold Z . Observed response times are obtained by an addition of constant non-decision time NDT prior to evidence accumulation and following the reaching of a threshold by an accumulator. The leakage and inhibition parameters are unique to LCA, and they appear to be motivated by neuroscientific observations demonstrating parallels to evidence accumulation on a neural level. Leakage describes the amount of information lost at each time point by an accumulator. Inhibition describes the proportion of an accumulator's activation that suppresses the activation of competing accumulators.

We chose to apply our approach to the LCA model because the model parameters do not have a known likelihood function. Consequently, most parameter estimation procedures resort to a simulation-based approach (Miletić *et al.*, 2017; Turner & Sederberg, 2014).

For the following example, we simulate response times from the LCA model and attempt to recover the data-generating parameters. We place the same prior distributions over the model parameters as described in Miletić *et al.* (2017):

$$\begin{aligned}
 \Delta I &\sim u(0.05, 0.3) \\
 I &\sim u(0.8, 1.2) \\
 k &\sim u(1, 8) \\
 \beta &\sim u(1, 8) \\
 Z &\sim u(0.05, 0.25) \\
 \text{NDT} &\sim u(0.200, 0.500).
 \end{aligned} \tag{14}$$

Note that by choosing uniform priors, the Bayesian approximation of the posterior mean and variance would coincide with the maximum likelihood estimation of the mean and variance, as $p(\theta)$ acts solely as a constant multiplicative factor on the likelihood.

For the current example, we fixed the noise variance to $\xi = 0.1$. We generated a reference table with 500,000 simulated datasets, sampling parameters from equation (14) and generating 1,000 response times according to equation (13). In order to allow for a comparison between methods, we simulated the same number of response times as in Miletić *et al.* (2017). We held out 1,000 datasets for validation and generated 500 further independent datasets \mathbf{X}_{test} for evaluation. We implemented a deeper CNN architecture consisting of five FC convolutional layers with the following increasing number of kernels in each layer: (64, 64, 128, 128, 128).

With these settings, the model receives as input each simulated response time distribution $\{\tilde{\mathbf{x}}_k^{(i)}\}_{k=1}^{1,000}$ as well as the information about which accumulator ‘won’ the threshold competition, encoded as a one-hot vector. It is then necessary to predict each of the following LCA parameters: $(\Delta I^{(i)}, I^{(i)}, k^{(i)}, \beta^{(i)}, Z^{(i)}, \text{NDT}^{(i)})$. Figure 4 compares the data-generating parameters with the approximations obtained by the CNN model.

We observe that our model is able to recover the parameters ΔI , NDT and Z , reasonably well, indexed by a low RMSE and a high correlation between CNN estimates and actual data-generating parameters. The recovery of the leakage and inhibition parameters, k and β , respectively, is less accurate, but seems to be improvable with the addition of more simulated samples or a deeper architecture with more layers. The CNN is unable to recover the input parameter I , as it appears as if the data did not contain any useful information about this particular parameter.

Table 3 lists the RMSE values and correlations between true and recovered parameter values obtained by our method. It also lists the performance of methods based on probability density estimation (PDA), such as maximum likelihood estimation and Bayesian estimation performed via differential evolution (DE) MCMC, as observed in the recent parameter recovery study by Miletić *et al.* (2017). The PDA and DE–MCMC results are taken as reported by Miletić *et al.* (2017) on the 1,000-trial datasets. Note that these results are not to be regarded as a direct comparison between the two methods, as the methods were not applied to the same test data, and also not tuned for optimal hyperparameter settings. It is worth noting, however, that the relative estimation quality of our method appears similar to that of the PDA-based fitting methods but it entails no simulations at inference time.

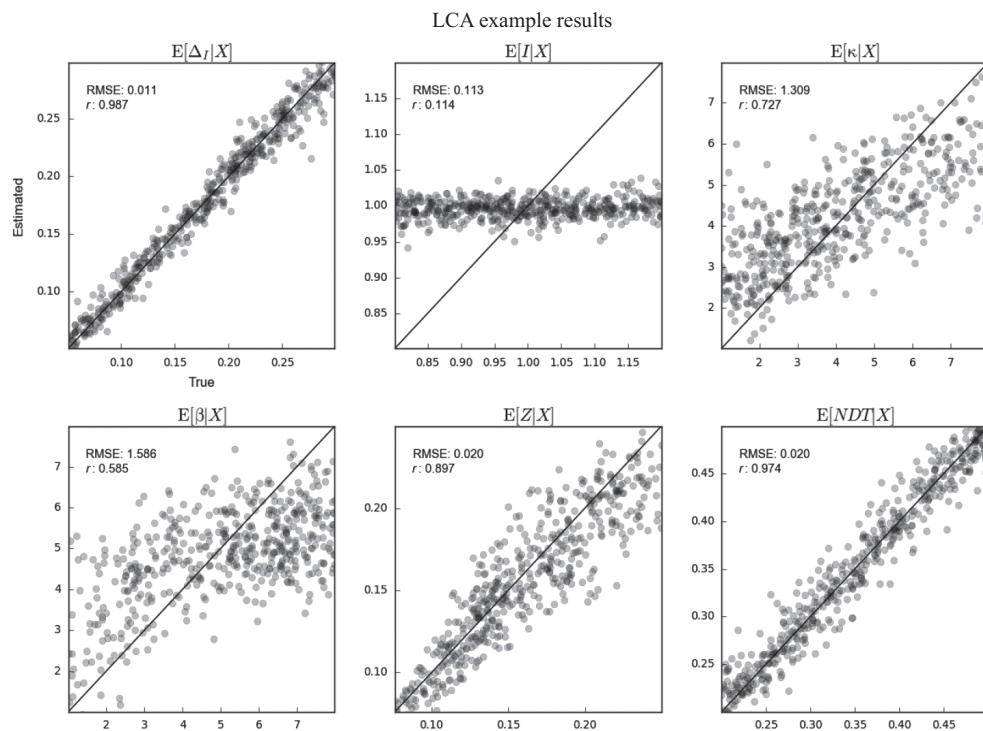


Figure 4. Depiction of parameter recovery of the LCA model. Values of the data-generating parameters are plotted on the x-axis, and corresponding CNN estimates on the y-axis. Both the RMSE and Pearson's correlation coefficient (r) are also reported.

4. Discussion

Developing effective procedures for estimating parameters of complex models is a vital endeavour of any principled science. In this paper, we have proposed a novel method for ABC parameter estimation using CNNs to simultaneously approximate the means and variances of different posterior parameter distributions. We have verified the usefulness of our approach on two toy examples and a 'real' example from the cognitive modelling literature.

Compared to previous ABC methods rooted in the supervised machine learning approach, our method is truly end-to-end, as it requires no computation or selection of summary statistics and no other feature-engineering activities. The possibility of representing raw simulated samples in a grid-like data structure allows us to exploit the advantages of modern CNNs over traditional DNNs, such as sparse interactions, parameter sharing and equivariant representations (Goodfellow *et al.*, 2016).

'Sparsity of interactions' refers to the fact that in place of matrix multiplication, CNNs utilize convolution with a small kernel size, which make CNNs much more efficient in terms of memory usage and computational efficiency. 'Parameter sharing' is directly related to the property of sparse interactions. It refers to the fact that the parameters of CNN layers are not tied to particular locations of the input, but rather used at every position of the input. Parameter sharing allows the CNN to learn a single set of parameters that capture this particular data point or series of data points. The third property,

Table 3. RMSE metric for the parameter estimates of the LCA model and Pearson’s correlation coefficient r between true and recovered parameter values. The performances of the ML-PDA and DE-MCMC methods on the same parameter recovery task are also included as a reference point for comparison

Method	RMSE					Correlation r				
	ΔI	I	k	β	Z	NDT	ΔI	I	k	β
DeepInference	0.011	0.113	1.309	1.586	0.02	0.02	0.99	0.11	0.73	0.59
ML-PDA	0.024	0.508	2.487	3.216	0.045	0.044	0.96	0.07	0.4	0.27
DE-MCMC	0.021	1.073	2.916	4.742	0.03	0.02	N/A	N/A	N/A	N/A

Note. RMSEs and correlations for the ML-PDA and DE-MCMC methods with sample of size $N = 1,000$ are reproduced from Miletic *et al.* (2017). Bold values indicate which algorithm performs best for the respective parameter.

'equivariance', is a consequence of the previous two. It implies that shifting the distribution (in the case of time-series), or shuffling the distribution (in the case of i.i.d. sequences), does not dramatically change the representation of the input learned by a particular kernel.

Taken together, these properties make CNNs an excellent choice for implicitly learning summary statistics from simulated samples. In addition, they enable the network to work with inputs of variable size, making the size of the simulated samples $\tilde{\mathbf{x}}$ independent of the size of the observed data \mathbf{x} , and thus overcoming a further huge limitation of all previous supervised machine learning approaches.

The use of CNNs in the way described in this paper allows for transfer learning to enter the field of likelihood-free inference. Transfer learning refers to the process of re-using a pre-trained network for the same or similar task. In other words, once trained, a CNN can be shared across and used as a fast, likelihood-free parameter estimator in a given cognitive domain (i.e., to estimate diffusion model parameters from reaction times in a single forward pass). Moreover, our approach is not confined to instances of models with intractable likelihoods. Even for models with known likelihood, our approach might offer a good alternative to computationally costly sampling procedures. It is worth noting that the architecture of the neural network depends on the structure of the experimental task or modelling domain. Different tasks/domains might require different input layers or even different types of neural networks (e.g., recurrent neural networks for time-series data). However, within a single research domain (e.g., analysis of response times), the structure of the data is expected to remain largely the same, and therefore a single model can be reused with little or no modifications to address questions of parameter inference. Moreover, the application of our approach as a general black-box estimator of hard-to-estimate parameters of various white-box cognitive models is only limited by the capability of the white-box model to generate representative samples (i.e., a reference table). Needless to say, this transferability is a unique advantage of our deep learning model compared with different architectures. Further work should explore various different architectures for other models or research domains.

Throughout our experiments, the approximation of the posterior means was shown to be very accurate. The posterior variance tends to be slightly overestimated – a fact that requires further attention, but, at the very least, implies that our model does not suffer from overconfidence. One possibility for the overestimation of uncertainty is that, because of the finite sample, the estimated $\sigma^2(\tilde{\mathbf{x}})$ term fails to capture the true posterior variance. Another possibility is that, because the heteroscedastic loss criterion implies an underlying Gaussian model, deviations of the true posterior from the Gaussian model will cause $\sigma^2(\tilde{\mathbf{x}})$ to deviate from the true posterior variance. Nevertheless, the deviations we observed are not dramatic, and further experiments on models with known posteriors should reveal the true extent of the approximation error.

As an attempt to circumvent the shortcomings of ABC rejection algorithms (curse of dimensionality, tolerance level tuning and selection of distance function), our approach does not provide a way to sample from the full joint posterior distribution, but rather summarizes the distribution with its first two moments. However, if one wishes, one can still use the estimates obtained by the CNN as summary statistics in a rejection procedure (Jiang *et al.*, 2015), with the added benefit of having the approximation of the second posterior moment beside the first moment. Furthermore, a full approximate distribution can be computed either by simply sampling from a normal distribution parametrized by the predicted means and variances or by using this distribution as a more informative prior distribution in a usual rejection sampling approach to drastically improve the efficiency.

It is also straightforward to extend our approach to the problem of ABC model selection. In machine-learning terms, instead of performing regression, we simply need to frame the problem as a classification task, in which the network is trained to classify the label of the model used to generate a particular data sample (see Pudlo *et al.*, 2015 for a random-forest based approach). In future work, we intend to extend our open-source ABC software ABrox (see Mertens *et al.*, 2018) with DeepInference capabilities.

4.1. Conclusions

In this paper, we have proposed a method for reliable likelihood-free inference using FCNs to automatically learn optimal summary statistics from raw samples. We have obtained accurate estimates of the first two moments of the posterior parameter distribution on two toy examples and a cognitive model of choice reaction times. Furthermore, our method exhibits the following unique features compared with previous supervised learning approaches to ABC:

1. an end-to-end architecture eliminating the need to manually hand-craft summary statistics of the data distribution;
2. fast training due to simultaneous parameter estimation;
3. re-usable models due to the possibility of working with variable-sized inputs.

Further research should test the utility of our approach on various real-world examples from the cognitive modelling literature, which require non-standard solutions or are too expensive to compute with sampling-based procedures. Another future avenue for our approach is likelihood-free model choice, which would require a different optimization procedure for ‘labelling’ the data with the correct model and quantifying uncertainty at the same time.

Acknowledgements

We thank Steven Miletic and Leendert van Maanen, as well as Louis Raynal and Jean-Michel Marin for providing us with their R code for simulating data from the LCA and regression models, respectively. We also thank Alica Bucher for helping us with the visualization of our model architecture. S. T. Radev and U. K. Mertens contributed equally to this work.

References

- Blum, M. G. (2010). Approximate Bayesian computation: A nonparametric perspective. *Journal of the American Statistical Association*, 105, 1178–1187. <https://doi.org/10.1198/jasa.2010.tm09448>
- Blum, M. G., & François, O. (2010). Non-linear regression models for Approximate Bayesian Computation. *Statistics and Computing*, 20, 63–73. <https://doi.org/10.1007/s11222-009-9116-0>
- Breiman, L. (2001). Random forests. *Machine Learning*, 45, 5–32. <https://doi.org/10.1023/a:1010933404324>
- Chollet, F. (2017). *Deep learning with Python*. Shelter Island, New York: Manning Publications Co.
- Cybenko, G. (1989). 1Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals and Systems*, 2, 303–314. <https://doi.org/10.1007/bf02551274>
- Frazier, D. T., Martin, G., Robert, C. P., & Rousseau, J. (2016). Asymptotic properties of approximate Bayesian computation. *Biometrika*.

- Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., & Rubin, D. B. (2013). *Bayesian data analysis*. Boca Raton, FL: CRC Press.
- Goodfellow, I., Bengio, Y., Courville, A., & Bengio, Y. (2016). *Deep learning* (Vol. 1). Cambridge, UK: MIT Press.
- Hornik, K. (1991). Approximation capabilities of multilayer feedforward networks. *Neural Networks*, 4, 251–257. [https://doi.org/10.1016/0893-6080\(91\)90009-t](https://doi.org/10.1016/0893-6080(91)90009-t)
- Jiang, B., Wu, T.-y., Zheng, C., & Wong, W. H. (2015). Learning summary statistic for approximate Bayesian computation via deep neural network. *Statistica Sinica*, 1595–1618.
- Kendall, A., & Gal, Y. (2017). *What uncertainties do we need in Bayesian deep learning for computer vision?* Paper presented at the Advances in Neural Information Processing Systems (NIPS 2017), Long Beach, CA. *arXiv preprint arXiv:1703.04977*.
- Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic optimization. In *Proc. 3rd Int. Conf. Learn. Representations*.
- Lee, M. D. (2008). Three case studies in the Bayesian analysis of cognitive models. *Psychonomic Bulletin & Review*, 15, 1–15. <https://doi.org/10.3758/pbr.15.1.1>
- Lin, M., Chen, Q., & Yan, S. (2013). Network in network. *arXiv preprint arXiv:1312.4400*.
- Marin, J.-M., Pudlo, P., Robert, C. P., & Ryder, R. J. (2012). Approximate Bayesian computational methods. *Statistics and Computing*, 22, 1167–1180. <https://doi.org/10.1007/s11222-011-9288-2>
- Marin, J.-M., & Robert, C. P. (2014). *Bayesian essentials with R* (Vol. 48). Berlin: Springer.
- Mertens, U. K., Voss, A., & Radev, S. (2018). ABrox—A user-friendly Python module for approximate Bayesian computation with a focus on model comparison. *PLoS ONE*, 13, e0193981. <https://doi.org/10.1371/journal.pone.0193981>
- Miletić, S., Turner, B. M., Forstmann, B. U., & van Maanen, L. (2017). Parameter recovery for the leaky competing accumulator model. *Journal of Mathematical Psychology*, 76, 25–50. <https://doi.org/10.1016/j.jmp.2016.12.001>
- Nix, D. A., & Weigend, A. S. (1994, June). Estimating the mean and variance of the target probability distribution. In *Neural Networks, 1994. IEEE World Congress on Computational Intelligence., 1994 IEEE International Conference (1, pp. 55–60)*.
- O'Shea, K., & Nash, R. (2015). An introduction to convolutional neural networks. *arXiv preprint arXiv:1511.08458*.
- Pritchard, J. K., Seielstad, M. T., Perez-Lezaun, A., & Feldman, M. W. (1999). Population growth of human Y chromosomes: A study of Y chromosome microsatellites. *Molecular Biology and Evolution*, 16, 1791–1798. <https://doi.org/10.1093/oxfordjournals.molbev.a026091>
- Pudlo, P., Marin, J.-M., Estoup, A., Cornuet, J.-M., Gautier, M., & Robert, C. P. (2015). Reliable ABC model choice via random forests. *Bioinformatics*, 32, 859–866. <https://doi.org/10.1093/bioinformatics/btv684>
- Ratcliff, R., & McKoon, G. (2008). The diffusion decision model: Theory and data for two-choice decision tasks. *Neural Computation*, 20, 873–922. <https://doi.org/10.1162/neco.2008.12.06.420>
- Raynal, L., Marin, J.-M., Pudlo, P., Ribatet, M., Robert, C. P., & Estoup, A. (2016). ABC random forests for Bayesian parameter inference. *arXiv preprint arXiv:1605.05537*.
- Rubin, D. B. (1984). Bayesianly justifiable and relevant frequency calculations for the applied statistician. *Annals of Statistics*, 12, 1151–1172. <https://doi.org/10.1214/aos/1176346785>
- Schoot, R., Kaplan, D., Denissen, J., Asendorpf, J. B., Neyer, F. J., & Aken, M. A. (2014). A gentle introduction to Bayesian analysis: Applications to developmental research. *Child Development*, 85, 842–860. <https://doi.org/10.1111/cdev.12169>
- Sheehan, S., & Song, Y. S. (2016). Deep learning for population genetic inference. *PLoS Computational Biology*, 12, e1004845. <https://doi.org/10.1371/journal.pcbi.1004845>
- Tavaré, S., Balding, D. J., Griffiths, R. C., & Donnelly, P. (1997). Inferring coalescence times from DNA sequence data. *Genetics*, 145, 505–518.

- Toni, T., Welch, D., Strelkowa, N., Ipsen, A., & Stumpf, M. P. (2009). Approximate Bayesian computation scheme for parameter inference and model selection in dynamical systems. *Journal of the Royal Society Interface*, 6(31), 187–202. <https://doi.org/10.1098/rsif.2008.0172>
- Turner, B. M., & Sederberg, P. B. (2014). A generalized, likelihood-free method for posterior estimation. *Psychonomic Bulletin and Review*, 21, 227–250. <https://doi.org/10.3758/s13423-013-0530-0>
- Turner, B. M., & Van Zandt, T. (2012). A tutorial on approximate Bayesian computation. *Journal of Mathematical Psychology*, 56, 69–85. <https://doi.org/10.1016/j.jmp.2012.02.005>
- Usher, M., & McClelland, J. L. (2001). The time course of perceptual choice: The leaky, competing accumulator model. *Psychological Review*, 108, 550. <https://doi.org/10.1037/0033-295x.108.3.550>
- Voss, A., Nagler, M., & Lerche, V. (2013). Diffusion models in experimental psychology: A practical introduction. *Experimental Psychology*, 60, 385. <https://doi.org/10.1027/1618-3169/a000218>

Received 7 July 2018; revised version received 11 December 2018