

# JGR Solid Earth

## RESEARCH ARTICLE

10.1029/2022JB024703

### Key Points:

- A machine learning approach for the probabilistic solution of inverse problems by directly estimating posterior statistics of any continuous or discrete feature of the posterior distribution
- Allows the use of complex prior information and noise models
- Demonstrated non-linear probabilistic inversion of airborne electromagnetic; enables analysis of more than  $10^5$  1D soundings per second

### Correspondence to:

T. M. Hansen,  
tmeha@geo.au.dk

### Citation:

Hansen, T. M., & Finlay, C. C. (2022). Use of machine learning to estimate statistics of the posterior distribution in probabilistic inverse problems—An application to airborne EM data. *Journal of Geophysical Research: Solid Earth*, 127, e2022JB024703. <https://doi.org/10.1029/2022JB024703>

Received 3 MAY 2022  
Accepted 17 OCT 2022

### Author Contributions:

**Conceptualization:** T. M. Hansen, C. C. Finlay  
**Formal analysis:** T. M. Hansen, C. C. Finlay  
**Funding acquisition:** T. M. Hansen  
**Investigation:** T. M. Hansen  
**Methodology:** T. M. Hansen, C. C. Finlay  
**Resources:** T. M. Hansen  
**Software:** T. M. Hansen  
**Validation:** T. M. Hansen, C. C. Finlay  
**Visualization:** T. M. Hansen  
**Writing – original draft:** T. M. Hansen  
**Writing – review & editing:** T. M. Hansen, C. C. Finlay

© 2022 The Authors.

This is an open access article under the terms of the [Creative Commons Attribution-NonCommercial License](https://creativecommons.org/licenses/by-nc/4.0/), which permits use, distribution and reproduction in any medium, provided the original work is properly cited and is not used for commercial purposes.

# Use of Machine Learning to Estimate Statistics of the Posterior Distribution in Probabilistic Inverse Problems—An Application to Airborne EM Data

T. M. Hansen<sup>1</sup>  and C. C. Finlay<sup>2</sup> 

<sup>1</sup>Department of Geoscience, Aarhus University, Aarhus C, Denmark, <sup>2</sup>DTU Space, Technical University of Denmark, Lyngby, Denmark

**Abstract** The solution to a probabilistic inverse problem is the posterior probability distribution for which a full analytic expression is rarely possible. Sampling methods are therefore often used to generate a sample from the posterior. Decision-makers may be interested in the probability of features related to model parameters (e.g., existence of pollution or the cumulative clay thickness) rather than the individual realizations themselves. Such features and their associated uncertainty, are simple to compute once a sample from the posterior distribution has been generated. However, sampling methods are often associated with high computational costs, especially when the prior and posterior distribution is non-trivial (non-Gaussian), and when the inverse problem is non-linear. Here we demonstrate how to use a neural network to directly estimate posterior statistics of continuous or discrete features of the posterior distribution. The method is illustrated on a probabilistic inversion of airborne EM data from Morrill Nebraska, where the forward problem is nonlinear and the prior information is non-Gaussian. Once trained the application of the network is fast, with results similar to those obtained using much slower sampling methods.

**Plain Language Summary** Probabilistic inversion is in principle ideal as a method for combining available information about geo-models, in a way that will allow detailed risk analysis and hypothesis testing based on available information. However, practical use of such methods has historically been limited because they (a) require computationally expensive numerical algorithms and (b) typically rely on relatively simple assumptions about the model. Machine learning-based methods, such as neural networks, provide an alternative approach to probabilistic inversion. We discuss using neural networks to estimate in principle any statistics about continuous or discrete statistical features of the combined information, such as “What is the probability that a specific lithology exists below a certain depth?”. A focus is on the use of realistic/complex assumptions. As an example, the method is demonstrated on a probabilistic inversion of airborne electromagnetic data, and it is demonstrated as accurate, fast, and allows analysis of many (>100,000) 1D soundings per second, making it applicable to very large data sets.

## 1. Introduction

A key challenge in geoscience is that of combining different kinds of geo-information into one geo-model, typically describing the subsurface. This information can be direct information about geological processes, and spatial variability, or it can be indirect information from measurements of properties related to the subsurface, such as geophysical data. Ideally, when such a geo-model has been established, one should be able to quantify information about specific features related to the geo-model, consistent with all information.

This integration of geo-information is typically solved using inverse problem theory (Menke, 2012; Tarantola & Valette, 1982a). Fast deterministic methods exist and have been widely used. For such methods, the goal is to obtain one optimal model, such as the simplest possible model, consistent with available information, typically in the form of observed data (Constable et al., 1987; Menke, 2012; Tikhonov, 1963). In practice, in part due to noise on data and model nonlinearities and imperfections, infinitely many models exist that will be consistent with data, and the deterministic approach can in general not account properly for such uncertainty.

Probabilistic inversion methods can, in principle, take into account arbitrarily complex information, and integrate the information into one consistent model, as given by the posterior probability distribution. A full analytic expression of the posterior distribution is rarely possible. Instead, sampling methods can be used to generate a

sample of the posterior, which is a collection of realizations drawn from the posterior distribution. From such a sample, the posterior statistics of any feature related to the model parameters can be computed. The probabilistic approach is therefore ideal for decision-makers for uncertainty quantification, as it allows probabilistic analysis and risk assessment consistent with available information.

The main obstacle to applying the probabilistic methodology in practice is that sampling methods are computationally very demanding (Hastings, 1970; Mosegaard & Tarantola, 1995). In some cases information about the posterior distribution can be used, for example, to construct a proposal distribution similar to the posterior distribution (Khoshkholgh et al., 2022), or in the form of information about the gradient of the posterior distribution (Fichtner et al., 2018), which can lead to more efficient sampling algorithms. Such cases are however often based on rather simplistic choices of prior models. In general, sampling-based methods typically require sampling or evaluation of a prior model, evaluation of the physical forward response(s), and evaluation of a noise model, many times.

One approach for reducing the computational requirements is to make use of fast approximate forward modeling. This can be related to using simplified 1D forward modeling as opposed to 3D forward models, or by using approximate physical models, which leads to modeling errors that should be accounted for (Hansen et al., 2014; Köpke et al., 2018; Madsen & Hansen, 2018).

Machine learning algorithms, which are fast to evaluate once trained, have also been used to approximate the forward modeling (Bording et al., 2021; Conway et al., 2019; Hansen & Cordua, 2017; Moghadas et al., 2020). Unsupervised machine learning methods, for example, Generative adversarial neural networks, have been used more generally as a means of representing features in a prior dataset; once trained, these provide an efficient means of rapidly generating many prior realizations (Laloy et al., 2018; Mosser et al., 2017, 2020).

Attempts have also been made to use machine learning methods to learn a mapping from data to a model that can directly solve the inverse problem. Röth and Tarantola (1994) were amongst the first to solve an inverse problem in this way using a multilayer perceptron neural network and demonstrated an application of inversion of reflection seismic data to obtain single estimates of 1D velocity profiles. Recently, several authors have further explored this approach for directly solving a geophysical inverse problem, making use of convolutional neural networks (Bai et al., 2020; Moghadas, 2020; Puzryev & Swidinsky, 2021). A drawback of such methods is that, as in the deterministic solution of an inverse problem, they estimate only a single model, typically without accounting for uncertainty in geophysical data, and do not quantify the uncertainty on the predicted model parameters.

An important step toward finding probabilistic solutions to inverse problems using neural networks was made by Devilee et al. (1999) who considered training data sets consisting of realizations from the prior distribution and the corresponding forward simulated data with and without noise. They then used neural networks to learn a set of statistics about each model parameter, including median and equidistant histogram estimators. Meier et al. (2007) extended this work and used a mixture density network (MDN) to estimate the parameters of a Gaussian mixture model representing a parametric distribution that approximated the 1D marginal posterior distribution, and applied it to the problem of estimating global crustal thickness maps, comparing to results obtained using a Monte Carlo based sampling method. Several other applications of MDN to approximate the posterior distribution, for different geophysical problems, have followed Earp and Curtis (2020); Earp et al. (2020); Shahraeeni and Curtis (2011); de Wit et al. (2013).

Zhang and Curtis (2020a) argue that it may be problematic to apply such MDN's for higher dimensional inverse problems, and suggest using variational inference (Blei et al., 2017) to estimate the mean and standard deviation of the (non-Gaussian) posterior distribution in an example of a seismic tomographic inverse problem. This method has been developed further for variational full waveform inversion and tomographic inversion using normalizing flows (Zhang & Curtis, 2020b; Zhao et al., 2022). In all these cases a uniform prior was assumed.

Attempts have also been made to use so-called invertible neural networks to simultaneously estimate both the forward and inverse mapping between data and model parameters (Ardizzone et al., 2018). This approach, which has recently been applied to geophysical data by Zhang and Curtis (2021), allows the generation of multiple realizations of the posterior distribution, from which properties of the posterior distribution can be estimated, although constructing invertible neural networks involves more work than traditional neural networks and involves compromises related to the flexibility of the network.

Here we present a method where the goal is not primarily to estimate the marginal 1D posterior distribution (as in works based on Earp et al. [2020]; Meier et al. [2007]; Shahraneini & Curtis [2011]; de Wit et al. [2013]). Instead, we propose and demonstrate a machine learning-based method that provides direct estimates of any desired statistical property (continuous or discrete) of the posterior distribution, including any feature or property that can be computed from realizations of an, in principle, arbitrarily complex, prior model. This is achieved without generating realizations of the posterior distribution.

Following Devilee et al. (1999) and Meier et al. (2007), we construct a finite-size training data set, representing the information available in the probabilistic formulation of the inverse problem, namely prior information and information about the forward model and the noise. This is then used to train a neural network whose output parameterizes any desired statistical property of the posterior distribution for which a log-likelihood can be computed. These properties can, for example, represent a Gaussian, generalized Gaussian, log-normal, or a mixture model distribution, representing continuous model parameters. The output can also refer to the posterior probability of defined classes of model features or discrete model parameters. The neural network is designed to ensure that the estimated statistical properties of the posterior are similar to the same statistics derived from a sample of the posterior. Given a suitable training set the method provides accurate information regarding properties of the posterior distribution of interest in a given problem at a fraction of the computational cost of traditional sampling-based approaches.

The method is first presented for probabilistic inverse problems in general; this can be considered a generalization of the ideas proposed by Devilee et al. (1999) and followed up by for example, Meier et al. (2007); Earp et al. (2020). Next, we demonstrate the method, applying it to non-linear probabilistic inversion of airborne electromagnetic data using non-Gaussian prior models of varying complexity. We show the neural network approach can be used to accurately estimate statistical properties of the posterior, related to both discrete and continuous model parameters, using regression and classification networks. The results are compared to results obtained by calculating the same statistical properties from a sample of the posterior obtained using the extended rejection sampler (Hansen, 2021).

## 2. Method

Let  $\mathbf{m} = [m_1, m_2, \dots, m_{N_M}]$  represent  $N_M$  model parameters that define some properties of a system, such as for example, physical properties of a geo-model.  $\mathbf{m}$  is typically represented on a grid in a Cartesian or spherical coordinate system. For example,  $\mathbf{m}$  might represent geophysical properties such as resistivity, velocity, or any other geological/geophysical/geochemical parameter.

A key issue in geosciences is how to infer information about  $\mathbf{m}$  from different types of available information, such as geological expert knowledge, geophysical data, well log data, etc. This is generally referred to as an inverse problem. Tarantola and Valette (1982b) describe the inverse problem as a problem of probabilistic conjunction of information. Available information about  $\mathbf{m}$  is described in the form of probability densities and then combined using a conjunction of information to obtain a single probability density that describes the combined information. For example, consider a case when a specific type of information about structural properties is quantified by  $\rho(\mathbf{m})$ , and that information from observed electromagnetic (EM) data and well logs is quantified through  $L(\mathbf{m})$ . Then the conjunction of this information is given by the posterior probability distribution  $\sigma(\mathbf{m})$ , which, under the assumption that the individual types of information have been obtained independently, is given by

$$\sigma(\mathbf{m}) \propto \rho(\mathbf{m}) \cdot L(\mathbf{m}). \quad (1)$$

In other words, the conjunction of the independent information is proportional to the product of probability densities describing each independent set of information. The likelihood  $L(\mathbf{m})$  is a measure of how well the data  $\mathbf{d}$  computed from a specific model matches observed data  $\mathbf{d}_{obs}$  given to noise with a specified probability distribution. Noise-free data can be computed by evaluating the forward model

$$\mathbf{d} = g(\mathbf{m}), \quad (2)$$

where  $g$  is a non-linear operator that maps the model parameters into data.  $g$  typically refers to some numerical algorithm solving some physical equations (such as Maxwell's equations).

The probabilistic inverse problem is then to infer information about  $\sigma(\mathbf{m})$ , which contains the combined information of, for example, both structural prior information, through the prior  $\rho(\mathbf{m})$ , and information from observed geophysical data, through  $L(\mathbf{m})$ .

A general approach (that allows using a non-linear forward model and non-Gaussian prior) for solving probabilistic formulated inverse problems is the use of sampling methods to sample the posterior distribution, Equation 1 (Geman & Geman, 1984; Green, 1995; Hansen et al., 2013, 2016; Hastings, 1970; Laloy & Vrugt, 2012; Metropolis et al., 1953; Mosegaard & Tarantola, 1995). Unfortunately, such sampling methods can be extremely computationally demanding, to the point where they cannot be practically applied. They rely on solving the forward problem, Equation 2, many (often millions of) times.

Some algorithms make implicit assumptions about the prior model, such as a layered subsurface (Malinverno, 2002; Sambridge et al., 2013), while others, such as the classical rejection sampler and Metropolis algorithm (Hastings, 1970) require that both the prior and likelihood can be evaluated. This typically leads to using relatively simple prior models.

The extended variations of the Metropolis algorithm (Mosegaard & Tarantola, 1995) and the rejection sampler (Hansen, 2021; Hansen et al., 2016) do not require that an analytical description of the prior exists, as evaluation of the prior is not needed. It is sufficient that an algorithm exists that can generate a realization from the prior. This opens up the possibility of using a variety of more complex prior models, based on, for example, geostatistical simulation-based methods (Hansen et al., 2008, 2012).

### 2.1. Properties Related to Geophysical Model Parameters

The model parameters  $\mathbf{m}$  typically refer to physical parameters (e.g., resistivity when dealing with EM data, or elastic properties when dealing with seismic data). In practice, decision-makers may be more interested in related features, or specific questions, such as “What is the chance of penetrating a specific lithology when drilling?” (Scales & Snieder, 1997). Such features or occurrences of events will be referred to through  $\mathbf{n}$ .

In general, the relation between  $\mathbf{m}$  and  $\mathbf{n}$  can be complex and is formally described by a joint prior distribution  $\rho(\mathbf{m}, \mathbf{n})$ . This can, for example, be the case if  $\mathbf{n}$  refers to subsurface lithology, and  $\mathbf{m}$  to a geophysical property. This has been widely studied in the inversion of reflection seismic data, where information about geophysical properties is often assumed dependent on lithology, such that  $\rho(\mathbf{m}, \mathbf{n}) = \rho(\mathbf{n})\rho(\mathbf{m}|\mathbf{n})$  (Bosch et al., 2010; Grana & Della Rossa, 2010; Rimstad et al., 2012). A more general formulation of Equation 1, describing information on both  $\mathbf{m}$  and  $\mathbf{n}$  is then

$$\sigma(\mathbf{m}, \mathbf{n}) \propto \rho(\mathbf{m}, \mathbf{n}) \cdot L(\mathbf{m}, \mathbf{n}), \quad (3)$$

given the available joint prior information, the forward model, and the noise. The corresponding forward problem, generalizing Equation 2, takes the form

$$\mathbf{d} = g(\mathbf{m}, \mathbf{n}). \quad (4)$$

Sometimes the relation between  $\mathbf{m}$  and  $\mathbf{n}$  is so simple that  $\mathbf{n}$  can be computed from  $\mathbf{m}$  through a mapping function  $\mathbf{n} = h(\mathbf{m})$ . For example,  $\mathbf{n}$  can refer to the volume of a reservoir (a scalar) obtained from a high dimensional set of geophysical model parameters  $\mathbf{m}$ . Or,  $\mathbf{n}$  can refer to the cumulative thickness of layers with a resistivity ( $\mathbf{m}$ ) above some threshold. Another example is when  $\mathbf{m}$  refers to the properties of a groundwater model. Then flow modeling based on a set of realizations from the posterior, can be used to propagate uncertainties into, for example, the arrival time of polluted groundwater ( $\mathbf{n}$ ) at a specific location (Vilhelmsen et al., 2019). Such a focus on related features and properties derived from the posterior distribution, rather than the posterior distribution over the geophysical parameter  $\sigma(\mathbf{m})$  itself, is discussed by Scheidt et al. (2015).

The sampling algorithms described above can be used to generate a sample from  $\sigma(\mathbf{m}, \mathbf{n})$  from which statistical analysis of any feature related to  $\sigma(\mathbf{m}, \mathbf{n})$  can be computed.

The goal here is however not to generate realizations of the posterior distribution, but instead to compute directly statistical properties of the posterior distribution similar to those that would be obtained by computing it directly from a sample of the posterior distribution. In other words, given a sample  $\hat{\mathbf{n}}$  of the posterior,  $\sigma(\mathbf{n})$ , the goal is to compute parameters  $\Theta$  that define a desired statistical property of  $\sigma(\mathbf{n})$ . For example, if  $\mathbf{n}$  refers to a discrete

parameter with  $N_o$  possible outcomes, then  $\Theta = [\theta_1, \dots, \theta_{N_o}]$  could refer to the probability of realizing each possible outcome. If  $\mathbf{n}$  refers to a continuous parameter,  $\Theta = [\theta_0, \mathbf{C}_\theta]$  could represent the mean and covariance of a multivariate Gaussian distribution.  $\Theta = [\theta_0, \theta_1, \theta_2]$  could represent the mean, variance, and power of a generalized 1D Gaussian distribution.  $\Theta = [\theta_0]$  could represent the rate of a Poisson distribution.  $\Theta = [\theta_0, \theta_1]$  could represent a Binomial distribution.

Assume that a sample  $\hat{\mathbf{n}}$  of  $\sigma(\mathbf{n})$  is available. The optimal values of  $\Theta$  can be found maximizing the likelihood,  $L_\Theta$ , that each realization of the posterior,  $\hat{\mathbf{n}}^{i*}$ , is a realization of the probability distribution (described by the parameter(s)  $\Theta$ )  $f(\hat{\mathbf{n}}^{i*}|\Theta)$ , given as

$$L_\Theta = f(\hat{\mathbf{n}}|\Theta) = \prod_{i=1}^{N_\sigma} f(\hat{\mathbf{n}}^{i*}|\Theta), \quad (5)$$

where  $N_\sigma$  is the number of independent realizations of  $\hat{\mathbf{n}}$ . The specific choice of  $f(\hat{\mathbf{n}}^{i*}|\Theta)$  depends on the type of statistical parameters to be estimated. Examples will be given below. Maximization of Equation 5 is equivalent to minimizing the negative log-likelihood (which we refer to as the loss,  $J_\Theta$ ):

$$J_\Theta = -\log\left(\prod_{i=1}^{N_\sigma} f(\hat{\mathbf{n}}^{i*}|\Theta)\right) \quad (6)$$

$$= -\sum_{i=1}^{N_\sigma} \log(f(\hat{\mathbf{n}}^{i*}|\Theta)). \quad (7)$$

Minimization of the loss function, Equation 7, can be used to obtain estimates of the parameters  $\Theta$  representing statistical properties of  $\sigma(\mathbf{n})$ .

Here a method is proposed that allows direct computation of the parameters,  $\Theta$ , that describe statistical properties of  $\sigma(\mathbf{m}, \mathbf{n})$ , using a neural network trained on a data set containing a sample of the known information (including the prior, forward, noise and modeling errors), without ever generating realizations from  $\sigma(\mathbf{m}, \mathbf{n})$ . The approach follows the basic strategy proposed by Devilee et al. (1999) and consists of two steps: (a) construction of a training data set and (b) construction and training of a neural network. This is done once. Then, the trained machine learning algorithm can be applied, very efficiently to compute desired properties of the posterior distribution, for potentially many sets of observed data.

## 2.2. A: Construction of Training Data Set

Equation 4 describes the forward problem of computing noise-free data. The forward problem describing simulation of data including noise,  $\mathbf{d}_{sim}$ , is

$$\mathbf{d}_{sim} = g(\mathbf{m}, \mathbf{n}) + r(\mathbf{m}, \mathbf{n}) = \mathbf{d} + r(\mathbf{m}, \mathbf{n}), \quad (8)$$

where  $r(\mathbf{m}, \mathbf{n})$  represent noise. Often geophysical data  $\mathbf{d}$  depends only directly on the physical parameters, in which case  $g(\mathbf{m}, \mathbf{n}) = g(\mathbf{m})$ .

Let  $\mathbf{M}^* = [\mathbf{m}^{1*}, \mathbf{m}^{2*}, \dots, \mathbf{m}^{N_T*}]$  and  $\mathbf{N}^* = [\mathbf{n}^{1*}, \mathbf{n}^{2*}, \dots, \mathbf{n}^{N_T*}]$  represent  $N_T$  realizations of  $\rho(\mathbf{m}, \mathbf{n})$ . Let  $\mathbf{D}^* = [\mathbf{d}^{1*}, \mathbf{d}^{2*}, \dots, \mathbf{d}^{N_T*}]$  represent the corresponding  $N_T$  noise free data, obtained by evaluating Equation 4. Finally let  $\mathbf{D}_{sim}^* = [\mathbf{d}_{sim}^{1*}, \mathbf{d}_{sim}^{2*}, \dots, \mathbf{d}_{sim}^{N_T*}]$  represent  $N_T$  corresponding realizations of simulated noisy data, following Equation 8. This constitutes a training data set

$$\mathbf{T} = [\mathbf{N}^*; \mathbf{M}^*; \mathbf{D}^*; \mathbf{D}_{sim}^*], \quad (9)$$

that can be obtained by (a) sampling the prior, (b) solving the forward problem, and (c) simulation the noise.

The sample  $\mathbf{T}$  in Equation 9 represents the available information (prior, physics of the forward model, noise) in so far as it can be represented by a finite sample of size  $N_T$ . The larger the sample, the more complete the representation of the available information.

### 2.3. B: Construct and Train a Neural Network to Estimate Relevant Statistics of $\sigma(\mathbf{m}, \mathbf{n})$

The goal is to design and train a neural network to estimate  $\Theta$  directly from realizations of simulated data including noise  $\mathbf{d}_{sim}^{i*}$ . In principle any machine learning method capable of regression and/or classification, such as regression trees and support vector machines (Bishop et al., 1995), can be used to estimate the mapping  $\mathbf{d}_{sim}^{i*} \mapsto \Theta$  which after training can be used on real data to evaluate  $\mathbf{d}_{obs} \mapsto \Theta$ . Here we choose to make use of a fully connected artificial neural network. The presented approach builds on earlier work by Devilee et al. (1999), Meier et al. (2007), and Röth and Tarantola (1994).

#### 2.3.1. The Structure of the Neural Network

A neural network can be described in terms of an input layer, an inner part of the neural network (which can consist of many layers, referred to as hidden layers), and an output layer.

The input layer here represents the training data, which includes noise, and consists of  $N_d$  neurons. The output layer has  $N_\theta$  neurons representing the statistical parameters describing a distribution characterizing the features or properties of the posterior distribution that one wishes to predict.

The inner part of the network can be either simple or complex, and it can consist of either (fully) connected layers of neurons, convolutional layers, or combinations of these and other types of layers depending on the application. Here a fully connected neural network is considered as it has been demonstrated that such a neural network, with at least one hidden layer, can approximate any continuous function with arbitrary accuracy, when the number of hidden units is large enough (Hornik et al., 1990).

Each neuron has a number of adjustable parameters, the weights  $w_i$  (one for each neuron in the previous layer), and a bias  $b$ , as well as an activation function  $\Psi$ . All neurons in one layer are fully connected to all neurons in the following layer. The input for a neuron (except for the first layer where the input  $\mathbf{d}_{sim}^{i*}$ ) is the output of the neurons in the previous layer, and the output  $y$  of a neuron in response to inputs  $x_i$ , is given by

$$y = \Psi \left( \sum_i (w_i * x_i) + b \right). \quad (10)$$

For a specific network, with specified values for the weights and biases, one can compute the output, given some input, simply by evaluating the neurons layer by layer, starting from the input layer. See, for example, Bishop et al. (1995) for more details.

#### 2.3.2. The Loss Function

When a neural network is trained using the training data set, its free parameters (the weight and bias of each node for a fully connected network) are adjusted to minimize a specific loss function. In the present case, the training data set consists of (when properties of  $\sigma(\mathbf{n})$  are of interest)  $\mathbf{T} = [\mathbf{N}^*, \mathbf{D}_{sim}^*]$ . The goal is to estimate  $\mathbf{d}_{sim}^{i*} \mapsto \Theta$  rather than simply  $\mathbf{d}_{sim}^{i*} \mapsto \mathbf{n}$ .

This is achieved by constructing a loss function with unknown parameters  $\Theta$  that describe statistical properties of the desired probability distribution, Equation 5, and whose parameters can be found by minimizing the loss function, Equation 7. The key here is to choose a loss function that is the negative log-likelihood of the property of interest as described by the parameters  $\Theta$  one wishes to estimate.

At each iteration of training the neural network, the loss is computed by applying the following steps for each dataset  $T^i = [\mathbf{n}^{i*}, \mathbf{d}_{sim}^{i*}]$  in the training data set  $\mathbf{T}$ :

1. Evaluate the network using  $\mathbf{d}_{sim}^{i*}$  as input. This provides as output an estimate  $\hat{\Theta}_i$
2. Evaluate the corresponding loss,  $J^i$ , as  $J^i = -\log \left( f \left( \mathbf{n}^{i*} | \hat{\Theta}_i \right) \right)$ .

The total loss is then given by

$$\mathbf{J} = \sum_{i=1}^{N_T} J^i. \quad (11)$$

By construction, as  $\mathbf{d}_{sim}^{i*}$  has been computed from  $\mathbf{n}^{i*}$  using Equation 8,  $\mathbf{n}^{i*}$  can be considered a realization of  $\sigma(\mathbf{n})$ , given the data  $\mathbf{d}_{sim}^{i*}$ , and therefore, minimizing the loss in Equation 11 leads to estimates of statistical parameters  $\Theta$  that describe  $\sigma(\mathbf{n})$ , in the same manner as would minimizing Equation 7 given a sample,  $\hat{\mathbf{n}}$ , from  $\sigma(\mathbf{n})$ . The difference is that the proposed method achieves this without the need to first realize the sample  $\hat{\mathbf{n}}$ .

Minimizing the loss function thus maximizes the probability that each  $\mathbf{n}^{i*}$  can be seen as a realization of the probability distribution whose parameters  $\Theta_i$  are the result of evaluating the neural network  $\mathbf{d}_{sim}^{i*} \mapsto \Theta_i$ . Note that it is crucial that data with noise  $\mathbf{d}_{sim}^{i*}$  is used for training, as opposed to using noise-free data  $\mathbf{d}^{i*}$ , as this would imply ignoring noise completely, which would lead to overfitting.

In general,  $\mathbf{n}$  (and/or  $\mathbf{m}$ ) can refer to a continuous parameter (such as velocity, resistivity, temperature, or related properties) or a discrete parameter (such as lithology type and event type). Continuous model parameters lead to a regression-type problem, whereas discrete model parameters lead to a classification problem.

### 2.3.2.1. Continuous Model Parameters - Regression

We first consider the case when  $\mathbf{n}$  represents continuous parameters. Say we wish to estimate the mean and covariance,  $\hat{\Theta}_0$  and  $\hat{\mathbf{C}}_\theta$ , of the posterior distribution  $\sigma(\mathbf{n})$  given a set of observed data  $\mathbf{d}_{obs}$ . Assume a neural network exists that outputs a set of parameters describing  $\Theta = [\hat{\Theta}_0^i, \hat{\mathbf{C}}_\theta^i]$ , given the input  $\mathbf{d}_{sim}^i$ . The likelihood that a set of parameters from the training dataset  $\mathbf{n}^{i*}$  is a realization from the multivariate Gaussian distribution  $\mathcal{N}(\hat{\Theta}_0^i, \hat{\mathbf{C}}_\theta^i)$  as obtained from evaluating the neural network using  $\mathbf{d}_{sim}^{i*}$  as input, is given by

$$f(\mathbf{n}^{i*} | \hat{\Theta}_0^i, \hat{\mathbf{C}}_\theta^i) = k_C \exp\left(-0.5(\mathbf{n}^{i*} - \hat{\Theta}_0^i)^T \hat{\mathbf{C}}_\theta^{i-1} (\mathbf{n}^{i*} - \hat{\Theta}_0^i)\right), \quad (12)$$

where  $k_C = \left((2\pi)^{N_d} |\hat{\mathbf{C}}_\theta^i|\right)^{-\frac{1}{2}}$  is a normalization factor. The corresponding loss function  $J^i$  is

$$J^i = -\log\left(f(\mathbf{n}^{i*} | \hat{\Theta}_0^i, \hat{\mathbf{C}}_\theta^i)\right) \quad (13)$$

$$= 0.5(\mathbf{n}^{i*} - \hat{\Theta}_0^i)^T \hat{\mathbf{C}}_\theta^{i-1} (\mathbf{n}^{i*} - \hat{\Theta}_0^i) - \log(k_C) \quad (14)$$

The total loss is then given by Equation 11. Any neural network that minimizes this loss function, will lead to an estimate of the parameters of interest, here  $\Theta = [\hat{\Theta}_0, \hat{\mathbf{C}}_\theta]$ , that are computed directly without ever computing realizations of  $\sigma(\mathbf{n})$ .

To represent the posterior mean and full covariance, given  $N_m$  model parameters, an output layer of  $N_\Theta = N_m + N_m^2$  nodes must be used. If only the posterior mean and variance are estimated, an output layer of  $N_\Theta = N_m + N_m$  nodes is needed. If only the posterior mean is of interest an output layer of  $N_\Theta = N_m$  nodes is needed and minimizing Equation 14 is then similar to minimizing the widely used mean squared error loss function (Bishop et al., 1995), as utilized for example, in for example, Röth and Tarantola (1994). Recall, that the above scheme does not impose any assumptions on either the prior or the posterior distribution which may be complex. The estimated mean and covariance are simply statistical parameters of the posterior distribution, that may or may not be useful for a specific use case. The quality of the obtained estimate naturally depends on the complexity of the machine learning model used, and the size of the training data set, which will be considered in more detail in the application presented below.

Other statistical parameters of the posterior can be estimated by minimizing the appropriate negative log-likelihood function for the considered probability distribution. For example, a 1D generalized probability distribution is defined by three parameters  $\Theta = [\theta_1, \theta_2, \theta_3]$ , and its probability distribution given by Tarantola (2005)

$$f(n^i|\Theta) = \frac{1}{2\theta_2\Gamma(1+1/\theta_3)} \exp\left(-\left(\frac{|n^i - \theta_1|}{\theta_2}\right)^{\theta_3}\right). \quad (15)$$

A 1D Gaussian mixture model based on a mixture of  $N_c$  1D Gaussian distribution, as considered by for example, Meier et al. (2007), is defined by  $\Theta = [\mathbf{t}_1, \mathbf{t}_2, \mathbf{t}_3] = [t_1^1, \dots, t_1^{N_c}, t_2^1, \dots, t_2^{N_c}, t_3^1, \dots, t_3^{N_c}]$ , where  $\mathbf{t}_1$  refers to the mean,  $\mathbf{t}_2$  refers to the standard deviation of  $N_c$  Gaussian distribution, each with weight  $\mathbf{t}_3$ , and its probability distribution given by

$$f(n^i|\mathbf{t}_1, \mathbf{t}_2, \mathbf{t}_3) = \sum_{i=1}^{N_c} t_3^i (t_2^i \sqrt{2\pi})^{-1} \exp\left(-0.5\left(\frac{n^i - t_1^i}{t_2^i}\right)^2\right). \quad (16)$$

The corresponding negative log-likelihood for Equations 15 and 16 can trivially be obtained and used as a loss function in a neural network to estimate  $\Theta$ . In principle, any statistical parameter with a corresponding negative log-likelihood that can be computed, and used as a loss function, can be estimated using the proposed methodology.

### 2.3.2.2. Discrete Model Parameters - Classification

Say  $n_i$  represents a discrete parameter with  $N_o$  possible outcomes (classes). One's aim is then to estimate the posterior probability of each of the  $N_o$  classes given some data  $\mathbf{d}_{obs}$ .

Let  $\theta_i^* = [p_i^{1*}, p_i^{2*}, \dots, p_i^{N_o*}]$  represent the true probabilities of  $n_i^*$  belonging to a specific class. In practice, the true probability of one (the correct) class will be one, and the others zero. Further  $\hat{\theta}_i = [\hat{p}_i^1, \hat{p}_i^2, \dots, \hat{p}_i^{N_o}]$  represent the corresponding predictions by the neural network of the probabilities of each class for a specific model parameter,  $n_i$ .

The likelihood of observing  $\theta_i$  given  $\hat{\theta}_i$  is then

$$f(\theta_i|\hat{\theta}_i) = \prod_{j=1}^{N_o} (\hat{p}_i^j)^{p_i^{j*}}. \quad (17)$$

The corresponding loss function  $J^i$  is then

$$J^i = -\log(f(\theta_i|\hat{\theta}_i)) = -\sum_{j=1}^{N_o} p_i^{j*} \log(\hat{p}_i^j). \quad (18)$$

The choice of class probabilities  $\hat{\theta}_i$  that maximizes Equation 17 can be found by minimizing the negative log-likelihood given by the loss function, Equation 18, which is equivalent to the categorical cross-entropy between the two probability distributions (Bishop et al., 1995). Usually, the softmax activation is used for multi-class classification problems (while the sigmoid activation function is used for binary classification problems), as it forces all probabilities to be in the range 0–1, and ensures that  $\sum_{j=1}^{N_o} \hat{p}_i^j = 1$ , such that the output parameters can be interpreted as a probability. A neural network that estimates the mapping  $\mathbf{d}_{sim}^i \mapsto \hat{\theta}_i^j$  by minimizing Equation 18, using the softmax activation function in the output layer, therefore locates the maximum-likelihood of Equation 17, which directly estimates  $\sigma(p_i^*)$ , that is, the posterior class probability for a discrete model parameter.

To summarize, our proposed method involves first constructing a training data set (Equation 9) that represents (within the limits of the size of the training data set used) the known information (the prior, the forward, and the noise model), and specifically contains prior knowledge regarding any feature  $\mathbf{n}$ , that may be directly or indirectly related to the model parameters  $\mathbf{m}$ , about which one wishes to infer information. A neural network is then designed and trained by minimizing a specific loss function, that expresses the log-likelihood of the parameters  $\Theta$  describing the probability distribution of desired features  $\mathbf{n}$  that may be either continuous or discrete.



### 3. Application to Airborne EM Data From Morrill, Nebraska

The methodology described above is applied to the inversion of airborne electromagnetic (AEM) data. This inverse problem has been widely studied by deterministic linearized least-squares methods using both a 1D and 3D forward model (Auken & Christiansen, 2004; Auken et al., 2014; Christensen, 2002; Cox et al., 2010; Grayver et al., 2013; Viezzoli et al., 2008).

The full non-linear 1D inverse problem has also been addressed using Markov chain Monte Carlo (MCMC) sampling methods, based on, for example, the reversible-jump sampling method relying on a prior model representing a 1D layered subsurface (Brodie & Sambridge, 2012; Minsley, 2011; B. J. Minsley, Foks, & Bedrosian, 2021). Hansen and Minsley (2019) proposed the use of the extended Metropolis algorithm, also an MCMC method, that allows the use of any prior model that can be sampled. The 1D nonlinear inverse EM problem leads to a non-trivial sampling problem, due to the existence of model equivalences (significantly different models lead to the same forward response). A sufficient sampling of the 1D posterior distribution of resistivity values, to obtain a limited set of independent realizations, may require hundreds of thousands of MCMC iterations, and hence forward model evaluations. For a single sounding this may take at least 10 min per sounding, requiring access to supercomputers for the application of real-world data sets (Foks & Minsley, 2020). Hansen (2021) proposed 1D probabilistic inversion based on the extended rejection sampler (using lookup tables, similar to  $[N^*, M^*, D^*]$ ) that relies on the construction of a large sample for the prior along with the forward responses (generated once). This is then used to generate independent realizations of the posterior distribution numerically more efficiently than is possible using Markov Chain-based algorithms, and at the same time avoids issues related to model equivalences. This sampling approach is used for the comparison below.

The size of airborne EM surveys is becoming larger, so the use of any of the inversion methods discussed above will lead to considerable computational demands. Currently, two major airborne EM surveys are being carried out. The AusAEM20 project, by Geoscience Australia, is expected to collect around 65,000 flight-line-kilometers of data, leading to many hundreds of thousands of EM measurements (Howard, 2020). USGS has collected more than 43,000 flight-line-kilometer data in the Mississippi Alluvial Plain, and another 25,000 flight-line-kilometer is planned for 2021, leading to significantly more than 1,000,000 data points to be inverted in the Mississippi Alluvial Plain (Minsley, Rigby, et al., 2021).

As an example, we consider the inversion of airborne electromagnetic (AEM) data from Morrill, Nebraska (Abraham et al., 2012; Smith et al., 2010). We use data at 451 locations, at every 50m along a 22.5 km West-East profile, as also considered in Minsley (2011). Each observed data set consists of 13 measurements (in-phase and quadrature measurements from six pairs of transmitter and receiver coils, as well the measurement altitude).

Three different types of prior models will be defined, that represent different information about the subsurface resistivities ( $\mathbf{m}$ ) and related (both discrete and continuous) properties  $\mathbf{n}$  at Morrill. For each of the three prior models considered, a unique posterior probability distribution exists. Various properties of the posterior distribution will be computed using the proposed machine learning method and compared to results obtained from a finite sample of the posterior distributions obtained using the extended rejection sampler with a lookup table of size  $N_T = 2 \cdot 10^6$ .

#### 3.1. A Priori Models and Noise

##### 3.1.1. Parameterization

In this example, the subsurface is parameterized into 125 layers of  $dz = 1$  m thickness. Prior models based on up to four sets of parameters,  $\rho(\mathbf{m}, \mathbf{n}_1, \mathbf{n}_2, \mathbf{n}_3)$  are considered.

$\mathbf{m} = [m_1, m_2, \dots, m_{N_M}]$  represents the resistivity of each of the 125 layers.

$\mathbf{n}_1$  represents the existence of a sharp boundary between two neighboring layers ( $n_{1i} = 0$  when there is no boundary and  $n_{1i} = 1$  in case of a boundary). A sharp boundary is defined when two neighboring resistivity values differ by more than 20%.  $\mathbf{n}_1$  refers to 124 discrete parameters and can be directly computed from  $\mathbf{m}$ .

$\mathbf{n}_2$  represents the cumulative thickness of resistivity values above 225  $\Omega$ . It can be directly computed from  $\mathbf{m}$  using

$$\mathbf{n}_2 = \sum_i^{N_M} dz * I(m_i),$$

where  $I(m_i) = 1$  when  $m_i > 225 \Omega\text{m}$ , and  $I(m_i) = 0$  when  $m_i \leq 225 \Omega\text{m}$ .  $\mathbf{n}_2$  refers to a single continuous parameter.

$\mathbf{n}_3$  represents a category ('1', '2', and '3', representing three distinct lithologies) in each layer.  $\mathbf{n}_3$  cannot be computed from  $\mathbf{m}$ , but  $\mathbf{n}_3$  and  $\mathbf{m}$  are linked through a conditional prior distribution  $\rho(\mathbf{m}|\mathbf{n}_3)$  (see example below).  $\mathbf{n}_3$  refers to 125 discrete parameters with three possible outcomes.

For brevity, all model parameters combined will be referred to as  $\mathbf{p} = [\mathbf{m}, \mathbf{n}_1, \mathbf{n}_2, \mathbf{n}_3]$ . To illustrate the potential of the method three different non-Gaussian prior models are considered that vary in complexity and information content.

### 3.1.2. Prior Information

**$\rho_A(\mathbf{p}) = \rho_A(\mathbf{m}, \mathbf{n}_1, \mathbf{n}_2)$ , a uniform prior model.**  $\rho_A(\mathbf{p})$  represents a choice of independence between model parameters,  $\rho_A(m_i, m_j) = \rho_A(m_i)\rho_A(m_j) \forall (i, j)$ . The resistivity of each resistivity model parameter is assumed to be log-uniform distributed in the range  $\mathcal{U}[2, 280] \Omega$ . This is the least informative prior model considered. Eleven independent realizations of  $\rho_A(\mathbf{m}, \mathbf{n}_1)$  are shown in Figure 1a.

**$\rho_B(\mathbf{p}) = \rho_B(\mathbf{m}, \mathbf{n}_1, \mathbf{n}_2)$ , a discrete layered model.**  $\rho_B(\mathbf{p})$  represents a layered subsurface consisting of 1–8 layers (uniformly distributed), each with a constant resistivity. The resistivity in a specific layer is assumed to be log-uniform distributed in the range  $\mathcal{U}[2, 2800] \Omega$ .

A realization  $\mathbf{p}^*$  of  $\rho_B(\mathbf{p})$  is generated by first choosing the number of layers as a random number,  $Nl$ , between 1 and 8. Then  $Nl - 1$  layer interfaces are randomly selected from a uniform distribution of  $\mathcal{U}[0, 125]$  m. Then the resistivity within each layer is realized from a uniform distribution  $\mathcal{U}[2, 280] \Omega$ . This type of prior model is similar to the transdimensional prior considered by Minsley (2011). Eleven independent realizations of  $\rho_B(\mathbf{m}, \mathbf{n}_1)$  are shown in Figure 1c.

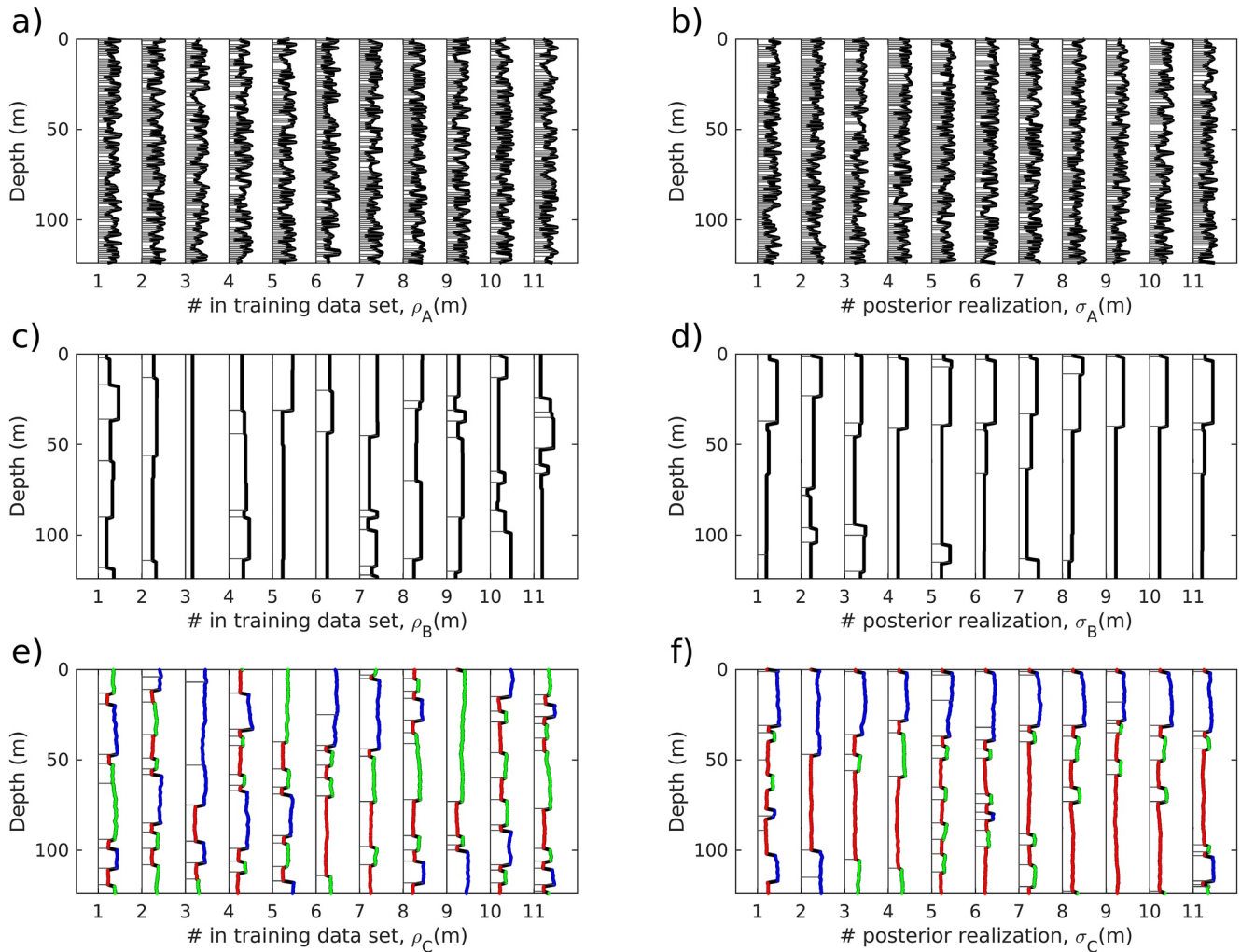
**$\rho_C(\mathbf{p})$ , a trimodal mixture Gaussian model.**  $\rho_C(\mathbf{p})$  represents a subsurface with three possible lithologies ('1', '2', and '3') each with a distinct resistivity distribution. See the discussion about the prior geological knowledge in Morrill in Abraham et al. (2012) and Hansen and Minsley (2019).

To sample  $\rho_C(\mathbf{p}) = \rho_C(\mathbf{m}, \mathbf{n}_1, \mathbf{n}_2, \mathbf{n}_3)$ , first a realization of  $\rho_C(\mathbf{n}_3)$  is generated as  $\rho_C(\mathbf{n}_3^*)$ , which represents an example of the distribution of the lithologies. This is achieved by generating a realization of a multivariate normal distribution with a Gaussian-type covariance model with a range of 30 m, followed by a simple truncation to obtain 40% of lithology A, 40% of lithology B, and 20% of lithology C. Then a realization of the resistivity  $\mathbf{m}^*$  is generated, conditional to the lithology type from  $\rho_C(\mathbf{m}|\mathbf{n}_3^*)$ . The resistivity, within each lithology, is generated as a realization of a multivariate normal distribution in  $\log_{10}$ -resistivity space with a range of 30 m, a specific mean,  $m_0$  and standard deviation,  $m_{std}$ . For lithology '1',  $m_0 = 1.1$  and  $m_{std} = 0.14$ . For lithology '2',  $m_0 = 2$  and  $m_{std} = 0.2$ . For lithology '3',  $m_0 = 2.75$  and  $m_{std} = 0.25$ . Finally,  $\mathbf{n}_1^*$  and  $\mathbf{n}_2^*$  are computed from  $\mathbf{m}^*$ . In this way a realization  $\mathbf{p}^* = [\mathbf{m}^*, \mathbf{n}_1^*, \mathbf{n}_2^*, \mathbf{n}_3^*]$  of  $\rho_C(\mathbf{p})$  is generated. Eleven independent realizations of  $\rho_C(\mathbf{p})$  are shown in Figure 1e.

$\rho_C(\mathbf{p})$  is designed to reflect available information related to the subsurface at Morrill (Abraham et al., 2012; Hansen & Minsley, 2019).  $\rho_A(\mathbf{p})$  and  $\rho_B(\mathbf{p})$  are considered here to investigate how the proposed methodology reacts to a uniform (maximum entropy) prior such as  $\rho_A(\mathbf{p})$ , and a simple (low entropy) prior as  $\rho_B(\mathbf{p})$ .

### 3.1.3. Noise

The noise of the EM data is assumed to be independent uncorrelated zero-mean Gaussian noise, with a standard deviation of 5 ppm (parts per million) plus 5% noise relative to the noise-free data value, which means the noise depends implicitly on the model. This is the same noise model as considered in previous works on the EM data from Morrill (Hansen & Minsley, 2019; Hansen, 2021; Minsley, 2011).



**Figure 1.** First 11 models from the training data,  $\mathbf{T}$ , for three prior models (a)  $\rho_a(\mathbf{m}, \mathbf{n}_1)$ , (c)  $\rho_b(\mathbf{m}, \mathbf{n}_1)$ , and (e)  $\rho_c(\mathbf{m}, \mathbf{n}_1, \mathbf{n}_2)$ , as well as 11 independent realizations from the posterior distribution obtained for the data at  $x = 15$  km for (b)  $\sigma_b(\mathbf{m}, \mathbf{n}_1)$ , (d)  $\sigma_d(\mathbf{m}, \mathbf{n}_1)$ , and (f)  $\sigma_c(\mathbf{m}, \mathbf{n}_1, \mathbf{n}_2)$ . Thin horizontal black lines indicate the existence of a layer interface ( $\mathbf{n}_1$ ). The thick line indicates variation in resistivity ( $\mathbf{m}$ ). In (e) and (f) the colors of the thick line represent lithology A (red), B (blue), and C (green) when defined.

### 3.2. Sampling of the Posterior Distribution

For reference, the extended rejection sampler, with a lookup table of size  $N_r = 5 \cdot 10^6$ , is used to sample the posterior distribution, as detailed in Hansen (2021). Eleven independent realizations of the three posterior distributions ( $\sigma_A(\mathbf{p})$ ,  $\sigma_B(\mathbf{p})$ , and  $\sigma_C(\mathbf{p})$ ) are shown in Figures 1b, 1d and 1f.

The goal of the proposed machine learning approach is to directly compute statistical properties of the posterior distribution similar to the same statistical properties obtained from a sample of the posterior using sampling, such as shown in Figures 1b, 1d and 1f.

### 3.3. Neural Network Design

Two fully connected neural networks are designed to allow the characterization of the 1D marginal posterior distribution of continuous and discrete parameters. The input layer, in both cases, consists of the observed data  $\mathbf{d}_{obs}$ , or simulated data with noise. For this specific case, it consists of 13 neurons. Twelve neurons refer to the 12 data measurements, and 1 neuron to the altitude measured during the recording of data.

The inner network is designed using either four or eight hidden layers depending on the application, each with 40 neurons using the rectified linear activation function (Bishop et al., 1995). This inner part of the network needs to be complex enough that the desired mapping can be represented, but simple enough to avoid overfitting, as discussed also by Meier et al. (2007). Network design is highly problem-dependent, and for the present problem, we found this network design provides results on par with, and in some cases better than, sampling-based approaches, while at the same time being relatively easy to optimize.

As discussed, the choice of the loss function, and to some extent the activation function, are set by the specific property of the posterior distribution that will be estimated. This leads to two specific types of output layers for regression and classification-type problems.

### 3.3.1. Regression Type Neural Network

The first neural network type is designed to estimate parameters  $\theta$  of a probability distribution describing the 1D marginal posterior distribution of continuous parameters (such as  $\mathbf{m}$  and  $\mathbf{n}_2$ ). If  $N_\theta$  is the number of parameters needed to describe a specific 1D distribution, then in total  $N_{out} = N_\theta N_m$  neurons are needed in the output layer if the target is properties of  $\sigma(\mathbf{m})$ , and  $N_{out} = N_\theta$  if the target is  $\sigma(\mathbf{n}_2)$ .

### 3.3.2. Classification Type Neural Network

The second neural network type is designed to estimate the posterior probability of possible classes for the discrete type model parameters  $\mathbf{n}_1$  and  $\mathbf{n}_3$ , that is, of  $\sigma(\mathbf{n}_1)$   $\sigma(\mathbf{n}_3)$ .

If the goal is to estimate the 1D marginal distribution of a discrete parameter with  $N_{cat}$  possible outcomes, this can be achieved by selecting an output layer with  $N_{out} = N_m$  when  $N_{cat} = 2$  (using a sigmoid activation function), and  $N_{out} = N_{cat} N_m$  when  $N_{cat} > 2$  (using the softmax activation function). As discussed above, using the cross-entropy loss function, Equation 18, will lead to direct estimation of the 1D posterior marginal probabilities in this case.

## 3.4. Network Training

Using the prior models, the nonlinear forward model, and the noise model, a training data set of size  $N_T = 5 \cdot 10^6$  is constructed (one for each type of prior model) and used for training. Both networks are trained using 67% of the training data set, while 33% is reserved for validation. In both cases, the loss function is minimized using the Adam optimizer (Kingma & Ba, 2014) using a learning rate of 0.001, for a maximum of 2000 epochs. Early stopping is utilized which stops the training if the loss function evaluated on the validation data does not decrease for 50 epochs. This is done to avoid over-fitting, where the loss on the training data will decrease, but where the loss on the validation data increases. TensorFlow with Keras and TensorFlow-probability have been used to implement and train the neural networks (Abadi et al., 2015; Chollet, 2015; Dillon et al., 2017).

The two considered networks, and the training of the networks, only differ concerning the definition of the output layer (the number of nodes and activation function), the choice of the loss function, and the chosen number of hidden layers.

## 3.5. Estimation of Properties of $\sigma(\mathbf{m})$

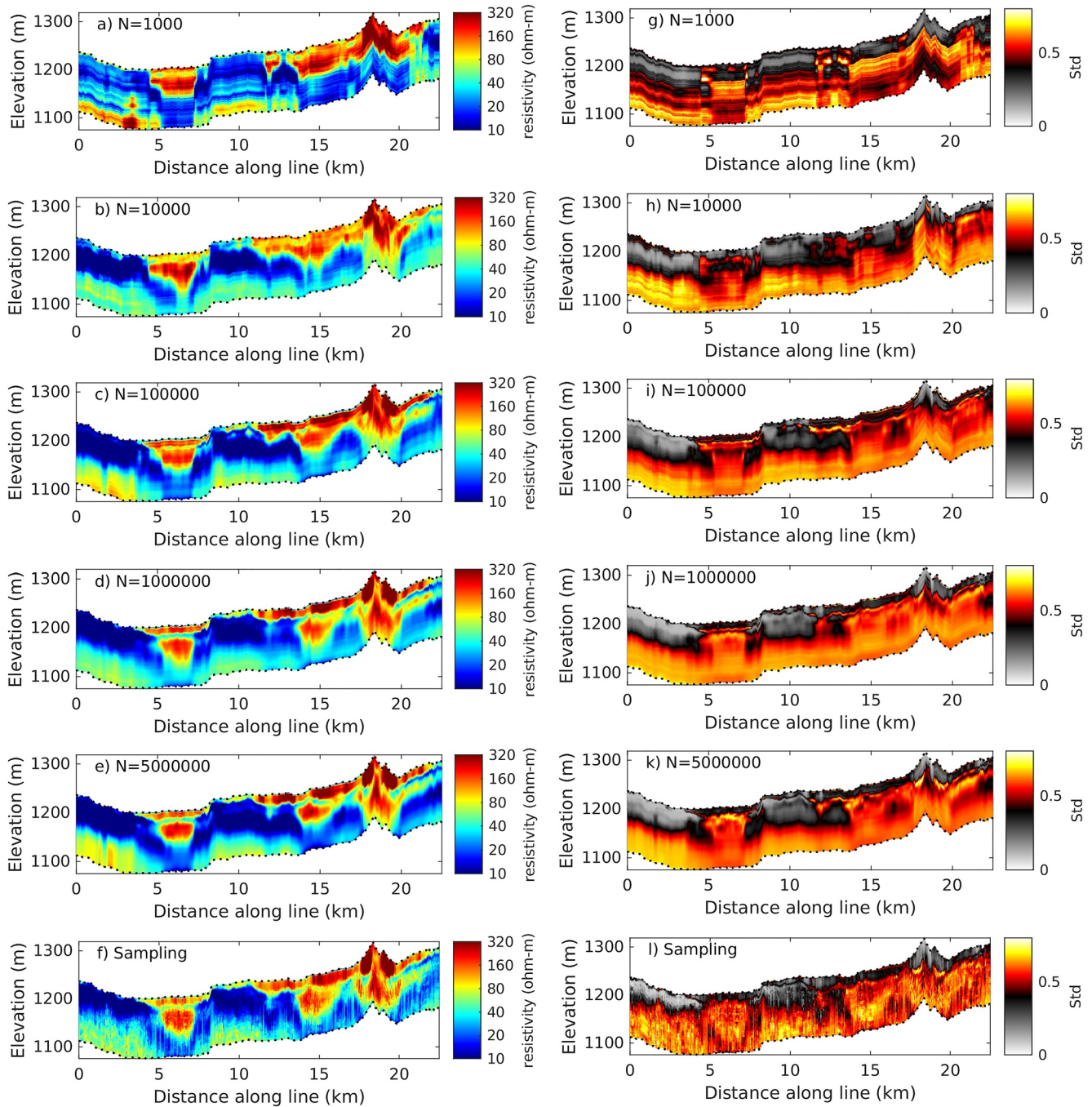
First, properties related to the posterior distribution of resistivity,  $\sigma(\mathbf{m})$ , are considered.

### 3.5.1. Estimation of Mean and Standard Deviation of $\sigma(\mathbf{m})$

A neural network is set up and trained to estimate the pointwise mean and standard deviation of  $\sigma(\mathbf{m})$ , using eight hidden layers, by minimizing the loss function in Equation 14.

Figures 2a–2e show the pointwise mean of the posterior distribution  $\sigma_c(\mathbf{m})$  obtained using the machine learning approach with a training data set of size  $N_T = [10^3, 10^4, 10^5, 10^6, 5 \cdot 10^6]$ , compared to the same statistics computed from a sample of the posterior obtained using the sampling method, Figure 2f. The corresponding standard deviation is shown in Figures 2g–2l.

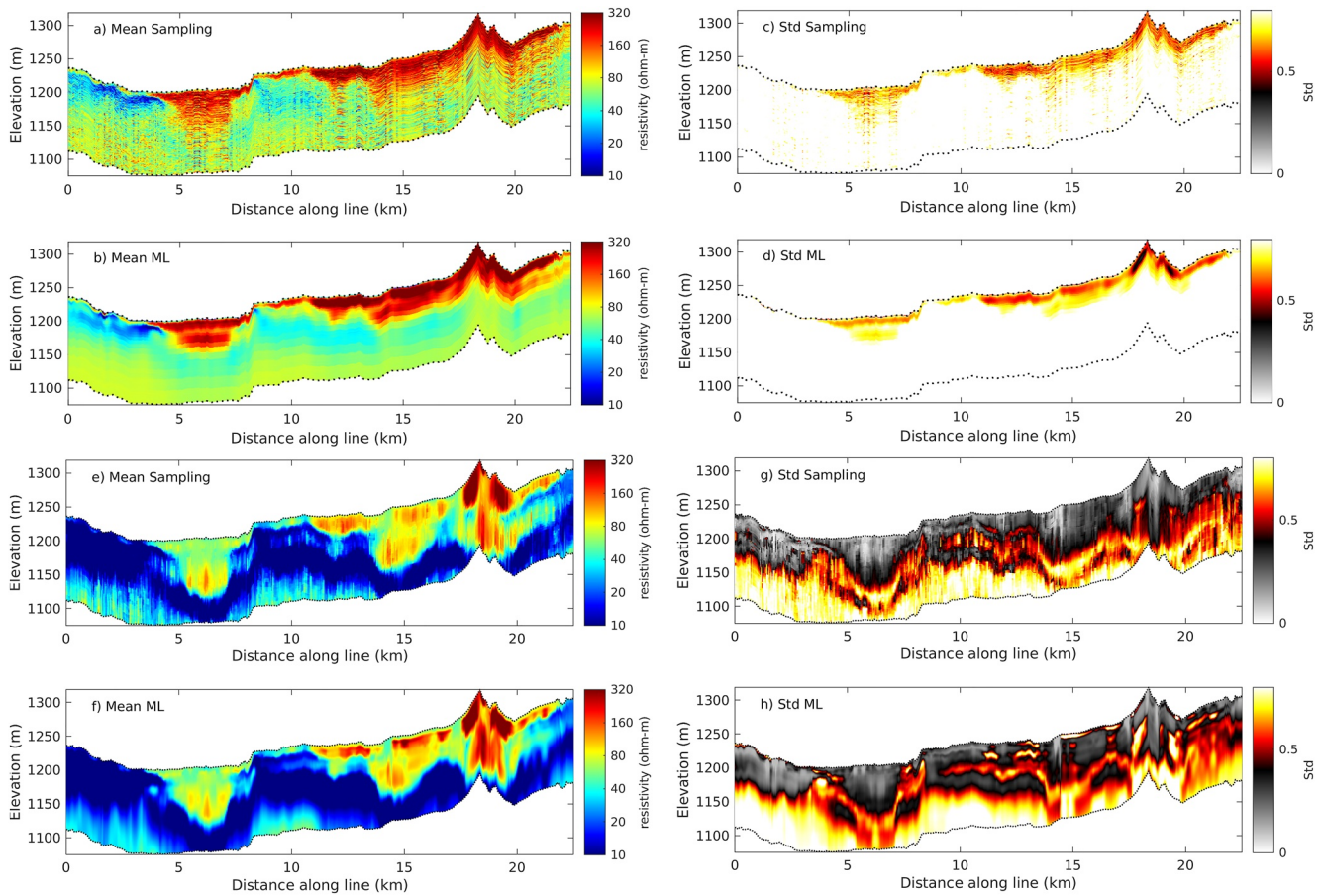
It is clear from Figures 2a and 2g that using  $N_T = 10^3$  provides very poor results, as compared to the results obtained using sampling, Figures 2e and 2j. But even using  $N_T = 10^4$  leads to results close to the sampling-based



**Figure 2.** Pointwise mean (a)–(f) and standard deviation (g)–(l) obtained from  $\sigma_c(\mathbf{m})$  obtained using machine learning based on a training data set of size  $10^3$  (a),(g),  $10^4$  (b),(h),  $10^5$  (c),(i),  $10^6$  (d),(j),  $5 \cdot 10^6$  (e),(k), and using the extended rejection sampler (f,l).

results. The changes in predicted mean and standard deviation become smaller as  $N_T$  increases, with only very subtle changes between the use of  $N_T = 10^6$  and  $N_T = 5 \cdot 10^6$ .

One notable difference when comparing Figures 2e and 2k ( $N_T = 5 \cdot 10^6$ ) and Figures 2f and 2l (sampling), is that sampling results in more small scale variability in the estimated parameters, as opposed to the more smooth result obtained using machine learning. The reason is simply that the sampling-based approach is based on inferring the statistics from a finite-sized sample of the posterior, whereas in the machine learning approach these statistics are estimated directly.



**Figure 3.** Pointwise mean (a),(b) and standard deviation (c),(d) obtained from  $\sigma_A(\mathbf{m})$  obtained using the extended rejection sampler (a),(c) and machine learning (b),(d) based on a training data set of size  $N_T = 5 \cdot 10^6$ . (e)-(h) As (a)-(d) but for  $\sigma_B(\mathbf{m})$ .

Figure 3 shows a comparison between the posterior mean and standard deviation obtained using the sampling approach and using the machine learning approach ( $N_T = 5 \cdot 10^6$ ), for  $\sigma_A(\mathbf{m})$  (Figures 3a–3d) and  $\sigma_B(\mathbf{m})$  (Figures 3e–3h), respectively.

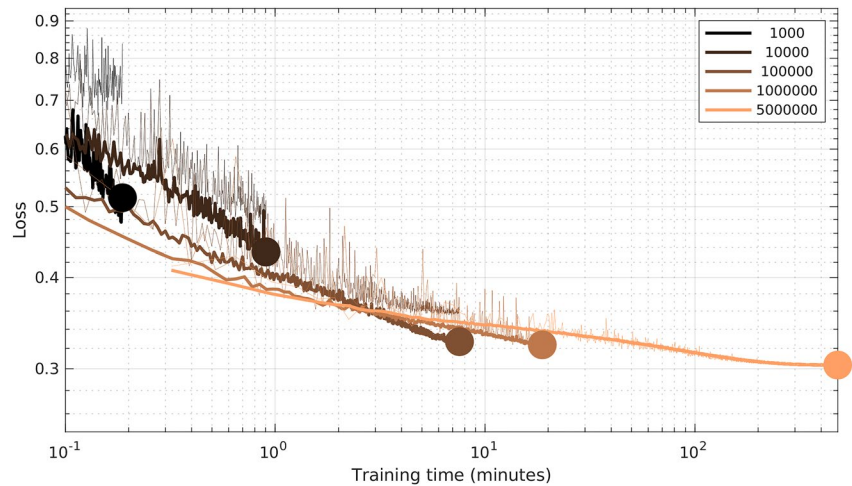
$\rho_A(\mathbf{m})$  refers to the least informed prior model, and hence one should expect the least resolution in the corresponding posterior distribution. This is what can be seen in the results from both the machine learning and the sampling approach, Figures 3a–3d, where only the resistivities at the top of the model are resolved.

While  $\rho_B(\mathbf{m})$  is somewhat simpler than  $\rho_C(\mathbf{m})$ , the mean of the corresponding posterior distribution is rather similar, Figures 2e and 3f, with the largest difference related to the posterior standard deviation, Figures 2k and 3h.

A key point from Figures 2 and 3 is that the use of the machine learning-based approach seems to provide results at least on par with the results obtained using sampling when the goal is to estimate the mean and standard deviation of the (non-Gaussian) posterior distribution. This is the case using both informed and uninformed prior models.

### 3.5.1.1. Computational Efficiency

Figure 4 shows the training and validation loss, and computation time needed to train the neural networks for the results presented in Figure 2, obtained using a workstation with an Intel Core(TM) i7-8700K CPU, Nvidia RTX 3090 GPU, and 64 Gb RAM was used. The training time increases with the size of the training data set,  $N_T$ . Both training and validation loss is reduced when  $N_T$  increases. It is also clear that the relative difference in loss decreases when comparing the use of  $N_T = 10^5$  to  $N_T = 5 \cdot 10^6$ , to when comparing the use of  $N_T = 10^3$  to  $N_T = 10^5$ . Hence, using  $N_T > 10^5$  leads to a substantially longer training time, but only to a minor loss reduction.



**Figure 4.** Training (thick lines) and validation (thin lines) loss as a function of training time for  $N_t = [10^3, 10^4, 10^5, 10^6, 5 \cdot 10^6]$ .

In addition, and as expected, when  $N_t$  increases the validation loss seems to match the training loss increasingly well, which indicates that there is no problem with overfitting.

Once set up and trained, the prediction of the network is very fast. For all the networks presented above, the prediction time for all 451 data locations was less than 5 ms. This means that more than 100,000 soundings can be analyzed per second.

### 3.5.2. Estimation of Multiple 1D Properties of $\sigma(m_i)$

As described above, any parameter of a probability distribution for which a loss function can be described through Equation 7 can be estimated using the machine learning method. To demonstrate this, four independent networks have been trained to estimate properties ( $\Theta$ ) of the 1D marginal posterior distribution  $\sigma(m_i)$  given by (a) a normal distribution (Equation 12, as in Figure 2), (b) a generalized normal distribution (Equation 15), (c) a mixture distribution based on two Gaussian distributions (Equation 16), and (d) a mixture distribution based on three Gaussian distributions (Equation 16). The loss functions used are the negative log-likelihood of the probability distribution in Equations 12, 15 and 16, respectively.

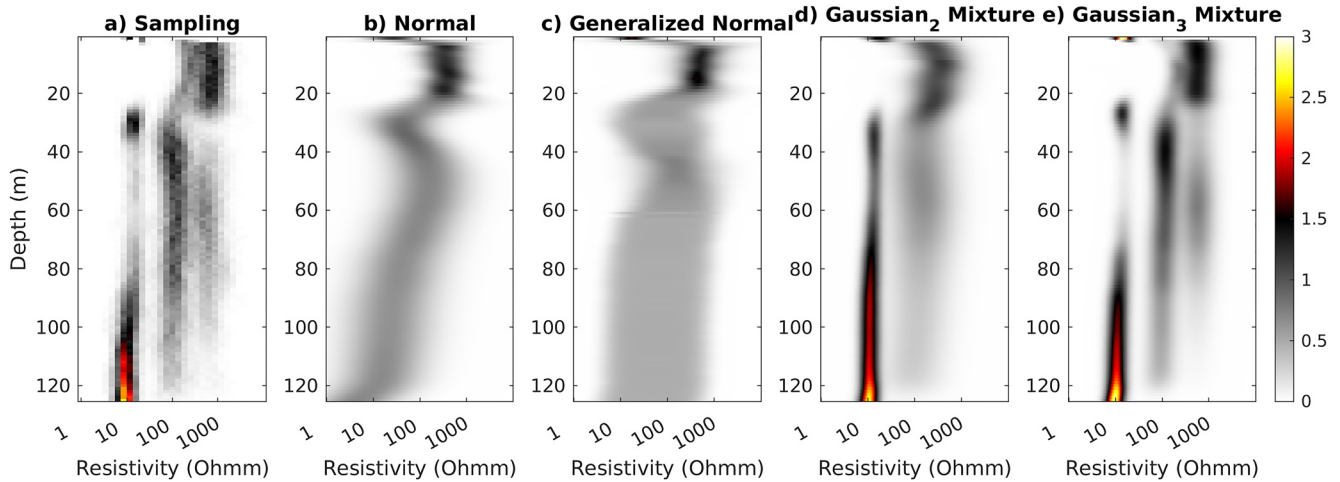
The number of parameters to estimate for the four cases, and hence neurons in the output layer, are  $N_\theta = [2*N_m, 3*N_m, 2*N_m*N_c, 3*N_m*N_c] = [250, 375, 750, 1125]$ , where  $N_c$  is the number of distributions in the mixture model.

Figure 5a shows the posterior 1D marginal distribution of resistivity values obtained using sampling, based on a finite set of realizations, obtained at  $x = 15$  km. One can clearly identify a bimodal to trimodal distribution at depth representing the three possible lithologies from the prior model  $\rho_c(\mathbf{m})$  with different resistivity values.

Figures 5b–5e show the probability distributions representing the estimated statistical properties of the 4 considered distributions. These distributions do not represent assumptions about the posterior distribution (which can be arbitrarily complex) but reflect example statistical properties that one might calculate from a sample of the posterior.

If the goal is to compute a representation of the 1D posterior marginal distribution, as considered by Meier et al. (2007); Shahraneini & Curtis (2011), then care should be taken to use a parameterization for the chosen 1D distribution complex enough to describe the variability of the posterior. From Figure 5 it is evident that only in case using the mixture model with 3 Gaussian distributions, does the estimated marginal probability density represents the actual 1D marginal posterior distribution well.

The statistical properties of the posterior distribution which is relevant to compute for a specific inverse problem, are naturally problem-dependent. This example nonetheless demonstrates that the machine learning methodology is capable of estimating parameters of different types of probability distributions, for which a probability density, and hence the corresponding loss function, can be computed.



**Figure 5.** 1D posterior probability density with depth using data at  $X = 6.2$  km (a) obtained using sampling followed by computation of the marginal posterior probability, and constructed from statistical properties inferred for a (b) normal distribution, (c) generalized normal distribution, (d-e) a mixture model based on 2 and 3 1D normal distributions. Obtained using  $N_T = [5 \cdot 10^6]$ .

### 3.6. Estimation of Properties of $\sigma(\mathbf{n}_1)$

$\sigma(\mathbf{n}_1)$  refers to the existence (or lack of) of a layer interface, which can be formulated as a binary classification problem. Therefore, a classification-type network is constructed using a sigmoid activation function, and Equation 18 as the loss function. 4 hidden layers are used.

Figures 6a and 6c refer to the pointwise posterior probability of locating a layer interface, as computed from a sample from the posterior distribution of  $\sigma_B(\mathbf{n}_1)$  and  $\sigma_C(\mathbf{n}_1)$ . The corresponding results obtained as the output of a trained neural network based on a training data set of size  $N_T = 5 \cdot 10^6$  are shown in Figures 6b and 6d. The prior probability of a layer interface is around 0.1, and hence a posterior probability of 0.25 is indicative of a layer interface.

The results using sampling and the machine learning approach are in both cases very similar with a bit more variability in the results obtained using sampling, due to the use of a finite-sized sample of the posterior distribution.

### 3.7. Estimation of Properties of $\sigma(\mathbf{n}_2)$

We consider the simpler problem of inferring information about a single continuous parameter,  $\mathbf{n}_2$ , representing the cumulative thickness of layers with a resistivity above 225  $\Omega\text{m}$ . The same neural network as considered above to estimate properties related to  $\mathbf{m}$  is used here, except that here only 4 hidden layers are used.

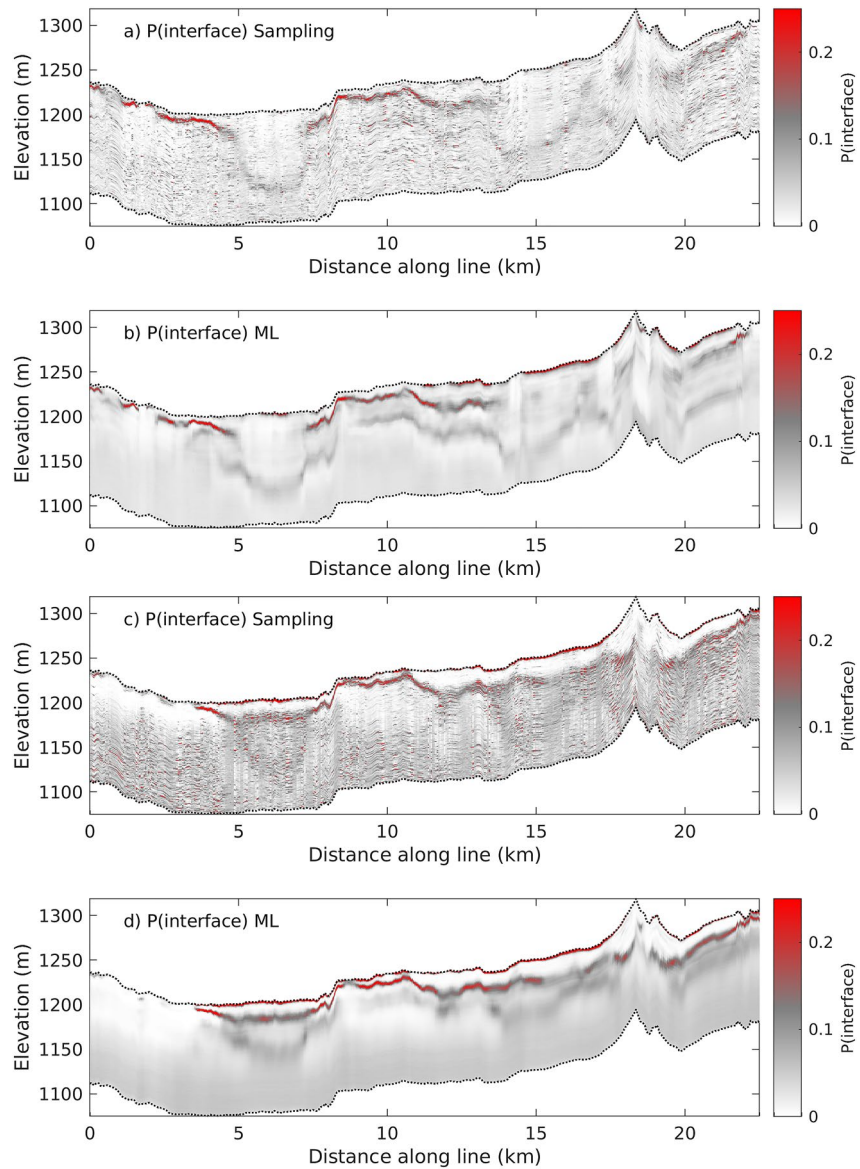
Figure 7 shows the mean of  $\sigma_C(\mathbf{n}_2)$  (black line), as well as the probability distribution reflecting the mean and standard deviation estimated using the machine learning approach for  $N_T = [10^3, 10^4, 10^5, 10^6, 5 \cdot 10^6]$  in Figures 7a–7e. The mean computed using the machine learning approach compares well to the mean obtained using sampling methods for  $N_T \geq 10^5$ .

### 3.8. Estimation of Properties of $\sigma(\mathbf{n}_3)$

Finally, we consider the discrete parameter  $\mathbf{n}_3$  which refers to lithology type, which can be of type '1', '2', and '3'. The outcome for each model parameter is then a multi-class (three classes) classification problem. Therefore, a classification-type network is constructed using a softmax activation function, and the loss function in Equation 18 four hidden layers are used.

Figures 8a, 8c and 8e show the posterior probability for each of the three classes obtained using sampling, while Figures 8b, 8d and 8f show the corresponding results obtained by evaluating the trained network. Except for some small-scale variations in the sampling results, due to using finite sample size, the obtained posterior statistics are strikingly similar.





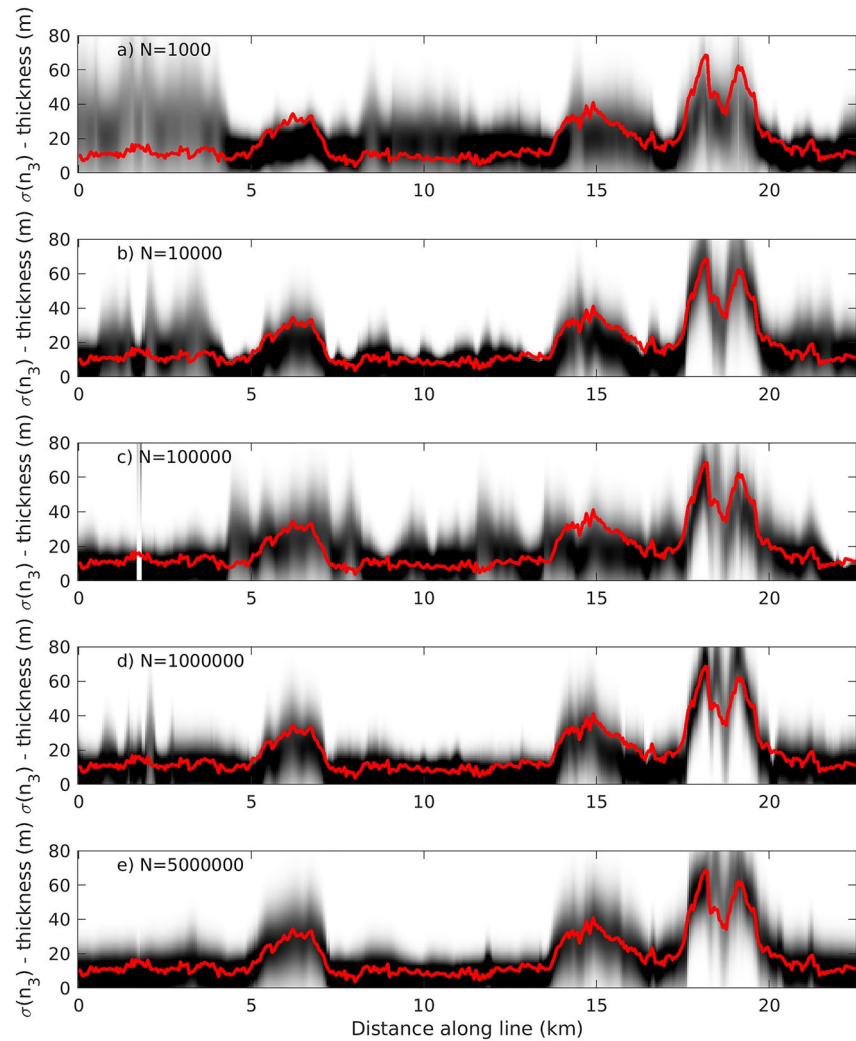
**Figure 6.** (a)-(b) Posterior probability of a layer interface obtained using extended rejection sampling (a), and machine learning (b), for  $\sigma_B(\mathbf{n}_1)$ . (c)-(d) Posterior probability of a layer interface obtained using extended rejection sampling (c), and machine learning (d), for  $\sigma_C(\mathbf{n}_1)$ . Obtained using  $N_T = [5 \cdot 10^6]$ .

#### 4. Discussion

A typical application of probabilistic inversion is to use some sampling method to generate a large sample from the posterior distribution. Then some appropriate statistic, computed from the sample of the posterior distribution, is chosen and visualized.

The theory presented above proposes how one can construct a neural network that can directly estimate any statistical property of the posterior distribution (for discrete or continuous parameters) for which a probability distribution can be evaluated, without ever generating realizations of the posterior distribution. This can be achieved by the following steps:

1. Construct a training data set, in the style of Devilee et al. (1999),  $\mathbf{T}^* = [\mathbf{N}^*, \mathbf{D}_{sim}^*]$ , where  $\mathbf{N}^*$  represents a set of features/properties of interest, and  $\mathbf{D}_{sim}^*$  represents a corresponding set of simulated data with noise, using both the forward and the noise model.

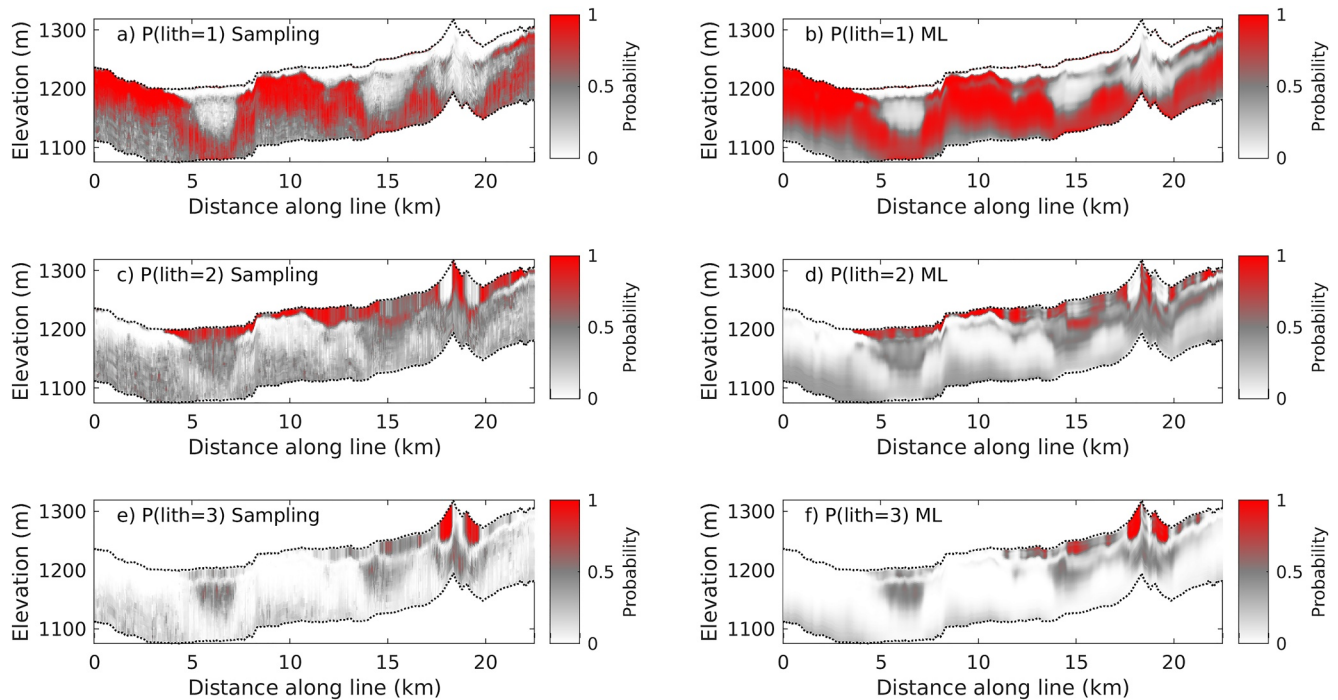


**Figure 7.** Mean of the posterior distributions  $\sigma_c(\mathbf{n}_3)$  estimated using sampling (red line) compared with the estimated 1d normal mean and standard deviation of  $\sigma_c(\mathbf{n}_3)$  plotted as probability density in grayscale, estimated using the machine learning approach using training data set of size  $N_T = [10^3, 10^4, 10^5, 10^6, 5 \cdot 10^6]$  in (a-e).

2. Design a neural network whose output layer represents the relevant statistical parameters  $\Theta$  of the posterior distribution  $\sigma(\mathbf{n})$  of interest.
3. Train the neural network by minimizing a loss function that is the negative log-likelihood of the probability density,  $f(\Theta)$ , whose properties one wishes to estimate.

Practical application of the methodology requires a) a neural network structure complex enough to be able to estimate the mapping  $\mathbf{d}_{sim}^* \mapsto \Theta$ , and b) a training data set large enough to allow the mapping to be inferred.

The methodology was applied and demonstrated in a case study using airborne EM data from Morrill, Nebraska. Several (uninformed to more informed) prior models were considered, describing both subsurface resistivity (a continuous parameter,  $\mathbf{m}$ ) and lithology (a discrete parameter,  $\mathbf{n}_3$ ) and the considered forward problem was nonlinear. In addition, the method was used to estimate posterior statistics of low-dimensional features of the prior models, such as the existence of a layer interface,  $\mathbf{n}_1$ , and the thickness of layers with resistivity above  $225 \Omega\text{m}$ ,  $\mathbf{n}_2$ . Results showed that using a training data set of size  $N_T > 10^5$ , in this case, leads to a trained neural network that provides estimates of posterior statistics similar to those obtained using sampling methods, using a fraction of the computational power (about 5 ms per sounding).



**Figure 8.** Posterior probability of lithology, using (a, c, e) sampling and (b, d, f) machine learning for  $\sigma_C(\mathbf{n}_3)$ .

#### 4.1. Limitations

The proposed method does not generate realizations of the posterior distribution, as do other sampling-based methods (B. J. Minsley, 2011; Brodie & Sambridge, 2012; Hansen & Minsley, 2019; Hansen, 2021). Instead, statistics of the posterior distribution for features of interest are estimated directly by applying a trained neural network.

In some use cases, one may actually need the realizations, for example, to propagate flow responses from a set of realizations from the posterior representing hydraulic parameters (Vilhelmsen et al., 2019). But, in many applications, where one is primarily interested in some statistical parameter describing the posterior, such as the posterior probability of a lithology type, the presented methodology may be useful.

The key practical difference to using sampling methods is that one has to identify the feature one is interested in and specify an appropriate loss function before running the inversion. Whereas using sampling methods to sample the posterior, one can convert the realizations of the posterior into a specific feature, and perform the posterior analysis, after the sampling algorithm has run.

In the example application, we adopted a widely used uncorrelated Gaussian noise model. In practice, real data are often affected by correlated noise (Bai et al., 2021; Hansen et al., 2014; Hauser et al., 2015). While in principle any noise model can be handled by the proposed methodology, as long as realizations of the noise can be generated, it remains to be tested how well the methodology works with more complex noise models.

The methodology is particularly promising for localized inverse problems, where the trained neural network can be set up and trained once, but applied many times. It is less obviously suited to 3D inversions with very large model dimensions because (a) construction of an adequately large training data set will be difficult and CPU intensive, (b) solving the 3D forward problem may be CPU intensive, and (c) it may be very difficult to train a neural network with millions of parameters in the output layer. While the method appears to work well for the AEM case considered, it remains to be seen how the method performs for other, possibly more nonlinear, inverse problems.

#### 4.2. Potential

The immediate appeal of the proposed methodology is that it leads to fast prediction times. One can get similar results, but much faster, compared with using sampling-based methods to analyze the posterior distribution. The presented method is faster than linearized least squares-based deterministic inversion of EM data (which uses less than a second CPU time per 4 soundings), which has been widely used for the inversion of large surveys (Auken et al., 2017; Minsley, Foks, & Bedrosian, 2021; Minsley, Rigby, et al., 2021) because they require much less computational resources than sampling-based methods. With the computational efficiency of the proposed method, the computational benefits of linearized methods are no longer so substantial that one should ignore the benefits of using the probabilistic methods that allow the use of site-specific prior information, a non-linear forward model, and full exploration of the space of uncertainty.

The more general appeal is that the proposed methodology allows the use of in principle arbitrarily complex prior models. The only requirement is that one must be able to generate independent realizations of the prior model. This allows an end-user to actively choose a prior model based on available information, as opposed to being forced to use the prior assumptions implicit in most available inversion algorithms, such as the assumptions of a layered subsurface (B. J. Minsley, Foks, & Bedrosian, 2021) or a Gaussian type smooth prior (Auken & Christiansen, 2004). The prior can be constructed according to site-specific information, and posterior statistics can be estimated for any parameter that can be computed from the prior model, as illustrated by the parameters  $\mathbf{n}_1$  and  $\mathbf{n}_2$  in the case study.

The main challenge then becomes the construction of realistic prior models that represent geological realistic information as well as realistic noise models.

#### 5. Conclusions

A simple, yet powerful, approach to probabilistic inversion has been proposed. Its application requires that one can simulate sets of examples capturing the known information. That is (a) sample from an arbitrarily complex prior model, (b) solving the forward problem, and (c) adding realistic noise to the simulated data. From each of these sets of models and data, a set of corresponding features related to the model parameters can be obtained. Together these represent, up to the limit of the finite set of models, all known information about these features of interest.

From such sets of features and corresponding noisy input data, a neural network can be used to estimate the statistical properties of the posterior distribution directly, by training the network to minimize an appropriate loss function. This provides the ability to carry out a fast and accurate estimation of relevant posterior statistics given an observed dataset.

A case study of the methodology applied to a nonlinear probabilistic inversion of EM data demonstrates it is possible to directly obtain posterior statistics similar to those obtained using sampling methods, using a fraction of the computation time. This approach allows the use and testing of multiple prior models, and to consider multiple features related to the prior distributions, in a fully probabilistic setting using only modest computational resources. The method has the most appeal for localized inverse problems, where the same trained neural network can be applied to many datasets with little computational effort.

#### Data Availability Statement

The airborne EM data used in this study is freely available and can be accessed at <https://doi.org/10.3133/ofr20101259> (Smith et al., 2010). Training data sets and python code for training and prediction will be made available upon publication at Zenodo <https://doi.org/10.5281/zenodo.7037407>.

#### References

- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., & Zheng, X. (2015). TensorFlow: Large-scale machine learning on heterogeneous systems. (Software available from tensorflow.org). Retrieved from <https://www.tensorflow.org/>
- Abraham, J. D., Cannia, J. C., Bedrosian, P. A., Johnson, M. R., Ball, L. B., & Sibray, S. S. (2012). Airborne electromagnetic mapping of the base of aquifer in areas of Western Nebraska. *U.S. Geological Survey Scientific Investigations Report, 2011-5219*, 38. <https://doi.org/10.3133/sir20115219>

#### Acknowledgments

This work is funded by the Danish Free Research Council, project 717-00160B.

- Ardizzone, L., Kruse, J., Wirkert, S., Rahner, D., Pellegrini, E. W., Klessen, R. S., et al. (2018). Analyzing inverse problems with invertible neural networks. arXiv preprint arXiv:1808.04730.
- Auken, E., Boesen, T., & Christiansen, A. V. (2017). A review of airborne electromagnetic methods with focus on geotechnical and hydrological applications from 2007 to 2017. *Advances in Geophysics*, 58, 47–93. <https://doi.org/10.1016/bs.agph.2017.10.002>
- Auken, E., & Christiansen, A. V. (2004). Layered and laterally constrained 2d inversion of resistivity data. *Geophysics*, 69(3), 752–761. <https://doi.org/10.1190/1.1759461>
- Auken, E., Christiansen, A. V., Kirkegaard, C., Fiandaca, G., Schamper, C., Behroozmand, A. A., et al. (2014). An overview of a highly versatile forward and stable inverse algorithm for airborne, ground-based and borehole electromagnetic and electric data. *Exploration Geophysics*, 46(3), 223–235.
- Bai, P., Vignoli, G., Andrea, V., Jouni, N., & Vacca, G. (2020). (quasi-) real-time inversion of airborne time-domain electromagnetic data via artificial neural network. *Remote Sensing*, 12(20), 3440. <https://doi.org/10.3390/rs12203440>
- Bai, P., Vignoli, G., & Hansen, T. M. (2021). 1d stochastic inversion of airborne time-domain electromagnetic data with realistic prior and accounting for the forward modeling error. *Remote Sensing*, 13(19), 3881. <https://doi.org/10.3390/rs13193881>
- Bishop, C. M. (1995). *Neural networks for pattern recognition*. Oxford university press.
- Blei, D. M., Kucukelbir, A., & McAuliffe, J. D. (2017). Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518), 859–877. <https://doi.org/10.1080/01621459.2017.1285773>
- Bording, T. S., Asif, M. R., Barfod, A. S., Larsen, J. J., Zhang, B., Grombacher, D. J., et al. (2021). Machine learning based fast forward modelling of ground-based time-domain electromagnetic data. *Journal of Applied Geophysics*, 187, 104290. <https://doi.org/10.1016/j.jappgeo.2021.104290>
- Bosch, M., Mukerji, T., & Gonzalez, E. F. (2010). Seismic inversion for reservoir properties combining statistical rock physics and geostatistics: A review. *Geophysics*, 75(5), 75A165–75A176. <https://doi.org/10.1190/1.3478209>
- Brodie, R. C., & Sambridge, M. (2012). Transdimensional Monte Carlo inversion of AEM data. *ASEG Extended Abstracts*, 2012(1), 1–4. <https://doi.org/10.1071/aseg2012ab095>
- Chollet, F. (2015). Keras. GitHub. Retrieved from <https://github.com/fchollet/keras>
- Christensen, N. B. (2002). A generic 1-D imaging method for transient electromagnetic data. *Geophysics*, 67(2), 438–447. <https://doi.org/10.1190/1.1468603>
- Constable, S. C., Parker, R. L., & Constable, C. G. (1987). Occam's inversion: A practical algorithm for generating smooth models from electromagnetic sounding data. *Geophysics*, 52(3), 289–300. <https://doi.org/10.1190/1.1442303>
- Conway, D., Alexander, B., King, M., Heinson, G., & Kee, Y. (2019). Inverting magnetotelluric responses in a three-dimensional Earth using fast forward approximations based on artificial neural networks. *Computers & Geosciences*, 127, 44–52. <https://doi.org/10.1016/j.cageo.2019.03.002>
- Cox, L. H., Wilson, G. A., & Zhdanov, M. S. (2010). 3D inversion of airborne electromagnetic data using a moving footprint. *Exploration Geophysics*, 41(4), 250–259. <https://doi.org/10.1071/eg10003>
- Devilee, R., Curtis, A., & Roy-Chowdhury, K. (1999). An efficient, probabilistic neural network approach to solving inverse problems: Inverting surface wave velocities for Eurasian crustal thickness. *Journal of Geophysical Research*, 104(B12), 28841–28857. <https://doi.org/10.1029/1999jb900273>
- de Wit, R. W., Valentine, A. P., & Trampert, J. (2013). Bayesian inference of Earth's radial seismic structure from body-wave traveltimes using neural networks. *Geophysical Journal International*, 195(1), 408–422. <https://doi.org/10.1093/gji/ggt220>
- Dillon, J. V., Langmore, I., Tran, D., Brevdo, E., Vasudevan, S., Moore, D., et al. (2017). Tensorflow distributions. arXiv preprint arXiv:1711.10604.
- Earp, S., & Curtis, A. (2020). Probabilistic neural network-based 2d travel-time tomography. *Neural Computing & Applications*, 32(22), 17077–17095. <https://doi.org/10.1007/s00521-020-04921-8>
- Earp, S., Curtis, A., Zhang, X., & Hansteen, F. (2020). Probabilistic neural network tomography across Grane field (North Sea) from surface wave dispersion data. *Geophysical Journal International*, 223(3), 1741–1757. <https://doi.org/10.1093/gji/ggaa328>
- Fichtner, A., Zunino, A., & Gebraad, L. (2018). Hamiltonian Monte Carlo solution of tomographic inverse problems. *Geophysical Journal International*, 216(2), 1344–1363. <https://doi.org/10.1093/gji/ggy496>
- Foks, N., & Minsley, B. (2020). Geophysical Bayesian inference in Python (GeoBiPy). <https://doi.org/10.5066/P9K3YH90>
- Geman, S., & Geman, D. (1984). Stochastic relaxation, Gibbs distributions, and the bayesian restoration of images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, (6), 721–741. <https://doi.org/10.1109/tpami.1984.4767596>
- Grana, D., & Della Rossa, E. (2010). Probabilistic petrophysical-properties estimation integrating statistical rock physics with seismic inversion. *Geophysics*, 75(3), O21–O37. <https://doi.org/10.1190/1.3386676>
- Grayver, A. V., Streich, R., & Ritter, O. (2013). Three-dimensional parallel distributed inversion of CSEM data using a direct forward solver. *Geophysical Journal International*, 193(3), 1432–1446. <https://doi.org/10.1093/gji/ggt055>
- Green, P. J. (1995). Reversible jump Markov chain Monte Carlo computation and Bayesian model determination. *Biometrika*, 82(4), 711–732. <https://doi.org/10.1093/biomet/82.4.711>
- Hansen, T. M. (2021). Efficient probabilistic inversion using the rejection sampler—Exemplified on airborne EM data. *Geophysical Journal International*, 224(1), 543–557. <https://doi.org/10.1093/gji/ggaa491>
- Hansen, T. M., & Cordua, K. S. (2017). Efficient Monte Carlo sampling of inverse problems using a neural network-based forward—Applied to gpr crosshole traveltime inversion. *Geophysical Journal International*, 211(3), 1524–1533. <https://doi.org/10.1093/gji/ggx380>
- Hansen, T. M., Cordua, K. S., Jacobsen, B. H., & Mosegaard, K. (2014). Accounting for imperfect forward modeling in geophysical inverse problems—Exemplified for crosshole tomography. *Geophysics*, 79(3), H1–H21. <https://doi.org/10.1190/geo2013-0215.1>
- Hansen, T. M., Cordua, K. S., Looms, M. C., & Mosegaard, K. (2013). SIPP: A Matlab toolbox for sampling the solution to inverse problems with complex prior information: Part 1, methodology. *Computers & Geosciences*, 52, 470–480. <https://doi.org/10.1016/j.cageo.2012.09.004>
- Hansen, T. M., Cordua, K. C., & Mosegaard, K. (2012). Inverse problems with non-trivial priors - Efficient solution through sequential Gibbs sampling. *Computational Geosciences*, 16(3), 593–611. <https://doi.org/10.1007/s10596-011-9271-1>
- Hansen, T. M., Cordua, K. S., Zunino, A., & Mosegaard, K. (2016). Probabilistic integration of geo-information. In M. Moorekamp, P. G. Lelièvre, N. Linde, & A. Khan (Eds.), *Integrated imaging of the earth: Theory and applications* (Vol. 218). John Wiley & Sons.
- Hansen, T. M., & Minsley, B. J. (2019). Inversion of airborne EM data with an explicit choice of prior model. *Geophysical Journal International*, 218(2), 1348–1366. <https://doi.org/10.1093/gji/ggz230>
- Hansen, T. M., Mosegaard, K., & Cordua, K. C. (2008). Using geostatistics to describe complex a priori information for inverse problems. In J. M. Ortiz, & X. Emery (Eds.), (Vol. 1, pp. 329–338). Mining Engineering Department, University of Chile. VIII international geostatistics congress.

- Hastings, W. K. (1970). Monte Carlo sampling methods using Markov chains and their applications. *Biometrika*, 57(1), 97–109. <https://doi.org/10.1093/biomet/57.1.97>
- Hauser, J., Gunning, J., & Annetts, D. (2015). Probabilistic inversion of airborne electromagnetic data under spatial constraints. *Geophysics*, 80(2), E135–E146. <https://doi.org/10.1190/geo2014-0389.1>
- Hornik, K., Stinchcombe, M., & White, H. (1990). Universal approximation of an unknown mapping and its derivatives using multilayer feedforward networks. *Neural Networks*, 3(5), 551–560. [https://doi.org/10.1016/0893-6080\(90\)90005-6](https://doi.org/10.1016/0893-6080(90)90005-6)
- Howard, D. (2020). Geological survey of Western Australia: Ausaem20-wa project. *Preview*, 2020(205), 18. <https://doi.org/10.1080/14432471.2020.1751781>
- Khosshkolgh, S., Zunino, A., & Mosegaard, K. (2022). Full-waveform inversion by informed-proposal Monte Carlo. *Geophysical Journal International*, 230(3), 1824–1833. <https://doi.org/10.1093/gji/ggac150>
- Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic optimization. arXiv preprint [arXiv:1412.6980](https://arxiv.org/abs/1412.6980)
- Köpke, C., Irving, J., & Elsheikh, A. H. (2018). Accounting for model error in Bayesian solutions to hydrogeophysical inverse problems using a local basis approach. *Advances in Water Resources*, 116, 195–207. <https://doi.org/10.1016/j.advwatres.2017.11.013>
- Laloy, E., Héroult, R., Jacques, D., & Linde, N. (2018). Training-image based geostatistical inversion using a spatial generative adversarial neural network. *Water Resources Research*, 54(1), 381–406. <https://doi.org/10.1002/2017wr022148>
- Laloy, E., & Vrugt, J. A. (2012). High-dimensional posterior exploration of hydrologic models using multiple-try dream (zs) and high-performance computing. *Water Resources Research*, 48(1), W01526. <https://doi.org/10.1029/2011wr010608>
- Madsen, R. B., & Hansen, T. M. (2018). Estimation and accounting for the modeling error in probabilistic linearized amplitude variation with offset inversion. *Geophysics*, 83(2), N15–N30. <https://doi.org/10.1190/geo2017-0404.1>
- Malinverno, A. (2002). Parsimonious Bayesian Markov chain Monte Carlo inversion in a nonlinear geophysical problem. *Geophysical Journal International*, 151(3), 675–688. <https://doi.org/10.1046/j.1365-246x.2002.01847.x>
- Meier, U., Curtis, A., & Trampert, J. (2007). Global crustal thickness from neural network inversion of surface wave data. *Geophysical Journal International*, 169(2), 706–722. <https://doi.org/10.1111/j.1365-246x.2007.03373.x>
- Menke, W. (2012). *Geophysical data analysis: Discrete inverse theory* (Vol. 45). Academic Press.
- Metropolis, N., Rosenbluth, M., Rosenbluth, A., Teller, A., & Teller, E. (1953). Equation of state calculations by fast computing machines. *The Journal of Chemical Physics*, 21, 1087–1092. <https://doi.org/10.1063/1.1699114>
- Minsley, B. J. (2011). A trans-dimensional Bayesian Markov chain Monte Carlo algorithm for model assessment using frequency-domain electromagnetic data. *Geophysical Journal International*, 187(1), 252–272. <https://doi.org/10.1111/j.1365-246x.2011.05165.x>
- Minsley, B. J., Foks, N. L., & Bedrosian, P. A. (2021). Quantifying model structural uncertainty using airborne electromagnetic data. *Geophysical Journal International*, 224(1), 590–607. <https://doi.org/10.1093/gji/ggaa393>
- Minsley, B. J., Rigby, J. R., James, S. R., Burton, B. L., Knierim, K. J., Pace, M. D. M., et al. (2021). Airborne geophysical surveys of the lower Mississippi Valley demonstrate system-scale mapping of subsurface architecture. *Communications Earth & Environment*, 2, 131. <https://doi.org/10.1038/s43247-021-00200-z>
- Moghadam, D. (2020). One-dimensional deep learning inversion of electromagnetic induction data using convolutional neural network. *Geophysical Journal International*, 222(1), 247–259. <https://doi.org/10.1093/gji/ggaa161>
- Moghadam, D., Behroozmand, A. A., & Christiansen, A. V. (2020). Soil electrical conductivity imaging using a neural network-based forward solver: Applied to large-scale bayesian electromagnetic inversion. *Journal of Applied Geophysics*, 104012. <https://doi.org/10.1016/j.jappgeo.2020.104012>
- Mosegaard, K., & Tarantola, A. (1995). Monte Carlo sampling of solutions to inverse problems. *Journal of Geophysical Research*, 100(B7), 12431–12447. <https://doi.org/10.1029/94jb03097>
- Mosser, L., Dubrulle, O., & Blunt, M. J. (2017). Reconstruction of three-dimensional porous media using generative adversarial neural networks. *Physical Review E*, 96(4), 043309. <https://doi.org/10.1103/physreve.96.043309>
- Mosser, L., Dubrulle, O., & Blunt, M. J. (2020). Stochastic seismic waveform inversion using generative adversarial networks as a geological prior. *Mathematical Geosciences*, 52(1), 53–79. <https://doi.org/10.1007/s11004-019-09832-6>
- Puzrev, V., & Swidinsky, A. (2021). Inversion of 1D frequency-and time-domain electromagnetic data with convolutional neural networks. *Computers & Geosciences*, 149, 104681.
- Rimstad, K., Avseth, P., & Omre, H. (2012). Hierarchical bayesian lithology/fluid prediction: A North Sea case study. *Geophysics*, 77(2), B69–B85. <https://doi.org/10.1190/geo2011-0202.1>
- Röth, G., & Tarantola, A. (1994). Neural networks and inversion of seismic data. *Journal of Geophysical Research*, 99(B4), 6753–6768. <https://doi.org/10.1029/93jb01563>
- Sambridge, M., Bodin, T., Gallagher, K., & Tkalčić, H. (2013). Transdimensional inference in the geosciences. *Philosophical Transactions of the Royal Society A: Mathematical, Physical & Engineering Sciences*, 371(1984), 20110547. <https://doi.org/10.1098/rsta.2011.0547>
- Scales, J. A., & Snieder, R. (1997). To [b]ayes or not to [b]ayes? *Geophysics*, 62(4), 1045–1046. <https://doi.org/10.1190/1.6241045.1>
- Scheidt, C., Renard, P., & Caers, J. (2015). Prediction-focused subsurface modeling: Investigating the need for accuracy in flow-based inverse modeling. *Mathematical Geosciences*, 47(2), 173–191. <https://doi.org/10.1007/s11004-014-9521-6>
- Shahraeeni, M. S., & Curtis, A. (2011). Fast probabilistic nonlinear petrophysical inversion. *Geophysics*, 76(2), E45–E58. <https://doi.org/10.1190/1.3540628>
- Smith, B. D., Abraham, J. D., Cannia, J. C., Minsley, B. J., Deszcz-Pan, M., & Ball, L. B. (2010). Helicopter electromagnetic and magnetic geophysical survey data, portions of the North Platte and South Platte natural resources Districts, Western Nebraska (Tech. Rep. No. 2010-1259). U.S. Geological Survey Scientific Investigations Report. <https://doi.org/10.3133/ofr20101259>
- Tarantola, A. (2005). *Inverse problem theory and methods for model parameter estimation*. SIAM.
- Tarantola, A., & Valette, B. (1982a). Generalized nonlinear inverse problems solved using the least squares criterion. *Reviews of Geophysics and Space Physics*, 20(2), 219–232. <https://doi.org/10.1029/rg020i002p00219>
- Tarantola, A., & Valette, B. (1982b). Inverse problems= quest for information. *Journal of Geophysics*, 50(3), 150–170.
- Tikhonov, A. N. (1963). On the solution of ill-posed problems and the method of regularization. *Doklady akademii nauk*, 151, 501–504.
- Viezzoli, A., Christiansen, A. V., Auken, E., & Sørensen, K. (2008). Quasi-3D modeling of airborne TEM data by spatially constrained inversion. *Geophysics*, 73(3), F105–F113. <https://doi.org/10.1190/1.2895521>
- Vilhelmsen, T. N., Auken, E., Christiansen, A. V., Barfod, A. S., Marker, P. A., & Bauer-Gottwein, P. (2019). Combining clustering methods with mps to estimate structural uncertainty for hydrological models. *Frontiers of Earth Science*, 7, 181. <https://doi.org/10.3389/feart.2019.00181>
- Zhang, X., & Curtis, A. (2020a). Seismic tomography using variational inference methods. *Journal of Geophysical Research: Solid Earth*, 125(4), e2019JB018589. <https://doi.org/10.1029/2019jb018589>

- Zhang, X., & Curtis, A. (2020b). Variational full-waveform inversion. *Geophysical Journal International*, 222(1), 406–411. <https://doi.org/10.1093/gji/ggaa170>
- Zhang, X., & Curtis, A. (2021). Bayesian geophysical inversion using invertible neural networks. *Journal of Geophysical Research: Solid Earth*, 126(7), e2021JB022320. <https://doi.org/10.1029/2021jb022320>
- Zhao, X., Curtis, A., & Zhang, X. (2022). Bayesian seismic tomography using normalizing flows. *Geophysical Journal International*, 228(1), 213–239. <https://doi.org/10.1093/gji/ggab298>