



GemPy 1.0: open-source stochastic geological modeling and inversion

Miguel de la Varga, Alexander Schaaf, and Florian Wellmann

Institute for Computational Geoscience and Reservoir Engineering, RWTH Aachen University, Aachen, Germany

Correspondence: Miguel de la Varga (varga@aices.rwth-aachen.de)

Received: 27 February 2018 – Discussion started: 6 March 2018

Revised: 15 November 2018 – Accepted: 16 November 2018 – Published: 2 January 2019

Abstract. The representation of subsurface structures is an essential aspect of a wide variety of geoscientific investigations and applications, ranging from geofluid reservoir studies, over raw material investigations, to geosequestration, as well as many branches of geoscientific research and applications in geological surveys. A wide range of methods exist to generate geological models. However, the powerful methods are behind a paywall in expensive commercial packages. We present here a full open-source geomodeling method, based on an implicit potential-field interpolation approach. The interpolation algorithm is comparable to implementations in commercial packages and capable of constructing complex full 3-D geological models, including fault networks, fault–surface interactions, unconformities and dome structures. This algorithm is implemented in the programming language Python, making use of a highly efficient underlying library for efficient code generation (*Theano*) that enables a direct execution on GPUs. The functionality can be separated into the core aspects required to generate 3-D geological models and additional assets for advanced scientific investigations. These assets provide the full power behind our approach, as they enable the link to machine-learning and Bayesian inference frameworks and thus a path to stochastic geological modeling and inversions. In addition, we provide methods to analyze model topology and to compute gravity fields on the basis of the geological models and assigned density values. In summary, we provide a basis for open scientific research using geological models, with the aim to foster reproducible research in the field of geomodeling.

1 Introduction

We commonly capture our knowledge about relevant geological features in the subsurface in the form of geological models, as 3-D representations of the geometric structural setting. Computer-aided geological modeling methods have existed for decades, and many advanced and elaborate commercial packages exist to generate these models (e.g., GoCAD, Petrel, GeoModeller). But even though these packages partly enable an external access to the modeling functionality through implemented APIs or scripting interfaces, it is a significant disadvantage that the source code is not accessible, and therefore the true inner workings are not clear. More importantly still, the possibility to extend these methods is limited – and, especially with the current rapid development of highly efficient open-source libraries for machine-learning and computational inference (e.g., *TensorFlow*, *Stan*, *pymc*, *PyTorch*, *Infer.NET*), the integration into other computational frameworks is limited.

However, there is to date no fully flexible open-source project that integrates state-of-the-art geological modeling methods. Conventional 3-D construction tools (CAD, e.g., *pythonOCC*, *PyGem*) are only useful to a limited extent, as they do not consider the specific aspects of subsurface structures and the inherent sparsity of data. Open source GIS tools exist (e.g., QGIS, *gdal*), but they are typically limited to 2-D (or 2.5-D) structures and do not facilitate the modeling and representation of fault networks, complex structures like overturned folds or dome structures, or combined stratigraphic sequences.

With the aim to close this gap, we present *GemPy*, an open-source implementation of a modern and powerful implicit geological modeling method based on a potential-field approach. The method was first introduced by

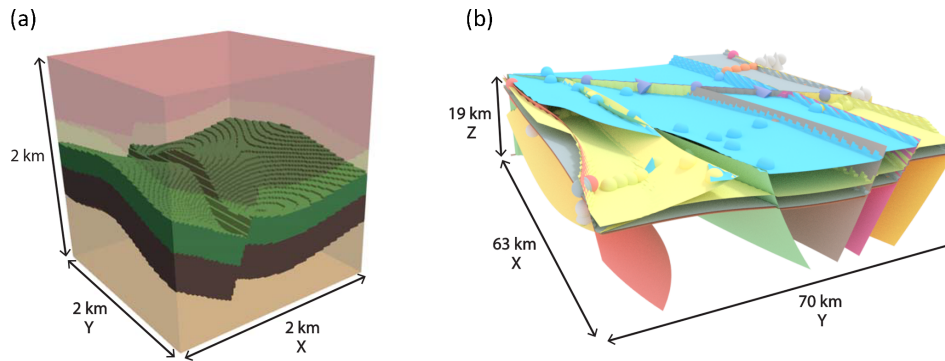


Figure 1. Example of models generated using *GemPy*. (a) Synthetic model representing a reservoir trap, visualized in Paraview (Stamm, 2017); (b) Geological model of the Perth basin (Australia) rendered using *GemPy* on the in-built Python in Blender (see Appendix F for more details), spheres and cones represent the input data.

Lajaunie et al. (1997) and it is grounded on the mathematical principles of universal cokriging. In distinction to surface-based modeling approaches (see ; Caumon et al., 2009, for a good overview), these approaches allow for the direct interpolation of multiple conformal sequences in a single scalar field, and the consideration of discontinuities (e.g., metamorphic contacts, unconformities) through the interaction of multiple sequences (Lajaunie et al., 1997; Mallet, 2004; Calcagno et al., 2008; Caumon, 2010; Hillier et al., 2014). Also, these methods allow for the construction of complex fault networks and enable, in addition, a direct global interpolation of all available geological data in a single step. This last aspect is relevant as it facilitates the integration of these methods into other diverse workflows. Most importantly, we show how we can integrate the method into novel and advanced machine-learning and Bayesian inference frameworks (Salvatier et al., 2016) for stochastic geomodeling and Bayesian inversion. Recent developments in this field have seen a surge in new methods and frameworks, e.g., using gradient-based Monte Carlo methods (Duane et al., 1987; Hoffman and Gelman, 2014) or variational inferences (Kucukelbir et al., 2017), making use of efficient implementations of automatic differentiation (Rall, 1981) in novel machine-learning frameworks. For this reason, *GemPy* is built on top of *Theano*, which provides not only the mentioned capacity to efficiently compute gradients but also provides optimized compiled code (for more details see Sect. 2.3.2). In addition, we utilize *pandas* for data storage and manipulation (McKinney, 2011), Visualization Toolkit (*vtk*) Python-bindings for interactive 3-D visualization (Schroeder et al., 2004), the de facto standard 2-D visualization library *Matplotlib* (Hunter, 2007) and *NumPy* for efficient numerical computations (Walt et al., 2011). Our implementation is specifically intended for combination with other packages to harvest efficient implementations in the best possible way.

Especially in this current time of rapid development of open-source scientific software packages and powerful

machine-learning frameworks, we consider an open-source implementation of a geological modeling tool as essential. We therefore aim to open up this possibility to a wide community, by combining state-of-the-art implicit geological modeling techniques with additional sophisticated Python packages for scientific programming and data analysis in an open-source ecosystem. The aim is explicitly not to rival the existing commercial packages with well-designed graphical user interfaces, underlying databases and highly advanced workflows for specific tasks in subsurface engineering, but to provide an environment to enhance existing methodologies as well as give access to an advanced modeling algorithm for scientific experiments in the field of geomodeling.

In the following, we will present the implementation of our code in the form of core modules, related to the task of geological modeling itself, and additional assets, which provide the link to external libraries, e.g., to facilitate stochastic geomodeling and the inversion of structural data. Each part is supported and supplemented with Jupyter notebooks that are available as additional online material and part of the package documentation, which enable the direct testing of our methods (see Sect. A3). These notebooks can also be directly executed in an online environment (Binder). We encourage the reader to use these Jupyter tutorial notebooks to follow along the steps explained in the following. Finally, we discuss our approach, specifically with respect to alternative modeling approaches in the field, and provide an outlook to our planned future work for this project.

2 CORE – geological modeling with *GemPy*

In this section, we describe the core functionality of *GemPy*: the construction of 3-D geological models from geological input data (surface contact points and orientation measurements) and defined topological relationships (stratigraphic sequences and fault networks). We begin with a brief review of the theory underlying the implemented interpolation algo-

rithm. We then describe the translation of this algorithm and the subsequent model generation and visualization using the Python front end of *GemPy* and how an entire model can be constructed by calling only a few functions. Throughout the text, we include code snippets with minimal working examples to demonstrate the use of the library.

After describing the simple functionality required to construct models, we go deeper into the underlying architecture of *GemPy*. This part is not only relevant for advanced users and potential developers but also highlights a key aspect: the link to *Theano* (Theano Development Team, 2016), a highly evolved Python library for efficient vector algebra and machine learning, which is an essential aspect required for making use of the more advanced aspects of stochastic geomodeling and Bayesian inversion, which will also be explained in the subsequent sections.

2.1 Geological modeling and the potential-field approach

2.1.1 Concept of the potential-field method

The potential-field method developed by Lajaunie et al. (1997) is the central method to generate the 3-D geological models in *GemPy*, which has already been successfully deployed in the modeling software GeoModeller 3-D (see Calcagno et al., 2008). The general idea is to construct an interpolation function $\mathbf{Z}(\mathbf{x}_0)$ where \mathbf{x} is any point in the continuous three-dimensional space $(x, y, z) \in \mathbb{R}^3$, which describes the domain \mathcal{D} as a scalar field. The gradient of the scalar field will follow the planar orientation of the stratigraphic structure throughout the volume or, in other words, every possible isosurface of the scalar field will represent every synchronous deposition of the layer (see Fig. 2).

Let's break down what we actually mean by this: imagine that a geological setting is formed by a perfect sequence of horizontal layers piled one above the other. If we know the exact timing of when one of these surfaces was deposited, we would know that any layer above had to occur afterwards while any layer below had to be deposited earlier in time. Obviously, we cannot have data for each of these infinitesimal synchronous layers, but we can interpolate the “date” between them. In reality, the exact year of the synchronous deposition is meaningless – as it is not possible to remotely obtain accurate estimates. What has value to generate a 3-D geomodel is the location of those synchronous layers and especially the lithological interfaces where the change of physical properties are notable. Because of this, instead of interpolating *time*, we use a simple dimensionless parameter that we simply refer to as *scalar field value*.

The advantages of using a global interpolator instead of interpolating each layer of interest independently are twofold: (i) the location of one layer affects the location of others in the same depositional environment, making it impossible for two layers in the same potential field to cross; and (ii) it en-

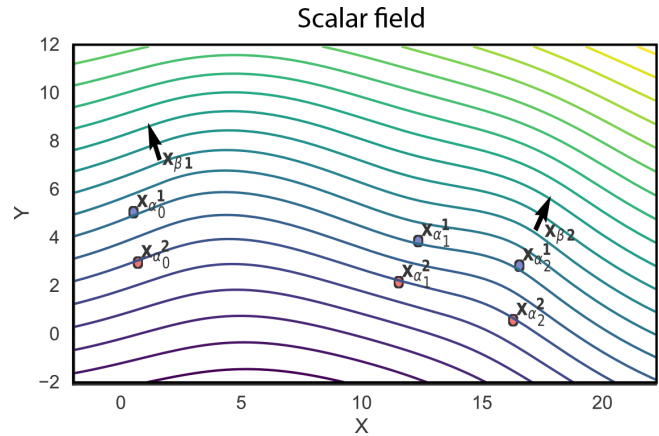


Figure 2. Example of scalar field. The input data are formed by six points distributed in two layers ($\mathbf{x}_{\alpha_i}^1$ and $\mathbf{x}_{\alpha_i}^2$) and two orientations (\mathbf{x}_{β_j}). An isosurface connects the interface points and the scalar field is perpendicular to the foliation gradient.

ables the use of data between the interfaces of interest, opening the range of possible measurements that can be used in the interpolation.

The interpolation function is obtained as a weighted interpolation based on universal cokriging (Chiles and Delfiner, 2009). Kriging or Gaussian process regression (Matheron, 1981) is a spatial interpolation that treats each input as a random variable, aiming to minimize the covariance function to obtain the best linear unbiased predictor (for a detailed description see Chap. 3 in Wackernagel, 2013). Furthermore, it is possible to combine more than one type of data, i.e., a multivariate case or cokriging, to increase the amount of information in the interpolator, as long as we capture their relation using a cross-covariance. The main advantage in our case is to be able to utilize orientations sampled from different locations in space for the estimation. Simple kriging, as a regression, only minimizes the second moment of the data (or variances). However, in most geological settings, we can expect linear trends in our data, i.e., the mean thickness of a layer varies across the region linearly. This trend is captured using polynomial drift functions to the system of equations in what is called universal kriging.

2.1.2 Adjustments to structural geological modeling

So far we have shown what we want to obtain and how universal cokriging is a suitable interpolation method to get there. In the following, we will describe the concrete steps from taking our input data to the final interpolation function $\mathbf{Z}(\mathbf{x}_0)$, where \mathbf{x}_0 refers to the estimated quantity for some integrable measure p_0 . Much of the complexity of the method comes from the difficulty of keeping highly nested nomenclature consistent across literature. For this reason, we will try to be especially verbose regarding the mathematical terminology based primarily on Chiles et al. (2004). The terms

of *potential field* (original coined by Lajaunie et al., 1997) and *scalar field* (preferred by the authors) are used interchangeably throughout the paper. The result of a kriging interpolation is a random function and hence both *interpolation function* and *random function* are used to refer to the function of interest, $\mathbf{Z}(\mathbf{x}_0)$. The cokriging nomenclature quickly grows convoluted, since it has to consider p random functions \mathbf{Z}_i , with p being the number of distinct parameters involved in the interpolation, sampled at different points \mathbf{x} of the three-dimensional domain \mathbb{R}^3 . Two types of parameters are used to characterize the *scalar field* in the interpolation: (i) layer interface points \mathbf{x}_α describing the respective isosurfaces of interest, usually the interface between two layers; and (ii) the gradients of the scalar field, \mathbf{x}_β , or in geological terms, poles of the layer, i.e., normal vectors to the dip plane. Therefore, gradients will be oriented perpendicular to the isosurfaces and can be located anywhere in space. We will refer to the main random function \mathbf{Z}_α – the scalar field itself – simply as \mathbf{Z} , and its set of samples as \mathbf{x}_α , while the second random function \mathbf{Z}_β – the gradient of the scalar field – will be referred to as $\partial\mathbf{Z}/\partial u$ with u being any unit vector and its samples as \mathbf{x}_β . We can capture the relationship between the scalar field \mathbf{Z} and its gradient as

$$\frac{\partial\mathbf{Z}}{\partial u}(\mathbf{x}) = \lim_{\rho \rightarrow 0} \frac{\mathbf{Z}(\mathbf{x} + u) - \mathbf{Z}(\mathbf{x})}{\rho}. \quad (1)$$

It is also important to keep the values of every individual synchronal layer identified since they have the same scalar field value. Therefore, samples that belong to a single layer k will be expressed as a subset denoted using superscript as \mathbf{x}_{α}^k and every individual point by a subscript, $\mathbf{x}_{\alpha i}^k$ (see Fig. 2).

Note that in this context the scalar field property α is dimensionless. The only mathematical constrain is that the value must increase in the direction of the gradient, which in turn describes the stratigraphic deposition. Therefore the two constraints we want to conserve in the interpolated scalar field are (i) all points belonging to a determined interface $\mathbf{x}_{\alpha i}^k$ must have the same scalar field value (i.e., there is an isosurface connecting all data points)

$$\mathbf{Z}(\mathbf{x}_{\alpha i}^k) - \mathbf{Z}(\mathbf{x}_{\alpha 0}^k) = 0, \quad (2)$$

where $\mathbf{x}_{\alpha 0}^k$ is a reference point of the interface and (ii) the scalar field will be perpendicular to the gradient (poles in geological nomenclature) \mathbf{x}_β anywhere in 3-D space. It is important to mention that the choice of the reference points $\mathbf{x}_{\alpha 0}^k$ has no effect on the results.

Considering Eq. (2), we do not care about the exact value at $\mathbf{Z}(\mathbf{x}_{\alpha i}^k)$ as long as it is constant at all points $\mathbf{x}_{\alpha i}^k$. Therefore, the random function \mathbf{Z} in the cokriging system (Eq. 4) can be substituted by Eq. (2). This formulation entails that the specific *scalar field values* will only depend on the gradients and hence at least one gradient is necessary to keep the system of equations defined. The advantage of this mathematical construction is that by not fixing the values of each interface

$\mathbf{Z}(\mathbf{x}_{\alpha}^k)$, the compression of layers – i.e., the rate of change of the scalar field – will only be defined by the gradients $\partial\mathbf{Z}/\partial u$. This allows us to propagate the effect of each gradient beyond the surrounding interfaces creating smoother formations.

The algebraic dependency between \mathbf{Z} and $\partial\mathbf{Z}/\partial u$ (Eq. 1) gives a mathematical definition of the relation between the two variables, avoiding the need of an empirical cross-variogram and instead enabling the use of the derivation of the covariance function. This dependency must be taken into consideration in the computation of the drift of the first moment as well having a different function for each of the variables,

$$\lambda F_1 + \lambda F_2 = f_{10}, \quad (3)$$

where F_1 is a polynomial of degree n , F_2 its derivative between the input data x_α and x_β , and f_{10} corresponds to the same polynomial to the objective point x_0 . Having taken this into consideration, the resulting cokriging system takes the following form:

$$\begin{bmatrix} \mathbf{C}_{\partial\mathbf{Z}/\partial u, \partial\mathbf{Z}/\partial v} & \mathbf{C}_{\partial\mathbf{Z}/\partial u, \mathbf{Z}} & \mathbf{U}_{\partial\mathbf{Z}/\partial u} \\ \mathbf{C}_{\mathbf{Z}, \partial\mathbf{Z}/\partial u} & \mathbf{C}_{\mathbf{Z}, \mathbf{Z}} & \mathbf{U}_{\mathbf{Z}} \\ \mathbf{U}'_{\partial\mathbf{Z}/\partial u} & \mathbf{U}'_{\mathbf{Z}} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \lambda_{\partial\mathbf{Z}/\partial u, \partial\mathbf{Z}/\partial v} & \lambda_{\partial\mathbf{Z}/\partial u, \mathbf{Z}} \\ \lambda_{\mathbf{Z}, \partial\mathbf{Z}/\partial u} & \lambda_{\mathbf{Z}, \mathbf{Z}} \\ \mu_{\partial u} & \mu_u \end{bmatrix} = \begin{bmatrix} \mathbf{c}_{\partial\mathbf{Z}/\partial u, \partial\mathbf{Z}/\partial v} & \mathbf{c}_{\partial\mathbf{Z}/\partial u, \mathbf{Z}} \\ \mathbf{c}_{\mathbf{Z}, \partial\mathbf{Z}/\partial u} & \mathbf{c}_{\mathbf{Z}, \mathbf{Z}} \\ \mathbf{f}_{10} & \mathbf{f}_{20} \end{bmatrix}, \quad (4)$$

where $\mathbf{C}_{\partial\mathbf{Z}/\partial u}$ is the gradient covariance matrix; $\mathbf{C}_{\mathbf{Z}, \mathbf{Z}}$ the covariance matrix of the differences between each interface point to reference points in each layer,

$$\mathbf{C}_{x_{\alpha i}^r, x_{\alpha j}^s} = \mathbf{C}_{x_{\alpha, i}^r, x_{\alpha, j}^s} - \mathbf{C}_{x_{\alpha, 0}^r, x_{\alpha, j}^s} - \mathbf{C}_{x_{\alpha, i}^r, x_{\alpha, 0}^s} + \mathbf{C}_{x_{\alpha, 0}^r, x_{\alpha, 0}^s} \quad (5)$$

(see Appendix B2 for further analysis); $\mathbf{C}_{\mathbf{Z}, \partial\mathbf{Z}/\partial u}$ encapsulates the cross-covariance function; and $\mathbf{U}_{\mathbf{Z}}$ and $\mathbf{U}'_{\partial\mathbf{Z}/\partial u}$ are the drift functions and their gradients, respectively. On the right-hand side we find the matrix of independent terms, $\mathbf{c}_{\partial\mathbf{Z}/\partial u, \partial\mathbf{Z}/\partial v}$ being the gradient of the covariance function to the point \mathbf{x} of interest; $\mathbf{c}_{\mathbf{Z}, \partial\mathbf{Z}/\partial u}$ the cross-covariance; $\mathbf{c}_{\mathbf{Z}, \mathbf{Z}}$ the actual covariance function; and \mathbf{f}_{10} and \mathbf{f}_{20} the gradient of the drift functions and the drift functions themselves, respectively. Lastly, the unknown vectors are formed by the corresponding weights, λ , and constants of the drift functions μ . A more detail inspection of this system of equations is carried out in Appendix B.

As we can see in Eq. (4), it is possible to solve the kriging system for the scalar field \mathbf{Z} (second column in the weights vector), as well as its derivative $\partial\mathbf{Z}/\partial u$ (first column in the weights vector). Even though the main goal is the segmentation of the layers, which is done using the value of \mathbf{Z}

(see Sect. 2.2.1), the gradient of the scalar field can be used for further mathematical applications, such as meshing, geophysical forward calculations or locating geological structures of interest (e.g., spill points of a hydrocarbon trap).

Furthermore, since the choice of covariance parameters is ad hoc (Appendix D shows the covariance function used in *GemPy*), the uncertainty derived by the kriging interpolation does not bear any physical meaning. This fact promotes the idea of only using the mean value of the kriging solution. For this reason it is recommended to solve the kriging system (Eq. 4) in its dual form (Matheron, 1981, see Appendix C).

2.2 Geological model interpolation using *GemPy*

2.2.1 From scalar field to geological block model

In most scenarios the goal of structural modeling is to define the spatial distribution of geological structures, such as layers, interfaces and faults. In practice, this segmentation is usually done either by using a volumetric discretization or by depicting the interfaces as surfaces.

The result of the kriging interpolation is the random function $\mathbf{Z}(x)$ (and its gradient $\partial\mathbf{Z}/\partial u(x)$, which we will omit in the following), which allows for the evaluation of the value of the scalar field at any given point x in space. From this point on, the easiest way to segment the domains is to discretize the 3-D space (e.g., we use a regular grid in Fig. 3). First, we need to calculate the scalar value at every interface by computing $\mathbf{Z}(\mathbf{x}_{\alpha,i}^k)$ for every interface k_i . Once we know the value of the scalar field at the interfaces, we evaluate every point of the mesh and compare their value to those at the interfaces, identifying every point of the mesh with a topological volume. Each of these compartmentalizations will represent each individual domain, i.e., each lithology of interest (see Fig. 3a).

At the time of this manuscript's preparation, *GemPy* only provides rectilinear grids but it is important to notice that the computation of the scalar field happens in continuous space, and therefore allows for the use of any type of mesh. The result of this type of segmentation is referred to in *GemPy* as a *lithology block*.

The second alternative segmentation consists of locating the layer isosurfaces. *GemPy* makes use of the marching cube algorithm (Lorensen and Cline, 1987) provided by the *scikit-image* library (van der Walt et al., 2014). The basics of the marching cube algorithm are quite intuitive. (i) First, we discretize the volume in 3-D voxels and by comparison we look to see if the value of the isosurface we want to extract falls within the boundary of every single voxel; (ii) if so, for each edge of the voxel, we interpolate the values at the corners of the cube to obtain the coordinates of the intersection between the edges of the voxels and the isosurface of interest, commonly referred to as vertices; (iii) those intersections are analyzed and compared against all possible configurations to define the simplices (i.e., the vertices that form an individual

polygon) of the triangles. Once we obtain the coordinates of vertices and their correspondent simplices, we can use them for visualization (see Sect. 3.1) or any subsequent computation that may make use of them (e.g., weighted voxels). For more information on meshing algorithms refer to Geuzaine and Remacle (2009).

2.2.2 Combining scalar fields: depositional series and faults

In reality, most geological settings are formed by a concatenation of depositional phases partitioned by unconformity boundaries and subjected to tectonic stresses that displace and deform the layers. While the interpolation is able to represent realistic folding – given enough data – the method fails to describe discontinuities. To overcome this limitation, it is possible to combine several scalar fields to recreate the desired result.

So far the implemented discontinuities in *GemPy* are unconformities and infinite faults. Both types are computed by specific combinations of independent scalar fields. We call these independent scalar fields *series* (from stratigraphic series in accordance to the use in GeoModeller 3-D; Calcagno et al., 2008), and in essence they represent a subset of grouped interfaces – either layers or fault planes – that are interpolated together and therefore their spatial location affect each other. To handle and visualize these relationships, we use a so-called sequential pile, representing the order – from the first scalar field to the last – and the grouping of the layers (see Fig. 3). It is interesting to point out that the sequential pile only controls the order of each individual series. Within each series, the stratigraphic sequence is strictly determined by the geometry and the interpolation algorithm.

Modeling unconformities is rather straightforward. Once we have grouped the layers into their respective series, younger series will overlay older ones beyond the unconformity. The scalar fields themselves, computed for each of these series, could be seen as a continuous depositional sequence in the absence of an unconformity.

```

import gempy as gp

# Main data management object containing
geo_data = gp.create_data(extent=[0, 20, 0, 10, -10, 0],
                          resolution=[100, 10, 100],
                          path_o="paper_Foliations.csv",
                          path_i="paper_Points.csv")

# Creating object with data prepared for interpolation and compiling
interp_data = gp.InterpolatorData(geo_data)

# Computing result
lith, fault = gp.compute_model(interp_data)

# Plotting result: scalar field
gp.plot_scalar_field(geo_data, lith[1], 5, plot_data=True)

# Plotting result: lithology block
gp.plot_section(geo_data, lith[0], 5, plot_data=True)

# Getting vertices and faces
vertices, simplexes = gp.get_surfaces(interp_data, lith[1], [fault[1]], original_scale=True)

```

Listing 1. Code to generate a single scalar field model (as seen in Fig. 2) and plotting a section of a regular grid (Fig. 3a) and extracting surfaces points at the interfaces.

Faults are modeled by the inclusion of an extra drift term into the kriging system (Marechal, 1984):

$$\begin{bmatrix} \mathbf{C}_{\partial\mathbf{Z}/\partial u, \partial\mathbf{Z}/\partial v} & \mathbf{C}_{\partial\mathbf{Z}/\partial u, Z} & \mathbf{U}_{\partial\mathbf{Z}/\partial u} & \mathbf{F}_{\partial\mathbf{Z}/\partial u} \\ \mathbf{C}_{Z, \partial\mathbf{Z}/\partial u} & \mathbf{C}_{Z, Z} & \mathbf{U}_Z & \mathbf{F}_Z \\ \mathbf{U}'_{\partial\mathbf{Z}/\partial u} & \mathbf{U}'_Z & \mathbf{0} & \mathbf{0} \\ \mathbf{F}'_{\partial\mathbf{Z}/\partial u} & \mathbf{F}'_Z & \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \lambda_{\partial\mathbf{Z}/\partial u, \partial\mathbf{Z}/\partial v} & \lambda_{\partial\mathbf{Z}/\partial u, Z} \\ \lambda_{Z, \partial\mathbf{Z}/\partial u} & \lambda_{Z, Z} \\ \mu_{\partial u} & \mu_u \\ \mu_{\partial f} & \mu_f \end{bmatrix} = \begin{bmatrix} \mathbf{c}_{\partial\mathbf{Z}/\partial u, \partial\mathbf{Z}/\partial v} & \mathbf{c}_{\partial\mathbf{Z}/\partial u, Z} \\ \mathbf{c}_{Z, \partial\mathbf{Z}/\partial u} & \mathbf{c}_{Z, Z} \\ \mathbf{f}_{10} & \mathbf{f}_{20} \\ \mathbf{f}_{10} & \mathbf{f}_{20} \end{bmatrix}, \quad (6)$$

which is a function of the faulting structure. This means that for every location \mathbf{x}_0 the drift function will take a value depending on the fault compartment – i.e., a segmented domain of the fault network – and other geometrical constrains such as spatial influence of a fault or variability in the offset. To obtain the offset effect of a fault, the value of the drift function has to be different at each of its sides. The level of complexity of the drift functions will determine the quality of the characterization as well as its robustness. Furthermore, finite or localized faults can be recreated by selecting an adequate function that describes those specific trends.

The computation of the segmentation of fault compartments (called *fault block* in *GemPy*) – prior to the inclusion of the fault drift functions that depend on this segmentation – can be performed with the potential-field method itself. In the case of multiple faults, individual drift functions have to be included in the kriging system for each fault, representing the

subdivision of space that they produce. Naturally, younger faults may offset older tectonic events. This behavior is replicated by recursively adding drift functions of younger faults to the computation of the older *fault blocks*. To date, the fault relations – i.e., which faults offset others – is described by the user in a Boolean matrix. An easy-to-use implementation to generate fault networks is being worked on at the time of the manuscript preparation.

An important detail to consider is that drift functions will bend the isosurfaces according to the given rules, but they will conserve their continuity. This differs from the intuitive idea of offset, where the interface presents a sharp jump. This fact has a direct impact on the geometry of the final model, and can, for example, affect certain meshing algorithms. Furthermore, in the ideal case of choosing the perfect drift function, the isosurface would bend exactly along the faulting plane. In the current state, *GemPy* only includes the addition of an arbitrary integer to each segmented volume. This limits the quality to a constant offset, decreasing the sharpness of the offset as data deviates from that constraint. Any deviation from this theoretical concept results in a bending of the layers as they approximate the fault plane to accommodate the data, potentially leading to overly smooth transitions around the discontinuity.

2.3 “Under the hood”: the *GemPy* architecture

2.3.1 The graph structure

The architecture of *GemPy* follows the Python Software Foundation recommendations of modularity and reusability (van Rossum et al., 2001). The aim is to divide all functionality into small independent logical units in order to avoid du-

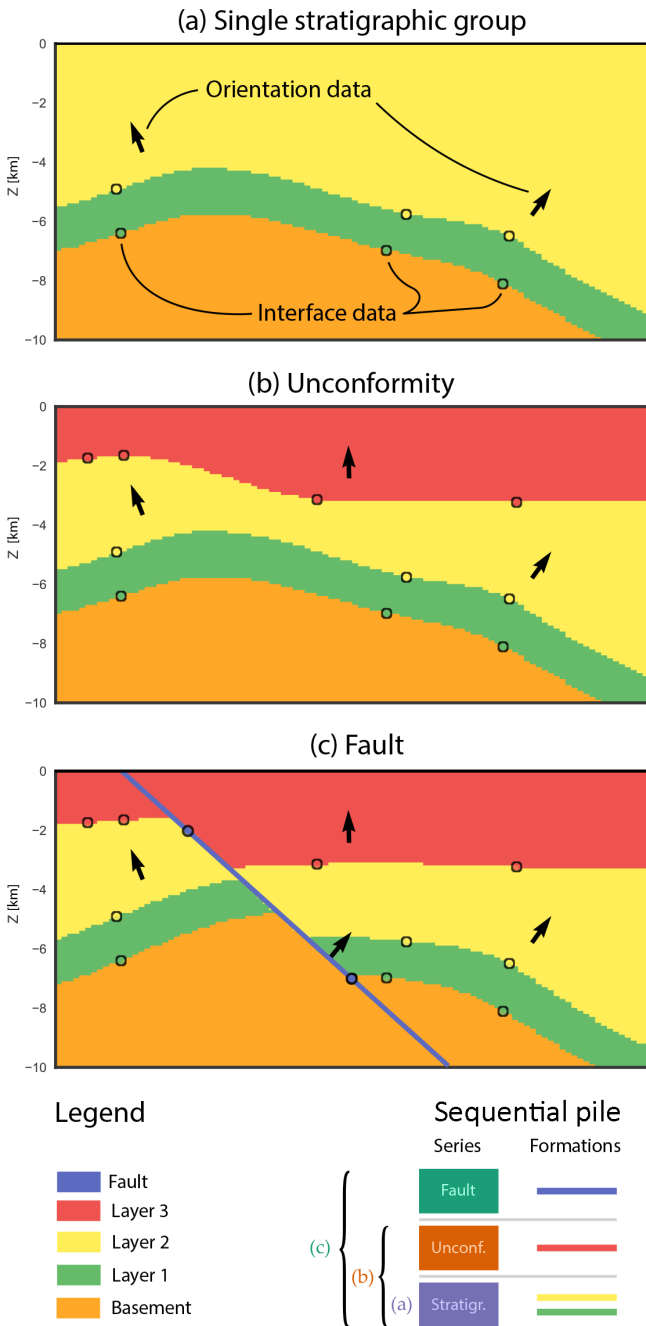


Figure 3. Example of different lithological units and their relation to scalar fields; (a) Simple stratigraphic sequence generated from a scalar field as product of the interpolation of interface points and orientation gradients; (b) Addition of an unconformity horizon from which the unconformity layer behaves independently from the older layers by overlying a second scalar field; (c) Combination of unconformity and faulting using three scalar fields.

plication, facilitate readability and make changes to the code base easier.

GemPy's architecture was designed from the ground up to accommodate an automatic differentiation (AD) library. The

main constraint is that the mathematical functions need to be continuous from the variables (in probabilistic jargon priors) to the cost function (or likelihoods), and therefore the code must be written in the same language (or at the very least compatible) to automatically compute the derivatives. In practice, this means that any operation involved in the AD must be coded symbolically using the library *Theano* (see Sect. 2.3.2 for further details). Writing symbolically requires a priori declaration of all algebra, from variables that will behave as latent parameters – i.e., the parameters we try to tune for optimization or uncertainty quantification – to all involved constants and the specific mathematical functions that relates them. These statements generate a so-called graph that symbolically encapsulates all the logic that enables us to perform further analysis on the logic itself (e.g., differentiation or optimization). However, the rigidity when constructing the graph dictates the whole design of input data management.

GemPy encapsulates this creation of the symbolic graph in its module `theanograph`. Due to the significant complexity in programming symbolically, features included in *GemPy* that heavily rely on external libraries are not written in *Theano* yet. The current functionality written in *Theano* can be seen in the Fig. 4 and it essentially encompasses all the interpolation of the geological modeling (Sect. 2.1) as well as forward calculation of the gravity (Sect. 3.2).

Regarding data structure, we make use of the Python package *pandas* (McKinney, 2011) to store and prepare the input data for the symbolic graph (red nodes in Fig. 4), or other processes such as visualization. All of the methodology to create, export and manipulate the original data is encapsulated in the class `DataManagement`. This class has several child classes to facilitate specific precomputation manipulations of data structures (e.g., for meshing). The aim is to have all constant data prepared before any inference or optimization is carried out to minimize the computational overhead of each iteration as much as possible.

It is important to keep in mind that, in this structure, once data enters the part of the symbolic graph, only algebraic operations are allowed. This limits the use of many high-level coding structures (e.g., dictionaries or undefined loops) and external dependencies. As a result of that, the preparation of data must be exhaustive before starting the computation. This includes ordering the data within the arrays and passing the exact lengths of the subsets we will need later on during the interpolation or the calculation of many necessary constant parameters. The preprocessing of data is done within the subclasses of `DataManagement`, the `InterpolatorData` class – of which an instance is used to call the *Theano* graph – and `InterpolatorClass`, which creates the *Theano* variables and compiles the symbolic graph.

The rest of the package is formed by a (always growing) series of modules that perform different tasks using the geological model as input (see Sect. 3 and the Assets area in Fig. 4).

```

import gempy as gp

# Main data management object containing
geo_data = gp.create_data(extent=[0, 20, 0, 10, -10, 0],
                          resolution=[100, 10, 100],
                          path_o="paper_Foliations.csv",
                          path_i="paper_Points.csv")

# Defining the series of the sequential pile
gp.set_series(geo_data,
              {'younger_serier' : 'Unconformity', 'older_serier': ('Layer1', 'Layer2')},
              order_formation= ['Unconformity', 'Layer2', 'Layer1'])

# Creating object with data prepared for interpolation and compiling
interp_data = gp.InterpolatorData(geo_data)

# Computing result
lith, fault = gp.compute_model(interp_data)

# Plotting result
gp.plot_section(geo_data, lith[0], 5, plot_data=True)

```

Listing 2. Extension of the code in Listing 1 to generate an unconformity by using two scalar fields. The corresponding model is shown in Fig. 3b).

```

import gempy as gp

# Main data management object containing the location of all interfaces and orientations
# parameters as well as the formation/fault to which belong
geo_data = gp.create_data(extent=[0,20,0,10,-10,0],
                          resolution=[100,10,100],
                          path_o = "paper_Foliations.csv",
                          path_i = "paper_Points.csv")

# Defining the series of the sequential pile. This is done by categorizing interfaces and
# orientations by its formation label
gp.set_series(geo_data, series_distribution={'fault_serier1': 'fault1',
                                           'younger_serier' : 'Unconformity',
                                           'older_serier': ('Layer1', 'Layer2')},
              order_formation= ['fault1', 'Unconformity', 'Layer2', 'Layer1'])

# Creating object with data prepared for interpolation and compiling
interp_data = gp.InterpolatorData(geo_data)

# Computing result
lith, fault = gp.compute_model(interp_data)

# Plotting result
gp.plot_section(geo_data, lith[0], 5, plot_data=True)

# Getting vertices and faces and plotting
vertices, simplexes = gp.get_surfaces(interp_data,lith[1], [fault[1]], original_scale=True)
gp.plot_surfaces_3D(geo_data, ver_s, sim_s)

```

Listing 3. Code to generate a model with an unconformity and a fault using the three scalar fields model (as seen in Fig. 3c) and the visualization 3-D using VTK (see Fig. 5).

2.3.2 Theano

Efficiently solving a large number of algebraic equations, and especially their derivatives, can easily get unmanageable in terms of both time and memory. Up to this point we have referenced *Theano* many times and its related terms such as AD or symbolic programming. In this section we will provide the motivation for its use and why its capabilities make all the

difference in making implicit geological modeling available for uncertainty analysis.

Theano is a Python package that takes over many of the optimization tasks in order to create a computationally feasible code implementation. *Theano* relies on the creation of symbolical graphs that represent the mathematical expressions to compute. Most of the extended programming paradigms (e.g., procedural languages and object-oriented

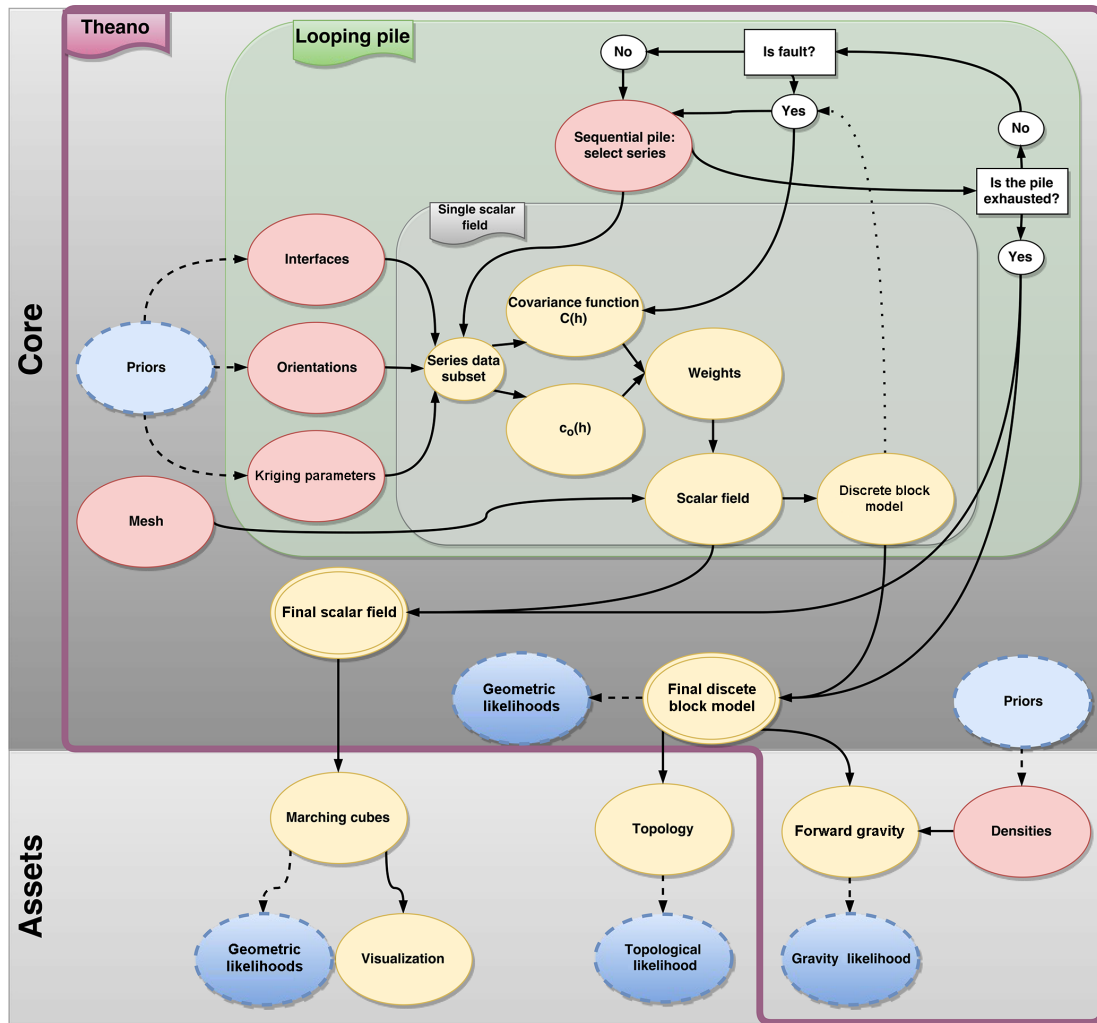


Figure 4. Graph of the logical structure of *GemPy* logic. There are several levels of abstraction represented. (i) The first division is between the implicit interpolation of the geological modeling (dark gray) and other subsequent operations for different objectives (light gray). (ii) All the logic required to perform automatic differentiation is presented under the *Theano* label (in purple). (iii) The parts under labels “Looping pile” (green) and “Single potential field” (gray) divide the logic to control the input data of each necessary scalar field and the operations within one of them. (iv) Finally, each superset of parameters is color coded according to their probabilistic nature and behavior in the graph: in blue, stochastic variables (priors or likelihoods); in yellow, deterministic functions; and in red, the inputs of the graph, which are either stochastic or constant depending on the problem.

programming) are sequentially executed without any interaction with the subsequent instructions. In other words, a later instruction has access to the memory states but is clueless about the previous instructions that have modified mentioned states. In contrast, symbolic programming defines, from the beginning to the end, not only the primary data structure but also the complete logic of a function, which in turn enables the optimization (e.g., redundancy) and manipulation (e.g., derivatives) of its logic.

Within the Python implementation, *Theano* creates an acyclic network graph where the parameters are represented by nodes, while the connections determine mathematical operators that relate them. The creation of the graph is done

in the class `theanograph`. Each individual method corresponds to a piece of the graph starting from the input data all the way to the geological model or the forward gravity (see Fig. 4, purple *Theano* area).

The symbolic graph is later analyzed to perform the optimization, the symbolic differentiation and the compilation to a faster language than Python (C or CUDA). This process is computationally demanding and therefore it must be avoided as much as possible.

Among the most outstanding optimizers included with *Theano* (for a detailed description see Theano Development Team, 2016), we can find (i) the canonicalization of the operations to reduce the number of duplicated computa-

tions, (ii) specialization of operations to improve consecutive element-wise operations, (iii) in-place operations to avoid duplications of memory or (iv) OpenMP parallelization for CPU computations. These optimizations and more can speed up the code by an order of magnitude.

However, although *Theano* code optimization is useful, the real game changer is its capability to perform AD. There is extensive literature explaining in detail the method and its related intuitions since it is a core algorithm to train neural networks (e.g., a detailed explanation is given by Baydin et al., 2015). Here, we will highlight the main differences with numerical approaches and how they can be used to improve the modeling process.

Many of the most advanced algorithms in computer science rely on an inverse framework, i.e., the result of a forward computation, $f(\mathbf{x})$, influences the value of one or many of the \mathbf{x} latent variables (e.g., neuronal networks, optimizations, inferences). The most emblematic example of this is the optimization of a cost function. All these problems can be described as an exploration of a multidimensional manifold $f: \mathbb{R}^N \rightarrow \mathbb{R}$. Hence the gradient of the function $\nabla f = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right)$ becomes key for an efficient analysis. In the case that the output is also multidimensional, i.e., $f: \mathbb{R}^N \rightarrow \mathbb{R}^M$, the entire manifold gradient can be expressed by the Jacobian matrix

$$Jf = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} \quad (7)$$

of dimension $N \cdot M$, where N is the number of variables and M the number of functions that depend on those variables. Now the question is how we compute the Jacobian matrix in a consistent and efficient manner. The most straightforward methodology consists of approximating the derivative by numerical differentiation and applying finite-difference approximations, e.g., a forward finite-difference scheme:

$$\frac{\partial f_i}{\partial x_j} = \lim_{h \rightarrow 0} \frac{f(x_j + h) - f(x_j)}{h}, \quad (8)$$

where h is a discrete increment. The main advantage of numerical differentiation is that it only computes f , evaluated for different values of x , which makes it very easy to implement it in any available code. By contrast, a drawback is that for every element of the Jacobian we are introducing an approximation error that can eventually lead to mathematical instabilities. But above all, the main limitation is the need of $2 \cdot M \cdot N$ evaluations of the function f , which quickly becomes prohibitively expensive to compute in high-dimensional problems.

The alternative is to create the symbolic differentiation of f . This encompasses decomposing f into its primal operators and applying the chain rule to the correspondent transformation by following the rules of differentiation to obtain

f' . However, symbolic differentiation is not enough since the application of the chain rule leads to exponentially large expressions of f' in what is known as “expression swell” (Cohen, 2003). Luckily, these large symbolic expressions have a high level of redundancy in their terms. This allows us to exploit this redundancy by storing the repeated intermediate steps in memory and simply invoking them when necessary, instead of computing the whole graph every time. This division of the program into subroutines to store the intermediate results – which are invoked several times – is called dynamic programming (Bellman, 2013). The simplified symbolic differentiation graph is ultimately what is called AD (Baydin et al., 2015). Additionally, in a multivariate/multi-objective case the benefits of using AD increase linearly as the difference between the number of parameters N and the number of objective functions M gets larger. By applying the same principle of redundancy explained above – this time between intermediate steps shared across multiple variables or multiple objective functions – it is possible to reduce the number of evaluations necessary to compute the Jacobian either to N in forward-propagation or to M in back-propagation, plus a small overhead on the evaluations (for a more detailed description of the two modes of AD see Cohen, 2003).

Theano provides a direct implementation of the back-propagation algorithm, which means in practice that a new graph of similar size is generated per cost function (or, in the probabilistic inference, per likelihood function; see Sect. 3.4 for more detail). Therefore, the computational time is independent of the number of input parameters, opening the door to solving high-dimensional problems.

3 Assets – model analysis and further use

In this second half of the paper we will explore different features that complement and expand the construction of the geological model itself. These extensions are just some examples of how *GemPy* can be used as a geological modeling engine for diverse research projects. The numerous libraries in the open-source ecosystem allow us to choose the best narrow-purpose tool for very specific tasks. Considering the visualization of *GemPy* for instance: *matplotlib* (Hunter, 2007) for 2-D visualization, *vtk* for fast and interactive 3-D visualization, *steno3D* for sharing block model visualizations online, or even the open-source 3-D modeling software *Blender* (Blender Online Community, 2017) for creating high-quality renderings and virtual reality are only some examples of the flexibility that the combination of *GemPy* with other open-source packages offers. In the same fashion, we can use the geological model as a basis for the subsequent geophysical simulations and process simulations. Because of Python’s modularity, combining distinct modules to extend the scope of a project to include the geological modeling process into a specific environment is effortless. In the next sections we will dive into some of the built-in functionality

implemented to date on top of the geological modeling core. Current assets are (i) 2-D and 3-D visualizations, (ii) forward calculation of gravity, (iii) topology analysis, (iv) uncertainty quantification (UQ), and (v) full Bayesian inference.

3.1 Visualization

The segmentation of meaningful units is the central task of geological modeling. It is often a prerequisite for engineering projects or process simulations. An intuitive 3-D visualization of a geological model is therefore a fundamental requirement.

For its data and model visualization, *GemPy* makes use of freely available tools in the Python module ecosystem to allow the user to inspect data and modeling results from all possible angles. The fundamental plotting library *matplotlib* (Hunter, 2007), enhanced by the statistical data visualization library *seaborn* (Waskom et al., 2017), provides the 2-D graphical interface to visualize input data and 2-D sections of scalar fields and geological models. In addition, making use of the capacities of *pyqt* implemented with *matplotlib*, we can generate interactive sequential piles, where the user can not only visualize the temporal relation of the different unconformities and faulting events but also modify it using intuitive drag-and-drop functionality (see Fig. 5).

On top of these features, *GemPy* offers in-built 3-D visualization based on the open-source Visualization Toolkit (VTK; Schroeder et al., 2004). It provides users with an interactive 3-D view of the geological model, as well as three additional orthogonal viewpoints (see Fig. 5). The user can decide to plot just the data, the geological surfaces or both. In addition to just visualizing the data in 3-D, *GemPy* makes use of the interaction capabilities provided by *vtk* to allow the user to move input data points on the fly via drag-and-drop functionality. Combined with *GemPy*'s optimized modeling process (and the ability to use GPUs for efficient model calculation), this feature allows for data modification with real-time updating of the geological model (in the order of milliseconds per scalar field). This functionality can not only improve the understanding of the model but can also help the user to obtain the desired outcome by working directly in 3-D space while getting direct visual feedback on the modeling results. However, due to the exponential increase in computational time with respect to the number of input data and the model resolution, very large and complex models may have difficulties to render fast enough to perceive continuity on conventional computer systems.

For additional high-quality visualization, we can generate *vtk* files using *pyevtk*. These files can later be loaded into external VTK viewer as Paraview (Ayachit, 2015) in order to take advantage of its intuitive interface and powerful visualization options. Another natural compatibility exists with Blender (Blender Online Community, 2017) due to its use of Python as front end. Using the Python distribution included within a Blender installation, it is possible to import, run and

automatically represent *GemPy*'s data and results (Fig. 1, see Appendix F for code extension). This not only allows users to render high-quality images and videos but also to visualize the models in virtual reality, making use of the Blender game engine and some of the plug-ins that enable this functionality.

For sharing models, *GemPy* also includes functionality to upload discretized models to the Steno 3-D platform (a freemium business model). Here, it is possible to visualize, manipulate and share the model with any number of people effortlessly by simple invitations or the distribution of a link.

In short, *GemPy* is not limited to a unique visualization library. Currently *GemPy* lends support to many of the available visualization options to fulfill the different needs of the developers accordingly. However, these are not by all means the only possible alternatives and in the future we expect that *GemPy* will be employed as the back end of other projects.

3.2 Gravity forward modeling

In recent years gravity measurements have increased in quality (Nabighian et al., 2005) and is by now a valuable additional geophysical data source to support geological modeling. There are different ways to include the new information into the modeling workflow, and one of the most common is via inversions (Tarantola, 2005). Geophysics can validate the quality of the model in a probabilistic or optimization framework, but also by back-propagating information geophysics can automatically improve the modeling process itself. As a drawback, simulating forward geophysics adds a significant computational cost and increases the uncertainty to the parametrization of the model. However, due to the amount of uncorrelated information – often continuous in space – the inclusion of geophysical data in the modeling process usually becomes significant to evaluate the quality of a given model.

GemPy includes built-in functionality to compute forward gravity conserving the AD of the package. It is calculated from the discretized block model applying the method of Nagy (1966) for rectangular prisms in the Z direction,

$$F_z = G_\rho [|x \ln(y+r) + y \ln(x+r) - z \arctan\left(\frac{xy}{zr}\right)|_{x_1}^{x_2} |_{y_1}^{y_2} |_{z_1}^{z_2}], \quad (9)$$

where x , y and z are the Cartesian components from the measuring point of the prism; r the Euclidean distance and G_ρ the average gravity pull of the prism. This integration provides the gravitational pull of every voxel for a given density and distance in the component z . Taking advantage of the immutability of the involved parameters, with the exception of density, allows us to precompute the decomposition of t_z – i.e., the distance-dependent side of Eq. (9) – leaving just its product with the weight G_ρ ,

$$F_z = G_\rho \cdot t_z, \quad (10)$$

as a recurrent operation.

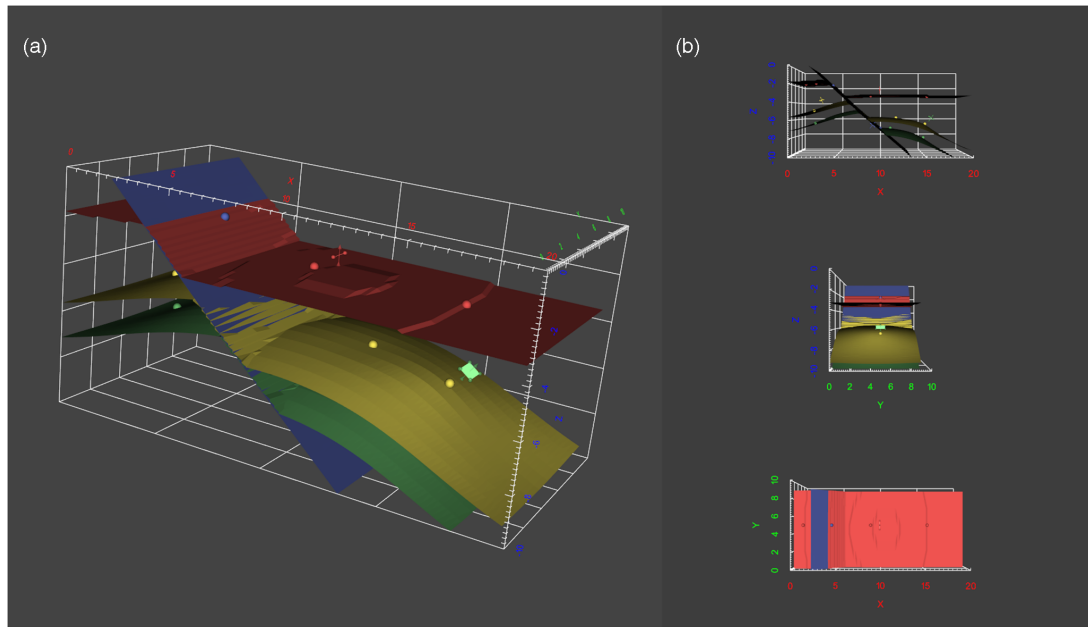


Figure 5. In-built *vtk* 3-D visualization of *GemPy* provides an interactive visualization of the geological model (a) and three additional orthogonal viewpoints (b) from different directions.

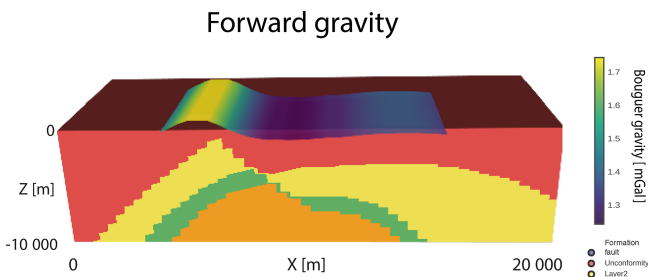


Figure 6. Forward gravity response overlaid on top of a 3-D lithology block sliced on the *Y* direction.

As an example, we show here the forward gravity response of the geological model in Fig. 3c. The first important detail is the increased extent of the interpolated model to avoid boundary errors. In general, a padding equal to the maximum distance used to compute the forward gravity computation would be the ideal value. In this example (Fig. 6) we add 10 km to the *X* and *Y* coordinates. The next step is to define the measurement 2-D grid – i.e., where to simulate the gravity response and the densities of each layer. The densities chosen are 2.92, 3.1, 2.61 and 2.92 kg m⁻³ for the basement, “Unconformity” layer (i.e., the layer on top of the unconformity), Layer 1 and Layer 2, respectively.

The computation of forward gravity is a required step towards a fully coupled gravity inversion. Embedding this step into a Bayesian inference allows us to condition the initial data used to create the model to the final gravity response. This idea will be further developed in Sect. 3.4.2.

3.3 Topology

The concept of topology provides a useful tool to describe adjacency relations in geomodels, such as stratigraphic contacts or across-fault connectivity (for a more detailed introduction see Thiele et al., 2016a, b). *GemPy* has in-built functionality to analyze the adjacency topology of its generated models as region adjacency graphs (RAGs), using the `topology_compute` method (see Listing 6). It can be directly visualized on top of model sections (see Fig. 7), where each unique topological region in the geomodel is represented by a graph node, and each connection as a graph edge. The function outputs the graph object *G*, the region centroid coordinates, a list of all the unique node labels and two look-up tables to conveniently reference node labels and lithologies

To analyze the model topology, *GemPy* makes use of a general connected component labeling (CCL) algorithm to uniquely label all separated geological entities in 3-D geomodels. The algorithm is provided via the widely used, open-source, Python-based image processing library *scikit-image* (van der Walt et al., 2014) by the function `skimage.measure.label`, which is based on the optimized algorithms of Fiorio and Gustedt (1996) and Wu et al. (2005). But just using CCL on a 3-D geomodel fails to discriminate a layer cut by a fault into two unique regions because in practice both sides of a fault are represented by the same label. To achieve the detection of edges across the fault, we need to precondition the 3-D geomodel matrix, which contains just the lithology information (layer id), with a 3-

```

import matplotlib.pyplot as plt
import gempy as gp

# Main data management object containing. The extent must be large enough respect the forward
# gravity plane to account the effect of all cells at a given distance, $d$ to any spatial
# direction $x, y, z$.
geo_data = gp.create_data(extent=[-10,30,-10,20,-10,0],
                          resolution=[50,50,50],
                          path_o = "paper_Foliations.csv",
                          path_i = "paper_Points.csv")

# Defining the series of the sequential pile
gp.set_series(geo_data, series_distribution={'fault_seriel': 'fault1',
                                           'younger_serie' : 'Unconformity',
                                           'older_serie' : ('Layer1', 'Layer2')},
              order_formation= ['fault1', 'Unconformity', 'Layer2', 'Layer1'])

# Creating object with data prepared for interpolation and compiling.
interp_data = gp.InterpolatorData(geo_data, output='gravity')

# Setting the 2D grid of the airborne where we want to compute the forward gravity
gp.set_geophysics_obj(interp_data_g, ai_extent = [0, 20, 0, 10, -10, 0],
                      ai_resolution = [30,10])

# Making all possible precomputations: Decomposing the value tz for every point of the 2D grid
# to each voxel
gp.precomputations_gravity(interp_data_g, 25, densities=[2.92, 3.1, 2.61, 2.92])

# Computing gravity (Eq. 10)
lith, fault, grav = gp.compute_model(interp_data_g, 'gravity')

# Plotting lithology section
gp.plot_section(geo_data, lith[0], 0, direction='z', plot_data=True)

# Plotting forward gravity
plt.imshow(grav.reshape(10,30), cmap='viridis', origin='lower', alpha=0.8, extent=[0,20,0,10])

```

Listing 4. Computing forward gravity of a *GemPy* model for a given 2-D grid (see Fig. 6).

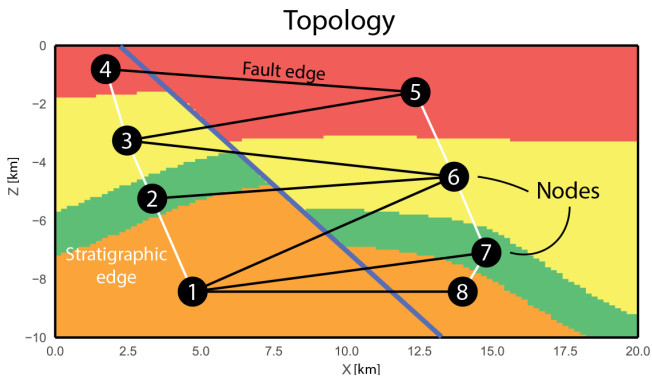


Figure 7. Section of the example geomodel with overlaid topology graph. The geomodel contains eight unique regions (graph nodes) and 13 unique connections (graph edges). White edges represent stratigraphic and unconformity connections, while black edges correspond to across-fault connections.

D matrix containing the information about the faults (fault block id). We multiply the binary fault array (0 for foot wall, 1 for hanging wall) with the maximum lithology value incremented by one. We then add it to the lithology array to make sure that layers that are in contact across faults are as-

signed a unique integer in the resulting array. This yields a 3-D matrix that combines the lithology information and the fault block information. This matrix can then be successfully labeled using CCL with a 2-connectivity stamp, resulting in a new matrix of uniquely labeled regions for the geomodel. From these, an adjacency graph is generated using `skimage.future.graph.RAG`, which created a RAG of all unique regions contained in a 2-D or 3-D matrix, representing each region with a node and their adjacency relations as edges, successfully capturing the topology information of our geomodel. The connections (edges) are then further classified into either stratigraphic or across-fault edges to provide further information. If the argument `compute_areas=True` was given, the contact area for the two regions of an edge is automatically calculated (number of voxels) and stored inside the adjacency graph.

3.4 Stochastic geomodeling and probabilistic programming

Raw geological data are noisy and measurements are usually sparse. As a result, geological models contain significant uncertainties (Wellmann et al., 2010; Bardossy and Fodor, 2004; Lark et al., 2013; Caers, 2011; McLane et al., 2008;

```

...
Add Listing 3
...

# Computing result
lith, fault = gp.compute_model(interp_data)

# Compute topology
G, centroids, labels_unique, labels_lot, lith_lot = gp.topology_compute(geo_data, lith[0],
    fault[0], compute_areas=True)

# Plotting topology network
gp.plot_section(geo_data, lith[0], 5)
gp.topology_plot(geo_data, G, centroids)

```

Listing 5. Topology analysis of a GemPy geomodel.

Chatfield, 1995) that must be addressed thoughtfully to reach a plausible level of confidence in the model. However, treating geological modeling stochastically implies many considerations. (i) From the tens or hundreds of variables involved in the mathematical equations, which ones should be latent? (ii) Can we filter all the possible outcomes that represent unreasonable geological settings? (iii) How can we use other sources of data – especially geophysics – to improve the accuracy of the inference itself?

The answers to these questions are still actively debated in research and are highly dependent on the type of mathematical and computational framework chosen. Uncertainty quantification and its logical extension into probabilistic machine learning will not be covered in the depth in this paper due to the broad scope of the subject. However, the main goal of *GemPy* is to serve as main generative model within these probabilistic approaches and as such we will provide a demonstration of how *GemPy* fits on the workflow of our previous work (de la Varga and Wellmann, 2016; Wellmann et al., 2017) as well as how this work may set the foundations for an easier expansion into the domain of probabilistic machine learning in the future.

As we have seen so far, the cokriging algorithm enables the construction of geological models for a wide range of geometric and topological settings with a limited number of parameters (Fig. 4, red):

- geometric parameters; interface points x_α , i.e., the three Cartesian coordinates x , y and z ; and orientations x_β , i.e., the three Cartesian coordinates x , y and z and the plane orientation normal Gx , Gy and Gz ;
- geophysical parameters, e.g., density;
- model parameters, e.g., covariance at distance zero C_0 (i.e., nugget effect) or the range of the covariance r (see Appendix D for an example of a covariance function).

Therefore, an implicit geological model is simply a graphical representation of a deterministic mathematical operation of these parameters and as such any of these parameters are suitable to behave as latent variables. From a probabilistic

point of view *GemPy* would act as the generative model that links two or more data sets.

GemPy is fully designed to be coupled with probabilistic frameworks, in particular with *pymc3* (Salvatier et al., 2016) as both libraries are based on *Theano*.

pymc is a series of Python libraries that provide intuitive tools to build and subsequently infer complex probabilistic graphical models (PGM; see Koller and Friedman, 2009, and Fig. 8 as an example of a PGM). These libraries offer expressive and clean syntax to write and use statistical distributions and different samplers. At the moment, two main libraries coexist due to their different strengths and weaknesses. On the one hand, we have *pymc2* (Patil et al., 2010) written in FORTRAN and Python. *pymc2* does not allow for gradient-based sampling methods, since it does not have AD capabilities. However, for that same reason, the model construction and debugging are more accessible. Furthermore, not computing gradients enables an easy integration with third party libraries and easy extensibility to other scientific libraries and languages. Therefore, for prototyping and lower dimensionality problems, where the posterior can be tracked by Metropolis–Hasting methods (Haario et al., 2001), *pymc2* is still the go-to choice.

On the other hand, the latest version, *pymc3* (Salvatier et al., 2016), allows for the use of next-generation gradient-based samplers such as No-U-Turn Sampler (Hoffman and Gelman, 2014) or Automatic Variational Inference (Kucukelbir et al., 2015). These sampling methods are proving to be a powerful tool to deal with multidimensional problems, i.e., models with a high number of uncertain parameters (Betancourt et al., 2017). The weakness of these methods are that they rely on the computation of gradients, which in many cases cannot be manually derived. To circumvent this limitation *pymc3* makes use of the AD capabilities of *Theano*. Being built on top of *Theano* confers the Bayesian inference process with all the capabilities discussed in Sect. 2.3.2 in exchange for the clarity and flexibility that pure Python provides.

In this context, the purpose of *GemPy* is to fill the gap of complex algebra between the prior data and observations,

such as geophysical responses (e.g., gravity or seismic inversions) or geological interpretations (e.g., tectonics, model topologies). Since *GemPy* is built on top of *Theano* as well, the compatibility with both libraries is relatively straightforward. However, being able to encode most of the conceivable probabilistic graphical models derived from, often, diverse and heterogeneous data would be a Herculean task. For this reason most of the construction of the PGM has to be coded by the user using the building blocks that the *pymc* packages offer (Bishop, 2013; Patil et al., 2010; Koller and Friedman, 2009, and see Listing 6). By doing so, we can guarantee full flexibility and adaptability to the necessities of every individual geological setting.

For this paper we will use *pymc2* for its higher readability and simplicity. *pymc3* architecture is analogous with the major difference that the PGM is constructed in *Theano* – and therefore symbolically coded (for examples using *pymc3* and *GemPy* check the online documentation detailed in Appendix A2).

3.4.1 Uncertainty quantification

An essential aspect of probabilistic programming is the inherent capability to quantify uncertainty. Monte Carlo error propagation (Ogilvie, 1984) has been introduced in the field of geological modeling a few years ago (Wellmann et al., 2010; Jessell et al., 2010; Lindsay et al., 2012), exploiting the automation of the model construction that implicit algorithms offer.

In this paper, for example Fig. 9 “priors”, we fit a normal distribution of standard deviation 300 m around the Z axis of the interface points in the initial model (Fig. 3c). In other words, we allow the interface points that define the model to oscillate independently along the Z axis randomly – using normal distributions – and we subsequently compute the geomodels that these new data describe. The choice of only perturbing the Z axis is merely due to computational limitations. Uncertainty tends to be higher in this direction (e.g., wells data or seismic velocity); however, there is a lot of room for further research on the definition of prior data – i.e., its choice and probabilistic description – in both directions to ensure that we properly explore the space of feasible models and to generate a parametric space as close as possible to the posterior.

The first step in the creation of a PGM is to define the parameters that are supposed to be stochastic and the probability functions that describe them. To do so, *pymc2* provides a large selection of distributions as well as a clear framework to create custom ones. Once we created the stochastic parameters we need to substitute the initial value in the *GemPy* database (`interp_data` in the snippets) for the corresponding *pymc2* objects. Next, we just need to follow the usual *GemPy* construction process, i.e., calling the `compute_model` function and wrapping it using a deterministic *pymc2* decorator to describe how this function is part

of the probabilistic model (Fig. 8). After creating the graphical model we can sample from the stochastic parameters using Monte Carlo sampling using *pymc2* methods.

The suite of possible realizations of the geological model are stored, as traces, in a database of choice (HDF5, SQL or Python pickles) for further analysis and visualization.

In 2-D we can display all possible locations of the interfaces on a cross section at the center of the model (see Fig. 9, priors 2-D representation); however, the extension of uncertainty visualization to 3D is not as trivial. *GemPy* makes use of the latest developments in uncertainty visualization for 3-D structural geological modeling (e.g., Lindsay et al., 2012, 2013a, b; Wellmann and Regenauer-Lieb, 2012). The first method consists of representing the probability of finding a given geological unit F at each discrete location in the model domain. This can be done by defining a probability function

$$p_F(x) = \sum_{k \in n} \frac{I_{F_k}(x)}{n}, \quad (11)$$

where n is the number of realizations and $I_{F_k}(x)$ is an indicator function of the mentioned geological unit (Fig. 9, probability shows the probability of finding Layer 1). However, this approach can only display each unit individually. A way to encapsulate geomodel uncertainty with a single parameter to quantify and visualize it is by applying the concept of information entropy (Wellmann and Regenauer-Lieb, 2012), based on the general concept developed by (Shannon, 1948). For a discretized geomodel the information entropy H (normalized by the total number of voxels n) can be defined as

$$H = - \sum_{i=1}^n p_i \log_2 p_i, \quad (12)$$

where p_i represents the probability of a layer at cell x . Therefore, we can use information entropy to reduce the dimensionality of probability fields into a single value at each voxel as an indication of uncertainty, reflecting the possible number of outcomes and their relative probability (see Fig. 9, “Entropy”).

3.4.2 Geological inversion: gravity and topology

Although computing the forward gravity has its own value for many applications, the main aim of *GemPy* is to integrate all possible sources of information into a single probabilistic framework. The use of likelihood functions in a Bayesian inference in comparison to simple forward simulation has been explored by the authors during recent years (de la Varga and Wellmann, 2016; Wellmann et al., 2017; Schaaf, 2017). This approach enables us to tune the conditioning of possible stochastic realizations by varying the probabilistic density function used as likelihoods. In addition, Bayesian networks allow us to combine several likelihood functions, generating a competition among the prior distribution of the input data and likelihood functions resulting in posterior distributions

```

...
Add Listing 3
...

# Copying the initial data
geo_data_stoch_init = deepcopy(interp_data.geo_data_res)
# MODEL CONSTRUCTION
# =====
# Positions (rows) of the data we want to make stochastic
ids = range(2,12)

# List with the stochastic parameters. pymc.Normal attributes: Name, mean, std
interface_Z_modifier = [pymc.Normal("interface_Z_mod_"+str(i), 0., 1./0.01**2) for i in ids]

# Modifying the input data at each iteration
@pymc.deterministic(trace=True)
def input_data(interface_Z_modifier = interface_Z_modifier,
               geo_data_stoch_init = geo_data_stoch_init,
               ids = ids, verbose=0):

# First we extract from our original interp_data object the numerical data that
# is necessary for the interpolation. geo_data_stoch is a pandas Dataframe
    geo_data_stoch = gp.get_data(geo_data_stoch_init, numeric=True)

# Now we loop each id which share the same uncertainty variable. In this case, each layer. We
# add the stochastic part to the initial value
    for num, i in enumerate(ids):
        interp_data.geo_data_res.interfaces.set_value(i, "Z",
            geo_data_stoch_init.interfaces.iloc[i]["Z"] + interface_Z_modifier[num])

# Return the input data to be input into the modeling function. Due to the way pymc2
# stores the traces we need to save the data as numpy arrays
    return interp_data.geo_data_res.interfaces[["X", "Y", "Z"]].values,
        interp_data.geo_data_res.orientations[["X", "Y", "Z", "dip", "azimuth",
            "polarity"]].values

# Computing the geological model
@pymc.deterministic(trace=True)
def gempy_model(value=0, input_data=input_data, verbose=False):

# modify input data values accordingly
    interp_data.geo_data_res.interfaces[["X", "Y", "Z"]] = input_data[0]

# Gx, Gy, Gz are just used for visualization. The Theano function gets azimuth dip and
# polarity
    interp_data.geo_data_res.orientations[["G_x", "G_y", "G_z", "X", "Y", "Z", 'dip',
        'azimuth', 'polarity']] = input_data[1]

# Some iterations will give a singular matrix, that's why we need to
# create a try to not break the code.
    try:
        lb, fb, grav = gp.compute_model(interp_data, outup='gravity')
        return lb, fb, grav

    except np.linalg.linalg.LinAlgError as err:
# If it fails (e.g. some input data combinations could lead to
# a singular matrix and thus break the chain) return an empty model
# with same dimensions (just zeros)
        if verbose:
            print("Exception_occured.")
        return np.zeros_like(lith_block), np.zeros_like(fault_block), np.zeros_like(
            grav_i)

# Extract the vertices in every iteration by applying the marching cube algorithm
@pymc.deterministic(trace=True)
def gempy_surfaces(value=0, gempy_model=gempy_model):
    vert, simp = gp.get_surfaces(interp_data, gempy_model[0][1], gempy_model[1][1],
                                original_scale=True)

    return vert

# We add all the pymc objects to a list
params = [input_data, gempy_model, gempy_surfaces, *interface_Z_modifier]

# We create the pymc model i.e. the probabilistic graph
model = pymc.Model(params)
runner = pymc.MCMC(model)

```

Listing 6. Probabilistic model construction and inference using *pymc2* and *GemPy*: Monte Carlo forward simulation (see Fig. 9, “priors” for the results).

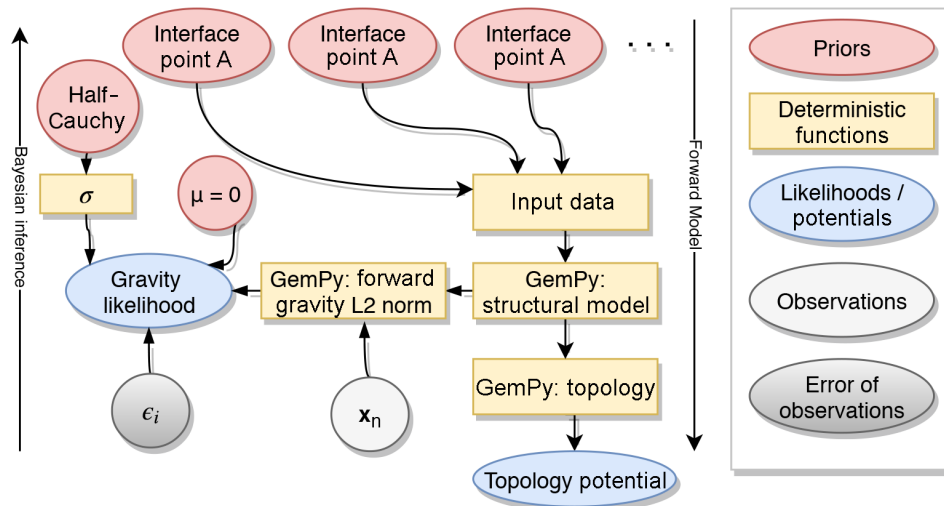


Figure 8. Probabilistic graphical model generated with *pymc2*. Ellipses represent stochastic parameters, while rectangles are deterministic functions that return intermediated states of the probabilistic model such as the GemPy model.

that best honor all the given information. To give a flavor of what is possible, we apply custom likelihoods to the previous example based on topology and gravity constrains in an inversion.

As we have shown above, topological graphs can represent the connectivity among the segmented areas of a geological model. As is expected, stochastic perturbations of the input data can rapidly alter the configuration of mentioned graphs. In order to preserve a given topological configuration partially or totally, we can construct specific likelihood functions. To exemplify the use of a topological likelihood function, we will use the topology computed in Sect. 3.3 derived from the initial model realization (Figs. 7 or 9, “Likelihoods”) as “ideal topology”. This can be based on an expert interpretation of kinematic data or deduced from auxiliary data.

The first challenge is to find a metric that captures the similarity of two graphs. As a graph is nothing but a set of nodes and their edges we can compare the intersection and union of two different sets using the Jaccard index (Jaccard, 1912; Thiele et al., 2016a). It calculates the ratio of intersection and union of two given graphs A and B :

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}. \quad (13)$$

The resulting ratio is zero for entirely different graphs, while the metric rises as the sets of edges and nodes become more similar between two graphs and reaches exactly one for an identical match. Therefore, the Jaccard index can be used to express the similarity of topology graphs as a single number we can evaluate using a probability density function. To evaluate the likelihood of the simulated model topology we use a factor potential with a half-Cauchy parametrization (rate parameter $\beta = 10^{-3}$) to constrain our model using the “soft

data” of our topological knowledge (Lauritzen et al., 1990; Jordan, 1998; Christakos, 2002). This specific parametrization was chosen due to empirical evidence from different model runs to allow for effective parameter space exploration in the used Markov chain Monte Carlo (MCMC) scheme.

Gravity likelihoods aim to exploit the spatial distribution of density, which can be related to different lithotypes (Dentith and Mudge, 2014). To test the likelihood function based on gravity data, we first generate the synthetic “measured” data. This was simply done by computing the forward gravity for one of the extreme models (to highlight the effect that a gravity likelihood can have) generated during the Monte Carlo error propagation in the previous section. This model is particularly characterized by its high dip values (Fig. 9, synthetic model to produce forward gravity). The construction of the likelihood function is done by applying an L2 norm between each “measured” data point and the forward computation and evaluating the result by a normal distribution of mean $\mu = 0$ and with the standard deviation, σ , as a half-Cauchy prior (rate parameter $\beta = 10^{-1}$). This likelihood function will push the model parameters (Fig. 4, red) in the direction to reduce the L2 norm as much as possible while keeping the standard deviation around the prior value (this prior value encapsulates the inherent measurement and model uncertainty).

Defining the topology potential and gravity likelihood on the same Bayesian network creates a joint likelihood value that will define the posterior space. To sample from the posterior we use adaptive Metropolis (Haario et al., 2001, for a more in depth explanation of samplers and their importance see de la Varga and Wellmann, 2016). This method varies the metropolis sampling size according to the covariance function that gets updated every n iterations. For the results exposed here, we performed 20 000 iterations, tuning the adap-

```

...
Add Listing 6
...

# Computation of topology
@pymc.deterministic(trace=True)
def gempy_topo(value=0, gm=gempy_model, verbose=False):
    G, c, lu, lot1, lot2 = gp.topology_compute(geo_data, gm[0][0], gm[1], cell_number=0,
        direction="y")

    if verbose:
        gp.plot_section(geo_data, gm[0][0], 0)
        gp.topology_plot(geo_data, G, c)

    return G, c, lu, lot1, lot2

# Computation of L2-Norm for the forward gravity
@pymc.deterministic
def e_sq(value = original_grav, model_grav = gempy_model[2], verbose = 0):
    square_error = np.sqrt(np.sum((value*10**-7 - (model_grav*10**-7))**2))
    return square_error

# Likelihoods
# =====
@pymc.stochastic
def like_topo_jaccard_cauchy(value=0, gempy_topo=gempy_topo, G=topo_G):
    """Compares the model output topology with a given topology graph G using an inverse Jaccard-
        index embedded in a half-cauchy likelihood."""
    # jaccard-index comparison
    j = gp.Topology.compare_graphs(G, gempy_topo[0])
    # the last parameter adjusts the "strength" of the likelihood
    return pymc.half_cauchy_like(1 - j, 0, 0.001)

@pymc.observed
def inversion(value = 1, e_sq = e_sq):
    return pymc.half_cauchy_like(e_sq, 0, 0.1)

# We add all the pymc objects to a list
params = [input_data, gempy_model, gempy_surfaces, gempy_topo, *interface_Z_modifier,
    like_topo_jaccard_cauchy, e_sq, inversion]

# We create the pymc model i.e. the probabilistic graph
model = pymc.Model(params)
runner = pymc.MCMC(model)

# BAYESIAN INFERENCE
# =====
# Number of iterations
iterations = 15000

# Inference. Adaptive Metropolis
runner.use_step_method(pymc.AdaptiveMetropolis, params, delay=1000)
runner.sample(iter = 20000, burn=1000, thin=20, tune_interval=1000, tune_throughout=True)

```

Listing 7. Probabilistic model construction and inference using *pymc2* and *GemPy*: Bayesian inference (see Fig. 9 for the results).

tive covariance every 1000 steps (a convergence analysis can be found in the Jupyter notebooks in the package repository).

As a result of applying likelihood functions we can appreciate a clear change in the posterior (i.e., the possible outcomes) of the inference. A closer look shows two main zones of influence, each of them related to one of the likelihood functions. On one hand, we observe a reduction of uncertainty along the fault plane due to the restrictions that the topology function imposes by conditioning the models to high Jaccard values. On the other hand, what in the first example – i.e., Monte Carlo error propagation, left in Fig. 9

– was just an outlier, due to the influence of the gravity inversion, now becomes the norm bending the layers pronouncedly. The purpose of this example is to highlight the functionality. For a realistic study, further detailed adjustments would have to be taken.

4 Discussion

We have introduced *GemPy*, a Python library for implicit geomodeling with special emphasis on the analysis of uncertainty. With the advent of powerful implicit methods to

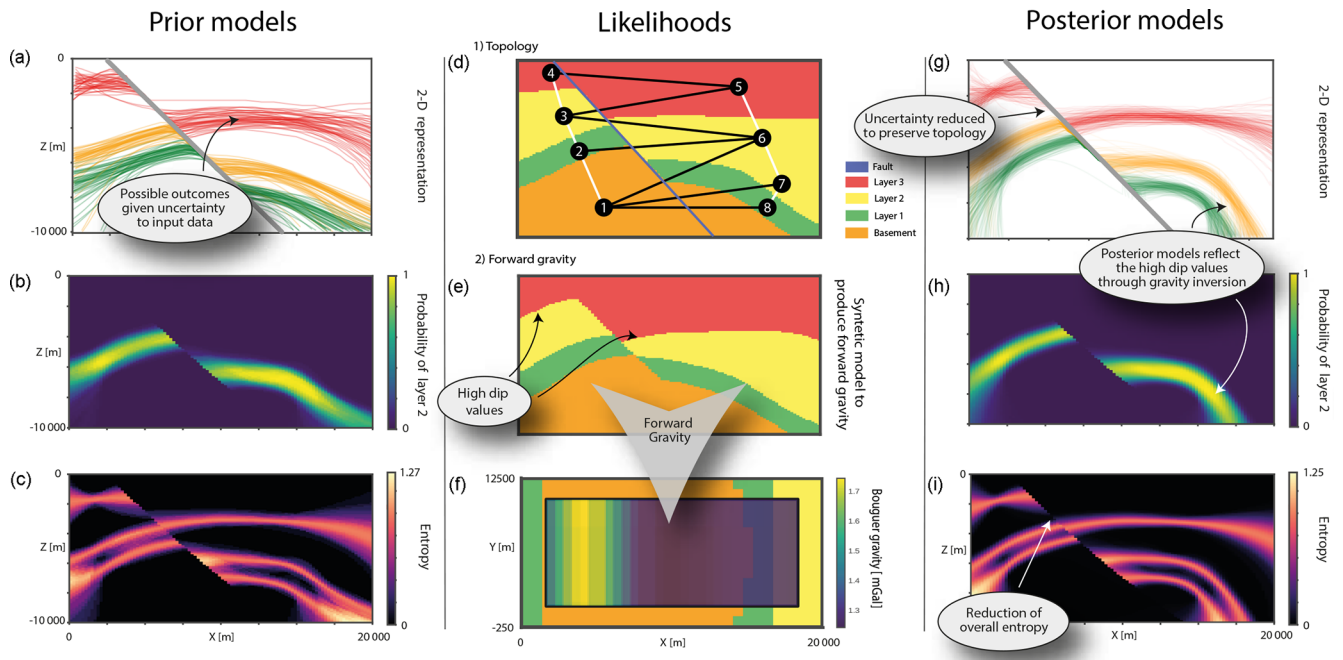


Figure 9. Probabilistic programming results on a cross section at the middle of the model ($Y = 10000$ m). (i) Priors UQ shows the uncertainty in geological models given stochastic values to the Z position of the input data (standard deviation, $\sigma = 300$): (a) 2-D interface representation; (b) probability of occurrence for Layer 2; (c) information entropy. (ii) Representation of data used as likelihood functions: (d) ideal topology graph; (e) synthetic model taken as reference for the gravity inversion; (f) reference forward gravity overlain on top of an XY cross section of the synthetic reference model. Posterior analysis after combining priors and likelihood in a Bayesian inference: (g) 2-D interface representation; (h) probability of occurrence for Layer 2; (i) information entropy.

automate many of the geological modeling steps, *GemPy* builds on these mathematical foundations to offer a reliable and easy-to-use technology to generate complex models with only a few lines of code. In many cases – and in research in particular – it is essential to have transparent software that allows for full manipulation and understanding of the logic beneath its front end to honor the scientific method and allow for reproducibility by using open-access software.

Up until now, implicit geological modeling was limited to proprietary software suites – for the petroleum industry (GoCad, Petrel, JewelSuite) or the mining sector (MicroMine, MIRA Geoscience, GeoModeller, Leapfrog) – with an important focus on industry needs and user experience (e.g., graphical user interfaces or data compatibilities). Despite access to the APIs of many of these software, their lack of transparency and the inability to fully manipulate any of the algorithms represents a serious obstacle for conducting appropriate *reproducible* research. To overcome these limitations, many scientific communities – e.g., *simpeg* in geophysics (Cockett et al., 2015), *astropy* in astronomy (Robitaille et al., 2013) or *pynoddy* in kinematic structural modeling (Wellmann et al., 2016) – are moving towards the open-source framework necessary for the full application of the scientific method. In this regard, the advent of open-source programming languages such as R or Python are playing a crucial role

in facilitating scientific programming and enabling the crucial reproducibility of simulations and script-based science. *GemPy* aims to fill the existing gap of implicit modeling in the open-source ecosystem in geosciences.

Implicit methods rely on interpolation functions to automate some or all the construction steps. Different mathematical approaches have been developed and improved in the recent years to tackle many of the challenges that particular geological settings pose (e.g., Lajaunie et al., 1997; Hillier et al., 2014; Calcagno et al., 2008; Caumon et al., 2013; Caumon, 2010). A significant advantage of some of these methods is that they directly enable the recomputation of the entire structure when input data are changed or added. Furthermore, they can provide a geologically meaningful interpolation function, e.g., considering deposition time (Caumon, 2010) or potential fields (Lajaunie et al., 1997) to encapsulate the essence of geological deposition in different environments. The creation of *GemPy* has been made possible at a moment when the automation of geological modeling via implicit algorithms, as well as the maturity of the Python open-source ecosystem, reached a point where a few thousand new lines of code are able to efficiently perform the millions of linear algebra operations and complex memory management necessary to create complex geological models. An important aspect in *GemPy*'s design has been the willingness to

allow users to simply use *GemPy* as a tool to construct geological models for different purposes as well as to encourage users to develop and expand the code base itself. With the purpose to facilitate a low-entry barrier we have taken two main structural decisions: (i) a clear separation between core features and extensible assets and (ii) a combination of functional and object-oriented programming. The aim of this dual design is to give a user-friendly, easy-to-use front end to the majority of users while keeping a modular structure of the code for future contributors.

Using *GemPy* requires a minimum familiarity with the Python syntax. The lack of an advanced graphical interface to place the input data interactively forces the user to provide data sets with the coordinates and angular data. For this reason, for complex initial models, *GemPy* could be seen more as a back-end library required to couple it with software providing 3-D graphical manipulation. Due to the development team's background, *GemPy* is fully integrated with GeoModeller through the built-in library *pygeomod*. *GemPy* is able to read and modify GeoModeller projects directly, allowing the user to take advantage of their respective unique features. All input data of *GemPy* itself is kept in open, standard Python formats, making use of the flexible *pandas* "DataFrames" and powerful *numpy* arrays. Hence, every user will be able to freely manipulate the data at any given point.

GemPy has built-in functionality to visualize results using the main visualization libraries offered in Python: *matplotlib* for 2-D and *vtk* for 3-D and allows the user to export *.vtk* files for later visualization in common open-source tools for scientific visualization such as ParaView. Although *GemPy* does not include an evolved user interface, we offer a certain level of interactivity using *GemPy*'s built-in 3-D visualization in VTK and interactive data frames through *qgrid*. Not only is the user able to move the input data via drag-and-drop or changing the data frame, but *GemPy* can immediately re-interpolate the perturbed model, enabling an extremely intuitive direct feedback on how the changes made affect the model. Visualization of vast model ensembles is also possible in 3-D using slider functionality. Future plans for the visualization of *GemPy* include virtual reality support to make data manipulation and model visualization more immersive and intuitive to use.

Another important feature of *GemPy* is the use of symbolic code. The lack of a domain-specific language allows for the compilation of the code to a highly efficient language. Furthermore, since all the logic has to be described prior to the compilation, memory allocation and parallelization of the code can be optimized. *Theano* uses BLAS (Lawson et al., 1979) to perform the algebraic operations with out-of-the-box OpenMP (Dagum and Menon, 1998) capabilities for multi-core operations. Additionally, parallel GPU computation is available and compatible with the use of CPUs, which allows the user to define certain operations to a specific device and even to split big arrays (e.g., grid) to multiple GPUs. In other words, the symbolic nature of the code

enables the separation of the logic according to the individual advantages of each device – i.e., sequential computations to CPUs and parallel calculations to the GPUs – allowing for better use of the available hardware. Hence, this scheme is portable to high-performance computing in the same fashion.

Up to now, structural geological models have significantly relied on the best deterministic and explicit realization that an expert is able to construct using often noisy and sparse data. Research into the interpretation uncertainty of geological data sets (e.g., seismic data) has recognized the significant impact of interpreter education and bias on the extracted input data for geological models (e.g., Bond et al., 2007; Bond, 2015). *GemPy*'s ability to be enveloped into probabilistic programming frameworks such as *pymc* allows for the consideration of input data uncertainties and could provide a free, open-source foundation for developing probabilistic geomodeling workflows that integrate uncertainties from the very beginning of data interpretation, through the geomodel interpolation and up to the geomodel application (e.g., flow simulations, economic estimations).

In the transition to a world dominated by data and optimization algorithms – e.g., deep neural networks or big data analytics – there are many attempts to apply those advances in geological modeling (Wang et al., 2017; Gonçalves et al., 2017). The biggest attempt to use data-driven models in geology comes from geophysical inversions (Tarantola and Valette, 1982; Mosegaard and Tarantola, 1995; Sambridge and Mosegaard, 2002; Tarantola, 2005). Their approaches consist of using the mismatch of one or many parameters, comparing model with reality and modifying them accordingly until reaching a given tolerance. However, since this solution is never unique, it is necessary to enclose the space of possibilities by some other means. This prior approach to the final solution is usually made using polygonal or elliptic bodies leading to oversimplified geometry of the distinct lithological units. Other researchers use additional data – geophysical data or other constrains (Jessell et al., 2010; Wellmann et al., 2014) – to validate multiple possible realizations of geological models generated either automatically by an interpolation function or manually. Here, the additional information is used as a hard deterministic filter for what is reasonable or not. The limitation of pure rejection filtering is that information does not propagate backward to modify the latent parameters that characterize the geological models, which makes computation infeasible to explore high-dimensional problems. In between these two approaches, we can find some attempts to reconcile both approaches meeting somewhere in the middle. An example of this is the approach followed in the software packages GeoModeller and SKUA. They optimize the layer densities, and when necessary the discretized model, to fit the geological model to the observed geophysical response. The consequence of only altering the discrete final model is that after optimization the original input data used for the construction of the geological model

(i.e., interface points and orientation data) gets overwritten and consequently hard to reproduce.

We propose a more general approach. By embedding the geological model construction into a probabilistic machine-learning framework (Bishop, 2013) – i.e., a Bayesian inference network. In short a Bayesian inference is a mathematical formulation to update beliefs in the light of new evidence. This statement applied to geological modeling is translated into keeping all or a subset of the parameters that generate the model uncertain and evaluate the quality of the model comparing its mismatch with additional data or geological knowledge encoded mathematically (de la Varga and Wellmann, 2016). In this way, we are able to utilize available information not only in a forward direction to construct models but also propagate information backwards in an inverse scheme to refine the probabilistic distributions that characterize the modeling parameters. Compared with previous approaches, we do not only use the inversion to improve a deterministic model but instead to learn about the parameters that define the model to begin with. In recent years, we have shown how this approach may help close the gap between geophysical inversions and geological modeling in an intuitive manner (Wellmann et al., 2017). At the end of the day, Bayesian inferences operate in a very similar way to how humans do, we create our best guess model; we compare it to the geophysical data or our geological knowledge, and in the case of disagreement we modify the input of the geological model in the direction we think is the best to honor the additional data.

Despite the convincing mathematical formulation of Bayesian inferences, there are caveats to be dealt with for practical applications. As mentioned in Sect. 3.4, the effective computational cost to perform such algorithms have prohibited its use beyond research and simplified models. However, recent developments in MCMC methods enable more efficient ways to explore the parametric space and hence open the door to a significant increase in the complexity of geological models. An in-depth study of the impact of gradient-based MCMC methods in geological modeling will be carried out in a future publication.

Nevertheless, performing AD does not come free of cost. The required code structures limit the use of libraries that do not perform AD themselves, which in essence imposes a requirement to rewrite most of the mathematical algorithms involved in the Bayesian network. Under these circumstances, we have rewritten in the potential field method *Theano* – with many of the add-ons developed in recent years (Calcagno et al., 2008) – and the computation of forward gravity responses for discrete rectangular prisms.

Currently *GemPy* is in active development moving towards three core topics: (i) increasing the probabilistic machine-learning capabilities by exploiting gradient-based methods and new types of likelihoods, (ii) post-processing of uncertainty quantification and its relation to decision theory and information theory, and (iii) exploring the new catalog of virtual reality and augmented reality solutions to improve

the visualization of both the final geological models and the building environment. Ideally *GemPy* will function as a platform to create a vibrant open-source community to push forward geological modeling into the new machine-learning era. Therefore, we hope to include functionality developed by other external users into the main package.

In conclusion, *GemPy* has evolved to a full approach for geological modeling in a probabilistic programming framework. The lack of available open-source tools in geological modeling and the necessity of writing all the logic symbolically has pushed the project to an unexpected stand-alone size. However, this would not have been possible without the immense, ever-growing, open-source community that provide numerous high-quality libraries that enable the creation of powerful software with relatively few new lines of code. And in the same fashion, we hope the community will make use of our library to perform geological modeling transparently and reproducibly, and incorporate the uncertainties inherent to earth sciences.

Code availability. GemPy is a free, open-source Python library licensed under the GNU Lesser General Public License v3.0 (GPLv3). It is hosted on the GitHub repository <https://github.com/cgre-aachen/gempy> (<https://doi.org/10.5281/zenodo.1186118>; de la Varga et al., 2018).

Appendix A: GemPy package information

A1 Installation

Installing GemPy can be done in two ways: (i) either by cloning the GitHub repository with `$ git clone https://github.com/cgre-aachen/gempy.git` (last access: 17 December 2018) and then manually installing it by running the Python setup script in the repository: `$ python install.py` (ii) or by using the Python Package Index (PyPI) with the command `$ pip install gempy`, which directly downloads and installs the library.

A2 Documentation

GemPy's documentation is hosted on <http://gempy.readthedocs.io/> (last access: 17 December 2018), which provides a general overview of the library and multiple in-depth tutorials. The tutorials are provided as Jupyter notebooks, which provide the convenient combination of documentation and executable script blocks in one document. The notebooks are part of the repository and located in the tutorials folder. See <http://jupyter.org/> (last access: 17 December 2018) for more information on installing and running Jupyter notebooks.

A3 Jupyter notebooks

We provide Jupyter notebooks as part of the online documentation. These notebooks can be executed in a local Python environment (if the required dependencies are correctly installed, see above). In addition, static versions of the notebooks can currently be inspected directly on the GitHub repository web page or through the use of nbviewer. In addition, it is possible to run interactive notebooks through the use of binder (provided through <https://mybinder.org> at the time of writing). For more details and up-to-date information, please refer to the repository page <https://github.com/cgre-aachen/gempy>.

A4 Unit tests

The *GemPy* package contains a set of tests, which can be executed in the standard Python testing environment. If you cloned or downloaded the repository, then these tests can be directly performed by going to the package folder and running the `pytest` command: `$ pytest`

If all tests are successful, you are ready to continue.

Appendix B: Kriging system expansion

The following equations have been derived from the work in Aug (2004); Lajaunie et al. (1997) and Chiles and Delfiner (2009).

B1 Gradient covariance matrix $C_{\partial\mathbf{Z}/\partial u}$

The gradient covariance matrix, $C_{\partial\mathbf{Z}/\partial u}$, is made up of as many variables as gradient directions that are taken into consideration. In 3-D, we would have the Cartesian coordinate dimensions – $\mathbf{Z}/\partial x$, $\mathbf{Z}/\partial y$ and $\mathbf{Z}/\partial z$ – and therefore they will derive from the partial differentiation of the covariance function $\sigma(x_i, x_j)$ of \mathbf{Z} .

As in our case the directional derivatives used are the three Cartesian directions; we can rewrite gradients covariance, $C_{\partial\mathbf{Z}/\partial u, \partial\mathbf{Z}/\partial v}$, for our specific case as

$$C_{\partial\mathbf{Z}/\partial x, \partial\mathbf{Z}/\partial y, \partial\mathbf{Z}/\partial z} = \begin{bmatrix} C_{\partial\mathbf{Z}/\partial x, \partial\mathbf{Z}/\partial x} & C_{\partial\mathbf{Z}/\partial x, \partial\mathbf{Z}/\partial y} & C_{\partial\mathbf{Z}/\partial x, \partial\mathbf{Z}/\partial z} \\ C_{\partial\mathbf{Z}/\partial y, \partial\mathbf{Z}/\partial x} & C_{\partial\mathbf{Z}/\partial y, \partial\mathbf{Z}/\partial y} & C_{\partial\mathbf{Z}/\partial y, \partial\mathbf{Z}/\partial z} \\ C_{\partial\mathbf{Z}/\partial z, \partial\mathbf{Z}/\partial x} & C_{\partial\mathbf{Z}/\partial z, \partial\mathbf{Z}/\partial y} & C_{\partial\mathbf{Z}/\partial z, \partial\mathbf{Z}/\partial z} \end{bmatrix}. \quad (\text{B1})$$

Notice, however, that covariance functions by definition are described in a polar coordinate system, and therefore it will be necessary to apply the chain rule for *directional derivatives*. Considering an isotropic and stationary covariance we can express the covariance function as

$$\sigma(x_i, x_j) = C(r) \quad (\text{B2})$$

with

$$r = \sqrt{h_x^2 + h_y^2 + h_z^2} \quad (\text{B3})$$

and h_u as the distance $u_i - u_j$ in the given direction (usually Cartesian directions). Therefore, since we aim to derive $C_Z(r)$ with respect to an arbitrary direction u , we must apply

the *directional derivative* rules as follows:

$$C_{\partial\mathbf{Z}/\partial u, \partial\mathbf{Z}/\partial u} = \frac{\partial^2 C_Z(r)}{\partial h_u^2} = \frac{\partial C_Z(r)}{\partial r} \frac{\partial}{\partial h_u} \left(\frac{\partial r}{\partial h_u} \right) + \frac{\partial}{\partial h_u} \left(\frac{\partial C_Z(r)}{\partial r} \right) \frac{\partial r}{\partial h_u}, \quad (\text{B4})$$

where

$$\frac{\partial C_Z(r)}{\partial r} = C'_Z(r), \quad (\text{B5})$$

$$\frac{\partial r}{\partial h_u} = \frac{h_u}{\sqrt{h_u^2 + h_v^2}} = -\frac{h_u}{r}, \quad (\text{B6})$$

$$\begin{aligned} \frac{\partial}{\partial h_u} \left(\frac{\partial r}{\partial h_u} \right) &= \frac{\partial}{\partial h_u} \left(\frac{h_u}{\sqrt{h_u^2 + h_v^2}} \right) \\ &= -\frac{2h_u^2}{2\sqrt{h_u^2 + h_v^2}} + \frac{1}{\sqrt{h_u^2 + h_v^2}} = -\frac{h_u^2}{r^3} + \frac{1}{r}, \end{aligned} \quad (\text{B7})$$

$$\begin{aligned} \frac{\partial}{\partial h_u} \left(\frac{\partial C_Z(r)}{\partial r} \right) &= \frac{\partial C'_Z(r)}{\partial h_u} = \frac{\partial C''_Z(r)}{\partial r} \frac{\partial r}{\partial h_u} \\ &= -\frac{h_u}{r} C''_Z. \end{aligned} \quad (\text{B8})$$

Substituting:

$$\begin{aligned} C_{\partial\mathbf{Z}/\partial u, \partial\mathbf{Z}/\partial u} &= C'_Z(r) \left(-\frac{h_u^2}{r^3} + \frac{1}{r} \right) - \frac{h_u}{r} C''_Z \frac{h_u}{r} \\ &= C'_Z(r) \left(-\frac{h_u^2}{r^3} + \frac{1}{r} \right) + \frac{h_u^2}{r^2} C''_Z. \end{aligned} \quad (\text{B9})$$

While in the case of two different directions the covariance will be

$$\begin{aligned} C_{\partial\mathbf{Z}/\partial u, \partial\mathbf{Z}/\partial v} &= \frac{\partial^2 C_Z(r)}{\partial h_u \partial h_v} = \frac{\partial C_Z(r)}{\partial r} \frac{\partial}{\partial h_v} \left(\frac{\partial r}{\partial h_u} \right) \\ &+ \frac{\partial}{\partial h_v} \left(\frac{\partial C_Z(r)}{\partial r} \right) \frac{\partial r}{\partial h_u} \end{aligned} \quad (\text{B10})$$

with

$$\frac{\partial}{\partial h_v} \left(\frac{\partial r}{\partial h_u} \right) = \frac{\partial}{\partial h_v} \left(\frac{h_u}{\sqrt{h_u^2 + h_v^2}} \right) = -\frac{h_u h_v}{r^3}, \quad (\text{B11})$$

$$\frac{\partial}{\partial h_v} \left(\frac{\partial C_Z(r)}{\partial r} \right) = \frac{\partial C'_Z(r)}{\partial h_v} = -C''_Z(r) \frac{h_v}{r}, \quad (\text{B12})$$

we have

$$\begin{aligned} C_{\partial\mathbf{Z}/\partial u, \partial\mathbf{Z}/\partial v} &= C'_Z(r) \left(-\frac{h_u h_v}{r^3} \right) + C''_Z(r) \frac{h_u h_v}{r^2} \\ &= \frac{h_u h_v}{r^2} \left(C''_Z(r) - \frac{C'_Z(r)}{r} \right). \end{aligned} \quad (\text{B13})$$

This derivation is independent to the covariance function of choice. However, some covariances may lead to mathematical indeterminations if they are not sufficiently differentiable.

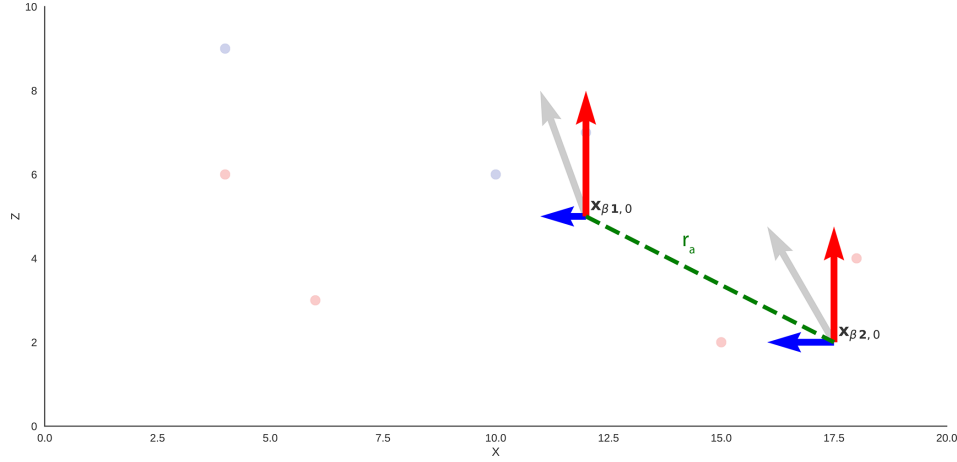


Figure B1. 2-D representation of the decomposition of the orientation vectors into Cartesian axes. Each Cartesian axis represents a variable of a sub-cokriging system. The dashed green line represents the covariance distance, r_a , for the covariance of the gradient.

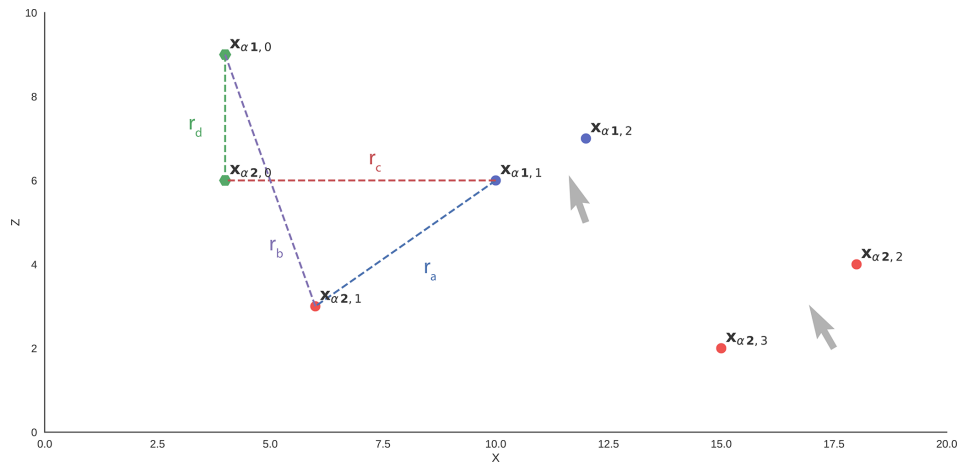


Figure B2. Distances r involved in the computation of the interface subsystem of the interpolation. Because all covariances are relative to a reference point $x_{\alpha,0}^i$, all four covariances with their respective distances, r_a , r_b , r_c and r_d , must be taken into account (Eq. B14).

B2 Interface covariance matrix

In a practical sense, keeping the value of the scalar field at every interface unfixed forces us to consider the covariance between the points within an interface as well as the covariance between different layers following equation

$$C_{x_{\alpha}^r, x_{\alpha}^s} = \overbrace{C_{x_{\alpha}^r, x_{\alpha}^s}}^{r_a} - \overbrace{C_{x_{\alpha}^r, x_{\alpha}^s}}^{r_b} - \overbrace{C_{x_{\alpha}^r, x_{\alpha}^s}}^{r_c} + \overbrace{C_{x_{\alpha}^r, x_{\alpha}^s}}^{r_d}. \quad (\text{B14})$$

This leads to the subdivision of the cokriging system respecting the interfaces:

$$C_{Z,Z} = \begin{bmatrix} C_{x_{\alpha}^1, x_{\alpha}^1} & C_{x_{\alpha}^1, x_{\alpha}^2} & \cdots & C_{x_{\alpha}^1, x_{\alpha}^s} \\ C_{x_{\alpha}^2, x_{\alpha}^1} & C_{x_{\alpha}^2, x_{\alpha}^2} & \cdots & C_{x_{\alpha}^2, x_{\alpha}^s} \\ \vdots & \vdots & \ddots & \vdots \\ C_{x_{\alpha}^r, x_{\alpha}^1} & C_{x_{\alpha}^r, x_{\alpha}^2} & \cdots & C_{x_{\alpha}^r, x_{\alpha}^s} \end{bmatrix}. \quad (\text{B15})$$

Combining Eqs. (5) and (B15) the covariance for the property *scalar field* will look like

$$C_{x_{\alpha}^r, x_{\alpha}^s} = \begin{bmatrix} C_{x_{\alpha}^1, x_{\alpha}^1} - C_{x_{\alpha}^1, x_{\alpha}^0} - C_{x_{\alpha}^1, x_{\alpha}^0} + C_{x_{\alpha}^0, x_{\alpha}^0} & C_{x_{\alpha}^1, x_{\alpha}^2} - C_{x_{\alpha}^1, x_{\alpha}^0} - C_{x_{\alpha}^2, x_{\alpha}^0} + C_{x_{\alpha}^0, x_{\alpha}^0} & \cdots & C_{x_{\alpha}^1, x_{\alpha}^s} - C_{x_{\alpha}^1, x_{\alpha}^0} - C_{x_{\alpha}^s, x_{\alpha}^0} + C_{x_{\alpha}^0, x_{\alpha}^0} \\ C_{x_{\alpha}^2, x_{\alpha}^1} - C_{x_{\alpha}^2, x_{\alpha}^0} - C_{x_{\alpha}^1, x_{\alpha}^0} + C_{x_{\alpha}^0, x_{\alpha}^0} & C_{x_{\alpha}^2, x_{\alpha}^2} - C_{x_{\alpha}^2, x_{\alpha}^0} - C_{x_{\alpha}^2, x_{\alpha}^0} + C_{x_{\alpha}^0, x_{\alpha}^0} & \cdots & C_{x_{\alpha}^2, x_{\alpha}^s} - C_{x_{\alpha}^2, x_{\alpha}^0} - C_{x_{\alpha}^s, x_{\alpha}^0} + C_{x_{\alpha}^0, x_{\alpha}^0} \\ \vdots & \vdots & \ddots & \vdots \\ C_{x_{\alpha}^r, x_{\alpha}^1} - C_{x_{\alpha}^r, x_{\alpha}^0} - C_{x_{\alpha}^1, x_{\alpha}^0} + C_{x_{\alpha}^0, x_{\alpha}^0} & C_{x_{\alpha}^r, x_{\alpha}^2} - C_{x_{\alpha}^r, x_{\alpha}^0} - C_{x_{\alpha}^2, x_{\alpha}^0} + C_{x_{\alpha}^0, x_{\alpha}^0} & \cdots & C_{x_{\alpha}^r, x_{\alpha}^s} - C_{x_{\alpha}^r, x_{\alpha}^0} - C_{x_{\alpha}^s, x_{\alpha}^0} + C_{x_{\alpha}^0, x_{\alpha}^0} \end{bmatrix}. \quad (\text{B16})$$

B3 Cross-covariance

In a cokriging system, the relation between the interpolated parameters is given by a cross-covariance function. As we saw above, the gradient covariance is subdivided into covariances with respect to the three Cartesian directions (Eq. B1), while the interface covariance is detached from the covariances matrices with respect to each individual interface (Eq. B15). In the same manner, the cross-covariance will reflect the relation of every interface to each gradient direction,

$$\mathbf{C}_{\mathbf{Z}, \partial \mathbf{Z} / \partial u} = \begin{bmatrix} \mathbf{C}_{x_{\alpha 1}^1, \partial \mathbf{Z}(x_{\beta 1}) / \partial x} & \mathbf{C}_{x_{\alpha 2}^1, \partial \mathbf{Z}(x_{\beta 1}) / \partial x} & \dots & \mathbf{C}_{x_{\alpha i}^r, \partial \mathbf{Z}(x_{\beta 1}) / \partial x} \\ \mathbf{C}_{x_{\alpha 1}^1, \partial \mathbf{Z}(x_{\beta 2}) / \partial x} & \mathbf{C}_{x_{\alpha 2}^1, \partial \mathbf{Z}(x_{\beta 2}) / \partial x} & \dots & \mathbf{C}_{x_{\alpha i}^r, \partial \mathbf{Z}(x_{\beta 2}) / \partial x} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{C}_{x_{\alpha 1}^1, \partial \mathbf{Z}(x_{\beta 1}) / \partial y} & \mathbf{C}_{x_{\alpha 2}^1, \partial \mathbf{Z}(x_{\beta 1}) / \partial y} & \dots & \mathbf{C}_{x_{\alpha i}^r, \partial \mathbf{Z}(x_{\beta 1}) / \partial y} \\ \mathbf{C}_{x_{\alpha 1}^1, \partial \mathbf{Z}(x_{\beta 2}) / \partial y} & \mathbf{C}_{x_{\alpha 2}^1, \partial \mathbf{Z}(x_{\beta 2}) / \partial y} & \dots & \mathbf{C}_{x_{\alpha i}^r, \partial \mathbf{Z}(x_{\beta 2}) / \partial y} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{C}_{x_{\alpha 1}^1, \partial \mathbf{Z}(x_{\beta j-1}) / \partial z} & \mathbf{C}_{x_{\alpha 2}^1, \partial \mathbf{Z}(x_{\beta j-1}) / \partial z} & \dots & \mathbf{C}_{x_{\alpha i}^r, \partial \mathbf{Z}(x_{\beta j-1}) / \partial z} \\ \mathbf{C}_{x_{\alpha 1}^1, \partial \mathbf{Z}(x_{\beta j}) / \partial z} & \mathbf{C}_{x_{\alpha 2}^1, \partial \mathbf{Z}(x_{\beta j}) / \partial z} & \dots & \mathbf{C}_{x_{\alpha i}^r, \partial \mathbf{Z}(x_{\beta j}) / \partial z} \end{bmatrix}. \quad (\text{B17})$$

As the interfaces are relative to a reference point per later $\mathbf{x}_{\alpha 0}^k$ the value of the covariance function will be the difference between this point and the rest on the same layer,

$$\mathbf{C}_{x_{\alpha i}^r, \partial \mathbf{Z}(x_{\beta j}) / \partial x} = \overbrace{\mathbf{C}_{\mathbf{Z}(x_{\alpha i}^r), \partial \mathbf{Z}(x_{\beta j}) / \partial x}}^{r_a} - \overbrace{\mathbf{C}_{\mathbf{Z}(x_{\alpha 0}^r), \partial \mathbf{Z}(x_{\beta j}) / \partial x}}^{r_b}, \quad (\text{B18})$$

with the covariance of the scalar field being a function of the vector \mathbf{r} , its directional derivative is analogous to the previous derivations:

$$\mathbf{C}_{\mathbf{Z}, \partial \mathbf{Z} / \partial u} = \frac{\partial \mathbf{C}_{\mathbf{Z}}(\mathbf{r})}{\partial \mathbf{r}} \frac{\partial \mathbf{r}}{\partial h_u} = -\frac{h_u}{r} \mathbf{C}'_{\mathbf{Z}}. \quad (\text{B19})$$

B4 Universal matrix

As the mean value of the scalar field is going to be always unknown, it needs to be estimated from data itself. The simplest approach is to consider the mean constant for the whole domain, i.e., ordinary kriging. However, in the *scalar field* case we can assume a certain drift towards the direction of the orientations. Therefore, the mean can be written as a function of known basis functions:

$$\mu(\mathbf{x}) = \sum_{l=0}^L a_l f^l(\mathbf{x}), \quad (\text{B20})$$

where l is the grade of the polynomials used to describe the drift. Because of the algebraic dependence of the variables, there is only one drift and therefore the unbiasedness can be expressed as

$$\mathbf{U}_{\mathbf{Z}} \lambda_1 + \mathbf{U}_{\partial \mathbf{Z} / \partial u} \lambda_2 = \mathbf{f}_{i0}. \quad (\text{B21})$$

Consequently, the number of equations are determined according to the grade of the polynomial and the number of equations forming the properties matrices, Eqs. (B22) and (B23):

$$\mathbf{U}_{\mathbf{Z}} = \begin{bmatrix} x_1^1 - x_0^1 & x_2^1 - x_0^1 & \dots & x_1^2 - x_0^2 & x_2^2 - x_0^2 & \dots & x_{i-1}^r - x_0^r & x_i^r - x_0^r \\ y_1^1 - y_0^1 & y_2^1 - y_0^1 & \dots & y_1^2 - y_0^2 & y_2^2 - y_0^2 & \dots & y_{i-1}^r - y_0^r & y_i^r - y_0^r \\ z_1^1 - z_0^1 & z_2^1 - z_0^1 & \dots & z_1^2 - z_0^2 & z_2^2 - z_0^2 & \dots & z_{i-1}^r - z_0^r & z_i^r - z_0^r \\ x_1^1 x_1^1 - x_0^1 x_0^1 & x_1^1 x_2^1 - x_0^1 x_0^1 & \dots & x_1^2 x_1^2 - x_0^2 x_0^2 & x_1^2 x_2^2 - x_0^2 x_0^2 & \dots & x_{i-1}^r x_{i-1}^r - x_0^r x_0^r & x_i^r x_i^r - x_0^r x_0^r \\ y_1^1 y_1^1 - y_0^1 y_0^1 & y_1^1 y_2^1 - y_0^1 y_0^1 & \dots & y_1^2 y_1^2 - y_0^2 y_0^2 & y_1^2 y_2^2 - y_0^2 y_0^2 & \dots & y_{i-1}^r y_{i-1}^r - y_0^r y_0^r & y_i^r y_i^r - y_0^r y_0^r \\ z_1^1 z_1^1 - z_0^1 z_0^1 & z_1^1 z_2^1 - z_0^1 z_0^1 & \dots & z_1^2 z_1^2 - z_0^2 z_0^2 & z_1^2 z_2^2 - z_0^2 z_0^2 & \dots & z_{i-1}^r z_{i-1}^r - z_0^r z_0^r & z_i^r z_i^r - z_0^r z_0^r \\ x_1^1 y_1^1 - x_0^1 y_0^1 & x_1^1 y_2^1 - x_0^1 y_0^1 & \dots & x_1^2 y_1^2 - x_0^2 y_0^2 & x_1^2 y_2^2 - x_0^2 y_0^2 & \dots & x_{i-1}^r y_{i-1}^r - x_0^r y_0^r & x_i^r y_i^r - x_0^r y_0^r \\ y_1^1 z_1^1 - y_0^1 z_0^1 & y_1^1 z_2^1 - y_0^1 z_0^1 & \dots & y_1^2 z_1^2 - y_0^2 z_0^2 & y_1^2 z_2^2 - y_0^2 z_0^2 & \dots & y_{i-1}^r z_{i-1}^r - y_0^r z_0^r & y_i^r z_i^r - y_0^r z_0^r \\ x_1^1 z_1^1 - x_0^1 z_0^1 & x_1^1 z_2^1 - x_0^1 z_0^1 & \dots & x_1^2 z_1^2 - x_0^2 z_0^2 & x_1^2 z_2^2 - x_0^2 z_0^2 & \dots & x_{i-1}^r z_{i-1}^r - x_0^r z_0^r & x_i^r z_i^r - x_0^r z_0^r \\ y_1^1 z_1^1 - y_0^1 z_0^1 & y_1^1 z_2^1 - y_0^1 z_0^1 & \dots & y_1^2 z_1^2 - y_0^2 z_0^2 & y_1^2 z_2^2 - y_0^2 z_0^2 & \dots & y_{i-1}^r z_{i-1}^r - y_0^r z_0^r & y_i^r z_i^r - y_0^r z_0^r \end{bmatrix} \quad (\text{B22})$$

$$\mathbf{U}_{\partial \mathbf{Z} / \partial u} = \begin{bmatrix} x_{\beta 1} & x_{\beta 2} & \dots & x_{\beta 1} & x_{\beta 2} & \dots & x_{\beta i-1} & x_{\beta i} \\ 1 & 1 & \dots & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & \dots & 1 & 1 & \dots & 0 & 0 \\ 0 & 0 & \dots & 0 & 0 & \dots & 1 & 1 \\ 2x_1 & 2x_2 & \dots & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & \dots & 2y_1 & 2y_2 & \dots & 0 & 0 \\ 0 & 0 & \dots & 0 & 0 & \dots & 2z_{i-1} & 2z_i \\ y_1 & y_2 & \dots & x_1 & x_2 & \dots & 0 & 0 \\ y_1 & y_2 & \dots & 0 & 0 & \dots & x_{i-1} & x_i \\ 0 & 0 & \dots & z_1 & z_2 & \dots & x_{i-1} & x_i \end{bmatrix} \left. \begin{array}{l} \partial x_{\beta i} / \partial x \\ \partial x_{\beta i} / \partial y \\ \partial x_{\beta i} / \partial z \\ \partial^2 x_{\beta i} / \partial x^2 \\ \partial^2 x_{\beta i} / \partial y^2 \\ \partial^2 x_{\beta i} / \partial z^2 \\ \partial^2 x_{\beta i} / \partial x \partial y \\ \partial^2 x_{\beta i} / \partial x \partial z \\ \partial^2 x_{\beta i} / \partial y \partial z \end{array} \right\} \quad (\text{B23})$$

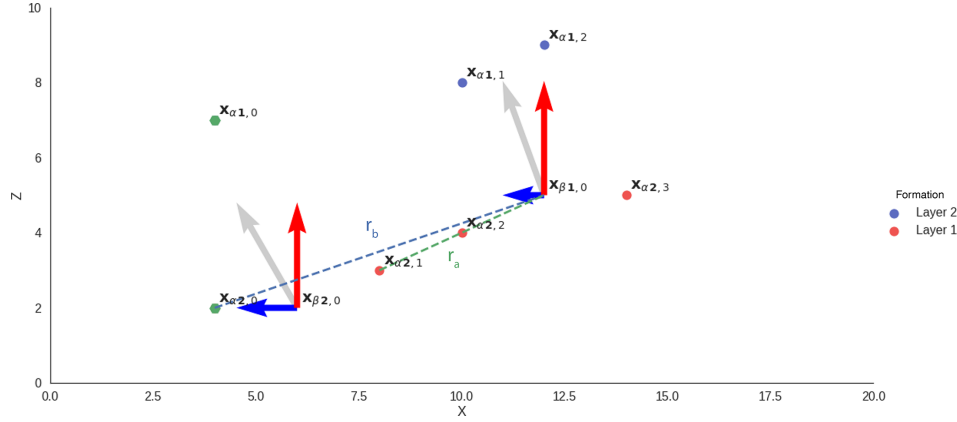


Figure B3. Distances r_a and r_b involved in the computation of the cross-covariance function. In a similar fashion as before, all interface covariance are computed relative to a reference point in each layer $x_{\alpha,0}^i$.

Appendix C: Kriging estimator

In normal kriging the right-hand term of the kriging system (Eq. 4) corresponds to covariances and drift matrices of dimensions $m \times n$, where m is the number of elements of the data sets – either x_α or x_β – and n the number of locations where the interpolation is performed, x_0 .

Since, in this case, the parameters of the variogram functions are arbitrarily chosen, the kriging variance does not hold any physical information of the domain. As a result of this, being interested in only the mean value, we can solve the kriging system in the dual form (Chiles and Delfiner, 2009; Matheron, 1981):

$$Z(x_0) = \begin{bmatrix} a'_{\partial Z/\partial u, \partial Z/\partial v} & b'_{Z, Z} & c' \\ \mathbf{c}_{\partial Z/\partial u, \partial Z/\partial v} & \mathbf{c}_{\partial Z/\partial u, Z} \\ \mathbf{c}_{Z, \partial Z/\partial u} & \mathbf{c}_{Z, Z} \\ f_{10} & f_{20} \end{bmatrix}, \quad (C1)$$

where

$$\begin{bmatrix} a_{\partial Z/\partial u, \partial Z/\partial v} \\ b_{Z, Z} \\ c \end{bmatrix} = \begin{bmatrix} \partial Z \\ 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} \mathbf{C}_{\partial Z/\partial u, \partial Z/\partial v} & \mathbf{C}_{\partial Z/\partial u, Z} & \mathbf{U}_{\partial Z/\partial u} \\ \mathbf{C}_{Z, \partial Z/\partial u} & \mathbf{C}_{Z, Z} & \mathbf{U}_Z \\ \mathbf{U}'_{\partial Z/\partial u} & \mathbf{U}'_Z & \mathbf{0} \end{bmatrix}^{-1} \quad (C2)$$

noticing that the 0 on the second row appears due to the fact that we are interpolating the difference of scalar fields instead of the scalar field itself Eq. (2).

Appendix D: Example of covariance function: cubic

The choice of the covariance function will govern the shape of the isosurfaces of the scalar field. As opposed to other

kriging uses, here the choice cannot be based on empirical measurements. Therefore, the choice of the covariance function is merely arbitrary trying to mimic, as far as possible, coherent geological structures.

The main requirement to take into consideration when the time comes to choose a covariance function is that it has to be twice differentiable, h^2 , at the origin in order to be able to calculate $\mathbf{C}_{\partial Z/\partial u, \partial Z/\partial v}$ as we saw in Eq. (B13). The use of a Gaussian model $C(r) = \exp(-r/a)^2$ and the non-divergent spline $C(r) = r^4 \text{Log}(r)$ and their correspondent flaws are explored in Lajaunie et al. (1997).

The most widely used function in the potential field method is the cubic covariance due to mathematical robustness and its coherent geological description of the space.

$$C(r) = \begin{cases} C_0(1 - 7(\frac{r}{a})^2 + \frac{35}{4}(\frac{r}{a})^3 - \frac{7}{2}(\frac{r}{a})^5 + \frac{3}{4}(\frac{r}{a})^7) & \text{for } 0 \leq r \leq a, \\ 0 & \text{for } r \geq a. \end{cases} \quad (D1)$$

with a being the range and C_0 the variance in the data. The value of a determines the maximum distance that a data point influences another. Since, we assume that all data belong to the same depositional phase, it is recommended to choose values close to the maximum extent to interpolate in order to avoid mathematical artifacts. For the values of the covariance at 0 and nuggets effects, so far only ad hoc values have been used. It is important to notice that the only effect that the values of the covariance in the potential-field method has is to weight the relative influence of both cokriging parameters (interfaces and orientations) since the absolute value of the field is meaningless. Regarding the nugget effect, the authors' recommendation is to use fairly small nugget effects to give stability to the computation – since we normally use the kriging mean, it should not have further impact to the result.

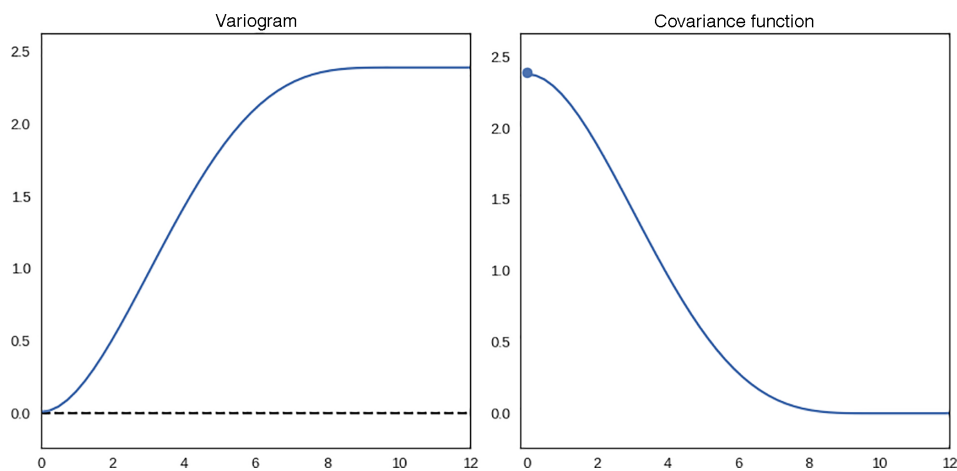


Figure D1. Representation of a cubic variogram and covariance for an arbitrary range and nugget effect.

Appendix E: Model checking and diagnostics

Estimating convergence of MCMC simulations. Here we show selected plots to evaluate convergence of the MCMC process. In Fig. E2, we present trace plots for selected parameters, and, in Fig. E1, the corresponding plots of the calculated Geweke statistics (Geweke, 1991). These parameters have been selected to show the overall range of convergence and trace behavior. Overall, the sampling algorithm performs well, although in some cases the step length could be adjusted further. Also, the Geweke statistics for some parameters fall partly outside of 2 standard deviations, indicating that the chain may not have fully converged. As described in the Discussion, we will attempt to address these issues with the use of a faster implementation of the modeling algorithm and by considering better sampling strategies in future work.

Appendix F: Blender integration

Throughout the paper we have mentioned and show Blender visualizations (Fig. 1b). The first step to obtain them is to be able to run *GemPy* in Blender's integrated Python code (there are several tutorials online to use external libraries in Blender). Once it is running, we can use Blender's library *bpy* to generate Blender's actors directly from code. Here we include the code listing with the extra functions necessary to automatically create *GemPy* models in Blender.

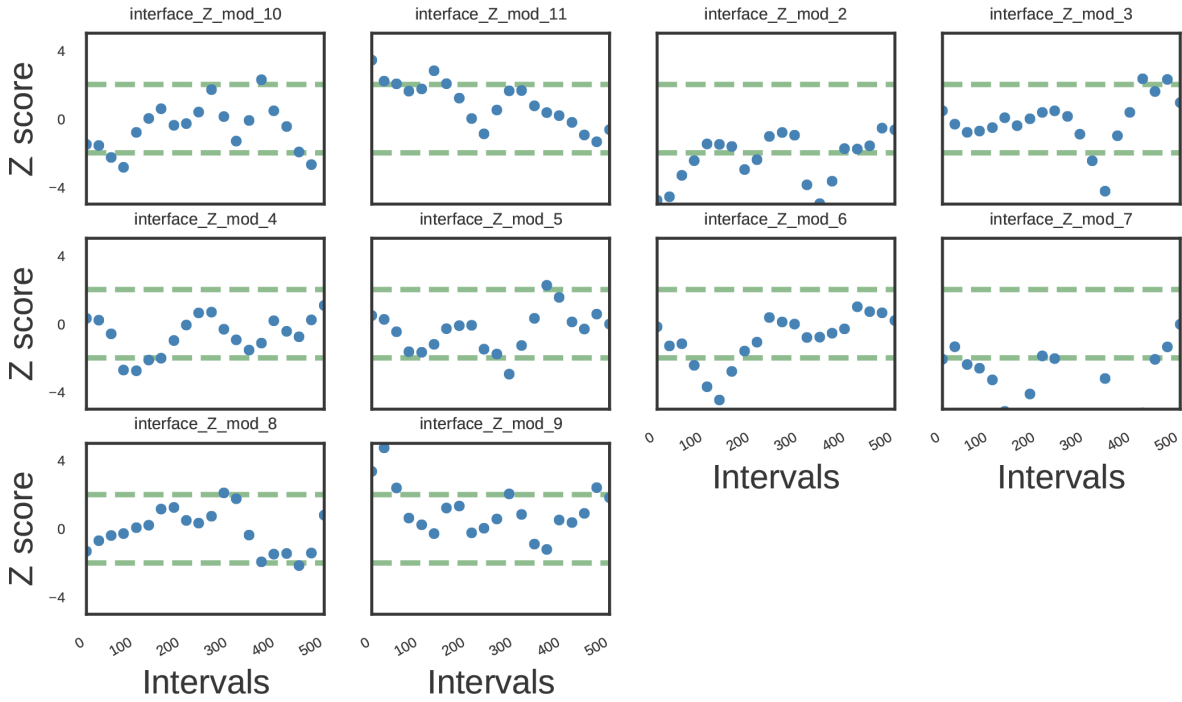


Figure E1. Geweke values of the parameters belonging to the inference of all likelihoods. Every point represents the mean of separated intervals of the chain. If interval *A* and interval *B* belong to the same distribution, most of the *Z* score should fall within 2 SD.

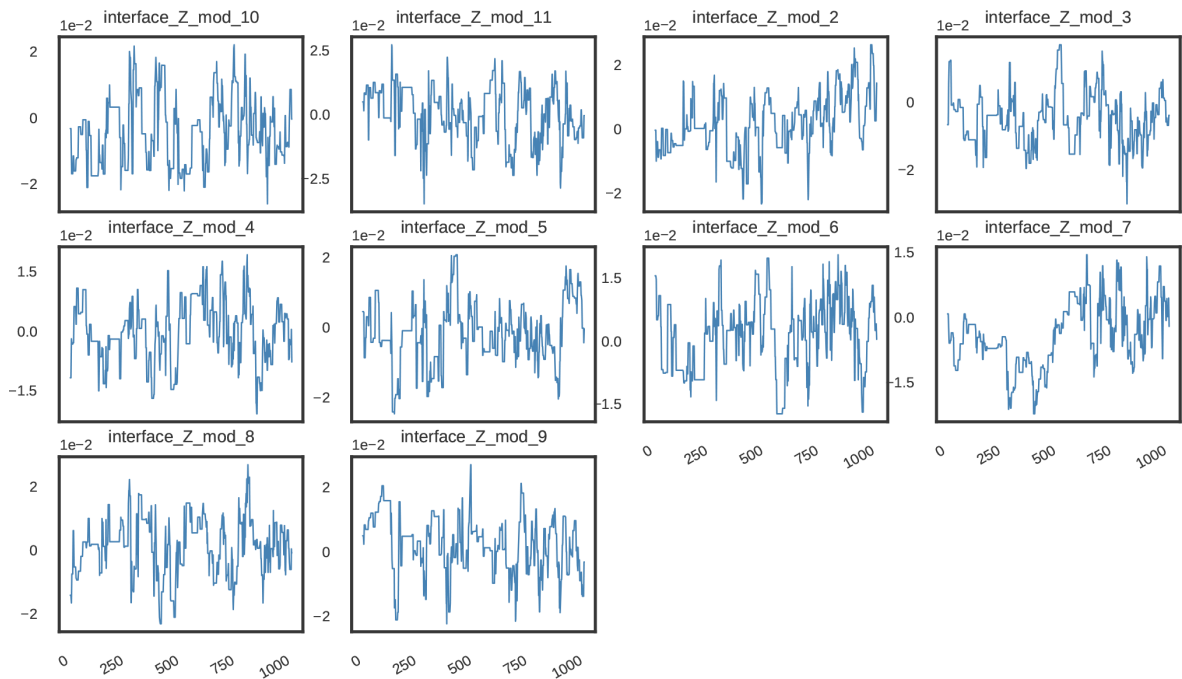


Figure E2. Traces of the parameters belonging to the inference with all likelihoods. The asymptotic behavior probes ergodicity. The first 1000 iterations (not displayed here) were used as burn-in and are not represented here.

```

# Import Blender library, GemPy and GemPy colors
import bpy
import gempy as gp
from gempy.colors import color_lot

# Delete previous objects
try:
    bpy.ops.object.mode_set(mode='OBJECT')
    bpy.ops.object.select_by_type(type='MESH')
    bpy.ops.object.delete(use_global=False)
    for item in bpy.data.meshes:
        bpy.data.meshes.remove(item)
except:
    pass

# Define function to create a Blender material, i.e. texture
def makeMaterial(name, diffuse, specular, alpha):
    mat = bpy.data.materials.new(name)
    mat.diffuse_color = diffuse
    mat.diffuse_shader = "LAMBERT"
    mat.diffuse_intensity = 1.0
    mat.specular_color = specular
    mat.specular_shader = "COOKTORR"
    mat.specular_intensity = 0.5
    mat.alpha = alpha
    mat.ambient = 1
    return mat

# Define function the assing material to object
def setMaterial(ob, mat):
    me = ob.data
    me.materials.append(mat)

...
add listing 1
...

# Create interface spheres and orientation copes out of a (rescaled) geodata
# Interfaces
for e, val in enumerate(interp_data.geo_data_res.interfaces.iterrows()):
    index = val[0]
    row = val[1]
    color = makeMaterial('color', color_lot[row['formation_number']], (1,1,1),1)
    origin = (row['X']*10, row['Y']*10, row['Z']*10)
    bpy.ops.mesh.primitive_uv_sphere_add(location=origin, size=0.1)
    setMaterial(bpy.context.object, color)

# Orientations
for e, val in enumerate(interp_data.geo_data_res.orientations.iterrows()):
    index = val[0]
    row = val[1]
    red = makeMaterial('Red', color_lot[row['formation_number']], (1,1,1),1)
    origin = (row['X']*10, row['Y']*10, row['Z']*10)
    rotation_p = (row['G_y'], row['G_x'], row['G_z'])
    bpy.ops.mesh.primitive_cone_add(location=origin, rotation=rotation_p)
    bpy.context.object.dimensions = [.3, .3, .3]
    #bpy.ops.transform.translate(value=(1,0,0))
    setMaterial(bpy.context.object, red)

# Create rescaled simpleces
verts, faces = gp.get_surfaces(interp_data, lith_block[1],
                               fault_block[1],
                               original_scale=False)

# or import them from a previous project
import numpy as np
verts = np.load('perth_ver.npy')
faces = np.load('perth_sim.npy')

# Create surfaces
for i in range(0, n_formation):
    mesh_data = bpy.data.meshes.new("cube_mesh_data")
    mesh_data.from_pydata(verts[i]*10, [], faces[i].tolist())
    mesh_data.update()

    obj = bpy.data.objects.new("My_Object", mesh_data)
    red = makeMaterial('Red', color_lot[i+1], (1,1,1),1)

    setMaterial(obj, red)
    scene = bpy.context.scene
    scene.objects.link(obj)
    obj.select = True

```

Listing F8. Extra functionality needed to create GemPy models in Blender.

Author contributions. MdIV and FW contributed to project conceptualization and method development. MdIV wrote and maintained the code with the help of AS (topology asset and the initial vtk-based visualization). MdIV prepared the manuscript with contributions of both co-authors in reviewing and editing. AS was involved in visualizing some results and writing some chapters (topology and visualization). FW provided overall project supervision and funding.

Competing interests. The authors declare that they have no conflict of interest.

Acknowledgements. The authors would like to acknowledge all of the people – all around the world – who have contributed to the final state of the library, either by stimulating mathematical discussions or finding bugs. This project could not be possible without their invaluable help.

Edited by: Lutz Gross

Reviewed by: Sally Cripps and one anonymous referee

References

- Aug, C.: Modélisation géologique 3D et caractérisation des incertitudes par la méthode du champ de potentiel: PhD thesis, PhD thesis, ENSMP, Paris, 2004.
- Ayachit, U.: The ParaView Guide: A Parallel Visualization Application, Kitware, Inc., USA, 2015.
- Bardossy, G. and Fodor, J.: Evaluation of Uncertainties and Risks in Geology: New Mathematical Approaches for their Handling, Springer, Berlin, Germany, 2004.
- Baydin, A. G., Pearlmutter, B. A., Radul, A. A., and Siskind, J. M.: Automatic differentiation in machine learning: a survey, arXiv preprint arXiv:1502.05767, 2015.
- Bellman, R.: Dynamic Programming, Courier Corporation, Dover Books on Computer Science, 366 pp., ISBN:9780486317199, 2013.
- Betancourt, M., Byrne, S., Livingstone, S., and Girolami, M.: The geometric foundations of Hamiltonian Monte Carlo, *Bernoulli*, 23, 2257–2298, 2017.
- Bishop, C. M.: Model-based machine learning, *Philos. T. R. Soc. A*, 371, <https://doi.org/10.1098/rsta.2012.0222>, 2013.
- Blender Online Community: Blender – a 3D modelling and rendering package, Blender Foundation, Blender Institute, Amsterdam, available at: <http://www.blender.org> (last access: 17 December 2018), 2017.
- Bond, C. E.: Uncertainty in structural interpretation: Lessons to be learnt, *J. Struct. Geol.*, 74, 185–200, 2015.
- Bond, C. E., Gibbs, A. D., Shipton, Z. K., and Jones, S.: What do you think this is? “Conceptual uncertainty” in geoscience interpretation, *GSA Today*, 17, <https://doi.org/10.1130/GSAT01711A.1>, 2007.
- Caers, J.: Introduction, in: Modeling Uncertainty in the Earth Sciences, John Wiley & Sons, Ltd, Chichester, UK, 1–8, 2011.
- Calcagno, P., Chiles, J.-P., Courrioux, G., and Guillen, A.: Geological modelling from field data and geological knowledge: Part I. Modelling method coupling 3D potential-field interpolation and geological rules: Recent Advances in Computational Geodynamics: Theory, Numerics and Applications, *Phys. Earth Planet. In.*, 171, 147–157, 2008.
- Caumon, G.: Towards Stochastic Time-Varying Geological Modelling, *Math. Geosci.*, 42, 555–569, 2010.
- Caumon, G., Collon-Drouaillet, P., Le Carlier de Veslud, C., Viseur, S., and Sausse, J.: Surface-Based 3D Modeling of Geological Structures, *Math. Geosci.*, 41, 927–945, 2009.
- Caumon, G., Gray, G., Antoine, C., and Titeux, M.-O.: Three-Dimensional Implicit Stratigraphic Model Building From Remote Sensing Data on Tetrahedral Meshes: Theory and Application to a Regional Model of La Popa Basin, NE Mexico, *IEEE T. Geosci. Remote*, 51, 1613–1621, 2013.
- Chatfield, C.: Model Uncertainty, Data Mining and Statistical Inference, *J. Roy. Stat. Soc. A*, 158, 419–466, 1995.
- Chiles, J.-P. and Delfiner, P.: Geostatistics: modeling spatial uncertainty, vol. 497, John Wiley & Sons, 2009.
- Chiles, J.-P., Aug, C., Guillen, A., and Lees, T.: Modelling of Geometry of Geological Units and its Uncertainty in 3D From Structural Data: The Potential-Field Method, Perth, 2004.
- Christakos, G.: On the assimilation of uncertain physical knowledge bases: Bayesian and non-Bayesian techniques, *Adv. Water Resour.*, 25, 1257–1274, 2002.
- Cockett, R., Kang, S., Heagy, L. J., Pidlisecky, A., and Oldenburg, D. W.: SimPEG: An open source framework for simulation and gradient based parameter estimation in geophysical applications, *Comput. Geosci.*, 85, 142–154, 2015.
- Cohen, J. S.: Computer algebra and symbolic computation: Mathematical methods, Universities Press, 2003.
- Dagum, L. and Menon, R.: OpenMP: an industry standard API for shared-memory programming, *IEEE Comput. Sci. Eng.*, 5, 46–55, 1998.
- de la Varga, M. and Wellmann, J. F.: Structural geologic modeling as an inference problem: A Bayesian perspective, *Interpretation*, 4, 1–16, 2016.
- de la Varga, M., Schaaf, A., Wellmann, F., Meeßen, F. C., and Wagner, F.: cgre-aachen/gempy: GemPy 1.0 pre-release (Version 1.0), Zenodo, <https://doi.org/10.5281/zenodo.1186118>, 2018.
- Dentith, M. and Mudge, S. T.: Geophysics for the mineral exploration geoscientist, Cambridge University Press, 2014.
- Duane, S., Kennedy, A. D., Pendleton, B. J., and Roweth, D.: Hybrid monte carlo, *Phys. Lett. B*, 195, 216–222, 1987.
- Fiorio, C. and Gustedt, J.: Two linear time union-find strategies for image processing, *Theor. Comput. Sci.*, 154, 165–181, 1996.
- Geuzaine, C. and Remacle, J.-F.: Gmsh: A 3-D finite element mesh generator with built-in pre-and post-processing facilities, *Int. J. Numer. Meth. Eng.*, 79, 1309–1331, 2009.
- Geweke, J.: Evaluating the accuracy of sampling-based approaches to the calculation of posterior moments, vol. 196, Federal Reserve Bank of Minneapolis, Research Department Minneapolis, MN, USA, 1991.
- Gonçalves, Í. G., Kumaira, S., and Guadagnin, F.: A machine learning approach to the potential-field method for implicit modeling of geological structures, *Comput. Geosci.*, 103, 173–182, 2017.
- Haario, H., Saksman, E., and Tamminen, J.: An adaptive metropolis algorithm, *Bernoulli*, 7, 223–242, 2001.
- Hillier, M. J., Schetselaar, E. M., de Kemp, E. A., and Perron, G.: Three-Dimensional Modelling of Geological Surfaces Us-

- ing Generalized Interpolation with Radial Basis Functions, *Math. Geosci.*, 46, 931–953, 2014.
- Hoffman, M. D. and Gelman, A.: The No-U-turn sampler: adaptively setting path lengths in Hamiltonian Monte Carlo, *J. Mach. Learn. Res.*, 15, 1593–1623, 2014.
- Hunter, J. D.: Matplotlib: A 2D graphics environment, *Comput. Sci. Eng.*, 9, 90–95, 2007.
- Jaccard, P.: The distribution of the flora in the alpine zone, *New Phytol.*, 11, 37–50, 1912.
- Jessell, M. W., Ailleres, L., and Kemp, A. E.: Towards an Integrated Inversion of Geoscientific data: what price of Geology?, *Tectonophysics*, 490, 294–306, 2010.
- Jordan, M. I.: Learning in graphical models, vol. 89, Springer Science & Business Media, 1998.
- Koller, D. and Friedman, N.: Probabilistic Graphical Models: Principles and Techniques – Adaptive Computation and Machine Learning, The MIT Press, 2009.
- Kucukelbir, A., Ranganath, R., Gelman, A., and Blei, D.: Automatic variational inference in Stan, Proceedings of the 28th International Conference on Neural Information Processing Systems, Montreal, Canada, 1, 568–576, 2015.
- Kucukelbir, A., Tran, D., Ranganath, R., Gelman, A., and Blei, D. M.: Automatic Differentiation Variational Inference, arXiv preprint arXiv:1603.00788, *J. Mach. Learn. Res.*, 18, 430–474, 2017.
- Lajaunie, C., Courrioux, G., and Manuel, L.: Foliation fields and 3D cartography in geology: Principles of a method based on potential interpolation, *Math. Geol.*, 29, 571–584, 1997.
- Lark, R. M., Mathers, S. J., Thorpe, S., Arkley, S. L. B., Morgan, D. J., and Lawrence, D. J. D.: A statistical assessment of the uncertainty in a 3-D geological framework model, *P. Geologists Assoc.*, 124, 946–958, 2013.
- Lauritzen, S. L., Dawid, A. P., Larsen, B. N., and Leimer, H.-G.: Independence properties of directed Markov fields, *Networks*, 20, 491–505, 1990.
- Lawson, C. L., Hanson, R. J., Kincaid, D. R., and Krogh, F. T.: Basic linear algebra subprograms for Fortran usage, *ACM T. Math. Software*, 5, 308–323, 1979.
- Lindsay, M., Ailleres, L., Jessell, M. W., de Kemp, E., and Betts, P. G.: Locating and quantifying geological uncertainty in three-dimensional models: Analysis of the Gippsland Basin, southeastern Australia, *Tectonophysics*, 546–547, 10–27, 2012.
- Lindsay, M. D., Jessell, M. W., Ailleres, L., Perrouty, S., de Kemp, E., and Betts, P. G.: Geodiversity: Exploration of 3D geological model space, *Tectonophysics*, 594, 27–37, 2013a.
- Lindsay, M. D., Perrouty, S., Jessell, M. W., and Ailleres, L.: Making the link between geological and geophysical uncertainty: geodiversity in the Ashanti Greenstone Belt, *Geophys. J. Int.*, 195, 903–922, 2013b.
- Lorensen, W. E. and Cline, H. E.: Marching cubes: A high resolution 3D surface construction algorithm, in: *ACM siggraph computer graphics*, 21, 163–169, 1987.
- Mallet, J.-L.: Space-time mathematical framework for sedimentary geology, *Math. Geol.*, 36, 1–32, 2004.
- Marechal, A.: Kriging seismic data in presence of faults, in: *Geostatistics for natural resources characterization*, Springer, 271–294, 1984.
- Matheron, G.: Splines and kriging: their formal equivalence, *Down-to-earth-statistics: Solutions looking for geological problems*, 77–95, 1981.
- McKinney, W.: pandas: a foundational Python library for data analysis and statistics, *Python for High Performance and Scientific Computing*, 1–9, 2011.
- McLane, M., Gouveia, J., Citron, G. P., MacKay, J., and Rose, P. R.: Responsible reporting of uncertain petroleum reserves, *AAPG Bull.*, 92, 1431–1452, 2008.
- Mosegaard, K. and Tarantola, A.: Monte Carlo sampling of solutions to inverse problems, *J. Geophys. Res.*, 100, 12–431, 1995.
- Nabighian, M. N., Ander, M. E., Grauch, V. J. S., Hansen, R. O., LaFehr, T. R., Li, Y., Pearson, W. C., Peirce, J. W., Phillips, J. D., and Ruder, M. E.: Historical development of the gravity method in exploration, *Geophysics*, 70, 63ND–89ND, <https://doi.org/10.1190/1.2133785>, 2005.
- Nagy, D.: The gravitational attraction of a right rectangular prism, *Geophysics*, 31, 362–371, 1966.
- Ogilvie, J. F.: A Monte-Carlo approach to error propagation, *Comput. Chem.*, 8, 205–207, 1984.
- Patil, A., Huard, D., and Fonnesbeck, C. J.: PyMC: Bayesian stochastic modelling in Python, *J. Stat. Softw.*, 35, 1–81, 2010.
- Rall, L. B.: Automatic differentiation: Techniques and Applications, *Lecture Notes in Computer Science*, Springer-Verlag Berlin Heidelberg, <https://doi.org/10.1007/3-540-10861-0>, 1981.
- Robitaille, T. P., Tollerud, E. J., Greenfield, P., Droettboom, M., Bray, E., Aldcroft, T., Davis, M., Ginsburg, A., Price-Whelan, A. M., Kerzendorf, W. E., Conley, A., Crighton, N., Barbary, K., Muna, D., Ferguson, H., Grollier, F., Parikh, M. M., Nair, P. H., Gunther, H. M., Deil, C., Woillez, J., Conseil, S., Kramer, R., Turner, J. E. H., Singer, L., Fox, R., Weaver, B. A., Zabalza, V., Edwards, Z. I., Azalee Bostroem, K., Burke, D. J., Casey, A. R., Crawford, S. M., Dencheva, N., Ely, J., Jenness, T., Labrie, K., Lim, P. L., Pierfederici, F., Pontzen, A., Ptak, A., Refsdal, B., Servillat, M., and Streicher, O.: Astropy: A community Python package for astronomy, *Astron. Astrophys.*, 558, <https://doi.org/10.1051/0004-6361/201322068>, 2013.
- Salvatier, J., Wiecki, T. V., and Fonnesbeck, C.: Probabilistic programming in Python using PyMC3, *PeerJ Computer Science*, 2, 55, <https://doi.org/10.7717/peerj-cs.55>, 2016.
- Sambridge, M. and Mosegaard, K.: Monte Carlo methods in geophysical inverse problems, *Rev. Geophys.*, 40, <https://doi.org/10.1029/2000RG000089>, 2002.
- Schaaf, A.: Geological Inference based on Kinematic Structural Models, Master’s thesis, RWTH Aachen University, Aachen, Germany, 2017.
- Schroeder, W. J., Lorensen, B., and Martin, K.: The visualization toolkit: an object-oriented approach to 3D graphics, *Kitware*, 2004.
- Shannon, E. C.: A mathematical theory of communication, *AT&T Tech. J.*, 27, <https://doi.org/10.1002/j.1538-7305.1948.tb01338.x>, 1948.
- Stamm, F. A.: Bayesian Decision Theory in Structural Geological Modeling – How Reducing Uncertainties Affects Reservoir Value Estimations, Master’s thesis, RWTH Aachen University, Aachen, Germany, 2017.
- Tarantola, A.: Inverse problem theory and methods for model parameter estimation, *Society for Industrial Mathematics*, 2005.

- Tarantola, A. and Valette, B.: Inverse Problems = Quest for Information, *J. Geophys.*, 50, 159–170, 1982.
- Theano Development Team: *Theano*: A Python framework for fast computation of mathematical expressions, arXiv e-prints, abs/1605.02688, <http://arxiv.org/abs/1605.02688> (last access: 17 December 2018), 2016.
- Thiele, S. T., Jessell, M. W., Lindsay, M., Ogarko, V., Wellmann, J. F., and Pakyuz-Charrier, E.: The topology of geology 1: Topological analysis, *J. Struct. Geol.*, 91, 27–38, 2016a.
- Thiele, S. T., Jessell, M. W., Lindsay, M., Wellmann, J. F., and Pakyuz-Charrier, E.: The topology of geology 2: Topological uncertainty, *J. Struct. Geol.*, 91, 74–87, 2016b.
- van der Walt, S., Schönberger, J. L., Nunez-Iglesias, J., Boulogne, F., Warner, J. D., Yager, N., Gouillart, E., and Yu, T.: scikit-image: image processing in Python, *PeerJ*, 2, 453, <https://doi.org/10.7717/peerj.453>, 2014.
- van der Walt, S., Schönberger, J. L., Nunez-Iglesias, J., Boulogne, F., Warner, J. D., Yager, N., Gouillart, E., Yu, T., and the scikit-image contributors: scikit-image: image processing in Python, *PeerJ*, 2, e453, <https://doi.org/10.7717/peerj.453>, 2014.
- van Rossum, G., Warsaw, B., and Coghlan, N.: PEP 8: style guide for Python code, available at: <http://www.python.org/dev/peps/pep-0008/> (last access: 17 December 2018), 2001.
- Wackernagel, H.: *Multivariate geostatistics: an introduction with applications*, Springer Science & Business Media, 2013.
- Walt, S. V. D., Colbert, S. C., and Varoquaux, G.: The NumPy array: a structure for efficient numerical computation, *Comput. Sci. Eng.*, 13, 22–30, 2011.
- Wang, H., Wellmann, J. F., Li, Z., Wang, X., and Liang, R. Y.: A Segmentation Approach for Stochastic Geological Modeling Using Hidden Markov Random Fields, *Math. Geosci.*, 49, 145–177, 2017.
- Waskom, M., Botvinnik, O., O’Kane, D., Hobson, P., Lukauskas, S., Gemperline, D. C., Augspurger, T., Halchenko, Y., Cole, J. B., Warmenhoven, J., de Ruyter, J., Pye, C., Hoyer, S., Vanderplas, J., Villalba, S., Kunter, G., Quintero, E., Bachant, P., Martin, M., Meyer, K., Miles, A., Ram, Y., Yarkoni, T., Williams, M. L., Evans, C., Fitzgerald, C. B., Fonnesbeck, C., Lee, A., and Qalieh, A.: *mwaskom/seaborn: v0.8.1* (September 2017), <https://doi.org/10.5281/zenodo.883859>, 2017.
- Wellmann, F. and Caumon, G.: 3-D Structural geological models: Concepts, methods, and uncertainties, *Cedric Schmelzbach, Adv. Geophys.*, 59, Elsevier, 1–121, <https://doi.org/10.1016/bs.agph.2018.09.001>, 2018.
- Wellmann, J. F. and Regenauer-Lieb, K.: Uncertainties have a meaning: Information entropy as a quality measure for 3-D geological models, *Tectonophysics*, 526–529, 207–216, 2012.
- Wellmann, J. F., Horowitz, F. G., Schill, E., and Regenauer-Lieb, K.: Towards incorporating uncertainty of structural data in 3D geological inversion, *Tectonophysics*, 490, 141–151, 2010.
- Wellmann, J. F., Lindsay, M., Poh, J., and Jessell, M. W.: Validating 3-D Structural Models with Geological Knowledge for Improved Uncertainty Evaluations, *Energy Proced.*, 59, 374–381, 2014.
- Wellmann, J. F., Thiele, S. T., Lindsay, M. D., and Jessell, M. W.: *pynoddy 1.0: an experimental platform for automated 3-D kinematic and potential field modelling*, *Geosci. Model Dev.*, 9, 1019–1035, <https://doi.org/10.5194/gmd-9-1019-2016>, 2016.
- Wellmann, J. F., de la Varga, M., Murdie, R. E., Gessner, K., and Jessell, M.: Uncertainty estimation for a geological model of the Sandstone greenstone belt, Western Australia – insights from integrated geological and geophysical inversion in a Bayesian inference framework, *Geol. Soc. Spec. Publ.*, 453, SP453–12, <https://doi.org/10.1144/SP453.12>, 2017.
- Wu, K., Otoo, E., and Shoshani, A.: *Optimizing connected component labeling algorithms*, Lawrence Berkeley National Laboratory, 2005.