



# EIS

Exploration  
Information  
System

## D 3.4: EIS Toolkit Beta Release

User Manual and Technical Specifications

Version 3.0

**Lead Beneficiary: GISPO**

04 / 2024

**Niko Aarnio<sup>1</sup>, Timo Aarnio<sup>1</sup>, Miikka Kallio<sup>1</sup>, Mika Sorvoja<sup>1</sup>, Johanna Torppa<sup>2</sup>, Bijal Chudasama<sup>2</sup>, Ina Storch<sup>3</sup>,  
Peggy Hielscher<sup>3</sup>, Michael Steffen<sup>3</sup>, Andreas Knobloch<sup>3</sup>**

<sup>1</sup>*GISPO*

<sup>2</sup>*GTK*

<sup>3</sup>*Beak Consultants GmbH*



Funded by  
the European Union

## Disclaimer

The content of this report reflects only the author's view. The European Commission is not responsible for any use that may be made of the information it contains.

## Document information

Grant Agreement / Proposal ID	101057357
Project Title	Exploration Information System
Project Acronym	EIS
Scientific Coordinator	Vesa Nykänen (vesa.nykanen@gtk.fi) – GTK
Project starting date (duration)	1 May 2022 (36 months)
Related Work Package	WP 3
Related Task(s)	Task 3.3
Lead Organisation	GISPO
Contributing Partner(s)	GTK, BEAK, UTU
Due Date	31.10.2023
Submission Date	09.11.2023, Revision: 06.12.2023, Revision: 05.04.2024
Dissemination level	PU

## History

Date	Version	Submitted by	Reviewed by	Comments
26.10.2023	0.1	Andreas Knobloch		Template
06.11.2023	0.2	Timo Aarnio		Draft
08.11.2023	0.9	Andreas Knobloch		Revision
09.11.2023	1.0	Hafsa Munia		Submission
06.12.2023	2.0	Andreas Knobloch		Revision
07.12.2023	2.0	Hafsa Munia		Submission
04.04.2024	3.0	Andreas Knobloch		Revision
05.04.2024	3.0	Hafsa Munia		Submission

## Table of Contents

1. Executive summary .....	6
2. Introduction .....	7
2.1. General remarks .....	7
2.2. EIS WP3 and Task 3.3.....	7
2.3. Purpose of the EIS Toolkit and user requirements.....	7
3. Components of the EIS Toolkit and related questions by the users to be answered .....	9
3.1. Pre-processing tools .....	11
3.2. Exploratory data analysis tools.....	11
3.3. Model data preparation tools .....	12
3.4. Prediction/modelling tools .....	13
3.5. Evaluation tools .....	13
3.6. Additional tools.....	14
4. User manual .....	15
4.1. System requirements.....	15
4.2. Installation guide .....	15
4.3. Jupyter Notebook description .....	15
5. Technical specifications .....	16
5.1. Conversions module .....	16
5.2. Exploratory analyses module .....	16
5.3. Prediction module .....	17
5.4. Raster processing module .....	18
5.5. Training data tools module.....	18
5.6. Transformations module .....	19
5.7. Validation module .....	19
5.8. Vector processing module .....	20
6. Initial testing .....	22
6.1. Listing of initial testers.....	22
6.2. Listing and description of initial test data for testers .....	22
6.3. Guidelines for testers to report bugs and problems .....	23
7. Listing of tools that will be implemented in next 6 months until final release .....	24
8. Conclusion.....	25
9. References .....	26

## List of figures

Figure 1. The mineral prospectivity modelling (MPM) workflow (blue boxes) and EIS toolkit functions related to each phase of the workflow (black boxes). Yellow boxes provide a general description of tools needed in each phase of MPM. ...10

## Abbreviations and Acronyms

Acronym	Description
WP	Work Package
GIS	Geographic Information System
EIS	Exploration Information System
GUI	Graphical User Interface

## Summary

The Deliverable D3.4 provides a description of the contents and general structure of the EIS Toolkit beta release.

## Keywords

Software Design, Mineral Prospectivity Modelling, Mineral Predictive Mapping, QGIS, Artificial Intelligence, Toolkit, Wizard, Python, Beta release

# 1. Executive summary

This document describes the first public beta release of the “EIS Toolkit”, which is a standalone Python library for conducting Mineral Prospectivity Mapping. The toolkit utilizes industry-standard libraries for well-defined tasks and provides an interface for integrations both with the upcoming “QGIS EIS Plugin” and other software.

A common way to use this kind of a library is via Jupyter Notebooks, which is an interactive environment for running Python software and a sample of such use is included with this deliverable in **Appendix 3** both as a runnable/interactive notebook file and a static export that visualizes the workflow and its results. The interactive version works as a user manual for end users.

The technical documentation for “EIS Toolkit” is automatically generated and is described in chapter 4 and an export of the documentation (also in beta stage) is included with this deliverable in the **Appendix 4**.

As the toolkit is still in development phase it is not feature complete meaning that new functionalities (see Deliverable D3.3) will be added and current implementations will be further optimized with the help of feedback from users.

This and future releases and related documentation can be downloaded from “EIS Toolkit” GitHub repository which is located at: [https://github.com/GispoCoding/eis\\_toolkit](https://github.com/GispoCoding/eis_toolkit).

## 2. Introduction

### 2.1. General remarks

This document will provide background information about the beta release of “EIS Toolkit”. The development is still in beta phase so not all functionality is included or fully optimized. The documentation for EIS Toolkit is currently in two forms:

- a Jupyter Notebook that guides the user through an example workflow (User Guide) (see **Appendix 3**) and
- a Technical Specification document that includes documentation for the actual EIS Toolkit functionality (see **Appendix 4**).

### 2.2. EIS WP3 and Task 3.3

The main objective of WP3 is the development of a GIS (Geographical Information System)-based Exploration Information System (EIS) for predictive mapping of mineral resources. EIS does not have a strict definition but can be characterized as an environment for performing data analysis and modelling, for managing data and other information, and for representing results in various forms.

In Task 3.3, a library of Python functions for EIS is implemented. We call this library the “EIS Toolkit”, and it is a comprehensive collection of independent functions relevant for performing mineral prospectivity analysis related tasks. These tasks include mainly predictive mapping and data integration via mathematical modelling, but also some general data processing and analysis. Also, tools for evaluation of the goodness of the models and modelling results are included for efficient decision making in the identification and prioritization of exploration targets. The tools in the EIS toolkit are not specific for any specific type of mineral system, as same the same data processing, modelling and model validation methods are applicable to a number of mineral system types. The idea of the EIS toolkit is to bring together special functions needed for performing, and also to facilitate, mineral prospectivity analysis.

The “EIS Toolkit” contains implementations of existing and new algorithms (from Deliverable D3.3), and new functions can be added even after the EIS project. Emphasis is given to exploring the applicability of modern machine learning methods, such as convolutional neural networks, the use of which in mineral prospectivity modelling is still at its infancy. Existing Python libraries have been reviewed and are used to minimize the coding effort.

The structure of and interfaces to the functions included in the library have been planned (in Deliverable D3.1) so that the functions can be smoothly integrated to other software, such as the “EIS QGIS Plugin. The EIS QGIS Plugin (will be submitted later in the project through Deliverables 3.6 and 3.7) provides the platform (or the exploration information system) for entire mineral prospectivity analysis workflow, taking also into account the specific features of different mineral system types.

### 2.3. Purpose of the EIS Toolkit and user requirements

The main purpose of EIS Toolkit is:

- to collect tools needed in different phases of mineral potential analysis into a single Python library, to avoid searching through different libraries for suitable functions.

The EIS Toolkit was designed in strong connection with WP2 and WP6 to answer the following specific needs of the end users:

- to process several typical types of input data to represent mineral system proxies
- to transform proxies into a unified form for modeling
- to prepare a training set for machine learning from different forms of input data
- to evaluate input data, generated proxy data and mineral potential results

The developed tools are based on the defined proxies by EIS WP2, as defined in the Deliverable D2.1 “Preliminary list of pathfinders and vectors for mineral systems”.

Furthermore, EIS WP6 conducted a market analysis of potential end users and how the EIS Toolkit and EIS Wizard can be exploited most effectively. The results have been presented in EIS Deliverable D6.3 “Updated Version of the Communication, Exploitation and Dissemination Plan”.

In addition, in the frame of the EIS Toolkit design phase at the beginning of the project, the EIS Deliverable D3.1 “Design Report – EIS Toolkit & QGIS EIS Wizard” by EIS WP3 includes a brief market analysis of existing software packages and the user requirements.



### 3. Components of the EIS Toolkit and related questions by the users to be answered

The functions within the EIS Toolkit are designed to be universally applicable across all mineral system types, without specific components tailored to individual types. Instead, the mineral system type is considered in selecting the data that the functions are applied to.

Together with WP2 - through their submitted Deliverable D2.1 (Preliminary list of pathfinders and vectors for mineral systems) - a proxy table was generated for each of the three types of mineral systems considered in the EIS project. A proxy is the geochemical, geophysical or geological feature that is related to the mineral system and is the input to the prospectivity modelling process.

Proxies are obtained from measurements and observations through several geo-processing steps using **pre-processing tools**. Example processing steps for the IOCG type mineral system are shown in Appendices 1 and 2. **Appendix 1** contains the four main workflows for processing the data. **Appendix 2** shows which workflow is connected with each of the proxies. In addition to the main workflows outlined in Appendix 1, certain proxies may require additional attributes such as density or distance to neighboring features to be calculated by the user. Alternatively, derivatives of raster data may need to be computed or raster files may need to be re-projected.

**Exploratory analysis** tools are available to assist in selecting representative threshold values for these tasks.

Once proxies are generated, further **data transformations** are necessary based on the chosen prediction tool for mineral prospectivity modeling. The **training data tools** module facilitates the generation or processing of training data, if required, while **validation tools** allow for the evaluation of the final prediction results.

All these functionalities are encompassed within the EIS Toolkit and are organized into highlighted modules to cover the entire workflow for mineral prospectivity modeling.



Figure 1. The mineral prospectivity modelling (MPM) workflow (blue boxes) and EIS toolkit functions related to each phase of the workflow (black boxes). Yellow boxes provide a general description of tools needed in each phase of MPM.

While the EIS Toolkit itself does not implement mineral system types or guide users through a Mineral Predictive Mapping (MPM) workflow, its contents have been carefully selected to support a comprehensive MPM workflow. This section aims to provide a structured approach to mineral prospectivity analysis and justify the necessity of each tool in satisfying user needs at each step. The tools are listed in an order that users are likely to follow, as defined within the EIS QGIS Plugin, although users may navigate between steps as needed.

Each tool included in the EIS Toolkit serves a specific purpose within the mineral prospectivity analysis workflow. The categories and tools presented here are intended to align with the concrete needs of users at each step. While the reasons for inclusion of each individual tool are not explicitly described in this section due to the beta release status of the EIS Toolkit, the tools are justified implicitly by addressing the need for each tool category. The categorization of tools may differ slightly within the toolkit itself, but this structure aims to provide clarity and guidance for users.

## 3.1. Pre-processing tools

*Related User Questions:*

- **How do I extract important information from my raw data?**
- **How do I produce datasets efficiently using a mineral system approach?**

*EIS Tools:*

- **Distance computation (Euclidean distance to vector features)**
- **Distance to anomaly (Euclidean distance to anomalous raster pixels)**
- **Density computation (density of vector features)**
- **Interpolation (vector features)**
  - **Inverse distance weighting (IDW)**
  - **Kriging interpolation (ordinary, universal)**

The resulting proxy datasets, defined in EIS as either distance raster or density raster, remain closely connected to the original phenomena or observations while being in a consistent format suitable for further processing in modeling tasks.

The Distance to anomaly workflow was identified as a common method for generating proxy data from raw raster data, particularly in geophysics. In this workflow, point data can be interpolated (optional if not already a raster file) and the raster is then reclassified with a user defined trash hold value defining anomalous values. Then, distances to anomalies are computed.

## 3.2. Exploratory data analysis tools

*Related User Questions:*

- **How do I explore the data I have?**
- **How do I select datasets for proxy processing or modeling?**
- **How can I define anomalous values in my data?**

*EIS Tools:*

- **Distribution plots: Histogram, KDE (kernel density estimation), ECDF (empirical cumulative distribution functions)**
- **Relational plots: Scatterplot, Lineplot, Pairplot**

- **Categorical plots: Boxplot, Barplot**
- **Other plots: Heatmap, Parallel coordinates plot**
- **Statistics: Descriptive statistics (Minimum, maximum, mean, quantiles, standard deviation, skewness), Chi-square test, Correlation matrix, Covariance matrix, Shapiro-Wilk Normality test**
- **Exploratory tools: PCA (Principal component analysis), Autoencoder model for image segmentation, Self-organizing map (SOM), K-means, DBSCAN (Density-Based Spatial Clustering of Applications with Noise), Feature importance**

The contents of exploratory data analysis have been selected with the mindset of providing a comprehensive set of different plots, statistics and methods to satisfy user needs. As such, not all tools will be needed for each analysis and some users might not be interested in using certain tools available. However, the categories and some of the tools can be highlighted here to be essential inclusions for EIS Toolkit.

Various plots are a powerful and standard way to inspect patterns in data. In practice, thresholds, correlations and magnitudes are often explored. Distribution plots are applicable for both vector and raster data, which makes them powerful and frequently used at all stages of analysis. Especially histograms are commonly used, as it is the most well-known distribution plot method. As an interesting and newer method to study data with many attributes or modalities, parallel coordinates plot was implemented.

Along with plotting, statistics are commonly used to understand the fundamental properties of data. Some statistics can be directly used for selecting a parameter value, such as threshold, for some tool such as Distance to anomaly in proxy processing phase, while other offer more interpretive – but necessary – information for decision-making. All included statistics were seen as commonly needed in MPM workflows.

The set of exploratory methods is a varied one and includes tools that might call for individual justification. Some of these tools can be used to support preparing proxies, while some offer an alternative approach to data processing. PCA is one of the most commonly used methods to reduce dimensionality of data, and in MPM can be used for example on a geochemical dataset with many different attributes. The output of PCA could be used to either inspect patterns in data, or be used directly in modeling later. Autoencoder for image segmentation, SOM, K-means clustering and DBSCAN are all primarily methods to inspect spatial patterns in the data, and together form a good supplementary set of data exploration. Feature importance is a method designed to support modeling by ranking the relevance of datasets in machine learning modeling.

### 3.3. Model data preparation tools

*Related User Questions:*

- **What data requirements does my chosen model have?**
- **What further proxy preprocessing is needed for optimal results?**
- **How do I address potential issues like data imbalance in my model?**

*EIS Tools:*

- **Machine learning data preparation: Balancing, Sampling, One-hot encoding**
- **General data preparation: Splitting datasets, Unifying raster, Resampling, Re-projecting, Clipping**
- **Data transformations: Reclassifying, Binarizing, Clip, Normalizing, Transforming, CoDa (compositional data transforms)**

Model data preparation consists of tools that are relevant for specific models and tools that the user might need or want to apply before proceeding to prediction phase. Data balancing, sampling and one-hot encoding are examples of model-specific methods that used to prepare data for machine learning methods. While data

balancing is not required but might improve the performance of the model, one-hot encoded data is strictly required for correct interpretation of categorical data. Scaling all datasets to same range, using same scaling method, is often necessary for modeling. The included data transformation functions satisfy this need in a very comprehensive way, offering multiple different methods. Perhaps the most commonly used one of these is min-max scaling data into range [0, 1], also called data normalization. To produce data that can be fed into Fuzzy overlay model, several of these methods can be utilized to produce the desired memberships (gaussian, large, linear, near, power, small). Of the listed tools under General data preparation, Unify raster is probably the most useful one to ensure all datasets are matching in their shape and size. Many of the other tools of the category are used internally by Unify raster, but made available for users that want to preprocess step-by-step.

## 3.4. Prediction/modelling tools

*Related User Questions:*

- **How can I predict mineral deposit locations?**
- **What machine learning models are suitable for my task?**
- **How do I evaluate and apply trained models?**

*EIS Tools:*

- **Train/Test various machine learning models: Random forest, Gradient boosting, Logistic regression, MLP, CNN**
- **Apply trained models**
- **Data-driven methods: Fuzzy overlay, Weights of evidence**

The prediction methods included in EIS Toolkit can be grouped into two categories: machine learning methods and data-driven methods. Fuzzy overlay and Weights of evidence constitute the data-driven methods and the rest are machine learning modeling related tools. As both Fuzzy overlay and Weights of evidence are commonly used methods in ML modeling, they were a natural inclusion. For all machine learning methods where applicable, both classifier and regressor models were included to cover all possible prediction tasks, even if classifiers are more direct to predict mineral occurrences. Of the machine learning methods, especially Random forest models have been widely applied in MPM contexts. On the other end, CNN models have only recently been built for mineral prediction and therefore represent the more innovative side of EIS Toolkit, bringing emerging and modern tools for advanced MPM. The rest of the machine learning models fall somewhere in-between these two in respect to their popularity, and were included to offer a comprehensive set of methods to choose from.

## 3.5. Evaluation tools

*Related User Questions:*

- **Is my model effectively predicting mineral deposits?**
- **How can I assess the performance of my model?**
- **Should I adjust my model's parameters or data?**

*EIS Tools:*

- **Classifier metrics (Accuracy, Precision, Recall, F-score)**
- **Regressor metrics (MAE, MSE, RMSE, R2)**
- **Area under curve (AUC)**

- **Prediction area (PA) curve**
- **Receiver Operating Characteristic (ROC) curve**
- **Confusion matrix**
- **Loss function (for neural networks)**

Evaluation methods allow users to assess the effectiveness of trained models using a variety of metrics and visualization tools. Both classifier and regressor metrics are available, along with ROC and PA curves for classification tasks and loss functions for neural networks. ROC curve has been included, since it can be a useful tool to when assessing one or more classifiers at varying classification thresholds and has been used in MPM contexts. For neural networks, classifier and regressor evaluation methods can be applied, but the loss function provides additional insight into the training performance.

### 3.6. Additional tools

- **Cell-based association**
- **Extract shared lines of vector features**
- **Rasterize vector**
- **Reproject vector**
- **Calculate vector feature geometry**
- **Create constant raster**
- **Extract values from raster**
- **Unique combinations in raster**
- **Surface derivatives**
- **Filters**
- **Replace data/no data**
- **Harmonize no data**

Additional tools not directly aligned with the previous steps are included in the toolkit to provide alternative or supplementary methods for data processing and analysis. These tools offer insights into data characteristics and address specific user needs, enhancing the toolkit's versatility and utility. Examples include tools for rasterizing vectors, reprojecting data, and computing surface derivatives.

By following the structured approach outlined in this section and leveraging the functionalities provided by the EIS Toolkit, users can efficiently preprocess data, generate proxies, train models, and evaluate predictions for mineral prospectivity mapping, ensuring a robust and informed decision-making process.

## 4. User manual

This chapter outlines how to use “EIS Toolkit”, what are the system requirements, how it can be installed and a Jupyter Notebook for interactive use with guidance. Additionally, a static export of the Jupyter Notebook is provided that visualizes an example workflow when done inside Jupyter Notebook.

### 4.1. System requirements

The only system requirement for installing the “EIS Toolkit” through the Python package is to have a supported version of Python installed on the system. The “EIS Toolkit” supports all Python 3.9.x and 3.10.x versions. Setting up a development environment requires cloning the source code from [GitHub](#), and additional requirements based on preferences. The environment can be set up using either Docker, Poetry or Conda. Note that setting up the development environment with Conda requires Python 3.9.

### 4.2. Installation guide

There are multiple ways to set up the “EIS Toolkit”. The simplest and most common way is to install the “EIS Toolkit” Python package through pip. This can be achieved with a single command. The package automatically installs all the necessary Python libraries, and the “EIS Toolkit” will be ready to use. The “EIS Toolkit” Python package can be downloaded from the GitHub repository under [releases](#).

The “EIS Toolkit” can also be installed by setting up the development environment. Clone the source code from GitHub, and follow the preferred instructions found under `instructions/` directory:

- Setting up with [Docker](#).
- Setting up with [Poetry](#).
- Setting up with [Conda](#).

### 4.3. Jupyter Notebook description

JupyterLab is included as a development dependency for testing purposes. It can be used for example to visualize data and test “EIS Toolkit” functions. The notebooks are found under the `notebooks/` directory. You can import and use the toolkit’s functions in these notebooks.

There is a notebook that contains general usage instructions for running and modifying Jupyter Notebooks and another one for testing that dependencies to other python packages work. At the time of writing (November 2023), there are a few notebooks to showcase some of the toolkit’s implemented functions and a notebook to demonstrate the workflow (please see Appendix 1).

Instructions to run JupyterLab can be found in [GitHub](#).

## 5. Technical specifications

The "EIS Toolkit" encompasses a comprehensive suite of functions designed to facilitate efficient and effective mineral exploration. These functions are meticulously organized into distinct modules, each serving a specific purpose within the exploration workflow. Technical specifications for the "EIS Toolkit" functions are automatically generated from docstrings, which are extensively documented within the source code. Below is an overview of the key modules.

For further details on the technical specifications and functionalities of each module, please refer to **Appendix 4**

In the following sections, the numbers in brackets refer to the chapters / section numbers in the Technical Specifications in Appendix 4.

### 5.1. Conversions module

The Conversions Module in the "EIS Toolkit" provides essential functionalities for converting different types of geospatial data, facilitating seamless integration and analysis within the exploration workflow.

#### **CSV to Geo Data Frame Conversion (3.1):**

This function enables users to read CSV files containing geospatial data and convert them into Geo Data Frames. It supports the transformation of point data represented by X and Y coordinates or geometric shapes specified in Well-Known Text (WKT) format. Users can also define the target Coordinate Reference System (CRS) for the resulting Geo Data Frame.

#### **Raster to Data Frame Conversion (3.2):**

The Raster to Data Frame conversion function allows users to convert raster datasets into Pandas Data Frames. Users can select specific bands from multi-band raster or utilize all bands for conversion. Additionally, pixel coordinates (row, col) can be optionally included in the Data Frame for each pixel of the raster, providing valuable spatial context to the data.

By utilizing these conversion functionalities, users can seamlessly transform and prepare diverse geospatial datasets for further analysis and exploration within the "EIS Toolkit."

### 5.2. Exploratory analyses module

The Exploratory Analyses Module within the "EIS Toolkit" encompasses a diverse range of functions aimed at uncovering insights and relationships within geospatial data, facilitating informed decision-making in mineral exploration endeavors.

#### **DBSCAN Clustering (4.1):**

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) is utilized to identify clusters within geospatial data based on density. By specifying parameters such as maximum distance and minimum samples, users can effectively identify spatial patterns indicative of mineralization occurrences.

#### **Descriptive Statistics (4.2):**

This function computes descriptive statistics from both vector and raster data, providing insights into the distribution and variability of geospatial attributes essential for exploration analysis and decision-making.



**Feature Importance Evaluation (4.3):**

Evaluate the importance of features derived from geospatial data using machine learning models. By assessing the contribution of each feature to model performance, users can prioritize exploration targets effectively.

**K-means Clustering (4.4):**

Perform K-means clustering on geospatial data to partition it into distinct clusters. This aids in identifying spatially coherent regions or anomalies associated with mineralization occurrences.

**Plot Parallel Coordinates (4.5):**

Generate parallel coordinates plots to visualize multivariate relationships within geospatial datasets. This aids in identifying patterns and trends across multiple variables, enhancing understanding of mineralization indicators and exploration targets.

**Principal Component Analysis (PCA) (4.6):**

Conduct PCA to reduce the dimensionality of geospatial datasets while preserving essential information. Visualization of principal components assists in identifying spatial patterns and relationships crucial for mineral exploration.

**Statistical (Hypothesis) Testing (4.7):**

Compute statistical tests such as the Chi-square test for independence to assess relationships between categorical variables. Additionally, functions are available for computing correlation matrices, covariance matrices, and testing for normality, providing further insights into the underlying structure and characteristics of geospatial data.

Through the utilization of these exploratory analyses functions, users can gain valuable insights into the spatial characteristics of geological features and mineral occurrences, facilitating more informed and effective decision-making in mineral exploration endeavors.

## 5.3. Prediction module

The Prediction Module in the "EIS Toolkit" equips users with powerful tools for predicting mineral occurrences and evaluating prospectivity.

**Fuzzy Overlay Operations (5.1):**

This suite of functions enables users to perform various fuzzy overlay operations on raster data. By utilizing fuzzy logic techniques such as AND, OR, product, sum, and gamma overlay, users can analyze multi-dimensional geospatial data to identify areas of interest and assess mineral prospectivity.

**Weights of Evidence Analysis (5.2):**

These functions facilitate the calculation of weights of evidence, allowing users to assess the spatial associations between input data and mineral deposits. By determining posterior probabilities and weights of spatial association, users can make informed decisions regarding mineral exploration targets and planning.

Through the capabilities offered by the Prediction Module, users can leverage fuzzy logic and weights of evidence

## 5.4. Raster processing module

The Raster Processing module offers a comprehensive suite of functions for efficiently manipulating raster data, ensuring its compatibility and suitability for various analytical tasks. These tools enable users to address diverse raster data processing needs, including data validation, spatial analysis, and format standardization. By providing a robust set of functionalities, the module empowers users to seamlessly manage and optimize raster data for subsequent analysis and modeling tasks, enhancing the efficiency and accuracy of geospatial workflows.

### **Check Raster Grids (6.1):**

This function ensures consistency in gridding among input raster, optionally checking for matching bounds.

### **Clipping (6.2):**

Clips a raster with polygon geometries, extracting specific regions of interest.

### **Create Constant Raster (6.3):**

Generates a constant raster based on user-defined parameters, offering flexibility in defining extent and coordinate system.

### **Extract Values from Raster (6.4):**

Extracts raster values using point data to a Data Frame, facilitating data extraction for analysis.

### **Re-projecting (6.5):**

Re-projects a raster to match a given coordinate reference system (CRS), ensuring compatibility between datasets.

### **Resampling (6.6):**

Resamples a raster according to a specified resolution, maintaining data integrity during transformation.

### **Snapping (6.7):**

Aligns a raster to a reference grid raster, ensuring alignment for consistent analysis.

### **Unifying (6.8):**

Unifies given raster relative to a base raster by re-projecting, resampling, aligning, and optionally clipping them. This operation ensures consistency in grid properties and extents across multiple raster.

### **Windowing (6.9):**

Extracts a window from a raster centered at specified coordinates, allowing for localized analysis. Padding with no data values is applied if the window extends beyond the raster bounds.

## 5.5. Training data tools module

The Training Data Tools module equips users with essential functionalities for preparing training datasets, particularly in the context of machine learning tasks. By offering these tools, the module facilitates the creation of balanced and well-prepared training datasets, crucial for achieving optimal model performance and accuracy in machine learning applications.

### **Class Balancing (7.1):**

With this function users can address class imbalance issues by applying the SMOTETomek resampling method, ensuring a more representative distribution of classes in the training data. It adjusts the class distribution by

oversampling the minority class and under-sampling the majority class to mitigate class imbalance issues. Users can specify various parameters such as the sampling strategy and random state to customize the resampling process. The function returns the resampled feature matrix and target labels, ensuring a more balanced dataset for training machine learning models

## 5.6. Transformations module

The Transformations module provides a suite of functions for preprocessing raster data, enabling users to apply various transformations to prepare their data for modeling and analysis. These functions include binarization, clipping, normalization, logarithmic transformation, sigmoid transformation, and winsorization, offering flexibility in data preprocessing to suit different modeling needs. Whether it's transforming data into binary format, clipping values based on specific thresholds, or normalizing data to a standard scale, these tools empower users to preprocess raster data effectively before further analysis or modeling tasks.

### **Binarize (8.1):**

This function binarizes raster data based on a given threshold. It replaces values less than or equal to the threshold with 0 and values greater than the threshold with 1.

### **Clip (8.2):**

The clip transform function clips raster data based on specified upper and lower limits. It replaces values below the lower limit and above the upper limit with provided values.

### **Linear (8.3):**

Min-max-scaling normalizes raster data based on a specified new range. It transforms the data into the new interval defined by the provided minimum and maximum values.

### **Z-Score Normalization (8.3):**

This function normalizes raster data based on the mean and standard deviation. The resulting data will have a mean of 0 and a standard deviation of 1.

### **Logarithmic (8.4):**

The log-transform function performs a logarithmic transformation on the provided raster data. It replaces negative values with a specified no data value and applies the logarithmic transformation based on the chosen base.

### **Sigmoid (8.5):**

The sigmoid-transform transforms data into a sigmoid shape based on a specified new range. It uses parameters such as minimum, maximum, slope, and center to perform the transformation.

### **Winsorize (8.6):**

The winsorize function winsorizes raster data based on specified percentile values. It replaces values outside the specified percentile ranges with the nearest values within the range.

## 5.7. Validation module

The Validation Module in the "EIS Toolkit" provides essential tools for evaluating the performance and reliability of predictive models and data analyses.

**Calculate AUC (9.1):**

This function computes the area under the curve (AUC), which is a crucial metric for assessing the performance of classification models. By providing X and Y values corresponding to false positive rate and true positive rate, respectively, users can obtain a quantitative measure of model performance.

**Calculate Base Metrics (9.2):**

The Calculate Base Metrics function computes true positive rate, proportion of area, and false positive rate values for different thresholds. It leverages mineral deposit locations and mineral prospectivity maps to evaluate model performance comprehensively.

**Get P-A Plot Intersection Point (9.3):**

This function calculates the intersection point for prediction rate and area curves in a prediction-area (P-A) plot. By analyzing threshold values along with true positive rate and proportion of area values, users can identify optimal thresholds for model evaluation.

**Plot Correlation Matrix (9.4):**

The Plot Correlation Matrix function generates a heatmap visualization of the correlation matrix. This visualization aids in understanding the relationships between variables, providing insights into potential multicollinearity and the overall structure of the data.

**Plot Prediction-Area (P-A) Curves (9.5):**

This function plots prediction-area (P-A) curves, which are valuable for evaluating mineral prospectivity maps and evidential layers. By visualizing true positive rate and proportion of area values against threshold values, users can assess the performance of predictive models.

**Plot Rate Curve (9.6):**

The Plot Rate Curve function generates success rate, prediction rate, or ROC curves based on provided X and Y values. This visualization helps users understand the trade-offs between true positive rate and false positive rate, aiding in model selection and evaluation.

## 5.8. Vector processing module

The Vector Processing Module in the "EIS Toolkit" offers a range of functionalities for handling and analyzing vector data, facilitating geospatial analysis tasks.

**Cell-Based Association (10.1):**

This function initializes a Cell-Based Association (CBA) matrix from vector files. It calculates the mesh based on geometries contained in the file and cell size specified by the user. Users can add multiple vector datasets to the matrix, incorporating targeted shapes and attributes for comprehensive analysis.

**Distance Computation (10.2):**

The Distance Computation function calculates the distance from raster cells to the nearest geometry in a given set of vector data. This tool is valuable for proximity analysis and understanding spatial relationships between raster and vector datasets.

**IDW (10.3):**

The Inverse Distance Weighted (IDW) interpolation function performs spatial interpolation on vector data, generating a raster output. Users can specify the target column containing values for interpolation, resolution of the output raster, and the power parameter to control the rate at which weights decrease with distance.

**Kriging Interpolation (10.4):**

This function implements Kriging interpolation on input vector data, producing a raster output with interpolated values. Users can select the variogram model, coordinates type, and method (ordinary or universal) for Kriging interpolation, providing flexibility in spatial analysis tasks.

**Rasterize Vector (10.5):**

Transforming vector data into raster format is facilitated by the Rasterize Vector function. Users can specify parameters such as resolution, value column, default value, and merge strategy to customize the rasterization process according to their analysis requirements.

**Re-project Vector (10.6):**

Re-projecting vector data to match a given coordinate reference system (CRS) is achieved with the Re-project Vector function. Users can specify the target CRS using its EPSG code, ensuring consistency and compatibility across geospatial datasets.

**Vector Density (10.7):**

The Vector Density function computes the density of geometries within a raster grid. Users can define parameters such as resolution, base raster profile, buffer value, and statistic (e.g., density).

## 6. Initial testing

### 6.1. Listing of initial testers

The initial testing is conducted by a small group of voluntary testers within the project who have not been actively involved in the development of “EIS Toolkit”. The initial testers are:

- Ina Storch (Beak),
- Roberto De La Rosa (Beak),
- Alex Vella (BRGM),
- Martiya Sadeghi (SGU),
- Patrick Casey (SGU),
- Sinovuyo Busakwe (UTU) and
- Gonzalo Ares (CSIC).

Testing is supported by other project members, primarily developers of “EIS Toolkit”.

### 6.2. Listing and description of initial test data for testers

In the first stage of testing, test data consists of two parts: small (< 1 MB) geospatial datasets uploaded in the GitHub repository and a set of IOCG data uploaded in cloud. The small datasets found in GitHub are meant only to test that the tools work without crashes and errors and cannot be used to perform a full mineral prospectivity mapping workflow. In contrast, the IOCG dataset was made available to perform full test analyses from beginning to end using “EIS Toolkit”, and to enable identifying tools needing performance optimization.

The IOCG dataset consists of the following files:

- **Deposits**
  - Deposit occurrences datasets
- **Geochemical data**
  - Geochemical dataset
  - Interpolated copper concentration dataset
  - Interpolated cobalt concentration dataset
  - Interpolated iron concentration dataset
- **Geological data**
  - Tectonic setting dataset
  - Structures dataset
  - Lithology dataset
  - Target area dataset
  - Distances to structures dataset
  - Tectonic setting dataset
- **Geophysical data**
  - AEM datasets

- Gravity datasets
- Magnetic dataset
- Radiometric dataset
- Discretized data (data prepared for Weights of evidence modeling)

## 6.3. Guidelines for testers to report bugs and problems

Development of “EIS Toolkit” is done in GitHub, and all initial testers with a GitHub account should use repository of “EIS Toolkit” to report found bugs and to send modification suggestions. The most practical way to do this is to use *GitHub issues*, which are widely used in programming projects for reporting bugs and development related tasks.

The “EIS Toolkit” GitHub page has basic instructions for developing, testing and using the toolkit, and the initial testers should refer to guidance found in there. In addition, a more detailed “Getting started with testing” document was sent to all initial testers.

In the initial testing phase, the possibility to report bugs and other findings via email was given to testers. However, this will not be possible in the future if the tester group expands.

## 7. Listing of tools that will be implemented in next 6 months until final release

The beta release of the “EIS Toolkit” includes most of the planned tools, but not all. Data preprocessing is almost complete, but some exploration and modeling functions are still being developed and refined.

The following planned tools are not included in the beta release:

- Surface derivatives (aspect, slope, curvature)
- Distance to anomaly
- Calculate vector geometry
- Extract shared lines
- Autoencoder
- Self-organizing map
- CoDa
- Training data sampler
- MLP
- CNN
- Random forests
- Model uncertainty estimation
- Model performance estimation
- Hyperparameter optimization.



## 8. Conclusion

The first beta release of the “EIS Toolkit” has been completed and the work continues with expanding and optimizing the toolkit’s functionality and with the implementation of “QGIS EIS Plugin”, which will include a graphical user-interface for conducting Mineral Prospectivity Mapping with QGIS.

Testing of the “EIS Toolkit” was already started before the beta release and revealed some challenges that require further analysing.

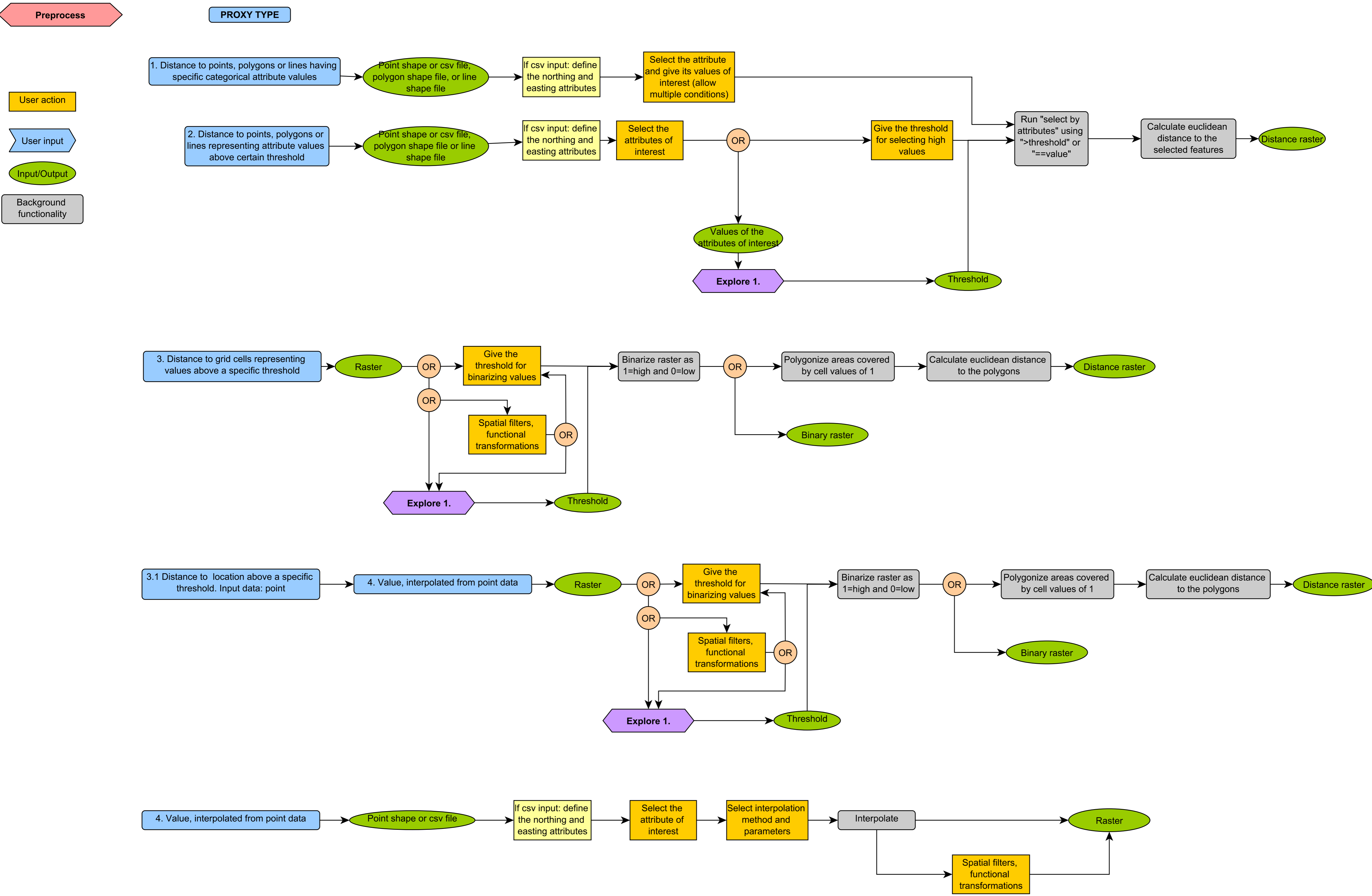
The role of testing will further increase after the first release of the “QGIS EIS Plugin”, and the plan is to gather a testing group and let them test the plugin individually.

After the testing phase, reported issues will be addressed before the final release, by the end of April 2024 – with the submission of Deliverable D3.5.

## 9. References

Nykänen, V., Lahti, I., Niiranen, T., & Korhonen, K. (2015). Receiver operating characteristics (ROC) as validation tool for prospectivity models — A magmatic Ni–Cu case study from the Central Lapland Greenstone Belt, Northern Finland. *Ore Geology Reviews*, 71, 853-860. doi:10.1016/j.oregeorev.2014.09.007

# Appendix 1: Main workflows for proxy processing



## Appendix 2: EIS proxy table for IOCG vectors

Mineral system Component	Mappable ingredient	Type of data	Format	Proxy	Processing workflow
Source	Albitized country rocks (Metal and ligand source)	Geological map	Polygon	Distance to albitized country rocks	1
		Mineralogy (outcrop observation data)	1) Point 2) Polygon		1
		Mineralogy (drilling data)	Point		1
	felsic (meta)volcanic rocks and subvolcanic rocks (Metal and fluid and Heat source);	Geological map	Polygon	Distance to felsic (meta)volcanic rocks and subvolcanic rocks	1
	Hydrothermal alteration of metavolcanic and granitoid rocks (Metal and ligand Source)	(1) Lithogeochemistry (2) Bottom of till geochemistry	Raster	Distance to high concentrations of Na, K or Mg±Fe	4 + 3
			Point		2
		EM	Raster	Distance to high conductivity anomalies	3
	Fe-rich rocks (magmatic hydrothermal fluids source and metal source)	(1) Lithogeochemistry	Point/ Raster	Distance to high Fe concentrations	2/ 4+3
			Raster Point		distance to Fe-oxides mapped from high magnetic anomalies
		(3) Gravity Data	Raster Point	distance to Fe-oxides mapped from high density anomalies	1 or 3, but processing is needed to get from density to the Fe oxides, if not directly proportional
		(4) Geological Map	Polygon	Distance to Fe-oxides (includes both magnetite and hematite minerals; or any other Fe-oxides related mineral of the interest of the users)	1
		(5) Mineralogy	Point		1
	Magmatic intrusions (Magmatic hydrothermal fluids, heat sources)	(1) Geological map	Polygon	Distance to magmatic intrusion of the relevant age (age from geology and radiometric data; distance from geophysical data)	1 or 3, if the polygons have the age information already
		(2) Radiometric ages	Point		
	Alkaline intrusions (Fluid Source)	(1) Lithogeochemistry	Point	Distance to rock units displaying alkaline magma signature	1&3 attribute value for alkaline magma signature needed
		(2) Geology	Polygon		1&3 attribute value for alkaline magma signature needed
		(3) Isotope composition	Point		1&3 attribute value for alkaline magma signature needed
	Metamorphic facies	1) Geology (metamorphic rock distributions) 2) Metamorphic indicator mineral	Polygon Point	Greenschist to mid-Amphibolite facies metamorphosed host rocks	1
	Tectonic setting (everything)	1) Geology 2) lithogeochemistry (lithotectonic rock associations)	Polygon Point	Syn- to late-orogenic back arc closure related intrusions.	1
	Magmatic rocks (Hydrothermal fluids Source)	Lithogeochemistry	Point	Distance to high Fe3O4 concentration	3.1
Co-REE anomaly (Metal Source)	Lithogeochemistry	Point	Distance to high concentrations of Co-REEs	3.1	
	mineralogy	Point			
I and A-type felsic igneous rocks; (Fluid and energy Source)	Lithogeochemistry	Point	Distance to magmatism contemporary with mineralisation (within 5 km buffer)	1 calculate buffer	
	Geology	1) Polygon 2) Lines			
Major structures	Geology	Lines	Distance to structures	1	
	Gravity worms	Lines	Distance to high order structures, interpreted from gravity worms	1	
Lower order structures	Geology	Lines	Distance to lower order structures	1	
	Gravity worms	Lines	Distance to lower order structures, interpreted from pseudo gravity worms	1	

Active/structural Pathways	Fluid flow (fault network) and deposition	Structure map	Poly line	Distance to kinks, offsets and breaks (interpretation from magnetic data, e.g. gradient)	1
	Fluid flow (fault network) and deposition	Structure map	Poly line	Distance to high density structure (interpretation of high density structures derived from geological map or airborne VLF)	2 attribute needed: density of structures
	Regional scale structure, crustal domain boundary; 75 km buffer	Structural Geology	Poly line	Distance to crustal boundary (interpretation from Solid geology, seismic reflection, passive seismic, gravity, magnetic data)	1
	Local scale structure; 2.5 km buffer	Geology	Shape file	Distance to intensity of structure	2
Depositional processes	Na-Ca Alteration	Bedrock & drill hole observations	Point	Distance to altered rocks	2
	K-alteration	Bedrock & drill hole observations	Point	Distance to altered rocks	2
		Till geochemistry	Point/raster	K-anomalies in till	3.1
	FeOx alteration	Bedrock & drill hole observations	Point	Distance to altered rocks	2
		Geophysics Mag	Raster Polygons	Magnetic highs (Or: distance to Fe-oxides mapped from high magnetic anomalies)	3.1 or 1
		Till geochemistry	Point Raster	Elevated V in till	2 3.1
	Geochemical signal	Till geochemistry	Point/raster	Elevated Ba, Co, Cu, P, Au in till Elevated Au, Ba, Co, Cu, La, P in till	3.1 or 3
		Radiometric U, Th	Raster	Radiometric U and/or Th anomalies	3
	Element deposition/mineralisation, marker horizon	Airborne Magnetic data	Raster	Distance to high-magnetic field anomalies	3
	Chemical trap, replacement	Geological map	Polygons (shp)	Presence of carbonate horizons in metavolcanic rocks	1
	Elemental deposition/mineralisation	Litho geochemistry	Csv (table)	Distance to high concentrations of REE, Fe	3.1
			Point		
	Magmatic crystallization	Geological map	Polygons (shp)	Massive magnetite	1
				Massive magnetite and vonsonite	
	Trap	Litho geochemistry	CSV	Ca-Fe metasomatism, Distance high Ca-Fe molar proportions relative to Na, K, and Mg	2 attribute needed: ratio of (Ca+Fe) to (Na + K + Mg)
Points (shp)					
Mineralization footprint	(1) Geology (Alteration minerals, e.g.; Albite) (2) Mineralogy (3) Litho geochemistry	Polygon	Distance to rock units displaying - sodic alteration signature - potassic and/or calcic alteration signature	1	
		Polygon			
		Point			
		Point			
		Polygons			
(1) Litho geochemistry (2) Petrography/Geology (3) Stable isotope ratios	Point	Low Ti signature in Fe-oxides	3.1		
	Point				
	Polygons				
(1) geology (2) Litho geochemistry	Points	Distance to (Sub) economic Cu-Au concentrations	1 3.1		
	Point				
Mineralisation, remobilisation, modification	Remobilisation	structural geology data	Polylines (shp)	Distance to kinks, offsets and breaks picked from magnetic data	2
	Remobilisation	Local-scale fault structure (large faults, high strain intensity)	Polylines (shp)	Distance to structures of high density/intensity Distance to structures	2 attribute for fault size, shear zone, strain, intensity needed
		Local-scale fault structure (shearzones with strain measured)			
Till geochemical anomalies	Till geochemistry	Csv (table) Points (shp)	Distance to elevated REE, Fe, Cu, Co, Au, Bi, Mo anomalies	3.1 repeat for each element	

## Appendix 3: EIS Toolkit – Workflow Demonstration



This Jupyter notebook serves as a first user guide to beta testers of EIS Toolkit. It shows how to import and call various tools of EIS Toolkit and includes steps of a simple MPM workflow.

Note that to run this notebook without modifications you need to have the test data (and under correct folder structure). One can always change the filepaths and use their own data, granted that the data is in right the format and is meaningful geospatial data. However, the primary function of this notebook is not to provide a template for conducting MPM workflows, but instead serve as example and showcase some of the tools of EIS Toolkit.

## 1. Imports and filepath definitions

```
In [2]: import rasterio
import geopandas as gpd
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from rasterio.io import MemoryFile
from rasterio.plot import show
import tempfile

import sys
sys.path.insert(0, "..")

from eis_toolkit.exploratory_analyses.basic_plots_seaborn import histogram,
from eis_toolkit.exploratory_analyses.descriptive_statistics import descript
from eis_toolkit.exploratory_analyses.pca import compute_pca

from eis_toolkit.prediction.fuzzy_overlay import gamma_overlay
from eis_toolkit.prediction.weights_of_evidence import weights_of_evidence_c

from eis_toolkit.raster_processing.extract_values_from_raster import extract
from eis_toolkit.raster_processing.unifying import unify_raster_grids

from eis_toolkit.transformations.sigmoid import _sigmoid_transform
from eis_toolkit.transformations.linear import _min_max_scaling

from eis_toolkit.vector_processing.idw_interpolation import idw
from eis_toolkit.vector_processing.distance_computation import distance_comp
from eis_toolkit.vector_processing.rasterize_vector import rasterize_vector
```

```
In [3]: # Filepaths
AEM_inphase_fp = "../tests/data/local/workflow_demo/IOCG_AEM_Inph_.tif"
AEM_quad_fp = "../tests/data/local/workflow_demo/IOCG_AEM_Quad.tif"
AEM_ratio_fp = "../tests/data/local/workflow_demo/IOCG_EM_ratio.tif"
Magn_AS_fp = "../tests/data/local/workflow_demo/Mag_DGRF_AS_FFT_ers_PCS_tif_

till_geochem_fp = "../tests/data/local/workflow_demo/Geochemical_Data/Vector
structures_fp = "../tests/data/local/workflow_demo/Geological_Data/IOCG_CLB_
```

```
lithology_fp = "../tests/data/local/workflow_demo/Geological_Data/IOCG_CLB_L  
known_occurrences_fp = "../tests/data/local/workflow_demo/Deposits_Occurrence
```

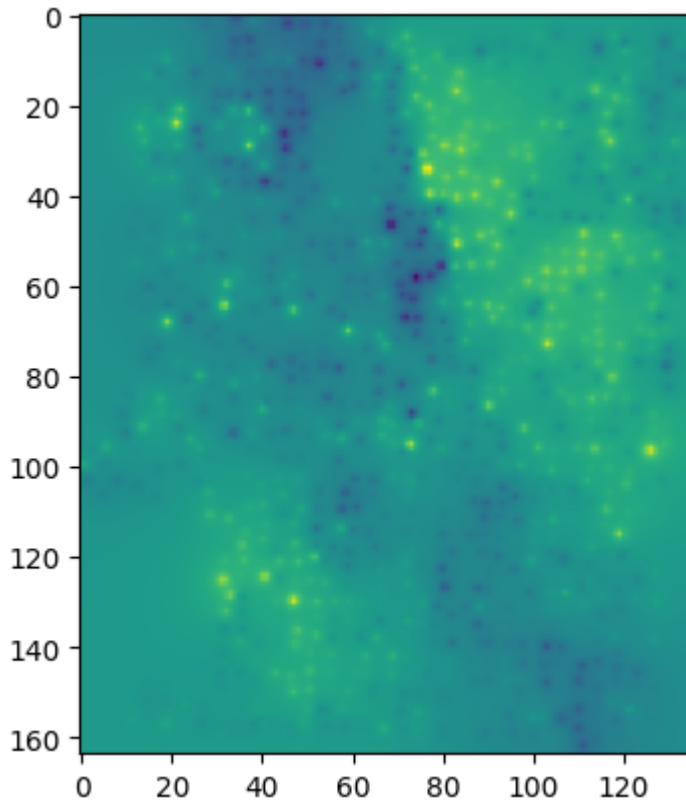
## 2. Preprocess data

```
In [4]: # Take AEM_inphase raster profile and grid shape as the base/target  
with rasterio.open(AEM_inphase_fp) as AEM_inphase:  
    raster_profile = AEM_inphase.profile  
    raster_extent = AEM_inphase.bounds  
    raster_pixel_size = AEM_inphase.transform[0]
```

```
In [ ]: # Preprocess geochemical data  
till_geochem = gpd.read_file(till_geochem_fp)  
  
columns_to_process = ["Fe_ppm_511", "Co_ppm_511", "Cu_ppm_511"]  
  
interpolation_results = {}  
interpolation_extent = (raster_extent[0], raster_extent[2], raster_extent[1])  
  
# Log transform concentrations and interpolate  
# NOTE: EIS Toolkit implements transforms only for rasters at the moment, th  
for column in columns_to_process:  
    new_col_name = column[:2] + "_log"  
    till_geochem[new_col_name] = np.log(till_geochem[column])  
  
    interpolation_results[new_col_name] = idw(till_geochem, new_col_name, (5
```

```
In [6]: # Visualize copper interpolated  
plt.imshow(interpolation_results["Cu_log"][0])
```

```
Out[6]: <matplotlib.image.AxesImage at 0x7f54f25a69e0>
```

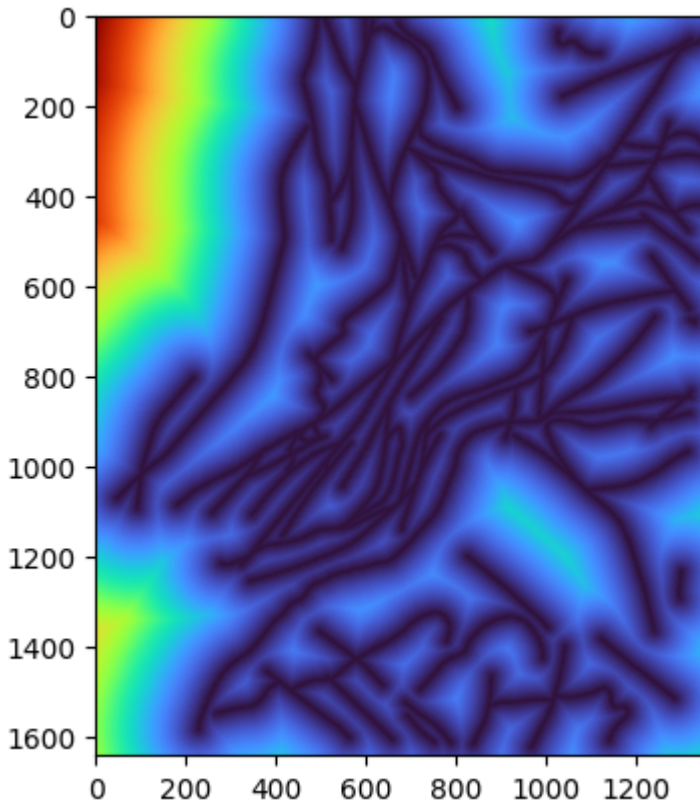


```
In [60]: # Preprocess geological data
structures = gpd.read_file(structures_fp)

distances_to_structures = distance_computation(raster_profile, structures)
```

```
In [7]: # Visualize distances to structures
plt.imshow(distances_to_structures, cmap="turbo")
```

```
Out[7]: <matplotlib.image.AxesImage at 0x7f12b8324400>
```



```
In [8]: # Rasterize lithology
lithology = gpd.read_file(lithology_fp)
rasterized_lithology, lithology_out_meta = rasterize_vector(lithology, base_
```

```
In [9]: # Extract raster values at deposits
known_occurrences = gpd.read_file(known_occurrences_fp)

# Open rasters with context manager, save lithology data to temporary raster
with \
    MemoryFile() as memfile, \
    rasterio.open(AEM_inphase_fp) as AEM_inphase, \
    rasterio.open(AEM_quad_fp) as AEM_quad, \
    rasterio.open(AEM_ratio_fp) as AEM_ratio, \
    rasterio.open(Magn_AS_fp) as Magn_AS:

    with tempfile.NamedTemporaryFile() as tmpfile:
        with rasterio.open(tmpfile.name, "w", **raster_profile) as dest:
            dest.write(rasterized_lithology, 1)
        with rasterio.open(tmpfile.name) as lithology_raster:
            raster_list = [AEM_inphase, AEM_quad, AEM_ratio, Magn_AS, lithol
            col_names = ["AEM_inphase", "AEM_quad", "AEM_ratio", "Magn_AS",
            raster_values_at_deposits = extract_values_from_raster(raster_li

            unified_rasters = unify_raster_grids(AEM_inphase, raster_list[1:

raster_values_at_deposits
```

Out[9]:

	x	y	AEM_inphase	AEM_quad	AEM_ratio	Magn_AS	Litr
0	370953.0	7502821.0	-408.969513	369.960022	-1.059486	17.158161	
1	371340.0	7496345.0	3889.745605	2384.806152	1.817769	21.545200	
2	385205.0	7497926.0	-339.839050	187.980026	-1.730280	14.112287	
3	369187.0	7490345.0	1468.369385	1285.413208	1.034502	20.598907	
4	368175.0	7489722.0	537.748718	324.102142	1.903474	43.523659	
5	350043.0	7489418.0	-6.511929	625.734375	0.025980	72.489517	
6	363838.0	7492354.0	-435.950928	248.234833	-1.672447	8.311082	
7	364672.0	7482684.0	-210.937698	1754.664917	-0.309453	31.408430	
8	396649.0	7501845.0	-373.599243	440.641846	-1.192332	12.663096	
9	372344.0	7492293.0	-885.269836	675.116272	-2.274024	8.680225	
10	371994.0	7491258.0	619.958496	578.203674	1.160073	25.673422	
11	368240.0	7487137.0	-259.313599	617.076233	-0.381119	58.577671	
12	373779.0	7504766.0	27.086342	577.757690	-0.123563	5.490668	

### 3. Explore data

#### Explore rasters

```
In [10]: # Calculate basic statistics of rasters at deposit locations
inphase_stats_at_points = descriptive_statistics_dataframe(raster_values_at_
quad_stats_at_points = descriptive_statistics_dataframe(raster_values_at_dep
ratio_stats_at_points = descriptive_statistics_dataframe(raster_values_at_de
magn_stats_at_points = descriptive_statistics_dataframe(raster_values_at_dep

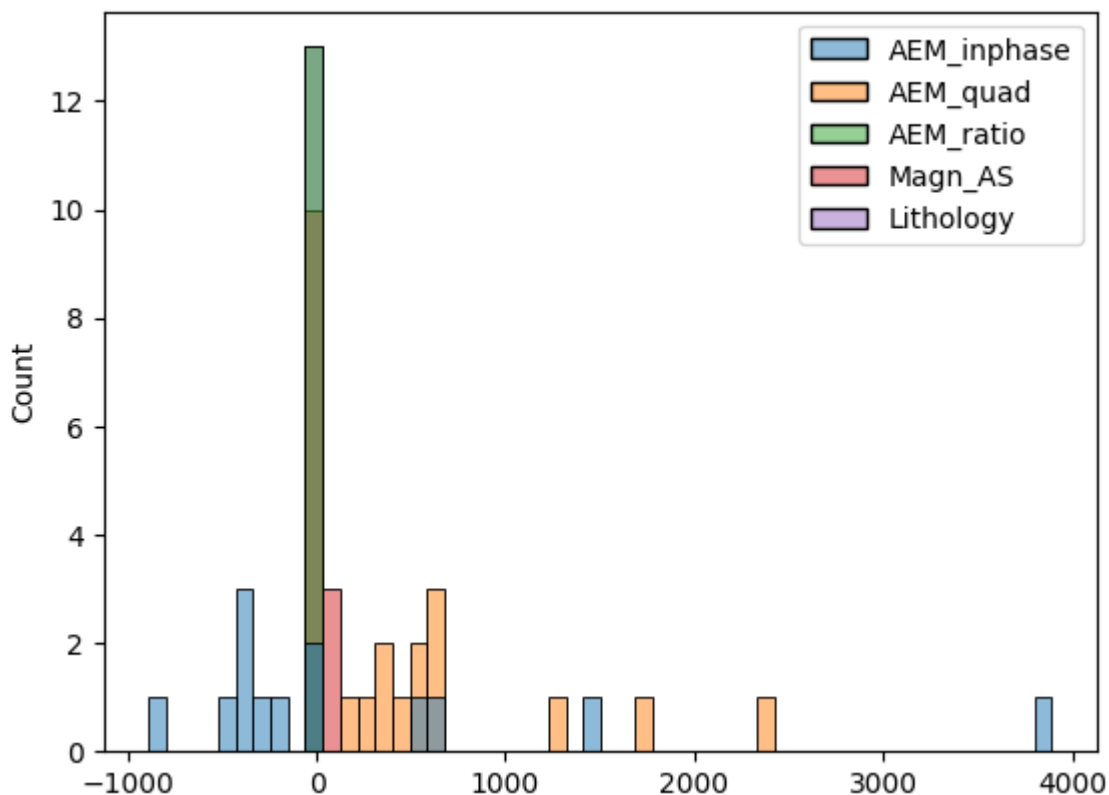
stats_dataframe = pd.DataFrame(
    {
        'inphase':pd.Series(inphase_stats_at_points),
        'quad':pd.Series(quad_stats_at_points),
        'ratio':pd.Series(ratio_stats_at_points),
        'magn':pd.Series(magn_stats_at_points),
    }
)
stats_dataframe
```

Out[10]:

	<b>inphase</b>	<b>quad</b>	<b>ratio</b>	<b>magn</b>
<b>min</b>	-885.269836	187.980026	-2.274024	5.490668
<b>max</b>	3889.745605	2384.806152	1.903474	72.489517
<b>mean</b>	278.655135	774.591645	-0.215454	26.171717
<b>25%</b>	-373.599243	369.960022	-1.192332	12.663096
<b>50%</b>	-210.937698	578.203674	-0.309453	20.598907
<b>75%</b>	537.748718	675.116272	1.034502	31.408430
<b>standard_deviation</b>	1191.602431	623.131335	1.315682	19.687020
<b>relative_standard_deviation</b>	4.276262	0.804464	-6.106545	0.752225
<b>skew</b>	2.079341	1.489026	0.191099	1.145725

```
In [11]: # Plot histogram of raster values at deposit locations
histogram(raster_values_at_deposits.drop(["x", "y"], axis=1))
```

Out[11]: <Axes: ylabel='Count'>

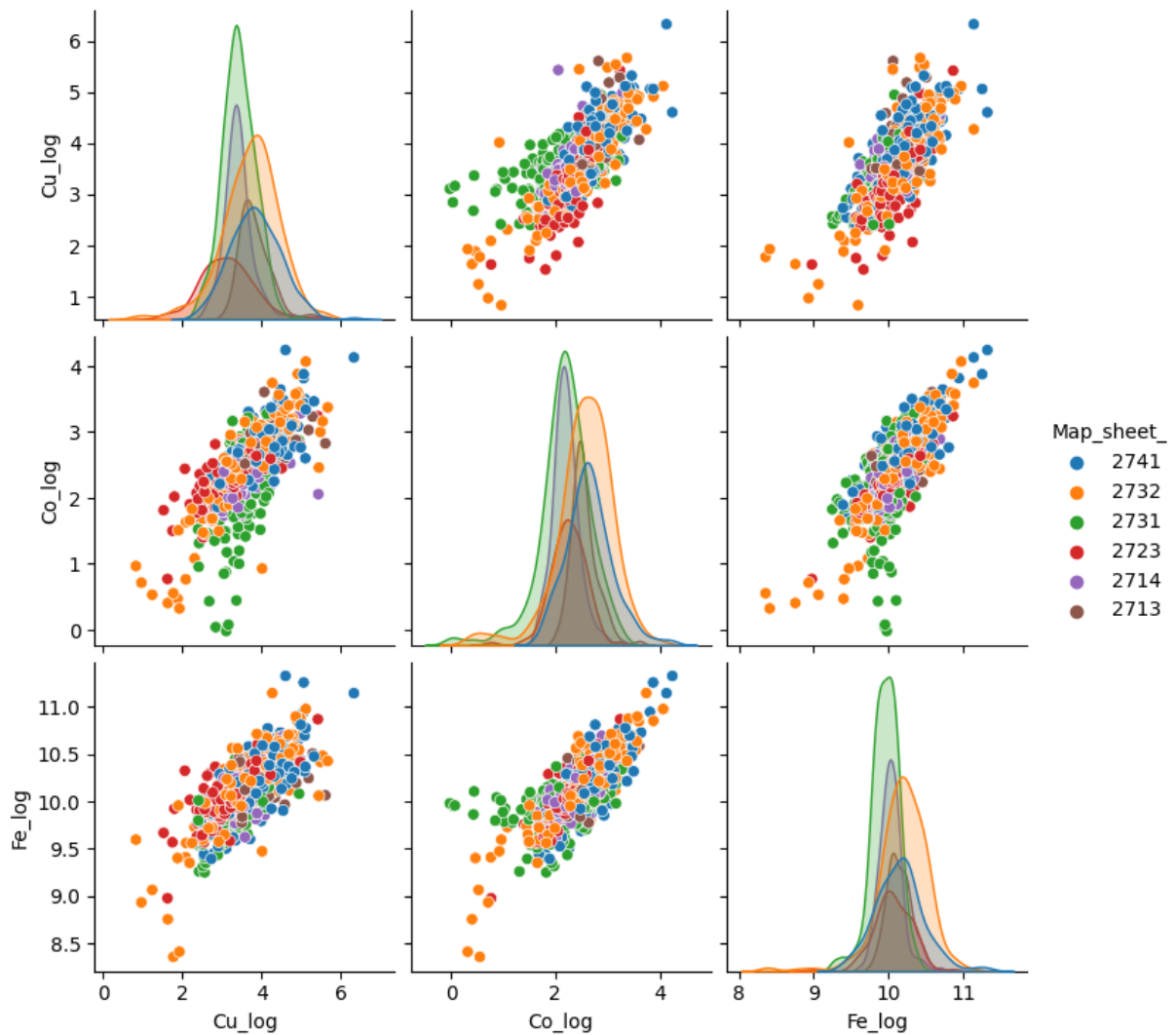


Explore geochemical data (vector data)

```
In [23]: # Plot pairplot of various log transformed concentrations
pairplot(till_geochem[["Cu_log", "Co_log", "Fe_log", "Map_sheet_"]], hue="Ma
```

```
/home/niko/.local/lib/python3.10/site-packages/seaborn/axisgrid.py:118: UserWarning: The figure layout has changed to tight
self._figure.tight_layout(*args, **kwargs)
```

```
Out[23]: <seaborn.axisgrid.PairGrid at 0x7f54e86b0820>
```



```
In [42]: # Define cmap to use in the following visualizations
cmap = plt.get_cmap('jet')
```

## Explore raster data with PCA

```
In [161].. # Preprocess nodata before modeling (for these rasters, min value is the nodata)
arrays_to_stack = []

for raster_array, meta in unified_rasters:
    raster_array[raster_array == np.nanmin(raster_array)] = np.nan
    arrays_to_stack.append(raster_array[0])

stacked_arrays = np.stack(arrays_to_stack)
```

```
In [170].. # Compute PCA for the input rasters
out_array, explained_variances = compute_pca(stacked_arrays, 3)
explained_variances
```

Out[170... array([0.41700809, 0.32194577, 0.16958726])

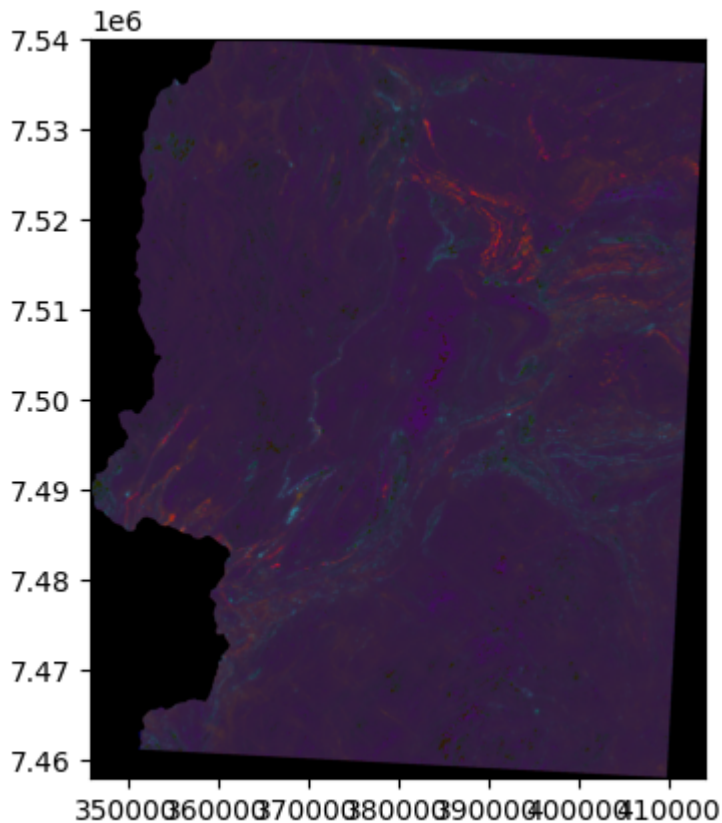
```
In [168... # Visualize PCA outputs

# Scale each band
scaled_bands = []
for band in range(out_array.shape[0]):
    scaled_band = _min_max_scaling(out_array[band], (0, 255))
    scaled_bands.append(scaled_band)

# Stack scaled bands back together
scaled_pca_output = np.stack(scaled_bands)

# Display the RGB image
show(scaled_pca_output.astype(np.uint8), transform=raster_profile["transform
```

```
/tmp/ipykernel_83370/1271570867.py:13: RuntimeWarning: invalid value encountered in cast
  show(scaled_pca_output.astype(np.uint8), transform=raster_profile["transform"])
```



Out[168... <Axes: >

## 4. Fuzzy logic modeling

```
In [53]: # Transform data before fuzzy overlay
arrays_to_stack = []
for raster_array, meta in unified_rasters:
    raster_array[raster_array == np.nanmin(raster_array)] = np.nan
```



```
out_array = _sigmoid_transform(raster_array, (0, 1), 1, True)
arrays_to_stack.append(out_array[0])
```

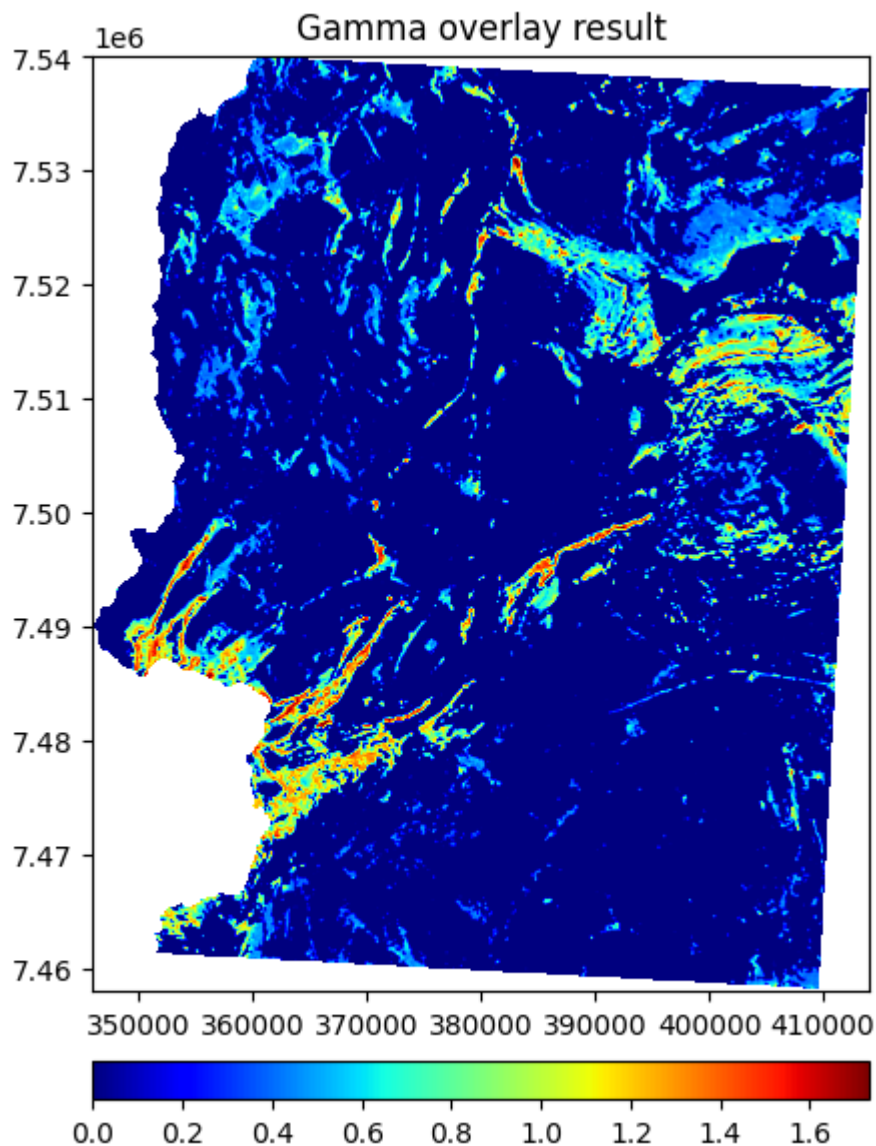
```
transformed_arrays_for_fuzzy_overlay = np.stack(arrays_to_stack)
```

```
In [54]: # Compute gamma overlay
overlay_result = gamma_overlay(transformed_arrays_for_fuzzy_overlay, 0.5)
```

```
In [55]: # Plot gamma overlay result
fig, ax = plt.subplots(1, 1, figsize = (5, 9))

ax.set_title("Gamma overlay result")
clrbar = plt.imshow(overlay_result, cmap=cmap)
plt.colorbar(clrbar, orientation="horizontal", pad = 0.05)
show(overlay_result, ax = ax, transform = raster_profile["transform"], cmap=
```

```
Out[55]: <Axes: title={'center': 'Gamma overlay result'}>
```



## 5. Weights of evidence modeling

```

In [56]: # Calculate weights
with rasterio.open("../tests/data/remote/wofe/wofe_evidence_raster.tif") as
    deposits = gpd.read_file("../tests/data/remote/wofe/wofe_deposits.shp")

weights_desc, arrays_desc, out_meta, deposit_pixels, evidence_pixels = w
    evidential_raster=evidence_raster,
    deposits=deposits,
    weights_type='descending',
    studentized_contrast_threshold=1
)

weights_asc, arrays_asc, _, _, _ = weights_of_evidence_calculate_weights
    evidential_raster=evidence_raster,
    deposits=deposits,
    weights_type='ascending',
    studentized_contrast_threshold=1
)

```

```

In [57]: # Plot descending weights
fig, axs = plt.subplots(3, 2, figsize = (14, 20))

axs[0, 0].set_title("Descending weights - Class")
clrbar = axs[0, 0].imshow(arrays_desc["Class"], cmap=cmap)
plt.colorbar(clrbar, orientation="horizontal", pad = 0.05)
show(arrays_desc["Class"], ax = axs[0, 0], transform = out_meta["transform"])

axs[1, 0].set_title("Descending weights - W+")
clrbar = axs[1, 0].imshow(arrays_desc["W+"], cmap=cmap)
plt.colorbar(clrbar, orientation="horizontal", pad = 0.05)
show(arrays_desc["W+"], ax = axs[1, 0], transform = out_meta["transform"], c

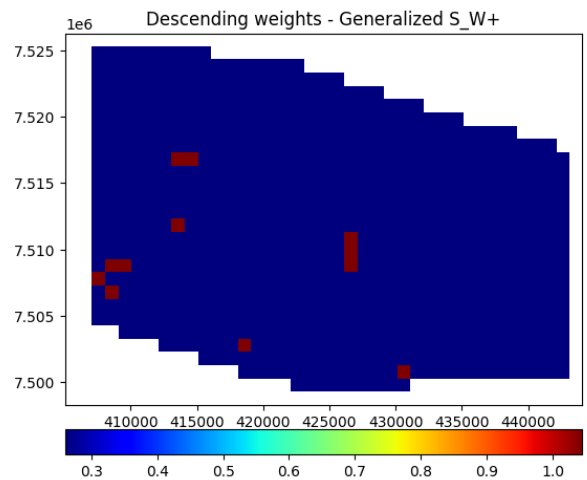
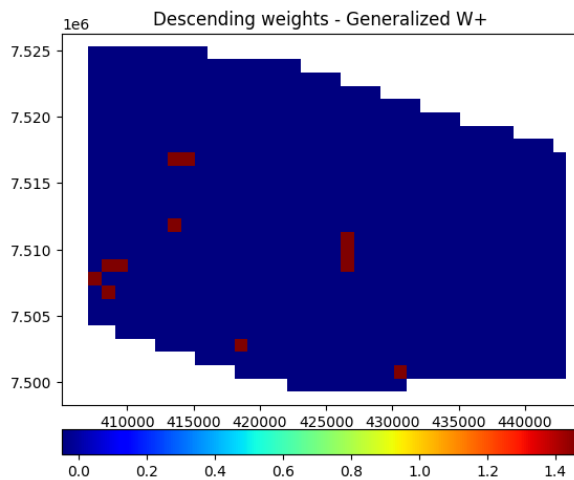
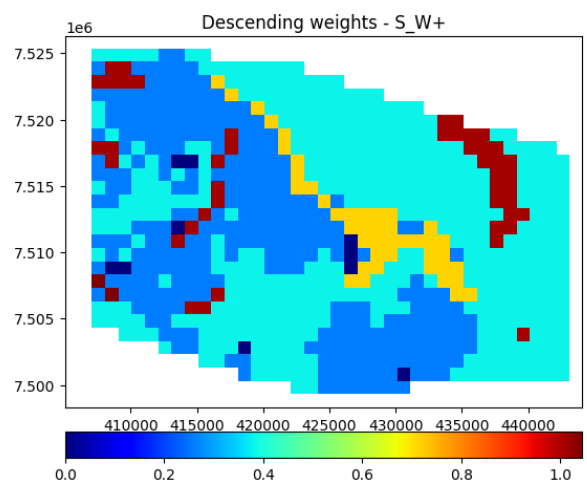
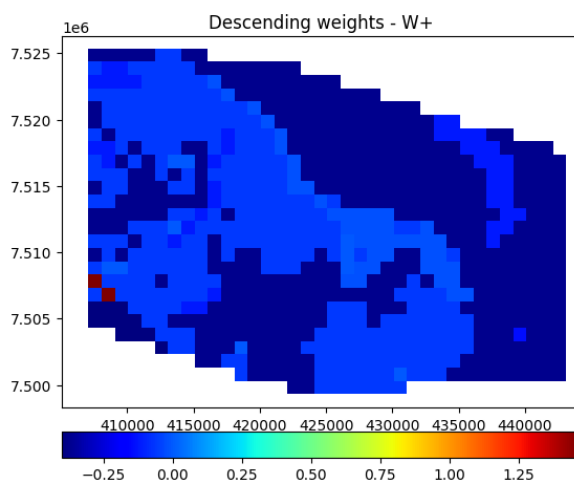
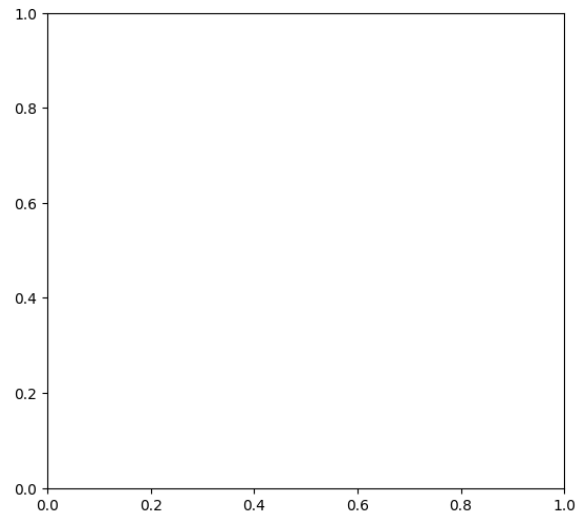
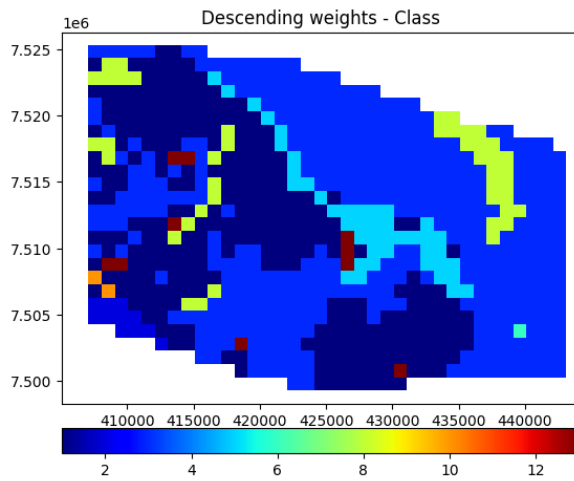
axs[1, 1].set_title("Descending weights - S_W+")
clrbar = axs[1, 1].imshow(arrays_desc["S_W+"], cmap=cmap)
plt.colorbar(clrbar, orientation="horizontal", pad = 0.05)
show(arrays_desc["S_W+"], ax = axs[1, 1], transform = out_meta["transform"],

axs[2, 0].set_title("Descending weights - Generalized W+")
clrbar = axs[2, 0].imshow(arrays_desc["Generalized W+"], cmap=cmap)
plt.colorbar(clrbar, orientation="horizontal", pad = 0.05)
show(arrays_desc["Generalized W+"], ax = axs[2, 0], transform = out_meta["tr

axs[2, 1].set_title("Descending weights - Generalized S_W+")
clrbar = axs[2, 1].imshow(arrays_desc["Generalized S_W+"], cmap=cmap)
plt.colorbar(clrbar, orientation="horizontal", pad = 0.05)
show(arrays_desc["Generalized S_W+"], ax = axs[2, 1], transform = out_meta["

```

Out[57]: <Axes: title={'center': 'Descending weights - Generalized S\_W+'}>



```
In [58]: # Calculate posterior probabilities / responses
posterior_array, posterior_array_std, posterior_confidence = weights_of_evid
```

```
In [59]: # Plot posterior probabilities weights

fig, axs = plt.subplots(2, 2, figsize = (14, 14))
```

```

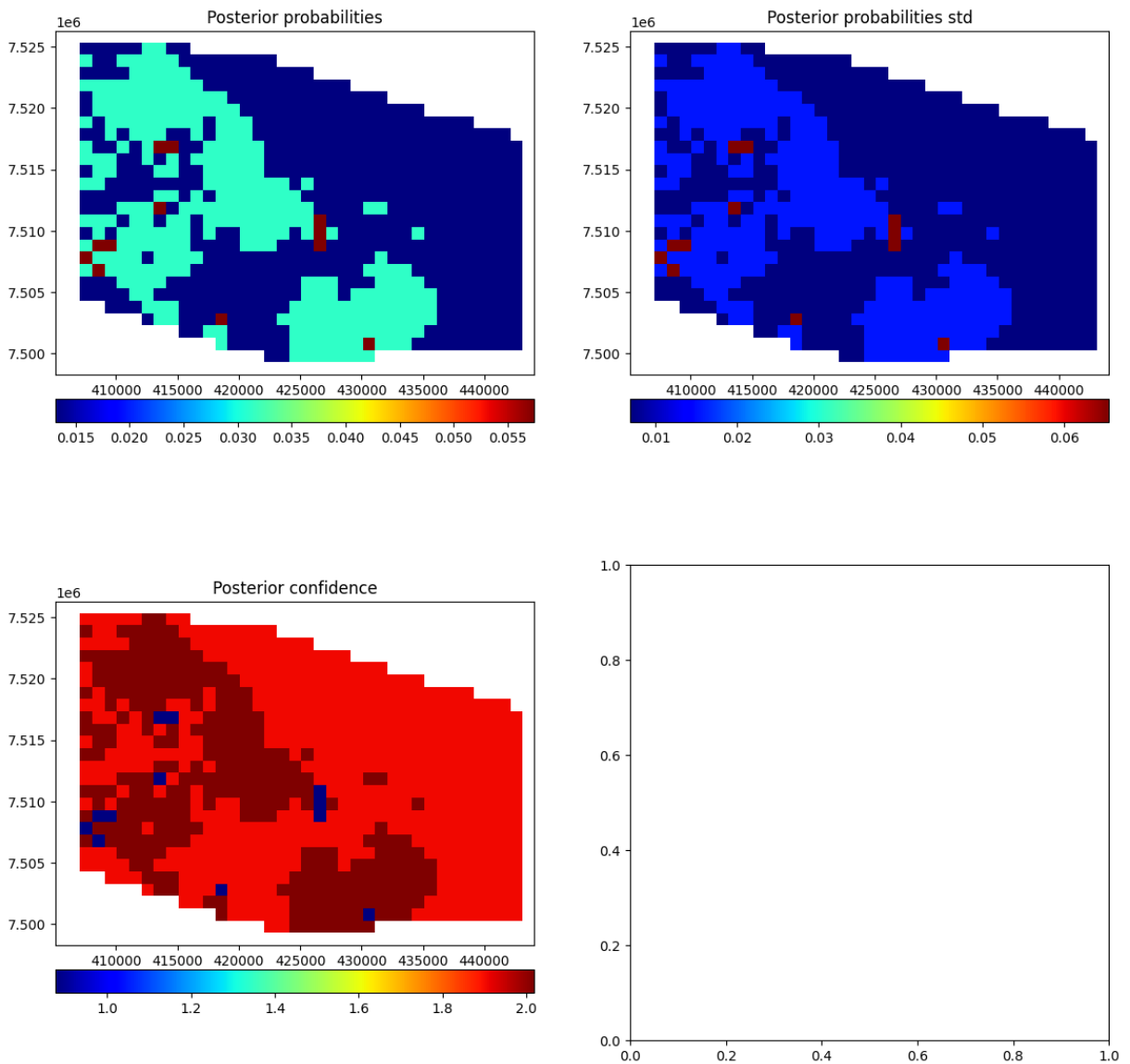
axs[0, 0].set_title("Posterior probabilities")
clrbar = axs[0, 0].imshow(posterior_array, cmap=cmap)
plt.colorbar(clrbar, orientation="horizontal", pad = 0.05)
show(posterior_array, ax = axs[0, 0], transform = out_meta["transform"], cma

axs[0, 1].set_title("Posterior probabilities std")
clrbar = axs[0, 1].imshow(posterior_array_std, cmap=cmap)
plt.colorbar(clrbar, orientation="horizontal", pad = 0.05)
show(posterior_array_std, ax = axs[0, 1], transform = out_meta["transform"],

axs[1, 0].set_title("Posterior confidence")
clrbar = axs[1, 0].imshow(posterior_confidence, cmap=cmap)
plt.colorbar(clrbar, orientation="horizontal", pad = 0.05)
show(posterior_confidence, ax = axs[1, 0], transform = out_meta["transform"])

```

Out[59]: <Axes: title={'center': 'Posterior confidence'}>



## Appendix 4: EIS Toolkit – Technical Specifications

# EIS Toolkit

---

**Technical specifications**

None

Made by <a href="mailto:info@gispo.fi">Gispo Ltd.</a>

# Table of contents

---

1. General	4
2. Dependency licenses	6
3. Conversions	10
3.1 Convert csv to geodataframe	10
3.2 Convert raster to dataframe	11
4. Exploratory analyses	12
4.1 DBSCAN	12
4.2 Descriptive statistics	14
4.3 Feature importance	16
4.4 K-means clustering	17
4.5 Plot parallel coordinates	18
4.6 PCA	20
4.7 Statistical (hypothesis) testing	29
5. Prediction	34
5.1 Fuzzy overlay	34
5.2 Weights of evidence	38
6. Raster processing	48
6.1 Check raster grids	48
6.2 Clipping	49
6.3 Create constant raster	50
6.4 Extract values from raster	52
6.5 Reprojecting	53
6.6 Resampling	54
6.7 Snapping	55
6.8 Unifying	56
6.9 Windowing	58
7. Training data tools	59
7.1 Class balancing	59
8. Transformations	61
8.1 Binarize	61
8.2 Clip	63
8.3 Linear	65
8.4 Logarithmic	69
8.5 Sigmoid	71
8.6 Winsorize	73

9. Validation	79
9.1 Calculate AUC	79
9.2 Calculate base metrics	80
9.3 Get P-A plot intersection point	82
9.4 Plot correlation matrix	84
9.5 Plot prediction-area (P-A) curves	86
9.6 Plot rate curve	88
10. Vector processing	90
10.1 Cell-Based Association	90
10.2 Distance computation	95
10.3 IDW	96
10.4 Kriging interpolation	98
10.5 Rasterize vector	101
10.6 Reproject vector	103
10.7 Vector density	104



# 1. General

---

This is the documentation site of the `eis_toolkit` python package. Here you can find documentation for each module. The documentation is automatically generated from docstrings.

Development of `eis_toolkit` is related to EIS Horizon EU project.



## 2. Dependency licenses

---

<b>Name</b>	<b>Version</b>	<b>License</b>
protobuf	3.19.4	3-Clause BSD License
tensorboard-plugin-wit	1.8.1	Apache 2.0
absl-py	1.2.0	Apache Software License
flatbuffers	1.12	Apache Software License
ghp-import	2.1.0	Apache Software License
google-auth	2.11.0	Apache Software License
google-auth-oauthlib	0.4.6	Apache Software License
google-pasta	0.2.0	Apache Software License
grpcio	1.48.1	Apache Software License
importlib-metadata	4.12.0	Apache Software License
keras	2.9.0	Apache Software License
libclang	14.0.6	Apache Software License
requests	2.28.1	Apache Software License
rsa	4.9	Apache Software License
tenacity	8.2.2	Apache Software License
tensorboard	2.9.1	Apache Software License
tensorboard-data-server	0.6.1	Apache Software License
tensorflow	2.9.2	Apache Software License
tensorflow-estimator	2.9.0	Apache Software License
tensorflow-io-gcs-filesystem	0.26.0	Apache Software License
watchdog	2.1.9	Apache Software License
packaging	21.3	Apache Software License; BSD License
python-dateutil	2.8.2	Apache Software License; BSD License
affine	2.3.1	BSD
cligj	0.7.2	BSD
geopandas	0.11.1	BSD
Fiona	1.8.21	BSD License
Jinja2	3.1.2	BSD License
Markdown	3.3.7	BSD License
MarkupSafe	2.1.1	BSD License
Pygments	2.13.0	BSD License
Shapely	1.8.4	BSD License
Werkzeug	2.2.2	BSD License
astunparse	1.6.3	BSD License
click	8.1.3	BSD License

<b>Name</b>	<b>Version</b>	<b>License</b>
click-plugins	1.1.1	BSD License
cycler	0.11.0	BSD License
gast	0.4.0	BSD License
h5py	3.7.0	BSD License
idna	3.3	BSD License
joblib	1.1.0	BSD License
kiwisolver	1.4.4	BSD License
mkdocs	1.3.1	BSD License
numpy	1.23.2	BSD License
oauthlib	3.2.0	BSD License
pandas	1.4.4	BSD License
patsy	0.5.2	BSD License
pyasn1	0.4.8	BSD License
pyasn1-modules	0.2.8	BSD License
rasterio	1.3.2	BSD License
requests-oauthlib	1.3.1	BSD License
scikit-learn	1.1.2	BSD License
scipy	1.9.1	BSD License
statsmodels	0.13.2	BSD License
threadpoolctl	3.1.0	BSD License
wrapt	1.14.1	BSD License
eis-toolkit	0.1.0	European Union Public Licence 1.2 (EUPL 1.2)
Pillow	9.2.0	Historical Permission Notice and Disclaimer (HPND)
opt-einsum	3.3.0	MIT
snuggs	1.4.7	MIT
GDAL	3.4.3	MIT License
Keras-Preprocessing	1.1.2	MIT License
PyYAML	6.0	MIT License
attrs	22.1.0	MIT License
cachetools	5.2.0	MIT License
charset-normalizer	2.1.1	MIT License
fonttools	4.37.1	MIT License
mergedeep	1.3.4	MIT License
mkdocs-material	8.4.2	MIT License
mkdocs-material-extensions	1.0.3	MIT License

<b>Name</b>	<b>Version</b>	<b>License</b>
munch	2.5.0	MIT License
plotly	5.14.0	MIT License
pydown-extensions	9.5	MIT License
pyparsing	3.0.9	MIT License
pyproj	3.3.1	MIT License
pytz	2022.2.1	MIT License
pyyaml_env_tag	0.1	MIT License
setuptools-scm	6.4.2	MIT License
six	1.16.0	MIT License
termcolor	1.1.0	MIT License
tomli	2.0.1	MIT License
urllib3	1.26.12	MIT License
zipp	3.8.1	MIT License
certifi	2022.6.15	Mozilla Public License 2.0 (MPL 2.0)
matplotlib	3.5.3	Python Software Foundation License
typing_extensions	4.3.0	Python Software Foundation License

## 3. Conversions

---

### 3.1 Convert csv to geodataframe

---

```
csv_to_geodataframe(csv, indexes, target_crs)
```

Read CSV file to a GeoDataFrame.

Usage of single index expects valid WKT geometry. Usage of two indexes expects POINT feature(s) X-coordinate as the first index and Y-coordinate as the second index.

#### Parameters:

Name	Type	Description	Default
csv	Path	Path to the .csv file to be read.	required
indexes	Sequence[int]	Index(es) of the geometry column(s).	required
target_crs	int	Target CRS as an EPSG code.	required

#### Returns:

Type	Description
GeoDataFrame	CSV file read to a GeoDataFrame.

#### Source code in `eis_toolkit/conversions/csv_to_geodataframe.py`

```

79 @beartype
80 def csv_to_geodataframe(
81     csv: Path,
82     indexes: Sequence[int],
83     target_crs: int,
84 ) -> geopandas.GeoDataFrame:
85     """
86     Read CSV file to a GeoDataFrame.
87
88     Usage of single index expects valid WKT geometry.
89     Usage of two indexes expects POINT feature(s) X-coordinate as the first index and Y-coordinate as the second index.
90
91     Args:
92         csv: Path to the .csv file to be read.
93         indexes: Index(es) of the geometry column(s).
94         target_crs: Target CRS as an EPSG code.
95
96     Returns:
97         CSV file read to a GeoDataFrame.
98     """
99
100     data_frame = _csv_to_geodataframe(
101         csv=csv,
102         indexes=indexes,
103         target_crs=target_crs,
104     )
105     return data_frame

```

## 3.2 Convert raster to dataframe

---

```
raster_to_dataframe(raster, bands=None, add_coordinates=False)
```

Convert raster to Pandas DataFrame.

If bands are not given, all bands are used for conversion. Selected bands are named based on their index e.g., band\_1, band\_2,...,band\_n. If wanted, image coordinates (row, col) for each pixel can be written to dataframe by setting add\_coordinates to True.

### Parameters:

Name	Type	Description	Default
raster	DatasetReader	Raster to be converted.	required
bands	Optional[Sequence[int]]	Selected bands from multiband raster. Indexing begins from one. Defaults to None.	None
add_coordinates	bool	Determines if pixel coordinates are written into dataframe. Defaults to False.	False

### Returns:

Type	Description
DataFrame	Raster converted to a DataFrame.

### Source code in eis\_toolkit/conversions/raster\_to\_dataframe.py

```

33 @beartype
34 def raster_to_dataframe(
35     raster: rasterio.io.DatasetReader,
36     bands: Optional[Sequence[int]] = None,
37     add_coordinates: bool = False,
38 ) -> pd.DataFrame:
39     """Convert raster to Pandas DataFrame.
40
41     If bands are not given, all bands are used for conversion. Selected bands are named based on their index e.g.,
42     band_1, band_2,...,band_n. If wanted, image coordinates (row, col) for each pixel can be written to
43     dataframe by setting add_coordinates to True.
44
45     Args:
46         raster: Raster to be converted.
47         bands: Selected bands from multiband raster. Indexing begins from one. Defaults to None.
48         add_coordinates: Determines if pixel coordinates are written into dataframe. Defaults to False.
49
50     Returns:
51         Raster converted to a DataFrame.
52     """
53
54     data_frame = _raster_to_dataframe(
55         raster=raster,
56         bands=bands,
57         add_coordinates=add_coordinates,
58     )
59     return data_frame

```



## 4. Exploratory analyses

---

### 4.1 DBSCAN

---

```
dbscan(data, max_distance=0.5, min_samples=5)
```

Perform DBSCAN clustering on the input data.

#### Parameters:


Name	Type	Description	Default
<code>data</code>	<code>GeoDataFrame</code>	GeoDataFrame containing the input data.	required
<code>max_distance</code>	<code>float</code>	The maximum distance between two samples for one to be considered as in the neighborhood of the other. Defaults to 0.5.	0.5
<code>min_samples</code>	<code>int</code>	The number of samples in a neighborhood for a point to be considered as a core point. Defaults to 5.	5

#### Returns:

Type	Description
<code>GeoDataFrame</code>	GeoDataFrame containing two new columns: one with assigned cluster labels and one indicating whether a point is a core point (1) or not (0).

#### Raises:

Type	Description
<code>EmptyDataFrameException</code>	The input GeoDataFrame is empty.
<code>InvalidParameterException</code>	The maximum distance between two samples in a neighborhood is not greater than zero or the number of samples in a neighborhood is not greater than one.

Source code in `eis_toolkit/exploratory_analyses/dbscan.py` 

```
25 @beartype
26 def dbscan(data: gdp.GeoDataFrame, max_distance: float = 0.5, min_samples: int = 5) -> gdp.GeoDataFrame:
27     """
28     Perform DBSCAN clustering on the input data.
29
30     Args:
31     data: GeoDataFrame containing the input data.
32     max_distance: The maximum distance between two samples for one to be considered as in the neighborhood of
33     the other. Defaults to 0.5.
34     min_samples: The number of samples in a neighborhood for a point to be considered as a core point.
35     Defaults to 5.
36
37     Returns:
38     GeoDataFrame containing two new columns: one with assigned cluster labels and one indicating whether a
39     point is a core point (1) or not (0).
40
41     Raises:
42     EmptyDataFrameException: The input GeoDataFrame is empty.
43     InvalidParameterException: The maximum distance between two samples in a neighborhood is not greater
44     than zero or the number of samples in a neighborhood is not greater than one.
45     """
46
47     if data.empty:
48         raise EmptyDataFrameException("The input GeoDataFrame is empty.")
49
50     if max_distance <= 0:
51         raise InvalidParameterValueException(
52             "The input value for the maximum distance between two samples in a neighborhood must be greater than zero."
53         )
54
55     if min_samples <= 1:
56         raise InvalidParameterValueException(
57             "The input value for the minimum number of samples in a neighborhood must be greater than one."
58         )
59
60     dbscan_gdf = _dbscan(data, max_distance, min_samples)
61
62     return dbscan_gdf
```

## 4.2 Descriptive statistics

---

### `descriptive_statistics_dataframe(input_data, column)`

Generate descriptive statistics from vector data.

Generates min, max, mean, quantiles(25%, 50% and 75%), standard deviation, relative standard deviation and skewness.

#### Parameters:

Name	Type	Description	Default
<code>input_data</code>	<code>Union[DataFrame, GeoDataFrame]</code>	Data to generate descriptive statistics from.	required
<code>column</code>	<code>str</code>	Specify the column to generate descriptive statistics from.	required

#### Returns:

Type	Description
<code>dict</code>	The descriptive statistics in previously described order.

#### Source code in `eis_toolkit/exploratory_analyses/descriptive_statistics.py`

```

42 @beartype
43 def descriptive_statistics_dataframe(input_data: Union[pd.DataFrame, gpd.GeoDataFrame], column: str) -> dict:
44     """Generate descriptive statistics from vector data.
45
46     Generates min, max, mean, quantiles(25%, 50% and 75%), standard deviation, relative standard deviation and skewness.
47
48     Args:
49         input_data: Data to generate descriptive statistics from.
50         column: Specify the column to generate descriptive statistics from.
51
52     Returns:
53         The descriptive statistics in previously described order.
54     """
55     if column not in input_data.columns:
56         raise InvalidColumnException
57     data = input_data[column]
58     statistics = _descriptive_statistics(data)
59     return statistics

```

### `descriptive_statistics_raster(input_data)`

Generate descriptive statistics from raster data.

Generates min, max, mean, quantiles(25%, 50% and 75%), standard deviation, relative standard deviation and skewness.

#### Parameters:

Name	Type	Description	Default
<code>input_data</code>	<code>DatasetReader</code>	Data to generate descriptive statistics from.	required

#### Returns:

Type	Description
<code>dict</code>	The descriptive statistics in previously described order.

Source code in [eis\\_toolkit/exploratory\\_analyses/descriptive\\_statistics.py](#) ✓

```
62 @beartype
63 def descriptive_statistics_raster(input_data: rasterio.io.DatasetReader) -> dict:
64     """Generate descriptive statistics from raster data.
65
66     Generates min, max, mean, quantiles(25%, 50% and 75%), standard deviation, relative standard deviation and skewness.
67
68     Args:
69         input_data: Data to generate descriptive statistics from.
70
71     Returns:
72         The descriptive statistics in previously described order.
73     """
74     data = input_data.read().flatten()
75     statistics = _descriptive_statistics(data)
76     return statistics
```

## 4.3 Feature importance

```
evaluate_feature_importance(classifier, x_test, y_test, feature_names,
                             number_of_repetition=50, random_state=0)
```

Evaluate the feature importance of a sklearn classifier or linear model.

### Parameters:

Name	Type	Description	Default
classifier	BaseEstimator	Trained classifier.	required
x_test	ndarray	Testing feature data (X data need to be normalized / standardized).	required
y_test	ndarray	Testing target data.	required
feature_names	Sequence[str]	Names of the feature columns.	required
number_of_repetition	int	Number of iteration used when calculate feature importance (default 50).	50
random_state	int	random state for repeatability of results (Default 0).	0

Return: A dataframe composed by features name and Importance value The resulted object with importance mean, importance std, and overall importance Raises: InvalidDatasetException: When the dataset is None.

### Source code in `eis_toolkit/exploratory_analyses/feature_importance.py`

```
12 @beartype
13 def evaluate_feature_importance(
14     classifier: sklearn.base.BaseEstimator,
15     x_test: np.ndarray,
16     y_test: np.ndarray,
17     feature_names: Sequence[str],
18     number_of_repetition: int = 50,
19     random_state: int = 0,
20 ) -> tuple[pd.DataFrame, dict]:
21     """
22     Evaluate the feature importance of a sklearn classifier or linear model.
23
24     Parameters:
25         classifier: Trained classifier.
26         x_test: Testing feature data (X data need to be normalized / standardized).
27         y_test: Testing target data.
28         feature_names: Names of the feature columns.
29         number_of_repetition: Number of iteration used when calculate feature importance (default 50).
30         random_state: random state for repeatability of results (Default 0).
31
32     Return:
33         A dataframe composed by features name and Importance value
34         The resulted object with importance mean, importance std, and overall importance
35
36     Raises:
37         InvalidDatasetException: When the dataset is None.
38     """
39
40     if x_test is None or y_test is None:
41         raise InvalidDatasetException
42
43     result = permutation_importance(
44         classifier, x_test, y_test.ravel(), n_repeats=number_of_repetition, random_state=random_state
45     )
46
47     feature_importance = pd.DataFrame({"Feature": feature_names, "Importance": result.importances_mean})
48
49     feature_importance["Importance"] = feature_importance["Importance"] * 100
50     feature_importance = feature_importance.sort_values(by="Importance", ascending=False)
51
52     return feature_importance, result
```

## 4.4 K-means clustering

---

```
k_means_clustering(data, number_of_clusters=None, random_state=None)
```

Perform k-means clustering on the input data.

### Parameters:

Name	Type	Description	Default
data	GeoDataFrame	A GeoDataFrame containing the input data.	required
number_of_clusters	Optional[int]	The number of clusters ( $\geq 1$ ) to form. Optional parameter. If not provided, optimal number of clusters is computed using the elbow method.	None
random_state	Optional[int]	A random number generation for centroid initialization to make the randomness deterministic. Optional parameter.	None

### Returns:

Type	Description
GeoDataFrame	GeoDataFrame containing assigned cluster labels.

### Raises:

Type	Description
EmptyDataFrameException	The input GeoDataFrame is empty.
InvalidParameterException	The number of clusters is less than one.

### Source code in `eis_toolkit/exploratory_analyses/k_means_cluster.py`

```

39 @beartype
40 def k_means_clustering(
41     data: gdp.GeoDataFrame, number_of_clusters: Optional[int] = None, random_state: Optional[int] = None
42 ) -> gdp.GeoDataFrame:
43     """
44     Perform k-means clustering on the input data.
45
46     Args:
47         data: A GeoDataFrame containing the input data.
48         number_of_clusters: The number of clusters ( $\geq 1$ ) to form. Optional parameter. If not provided,
49             optimal number of clusters is computed using the elbow method.
50         random_state: A random number generation for centroid initialization to make
51             the randomness deterministic. Optional parameter.
52
53     Returns:
54         GeoDataFrame containing assigned cluster labels.
55
56     Raises:
57         EmptyDataFrameException: The input GeoDataFrame is empty.
58         InvalidParameterException: The number of clusters is less than one.
59     """
60
61     if data.empty:
62         raise EmptyDataFrameException("The input GeoDataFrame is empty.")
63
64     if number_of_clusters is not None and number_of_clusters < 1:
65         raise InvalidParameterValueException("The input value for number of clusters must be at least one.")
66
67     k_means_gdf = _k_means_clustering(data, number_of_clusters, random_state)
68
69     return k_means_gdf

```

## 4.5 Plot parallel coordinates

---

```
plot_parallel_coordinates(df, color_column_name, plot_title=None,
palette_name=None, curved_lines=True)
```

Plot a parallel coordinates plot.

Automatically removes all rows containing null/nan values. Tries to convert columns to numeric to be able to plot them. If more than 8 columns are present (after numeric filtering), keeps only the first 8 to plot.

### Parameters:

Name	Type	Description	Default
df	DataFrame	The DataFrame to plot.	required
color_column_name	str	The name of the column in df to use for color encoding.	required
plot_title	Optional[str]	The title for the plot. Default is None.	None
palette_name	Optional[str]	The name of the color palette to use. Default is None.	None
curved_lines	bool	If True, the plot will have curved instead of straight lines. Default is True.	True

### Returns:

Type	Description
Figure	A matplotlib figure containing the parallel coordinates plot.

### Raises:

Type	Description
EmptyDataFrameException	Raised when the DataFrame is empty.
InvalidColumnException	Raised when the color column is not found in the DataFrame.
InconsistentDataTypesException	Raised when the color column has multiple data types.

Source code in `eis_toolkit/exploratory_analyses/parallel_coordinates.py` 

```

134 @beartype
135 def plot_parallel_coordinates(
136     df: pd.DataFrame,
137     color_column_name: str,
138     plot_title: Optional[str] = None,
139     palette_name: Optional[str] = None,
140     curved_lines: bool = True,
141 ) -> matplotlib.figure.Figure:
142     """Plot a parallel coordinates plot.
143
144     Automatically removes all rows containing null/nan values. Tries to convert columns to numeric
145     to be able to plot them. If more than 8 columns are present (after numeric filtering), keeps only
146     the first 8 to plot.
147
148     Args:
149         df: The DataFrame to plot.
150         color_column_name: The name of the column in df to use for color encoding.
151         plot_title: The title for the plot. Default is None.
152         palette_name: The name of the color palette to use. Default is None.
153         curved_lines: If True, the plot will have curved instead of straight lines. Default is True.
154
155     Returns:
156         A matplotlib figure containing the parallel coordinates plot.
157
158     Raises:
159         EmptyDataFrameException: Raised when the DataFrame is empty.
160         InvalidColumnException: Raised when the color column is not found in the DataFrame.
161         InconsistentDataTypesException: Raised when the color column has multiple data types.
162     """
163
164     if df.empty:
165         raise exceptions.EmptyDataFrameException("The input DataFrame is empty.")
166
167     if color_column_name not in df.columns:
168         raise exceptions.InvalidColumnException(
169             f"The provided color column {color_column_name} is not found in the DataFrame."
170         )
171
172     df = df.convert_dtypes()
173     df = df.apply(pd.to_numeric, errors="ignore")
174
175     color_data = df[color_column_name].to_numpy()
176     if len(set([type(elem) for elem in color_data])) != 1:
177         raise exceptions.InconsistentDataTypesException(
178             "The color column should have a consistent datatype. Multiple data types detected in the color column."
179         )
180
181     df = df.select_dtypes(include=np.number)
182
183     # Drop non-numeric columns and the column used for coloring
184     columns_to_drop = [color_column_name]
185     for column in df.columns.values:
186         if df[column].isnull().all():
187             columns_to_drop.append(column)
188     df = df.loc[:, ~df.columns.isin(columns_to_drop)]
189
190     # Keep only first 8 columns if more are still present
191     if len(df.columns.values) > 8:
192         df = df.iloc[:, :8]
193
194     data_labels = df.columns.values
195     data = df.to_numpy()
196
197     fig = _plot_parallel_coordinates(
198         data=data,
199         data_labels=data_labels,
200         color_data=color_data,
201         color_column_name=color_column_name,
202         plot_title=plot_title,
203         palette_name=palette_name,
204         curved_lines=curved_lines,
205     )
206     return fig

```



## 4.6 PCA

---

```
compute_pca(data, number_of_components, scaler_type='standard', nodata=None,  
color_column_name=None)
```

Compute given number of principal components for numeric input data.

Various input data formats are accepted and the output format depends on the input format. If input is (Geo)DataFrame, a pairplot is produced additionally. A column name used for coloring can be specified in this case.

**Parameters:**

<b>Name</b>	<b>Type</b>	<b>Description</b>	<b>Default</b>
<code>data</code>	<code>Union[ndarray, DataFrame, GeoDataFrame, DatasetReader]</code>	Input data for PCA.	required
<code>number_of_components</code>	<code>int</code>	The number of principal components to compute Should be $\geq 1$ and at most the number of numeric columns if input is (Geo)DataFrame or number of bands if input is raster.	required
<code>scaler_type</code>	<code>Literal['standard', 'min_max', 'robust']</code>	Transform data according to a specified Sklearn scaler. Options are "standard", "min_max" and "robust". Defaults to "standard".	'standard'
<code>nodata</code>	<code>Optional[Number]</code>	Define nodata value to be masked out. Optional parameter. If None and input is raster, looks for nodata value from raster metadata. Defaults to None.	None
<code>color_column_name</code>	<code>Optional[str]</code>	If input data is a DataFrame or a GeoDataFrame, column name used for coloring data points in the produced pairplot can be defined. Defaults to None.	None

**Returns:**

Type	Description
Union[ndarray, Tuple[DataFrame, PairGrid], Tuple[GeoDataFrame, PairGrid], Tuple[ndarray, Profile]]	The computed principal components in corresponding format as the input data (for raster, output is
ndarray	Numpy array containing the data and raster profile) and the explained variance ratios for each component.

**Raises:**

Type	Description
EmptyDataException	The input is empty.
InvalidNumberOfPrincipalComponents	The number of principal components is less than 1 or more than number of columns if input was (Geo)DataFrame.

Source code in [eis\\_toolkit/exploratory\\_analyses/pca.py](#) 

62  
63  
64  
65  
66  
67  
68  
69  
70  
71  
72  
73  
74  
75  
76  
77  
78  
79  
80  
81  
82  
83  
84  
85  
86  
87  
88  
89  
90  
91  
92  
93  
94  
95  
96  
97  
98  
99  
100  
101  
102  
103  
104  
105  
106  
107  
108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161  
162  
163  
164  
165  
166  
167  
168  
169

```

170 @beartype
171 def compute_pca(
172     data: Union[np.ndarray, pd.DataFrame, gpd.GeoDataFrame, rasterio.io.DatasetReader],
173     number_of_components: int,
174     scaler_type: Literal["standard", "min_max", "robust"] = "standard",
175     nodata: Optional[Number] = None,
176     color_column_name: Optional[str] = None,
177 ) -> Tuple[
178     Union[
179         np.ndarray,
180         Tuple[pd.DataFrame, sns.PairGrid],
181         Tuple[gpd.GeoDataFrame, sns.PairGrid],
182         Tuple[np.ndarray, rasterio.profiles.Profile],
183     ],
184     np.ndarray,
185 ]:
186     """
187     Compute given number of principal components for numeric input data.
188
189     Various input data formats are accepted and the output format depends on the input format. If
190     input is (Geo)DataFrame, a pairplot is produced additionally. A column name used for coloring can
191     be specified in this case.
192
193     Args:
194     data: Input data for PCA.
195     number_of_components: The number of principal components to compute Should be >= 1 and at most
196         the number of numeric columns if input is (Geo)DataFrame or number of bands if input is raster.
197     scaler_type: Transform data according to a specified Sklearn scaler.
198         Options are "standard", "min_max" and "robust". Defaults to "standard".
199     nodata: Define nodata value to be masked out. Optional parameter. If None and input is raster, looks
200         for nodata value from raster metadata. Defaults to None.
201     color_column_name: If input data is a DataFrame or a GeoDataFrame, column name used for
202         coloring data points in the produced pairplot can be defined. Defaults to None.
203
204     Returns:
205     The computed principal components in corresponding format as the input data (for raster, output is
206     Numpy array containing the data and raster profile) and the explained variance ratios for each component.
207
208     Raises:
209     EmptyDataException: The input is empty.
210     InvalidNumberOfPrincipalComponents: The number of principal components is less than 1 or more than
211         number of columns if input was (Geo)DataFrame.
212     """
213     if scaler_type not in SCALERS:
214         raise exceptions.InvalidParameterValueException(f"Invalid scaler. Choose from: {list(SCALERS.keys())}")
215
216     if number_of_components < 1:
217         raise exceptions.InvalidParameterValueException("The number of principal components should be >= 1.")
218
219     # Get feature matrix (Numpy array) from various input types
220     if isinstance(data, np.ndarray):
221         feature_matrix = data
222         if feature_matrix.ndim == 2: # Table-like data (assume it is a DataFrame transformed to Numpy array)
223             feature_matrix, nan_mask = _prepare_array_data(feature_matrix, nodata_value=nodata, reshape=False)
224         elif feature_matrix.ndim == 3: # Assume data represents multiband raster data
225             rows, cols = feature_matrix.shape[1], feature_matrix.shape[2]
226             feature_matrix, nan_mask = _prepare_array_data(feature_matrix, nodata_value=nodata, reshape=True)
227         else:
228             raise exceptions.InvalidParameterValueException(
229                 f"Unsupported input data format. {feature_matrix.ndim} dimensions detected."
230             )
231         if feature_matrix.size == 0:
232             raise exceptions.EmptyDataException("Input array is empty.")
233
234     elif isinstance(data, rasterio.io.DatasetReader):
235         feature_matrix = data.read()
236         if feature_matrix.ndim < 3:
237             raise exceptions.InvalidParameterValueException("Input raster should have multiple bands.")
238         rows, cols = feature_matrix.shape[1], feature_matrix.shape[2]
239         if nodata is None:
240             nodata = data.nodata
241         feature_matrix, nan_mask = _prepare_array_data(feature_matrix, nodata_value=nodata, reshape=True)
242
243     elif isinstance(data, pd.DataFrame):
244         df = data.copy()
245         if df.empty:
246             raise exceptions.EmptyDataException("Input DataFrame is empty.")
247         if number_of_components > len(df.columns):
248             raise exceptions.InvalidParameterValueException(
249                 "The number of principal should be at most the number of numeric columns in the input DataFrame."
250             )
251         if color_column_name is not None:
252             color_column_data = df[color_column_name]
253
254         if isinstance(data, gpd.GeoDataFrame):
255             geometries = data.geometry
256             crs = data.crs
257             df = df.drop(columns=["geometry"])
258
259         df = df.convert_dtypes()
260         df = df.apply(pd.to_numeric, errors="ignore")
261         df = df.select_dtypes(include=np.number)
262         df = df.astype(dtype=np.number)
263         feature_matrix = df.to_numpy()
264         feature_matrix = feature_matrix.astype(float)
265         feature_matrix, nan_mask = _handle_missing_values(feature_matrix, nodata)
266
267     # Core PCA computation
268     principal_components, explained_variances = _compute_pca(feature_matrix, number_of_components, scaler_type)
269
270     # Put nodata back in and consider new dimension of data
271     if nodata is not None:
272         principal_components[nan_mask[:, :number_of_components]] = nodata
273     else:
274         principal_components[nan_mask[:, :number_of_components]] = np.nan
275
276     # Convert PCA output to proper format
277     if isinstance(data, np.ndarray):

```

```
if data.ndim == 3:
    result_data = principal_components.reshape(rows, cols, -1).transpose(2, 0, 1)
else:
    result_data = principal_components

elif isinstance(data, rasterio.io.DatasetReader):
    principal_components = principal_components.reshape(rows, cols, -1).transpose(2, 0, 1)
    out_profile = data.profile.copy()
    out_profile["count"] = number_of_components
    out_profile["dtype"] = "float32"
    result_data = (principal_components, out_profile)

elif isinstance(data, pd.DataFrame):
    component_names = [f"principal_component_{i+1}" for i in range(number_of_components)]
    pca_df = pd.DataFrame(data=principal_components, columns=component_names)
    if color_column_name is not None:
        pca_df[color_column_name] = color_column_data
    sns_pair_grid = plot_pca(pca_df, explained_variances, color_column_name)
    if isinstance(data, gpd.GeoDataFrame):
        pca_df = gpd.GeoDataFrame(pca_df, geometry=geometries, crs=crs)
        result_data = (pca_df, sns_pair_grid)

return result_data, explained_variances
```

```
plot_pca(pca_df, explained_variances=None, color_column_name=None,
save_path=None)
```

Plot a scatter matrix of different principal component combinations.

#### Parameters:

Name	Type	Description	Default
<code>pca_df</code>	<code>DataFrame</code>	A <code>DataFrame</code> containing computed principal components.	required
<code>explained_variances</code>	<code>Optional[ndarray]</code>	The explained variance ratios for each principal component. Used for labeling axes in the plot. Optional parameter. Defaults to <code>None</code> .	<code>None</code>
<code>color_column_name</code>	<code>Optional[str]</code>	Name of the column that will be used for color-coding data points. Typically a categorical variable in the original data. Optional parameter, no colors if not provided. Defaults to <code>None</code> .	<code>None</code>
<code>save_path</code>	<code>Optional[str]</code>	The save path for the plot. Optional parameter, no saving if not provided. Defaults to <code>None</code> .	<code>None</code>

#### Returns:

Type	Description
<code>PairGrid</code>	A Seaborn pairgrid containing the PCA scatter matrix.

#### Raises:

Type	Description
<code>InvalidColumnException</code>	<code>DataFrame</code> does not contain the given color column.



Source code in `eis_toolkit/exploratory_analyses/pca.py` 

```

195 @beartype
196 def plot_pca(
197     pca_df: pd.DataFrame,
198     explained_variances: Optional[np.ndarray] = None,
199     color_column_name: Optional[str] = None,
200     save_path: Optional[str] = None,
201 ) -> sns.PairGrid:
202     """Plot a scatter matrix of different principal component combinations.
203
204     Args:
205         pca_df: A DataFrame containing computed principal components.
206         explained_variances: The explained variance ratios for each principal component. Used for labeling
207             axes in the plot. Optional parameter. Defaults to None.
208         color_column_name: Name of the column that will be used for color-coding data points. Typically a
209             categorical variable in the original data. Optional parameter, no colors if not provided.
210             Defaults to None.
211         save_path: The save path for the plot. Optional parameter, no saving if not provided. Defaults to None.
212
213     Returns:
214         A Seaborn pairgrid containing the PCA scatter matrix.
215
216     Raises:
217         InvalidColumnException: DataFrame does not contain the given color column.
218     """
219
220     if color_column_name and color_column_name not in pca_df.columns:
221         raise exceptions.InvalidColumnException("DataFrame does not contain the given color column.")
222
223     pair_grid = sns.pairplot(pca_df, hue=color_column_name)
224
225     # Add explained variances to axis labels if provided
226     if explained_variances is not None:
227         labels = [f"PC {i+1} ({var:.1f}%)" for i, var in enumerate(explained_variances * 100)]
228     else:
229         labels = [f"PC {i+1}" for i in range(len(pair_grid.axes))]
230
231     # Iterate over axes objects and set the labels
232     for i, ax_row in enumerate(pair_grid.axes):
233         for j, ax in enumerate(ax_row):
234             if j == 0: # Only the first column
235                 ax.set_ylabel(labels[i], fontsize="large")
236             if i == len(ax_row) - 1: # Only the last row
237                 ax.set_xlabel(labels[j], fontsize="large")
238
239     if save_path is not None:
240         plt.savefig(save_path)
241
242     return pair_grid

```

## 4.7 Statistical (hypothesis) testing

---

```
chi_square_test(data, target_column, columns=None)
```

Compute Chi-square test for independence on the input data.

It is assumed that the variables in the input data are independent and that they are categorical, i.e. strings, booleans or integers, but not floats.

### Parameters:

Name	Type	Description	Default
<code>data</code>	<code>DataFrame</code>	Dataframe containing the input data	required
<code>target_column</code>	<code>str</code>	Variable against which independence of other variables is tested.	required
<code>columns</code>	<code>Sequence[str]</code>	Variables that are tested against the variable in <code>target_column</code> . If <code>None</code> , every column is used.	<code>None</code>

### Raises:

Type	Description
<code>EmptyDataFrameException</code>	The input Dataframe is empty.
<code>InvalidParameterValueException</code>	The <code>target_column</code> is not in input Dataframe or invalid column is provided.

### Returns:

Type	Description
<code>dict</code>	Test statistics for each variable (except <code>target_column</code> ).

Source code in `eis_toolkit/exploratory_analyses/statistical_tests.py` 

```

10 @beartype
11 def chi_square_test(data: pd.DataFrame, target_column: str, columns: Sequence[str] = None) -> dict:
12     """Compute Chi-square test for independence on the input data.
13
14     It is assumed that the variables in the input data are independent and that they are categorical, i.e. strings,
15     booleans or integers, but not floats.
16
17     Args:
18         data: Dataframe containing the input data
19         target_column: Variable against which independence of other variables is tested.
20         columns: Variables that are tested against the variable in target_column. If None, every column is used.
21
22     Raises:
23         EmptyDataFrameException: The input Dataframe is empty.
24         InvalidParameterValueException: The target_column is not in input Dataframe or invalid column is provided.
25
26     Returns:
27         Test statistics for each variable (except target_column).
28     """
29     if check_empty_dataframe(data):
30         raise exceptions.EmptyDataFrameException("The input Dataframe is empty.")
31
32     if not check_columns_valid(data, target_column):
33         raise exceptions.InvalidParameterValueException("Target column not found in the Dataframe.")
34
35     if columns is not None:
36         invalid_columns = [column for column in columns if column not in data.columns]
37         if any(invalid_columns):
38             raise exceptions.InvalidParameterValueException(
39                 f"The following variables are not in the dataframe: {invalid_columns}"
40             )
41     else:
42         columns = data.columns
43
44     statistics = {}
45     for column in columns:
46         if column != target_column:
47             contingency_table = pd.crosstab(data[target_column], data[column])
48             chi_square, p_value, degrees_of_freedom, _ = chi2_contingency(contingency_table)
49             statistics[column] = (chi_square, p_value, degrees_of_freedom)
50
51     return statistics

```

```
correlation_matrix(data, correlation_method='pearson', min_periods=None)
```

Compute correlation matrix on the input data.

It is assumed that the data is numeric, i.e. integers or floats.

**Parameters:**

Name	Type	Description	Default
<code>data</code>	<code>DataFrame</code>	Dataframe containing the input data.	required
<code>correlation_method</code>	<code>Literal[pearson, kendall, spearman]</code>	'pearson', 'kendall', or 'spearman'. Defaults to 'pearson'.	'pearson'
<code>min_periods</code>	<code>Optional[int]</code>	Minimum number of observations required per pair of columns to have valid result. Optional.	None

**Raises:**

Type	Description
<code>EmptyDataFrameException</code>	The input Dataframe is empty.
<code>InvalidParameterValueException</code>	<code>min_periods</code> argument is used with method 'kendall'.

**Returns:**

Type	Description
DataFrame	Dataframe containing the correlation matrix

**Source code in eis\_toolkit/exploratory\_analyses/statistical\_tests.py** 

```

80 @beartype
81 def correlation_matrix(
82     data: pd.DataFrame,
83     correlation_method: Literal["pearson", "kendall", "spearman"] = "pearson",
84     min_periods: Optional[int] = None,
85 ) -> pd.DataFrame:
86     """Compute correlation matrix on the input data.
87
88     It is assumed that the data is numeric, i.e. integers or floats.
89
90     Args:
91         data: Dataframe containing the input data.
92         correlation_method: 'pearson', 'kendall', or 'spearman'. Defaults to 'pearson'.
93         min_periods: Minimum number of observations required per pair of columns to have valid result. Optional.
94
95     Raises:
96         EmptyDataFrameException: The input Dataframe is empty.
97         InvalidParameterValueException: min_periods argument is used with method 'kendall'.
98
99     Returns:
100         Dataframe containing the correlation matrix
101     """
102     if check_empty_dataframe(data):
103         raise exceptions.EmptyDataFrameException("The input Dataframe is empty.")
104
105     if correlation_method == "kendall" and min_periods is not None:
106         raise exceptions.InvalidParameterValueException(
107             "The argument min_periods is available only with correlation methods 'pearson' and 'spearman'."
108         )
109
110     matrix = data.corr(method=correlation_method, min_periods=min_periods)
111
112     return matrix

```

`covariance_matrix(data, min_periods=None, delta_degrees_of_freedom=1)`

Compute covariance matrix on the input data.

It is assumed that the data is numeric, i.e. integers or floats.

**Parameters:**

Name	Type	Description	Default
data	DataFrame	Dataframe containing the input data.	required
min_periods	Optional[int]	Minimum number of observations required per pair of columns to have valid result. Optional.	None
delta_degrees_of_freedom	int	Delta degrees of freedom used for computing covariance matrix. Defaults to 1.	1

**Raises:**

Type	Description
EmptyDataFrameException	The input Dataframe is empty.
InvalidParameterValueException	Provided value for delta_degrees_of_freedom is negative.

**Returns:**

Type	Description
DataFrame	Dataframe containing the covariance matrix

Source code in `eis_toolkit/exploratory_analyses/statistical_tests.py` 

```

115 @beartype
116 def covariance_matrix(
117     data: pd.DataFrame, min_periods: Optional[int] = None, delta_degrees_of_freedom: int = 1
118 ) -> pd.DataFrame:
119     """Compute covariance matrix on the input data.
120
121     It is assumed that the data is numeric, i.e. integers or floats.
122
123     Args:
124         data: Dataframe containing the input data.
125         min_periods: Minimum number of observations required per pair of columns to have valid result. Optional.
126         delta_degrees_of_freedom: Delta degrees of freedom used for computing covariance matrix. Defaults to 1.
127
128     Raises:
129         EmptyDataFrameException: The input Dataframe is empty.
130         InvalidParameterValueException: Provided value for delta_degrees_of_freedom is negative.
131
132     Returns:
133         Dataframe containing the covariance matrix
134     """
135     if check_empty_dataframe(data):
136         raise exceptions.EmptyDataFrameException("The input Dataframe is empty.")
137
138     if delta_degrees_of_freedom < 0:
139         raise exceptions.InvalidParameterValueException("Delta degrees of freedom must be non-negative.")
140
141     matrix = data.cov(min_periods=min_periods, ddof=delta_degrees_of_freedom)
142
143     return matrix

```

`normality_test(data)`

Compute Shapiro-Wilk test for normality on the input data.

It is assumed that the input data is normally distributed and numeric, i.e. integers or floats.

**Parameters:**

Name	Type	Description	Default
<code>data</code>	<code>DataFrame</code>	Dataframe containing the input data.	required

**Returns:**

Type	Description
<code>dict</code>	Test statistics for each variable.

**Raises:**

Type	Description
<code>EmptyDataFrameException</code>	The input Dataframe is empty.

Source code in `eis_toolkit/exploratory_analyses/statistical_tests.py` 

```
54 @beartype
55 def normality_test(data: pd.DataFrame) -> dict:
56     """Compute Shapiro-Wilk test for normality on the input data.
57
58     It is assumed that the input data is normally distributed and numeric, i.e. integers or floats.
59
60     Args:
61         data: Dataframe containing the input data.
62
63     Returns:
64         Test statistics for each variable.
65
66     Raises:
67         EmptyDataFrameException: The input Dataframe is empty.
68     """
69     if check_empty_dataframe(data):
70         raise exceptions.EmptyDataFrameException("The input Dataframe is empty.")
71
72     statistics = {}
73     for column in data.columns:
74         statistic, p_value = shapiro(data[column])
75         statistics[column] = (statistic, p_value)
76
77     return statistics
```

## 5. Prediction

---

### 5.1 Fuzzy overlay

---

#### `and_overlay(data)`

Compute an 'and' overlay operation with fuzzy logic.

#### Parameters:

Name	Type	Description	Default
<code>data</code>	<code>ndarray</code>	The input data as a 3D Numpy array. Each 2D array represents a raster band. Data points should be in the range [0, 1].	required

#### Returns:

Type	Description
<code>ndarray</code>	2D Numpy array with the result of the 'and' overlay operation. Values are in range [0, 1].

#### Raises:

Type	Description
<code>InvalidParameterValueException</code>	If data values are not in range [0, 1].

#### Source code in `eis_toolkit/prediction/fuzzy_overlay.py`

```

14 @beartype
15 def and_overlay(data: np.ndarray) -> np.ndarray:
16     """Compute an 'and' overlay operation with fuzzy logic.
17
18     Args:
19         data: The input data as a 3D Numpy array. Each 2D array represents a raster band.
20             Data points should be in the range [0, 1].
21
22     Returns:
23         2D Numpy array with the result of the 'and' overlay operation. Values are in range [0, 1].
24
25     Raises:
26         InvalidParameterValueException: If data values are not in range [0, 1].
27     """
28     _check_input_data(data=data)
29
30     return data.min(axis=0)

```

#### `gamma_overlay(data, gamma)`

Compute a 'gamma' overlay operation with fuzzy logic.

**Parameters:**

Name	Type	Description	Default
data	ndarray	The input data as a 3D Numpy array. Each 2D array represents a raster band. Data points should be in the range [0, 1].	required
gamma	float	The gamma parameter. With gamma value 0, result will be same as 'product'overlay. When gamma is closer to 1, the weight of 'sum' overlay is increased. Value must be in the range [0, 1].	required

**Returns:**

Type	Description
ndarray	2D Numpy array with the result of the 'gamma' overlay operation. Values are in range [0, 1].

**Raises:**

Type	Description
InvalidParameterValueException	If data values or gamma are not in range [0, 1].

**Source code in eis\_toolkit/prediction/fuzzy\_overlay.py** 

```

90 @beartype
91 def gamma_overlay(data: np.ndarray, gamma: float) -> np.ndarray:
92     """Compute a 'gamma' overlay operation with fuzzy logic.
93
94     Args:
95         data: The input data as a 3D Numpy array. Each 2D array represents a raster band.
96             Data points should be in the range [0, 1].
97         gamma: The gamma parameter. With gamma value 0, result will be same as 'product'overlay.
98             When gamma is closer to 1, the weight of 'sum' overlay is increased.
99             Value must be in the range [0, 1].
100
101     Returns:
102         2D Numpy array with the result of the 'gamma' overlay operation. Values are in range [0, 1].
103
104     Raises:
105         InvalidParameterValueException: If data values or gamma are not in range [0, 1].
106     """
107     if gamma < 0 or gamma > 1:
108         raise exceptions.InvalidParameterValueException("The gamma parameter must be in range [0, 1]")
109
110     sum = sum_overlay(data=data)
111     product = product_overlay(data=data)
112     return product ** (1 - gamma) * sum**gamma

```

**or\_overlay(data)**

Compute an 'or' overlay operation with fuzzy logic.

**Parameters:**

Name	Type	Description	Default
data	ndarray	The input data as a 3D Numpy array. Each 2D array represents a raster band. Data points should be in the range [0, 1].	required

**Returns:**

Type	Description
ndarray	2D Numpy array with the result of the 'or' overlay operation. Values are in range [0, 1].



**Raises:**

Type	Description
<code>InvalidParameterValueException</code>	If data values are not in range [0, 1].

**Source code in `eis_toolkit/prediction/fuzzy_overlay.py`**

```

33 @beartype
34 def or_overlay(data: np.ndarray) -> np.ndarray:
35     """Compute an 'or' overlay operation with fuzzy logic.
36
37     Args:
38         data: The input data as a 3D Numpy array. Each 2D array represents a raster band.
39             Data points should be in the range [0, 1].
40
41     Returns:
42         2D Numpy array with the result of the 'or' overlay operation. Values are in range [0, 1].
43
44     Raises:
45         InvalidParameterValueException: If data values are not in range [0, 1].
46     """
47     _check_input_data(data=data)
48
49     return data.max(axis=0)

```

**`product_overlay(data)`**

Compute a 'product' overlay operation with fuzzy logic.

**Parameters:**

Name	Type	Description	Default
<code>data</code>	<code>ndarray</code>	The input data as a 3D Numpy array. Each 2D array represents a raster band. Data points should be in the range [0, 1].	required

**Returns:**

Type	Description
<code>ndarray</code>	2D Numpy array with the result of the 'product' overlay operation. Values are in range [0, 1].

**Raises:**

Type	Description
<code>InvalidParameterValueException</code>	If data values are not in range [0, 1].

**Source code in `eis_toolkit/prediction/fuzzy_overlay.py`**

```

52 @beartype
53 def product_overlay(data: np.ndarray) -> np.ndarray:
54     """Compute a 'product' overlay operation with fuzzy logic.
55
56     Args:
57         data: The input data as a 3D Numpy array. Each 2D array represents a raster band.
58             Data points should be in the range [0, 1].
59
60     Returns:
61         2D Numpy array with the result of the 'product' overlay operation. Values are in range [0, 1].
62
63     Raises:
64         InvalidParameterValueException: If data values are not in range [0, 1].
65     """
66     _check_input_data(data=data)
67
68     return np.prod(data, axis=0)

```

## sum\_overlay(data)

Compute a 'sum' overlay operation with fuzzy logic.

### Parameters:

Name	Type	Description	Default
data	ndarray	The input data as a 3D Numpy array. Each 2D array represents a raster band. Data points should be in the range [0, 1].	required

### Returns:

Type	Description
ndarray	2D Numpy array with the result of the 'sum' overlay operation. Values are in range [0, 1].

### Raises:

Type	Description
InvalidParameterValueException	If data values are not in range [0, 1].

### Source code in eis\_toolkit/prediction/fuzzy\_overlay.py

```

71 @beartype
72 def sum_overlay(data: np.ndarray) -> np.ndarray:
73     """Compute a 'sum' overlay operation with fuzzy logic.
74
75     Args:
76         data: The input data as a 3D Numpy array. Each 2D array represents a raster band.
77             Data points should be in the range [0, 1].
78
79     Returns:
80         2D Numpy array with the result of the 'sum' overlay operation. Values are in range [0, 1].
81
82     Raises:
83         InvalidParameterValueException: If data values are not in range [0, 1].
84     """
85     _check_input_data(data=data)
86
87     return data.sum(axis=0) - np.prod(data, axis=0)

```

## 5.2 Weights of evidence

---

```
weights_of_evidence_calculate_responses(output_arrays, nr_of_deposits,
nr_of_pixels)
```

Calculate the posterior probabilities for the given generalized weight arrays.

### Parameters:

Name	Type	Description	Default
<code>output_arrays</code>	Sequence[Dict[str, ndarray]]	List of output array dictionaries returned by weights of evidence calculations. For each dictionary, generalized weight and generalized standard deviation arrays are used and summed together pixel-wise to calculate the posterior probabilities. If generalized arrays are not found, the W+ and S_W+ arrays are used (so if outputs from unique weight calculations are used for this function).	required
<code>nr_of_deposits</code>	int	Number of deposit pixels in the input data for weights of evidence calculations.	required
<code>nr_of_pixels</code>	int	Number of evidence pixels in the input data for weights of evidence calculations.	required

### Returns:

Type	Description
ndarray	Array of posterior probabilities.
ndarray	Array of standard deviations in the posterior probability calculations.
ndarray	Array of confidence of the prospectivity values obtained in the posterior probability array.

Source code in `eis_toolkit/prediction/weights_of_evidence.py` 

```

372 @beartype
373 def weights_of_evidence_calculate_responses(
374     output_arrays: Sequence[Dict[str, np.ndarray]], nr_of_deposits: int, nr_of_pixels: int
375 ) -> Tuple[np.ndarray, np.ndarray, np.ndarray]:
376     """Calculate the posterior probabilities for the given generalized weight arrays.
377
378     Args:
379         output_arrays: List of output array dictionaries returned by weights of evidence calculations.
380             For each dictionary, generalized weight and generalized standard deviation arrays are used and summed
381             together pixel-wise to calculate the posterior probabilities. If generalized arrays are not found,
382             the W+ and S W+ arrays are used (so if outputs from unique weight calculations are used for this function).
383         nr_of_deposits: Number of deposit pixels in the input data for weights of evidence calculations.
384         nr_of_pixels: Number of evidence pixels in the input data for weights of evidence calculations.
385
386     Returns:
387         Array of posterior probabilities.
388         Array of standard deviations in the posterior probability calculations.
389         Array of confidence of the prospectivity values obtained in the posterior probability array.
390     """
391     gen_weights_sum = sum(
392         [
393             item[GENERALIZED_WEIGHT_PLUS_COLUMN]
394             if GENERALIZED_WEIGHT_PLUS_COLUMN in item.keys()
395             else item[WEIGHT_PLUS_COLUMN]
396             for item in output_arrays
397         ]
398     )
399     gen_weights_variance_sum = sum(
400         [
401             np.square(item[GENERALIZED_S_WEIGHT_PLUS_COLUMN])
402             if GENERALIZED_S_WEIGHT_PLUS_COLUMN in item.keys()
403             else np.square(item[WEIGHT_S_PLUS_COLUMN])
404             for item in output_arrays
405         ]
406     )
407     prior_probabilities = nr_of_deposits / nr_of_pixels
408     prior_odds = np.log(prior_probabilities / (1 - prior_probabilities))
409     posterior_probabilities = np.exp(gen_weights_sum + prior_odds) / (1 + np.exp(gen_weights_sum + prior_odds))
410
411     posterior_probabilities_squared = np.square(posterior_probabilities)
412     posterior_probabilities_std = np.sqrt(
413         (1 / nr_of_deposits + gen_weights_variance_sum) * posterior_probabilities_squared
414     )
415
416     confidence_array = posterior_probabilities / posterior_probabilities_std
417     return posterior_probabilities, posterior_probabilities_std, confidence_array

```

```

weights_of_evidence_calculate_weights(evidential_raster, deposits,
raster_nodata=None, weights_type='unique', studentized_contrast_threshold=1,
arrays_to_generate=None)

```

Calculate weights of spatial associations.

**Parameters:**

<b>Name</b>	<b>Type</b>	<b>Description</b>	<b>Default</b>
<code>evidential_raster</code>	<code>DatasetReader</code>	The evidential raster.	required
<code>deposits</code>	<code>GeoDataFrame</code>	Vector data representing the mineral deposits or occurrences point data.	required
<code>raster_nodata</code>	<code>Optional[Number]</code>	If nodata value of raster is wanted to specify manually. Optional parameter, defaults to None (nodata from raster metadata is used).	None
<code>weights_type</code>	<code>Literal[unique, categorical, ascending, descending]</code>	Accepted values are 'unique', 'categorical', 'ascending' and 'descending'. Unique weights does not create generalized classes and does not use a studentized contrast threshold value while categorical, cumulative ascending and cumulative descending do. Categorical weights are calculated so that all classes with studentized contrast below the defined threshold are grouped into one generalized class. Cumulative	'unique'

Name	Type	Description	Default
studentized_contrast_threshold	Number	<p>ascending and descending weights find the class with max contrast and group classes above/below into generalized classes. Generalized weights are also calculated for generalized classes.</p>	1
arrays_to_generate	Optional[Sequence[str]]	<p>Studentized contrast threshold value used with 'categorical', 'ascending' and 'descending' weight types. Used either as reclassification threshold directly (categorical) or to check that class with max contrast has studentized contrast value at least the defined value (cumulative). Defaults to 1.</p>	None
		<p>Arrays to generate from the computed weight metrics. All column names in the produced weights_df are valid choices. Defaults to ["Class", "W+", "S_W+ ] for "unique" weights_type and ["Class", "W+", "S_W+",</p>	

Name	Type	Description	Default
		"Generalized W+", "Generalized S_W+" for the cumulative weight types.	

**Returns:**

Type	Description
DataFrame	Dataframe with weights of spatial association between the input data.
dict	Dictionary of arrays for specified metrics.
dict	Raster metadata.
int	Number of deposit pixels.
int	Number of all evidence pixels.



Source code in [eis\\_toolkit/prediction/weights\\_of\\_evidence.py](#) 

254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361

```

362 @beartype
363 def weights_of_evidence_calculate_weights(
364     evidential_raster: rasterio.io.DatasetReader,
365     deposits: gpd.GeoDataFrame,
366     raster_nodata: Optional[Number] = None,
367     weights_type: Literal["unique", "categorical", "ascending", "descending"] = "unique",
368     studentized_contrast_threshold: Number = 1,
369     arrays_to_generate: Optional[Sequence[str]] = None,
370 ) -> Tuple[pd.DataFrame, dict, dict, int, int]:
    """
    Calculate weights of spatial associations.

    Args:
        evidential_raster: The evidential raster.
        deposits: Vector data representing the mineral deposits or occurrences point data.
        raster_nodata: If nodata value of raster is wanted to specify manually. Optional parameter, defaults to None
            (nodata from raster metadata is used).
        weights_type: Accepted values are 'unique', 'categorical', 'ascending' and 'descending'.
            Unique weights does not create generalized classes and does not use a studentized contrast threshold value
            while categorical, cumulative ascending and cumulative descending do. Categorical weights are calculated so
            that all classes with studentized contrast below the defined threshold are grouped into one generalized
            class. Cumulative ascending and descending weights find the class with max contrast and group classes
            above/below into generalized classes. Generalized weights are also calculated for generalized classes.
        studentized_contrast_threshold: Studentized contrast threshold value used with 'categorical', 'ascending' and
            'descending' weight types. Used either as reclassification threshold directly (categorical) or to check
            that class with max contrast has studentized contrast value at least the defined value (cumulative).
            Defaults to 1.
        arrays_to_generate: Arrays to generate from the computed weight metrics. All column names
            in the produced weights_df are valid choices. Defaults to ["Class", "W+", "S_W+"]
            for "unique" weights type and ["Class", "W+", "S_W+", "Generalized W+", "Generalized S_W+"]
            for the cumulative weight types.

    Returns:
        Dataframe with weights of spatial association between the input data.
        Dictionary of arrays for specified metrics.
        Raster metadata.
        Number of deposit pixels.
        Number of all evidence pixels.
    """

    if arrays_to_generate is None:
        if weights_type == "unique":
            metrics_to_arrays = DEFAULT_METRICS_UNIQUE
        else:
            metrics_to_arrays = DEFAULT_METRICS_CUMULATIVE
    else:
        for col_name in arrays_to_generate:
            if col_name not in VALID_DF_COLUMNS:
                raise exceptions.InvalidColumnException(
                    f"Arrays to generate contains invalid metric / column name: {col_name}."
                )
            metrics_to_arrays = arrays_to_generate.copy()

    # 1. Preprocess data
    evidence_array = read_and_preprocess_evidence(evidential_raster, raster_nodata)
    raster_meta = evidential_raster.meta

    # Rasterize deposits
    deposit_array, _ = rasterize_vector(
        geodataframe=deposits, default_value=1.0, base_raster_profile=raster_meta, fill_value=0.0
    )

    # Mask NaN out of the array
    nodata_mask = np.isnan(evidence_array)
    masked_evidence_array = evidence_array[~nodata_mask]
    masked_deposit_array = deposit_array[~nodata_mask]

    # 2. WofE calculations
    if weights_type == "unique" or weights_type == "categorical":
        wofe_weights = _unique_weights(masked_deposit_array, masked_evidence_array)
    elif weights_type == "ascending":
        wofe_weights = _cumulative_weights(masked_deposit_array, masked_evidence_array, ascending=True)
    elif weights_type == "descending":
        wofe_weights = _cumulative_weights(masked_deposit_array, masked_evidence_array, ascending=False)
    else:
        raise exceptions.InvalidParameterValueException(
            "Expected weights_type to be one of unique, categorical, ascending or descending."
        )

    # 3. Create DataFrame based on calculated metrics
    df_entries = []
    for cls, metrics in wofe_weights.items():
        metrics = [round(metric, 4) if isinstance(metric, np.floating) else metric for metric in metrics]
        A, _, C, _, w_plus, s_w_plus, w_minus, s_w_minus, contrast, s_contrast, studentized_contrast = metrics
        df_entries.append(
            {
                CLASS_COLUMN: cls,
                PIXEL_COUNT_COLUMN: A + C,
                DEPOSIT_COUNT_COLUMN: A,
                WEIGHT_PLUS_COLUMN: w_plus,
                WEIGHT_S_PLUS_COLUMN: s_w_plus,
                WEIGHT_MINUS_COLUMN: w_minus,
                WEIGHT_S_MINUS_COLUMN: s_w_minus,
                CONTRAST_COLUMN: contrast,
                S_CONTRAST_COLUMN: s_contrast,
                STUDENTIZED_CONTRAST_COLUMN: studentized_contrast,
            }
        )
    weights_df = pd.DataFrame(df_entries)

    # 4. If we use cumulative weights type, calculate generalized classes and weights
    if weights_type == "categorical":
        weights_df = _generalized_classes_categorical(weights_df, studentized_contrast_threshold)
        weights_df = _generalized_weights_categorical(weights_df, masked_deposit_array)
    elif weights_type == "ascending" or weights_type == "descending":
        weights_df = _generalized_classes_cumulative(weights_df, studentized_contrast_threshold)
        weights_df = _generalized_weights_cumulative(weights_df, masked_deposit_array)

```

```
# 5. Generate arrays for desired metrics
arrays_dict = _generate_arrays_from_metrics(evidence_array, weights_df, metrics_to_arrays)

# Return nr. of deposit pixels and nr. of all evidence pixels for to be used in calculate responses
nr_of_deposits = int(np.sum(masked_deposit_array == 1))
nr_of_pixels = int(np.size(masked_evidence_array))

return weights_df, arrays_dict, raster_meta, nr_of_deposits, nr_of_pixels
```

## 6. Raster processing

---

### 6.1 Check raster grids

---

```
check_raster_grids(rasters, same_extent=False)
```

Check the set of input rasters for matching gridding and optionally matching bounds.

#### Parameters:

Name	Type	Description	Default
<code>rasters</code>	List[DatasetReader]	List of rasters to test for matching gridding.	required
<code>same_extent</code>	bool	optional boolean argument that determines if rasters are tested for matching bounds. Default set to False.	False

#### Returns:

Type	Description
bool	True if gridding and optionally bounds matches, False if not.

#### Source code in `eis_toolkit/raster_processing/check_raster_grids.py`

```

23 def check_raster_grids( # type: ignore[no-any-unimported]
24     rasters: List[rasterio.io.DatasetReader], same_extent: bool = False
25 ) -> bool:
26     """
27     Check the set of input rasters for matching gridding and optionally matching bounds.
28
29     Args:
30         rasters: List of rasters to test for matching gridding.
31         same_extent: optional boolean argument that determines if rasters are tested for matching bounds.
32                     Default set to False.
33
34     Returns:
35         True if gridding and optionally bounds matches, False if not.
36     """
37     check = _check_raster_grids(rasters=rasters, same_extent=same_extent)
38     return check

```

## 6.2 Clipping

### `clip_raster(raster, geodataframe)`

Clips a raster with polygon geometries.

#### Parameters:

Name	Type	Description	Default
<code>raster</code>	<code>DatasetReader</code>	The raster to be clipped.	required
<code>geodataframe</code>	<code>GeoDataFrame</code>	A geodataframe containing the geometries to do the clipping with. Should contain only polygon features.	required

#### Returns:

Type	Description
<code>ndarray</code>	The clipped raster data.
<code>dict</code>	The updated metadata.

#### Raises:

Type	Description
<code>NonMatchingCrsException</code>	The raster and geodataframe are not in the same CRS.
<code>NotApplicableGeometryTypeException</code>	The input geometries contain non-polygon features.

#### Source code in `eis_toolkit/raster_processing/clipping.py`

```

24 @beartype
25 def clip_raster(raster: rasterio.io.DatasetReader, geodataframe: geopandas.GeoDataFrame) -> Tuple[np.ndarray, dict]:
26     """Clips a raster with polygon geometries.
27
28     Args:
29         raster: The raster to be clipped.
30         geodataframe: A geodataframe containing the geometries to do the clipping with.
31             Should contain only polygon features.
32
33     Returns:
34         The clipped raster data.
35         The updated metadata.
36
37     Raises:
38         NonMatchingCrsException: The raster and geodataframe are not in the same CRS.
39         NotApplicableGeometryTypeException: The input geometries contain non-polygon features.
40     """
41     geometries = geodataframe["geometry"]
42
43     if not check_matching_crs(
44         objects=[raster, geometries],
45     ):
46         raise NonMatchingCrsException("The raster and geodataframe are not in the same CRS.")
47
48     if not check_geometry_types(
49         geometries=geometries,
50         allowed_types=["Polygon", "MultiPolygon"],
51     ):
52         raise NotApplicableGeometryTypeException("The input geometries contain non-polygon features.")
53
54     out_image, out_meta = _clip_raster(
55         raster=raster,
56         geometries=geometries,
57     )
58
59     return out_image, out_meta

```

## 6.3 Create constant raster

---

```
create_constant_raster(constant_value, template_raster=None, coord_west=None,
coord_north=None, coord_east=None, coord_south=None, target_epsg=None,
target_pixel_size=None, raster_width=None, raster_height=None,
nodata_value=None)
```

Create a constant raster based on a user-defined value.

Provide 3 methods for raster creation: 1. Set extent and coordinate system based on a template raster. 2. Set extent from origin, based on the western and northern coordinates and the pixel size. 3. Set extent from bounds, based on western, northern, eastern and southern points.

Always provide values for height and width for the last two options, which correspond to the desired number of pixels for rows and columns.

### Parameters:

Name	Type	Description	Default
constant_value	Number	The constant value to use in the raster.	required
template_raster	Optional[DatasetReader]	An optional raster to use as a template for the output.	None
coord_west	Optional[Number]	The western coordinate of the output raster in [m].	None
coord_east	Optional[Number]	The eastern coordinate of the output raster in [m].	None
coord_south	Optional[Number]	The southern coordinate of the output raster in [m].	None
coord_north	Optional[Number]	The northern coordinate of the output raster in [m].	None
target_epsg	Optional[int]	The EPSG code for the output raster.	None
target_pixel_size	Optional[int]	The pixel size of the output raster.	None
raster_width	Optional[int]	The width of the output raster.	None
raster_height	Optional[int]	The height of the output raster.	None
nodata_value	Optional[Number]	The nodata value of the output raster.	None

### Returns:

Type	Description
Tuple[ndarray, dict]	A tuple containing the output raster as a NumPy array and updated metadata.

### Raises:

Type	Description
InvalidParameterValueException	Provide invalid input parameter.

Source code in `eis_toolkit/raster_processing/create_constant_raster.py` 

```

88 @beartype
89 def create_constant_raster( # type: ignore[no-any-unimported]
90     constant_value: Number,
91     template_raster: Optional[rasterio.io.DatasetReader] = None,
92     coord_west: Optional[Number] = None,
93     coord_north: Optional[Number] = None,
94     coord_east: Optional[Number] = None,
95     coord_south: Optional[Number] = None,
96     target_epsg: Optional[int] = None,
97     target_pixel_size: Optional[int] = None,
98     raster_width: Optional[int] = None,
99     raster_height: Optional[int] = None,
100     nodata_value: Optional[Number] = None,
101 ) -> Tuple[np.ndarray, dict]:
102     """Create a constant raster based on a user-defined value.
103
104     Provide 3 methods for raster creation:
105     1. Set extent and coordinate system based on a template raster.
106     2. Set extent from origin, based on the western and northern coordinates and the pixel size.
107     3. Set extent from bounds, based on western, northern, eastern and southern points.
108
109     Always provide values for height and width for the last two options, which correspond to
110     the desired number of pixels for rows and columns.
111
112     Args:
113     constant_value: The constant value to use in the raster.
114     template_raster: An optional raster to use as a template for the output.
115     coord_west: The western coordinate of the output raster in [m].
116     coord_east: The eastern coordinate of the output raster in [m].
117     coord_south: The southern coordinate of the output raster in [m].
118     coord_north: The northern coordinate of the output raster in [m].
119     target_epsg: The EPSG code for the output raster.
120     target_pixel_size: The pixel size of the output raster.
121     raster_width: The width of the output raster.
122     raster_height: The height of the output raster.
123     nodata_value: The nodata value of the output raster.
124
125     Returns:
126     A tuple containing the output raster as a NumPy array and updated metadata.
127
128     Raises:
129     InvalidParameterValueException: Provide invalid input parameter.
130     """
131
132     if template_raster is not None:
133         out_array, out_meta = _create_constant_raster_from_template(constant_value, template_raster, nodata_value)
134
135     elif all(coords is not None for coords in [coord_west, coord_east, coord_south, coord_north]):
136         if raster_height <= 0 or raster_width <= 0:
137             raise InvalidParameterValueException("Invalid raster extent provided.")
138         if not check_minmax_position((coord_west, coord_east) or not check_minmax_position((coord_south, coord_north))):
139             raise InvalidParameterValueException("Invalid coordinate values provided.")
140
141         out_array, out_meta = _create_constant_raster_from_bounds(
142             constant_value,
143             coord_west,
144             coord_north,
145             coord_east,
146             coord_south,
147             target_epsg,
148             raster_width,
149             raster_height,
150             nodata_value,
151         )
152
153     elif all(coords is not None for coords in [coord_west, coord_north]) and all(
154         coords is None for coords in [coord_east, coord_south]
155     ):
156         if raster_height <= 0 or raster_width <= 0:
157             raise InvalidParameterValueException("Invalid raster extent provided.")
158         if target_pixel_size <= 0:
159             raise InvalidParameterValueException("Invalid pixel size.")
160
161         out_array, out_meta = _create_constant_raster_from_origin(
162             constant_value,
163             coord_west,
164             coord_north,
165             target_epsg,
166             target_pixel_size,
167             raster_width,
168             raster_height,
169             nodata_value,
170         )
171
172     else:
173         raise InvalidParameterValueException("Suitable parameter values were not provided for any of the 3 methods.")
174
175     constant_value = cast_scalar_to_int(constant_value)
176     nodata_value = cast_scalar_to_int(out_meta["nodata"])
177
178     if isinstance(constant_value, int) and isinstance(nodata_value, int):
179         target_dtype = np.result_type(get_min_int_type(constant_value), get_min_int_type(nodata_value))
180         out_array = out_array.astype(target_dtype)
181         out_meta["dtype"] = out_array.dtype
182     elif isinstance(constant_value, int) and isinstance(nodata_value, float):
183         out_array = out_array.astype(get_min_int_type(constant_value))
184         out_meta["dtype"] = np.float64.__name__
185     elif isinstance(constant_value, float):
186         out_array = out_array.astype(np.float64)
187         out_meta["dtype"] = out_array.dtype
188
189     return out_array, out_meta

```



## 6.4 Extract values from raster

```
extract_values_from_raster(raster_list, geodataframe, raster_column_names=None)
```

Extract raster values using point data to a DataFrame.

If custom column names are not given, column names are `file_name` for singleband files and `file_name_bandnumber` for multiband files. If custom column names are given, there should be column names for each raster provided in the raster list.

### Parameters:

Name	Type	Description	Default
<code>raster_list</code>	<code>Sequence[DatasetReader]</code>	List to extract values from.	required
<code>geodataframe</code>	<code>GeoDataFrame</code>	Object to extract values with.	required
<code>raster_column_names</code>	<code>Optional[Sequence[str]]</code>	List of optional column names for bands.	None

### Returns:

Type	Description
<code>DataFrame</code>	Dataframe with x & y coordinates and the values from the raster file(s) as columns.

### Raises:

Type	Description
<code>NonMatchingParameterLengthsException</code>	<code>raster_list</code> and <code>raster_columns_names</code> have different lengths.

### Source code in `eis_toolkit/raster_processing/extract_values_from_raster.py`

```

45 @beartype
46 def extract_values_from_raster(
47     raster_list: Sequence[rasterio.io.DatasetReader],
48     geodataframe: gpd.GeoDataFrame,
49     raster_column_names: Optional[Sequence[str]] = None,
50 ) -> pd.DataFrame:
51     """Extract raster values using point data to a DataFrame.
52
53     If custom column names are not given, column names are file_name for singleband files
54     and file_name_bandnumber for multiband files. If custom column names are given, there
55     should be column names for each raster provided in the raster list.
56
57     Args:
58         raster_list: List to extract values from.
59         geodataframe: Object to extract values with.
60         raster_column_names: List of optional column names for bands.
61
62     Returns:
63         Dataframe with x & y coordinates and the values from the raster file(s) as columns.
64
65     Raises:
66         NonMatchingParameterLengthsException: raster_list and raster_columns_names have different lengths.
67     """
68     if raster_column_names == []:
69         raster_column_names = None
70
71     if raster_column_names is not None and len(raster_list) != len(raster_column_names):
72         raise NonMatchingParameterLengthsException("Raster list and raster columns names have different lengths.")
73
74     data_frame = _extract_values_from_raster(
75         raster_list=raster_list, geodataframe=geodataframe, raster_column_names=raster_column_names
76     )
77
78     return data_frame

```

## 6.5 Reprojecting

---

```
reproject_raster(raster, target_crs, resampling_method=warp.Resampling.nearest)
```

Reprojects raster to match given coordinate reference system (EPSG).

### Parameters:

Name	Type	Description	Default
<code>raster</code>	<code>DatasetReader</code>	The raster to be reprojected.	required
<code>target_crs</code>	<code>int</code>	Target CRS as EPSG code.	required
<code>resampling_method</code>	<code>Resampling</code>	Resampling method. Most suitable method depends on the dataset and context. Nearest, bilinear and cubic are some common choices. This parameter defaults to nearest.	<code>nearest</code>

### Returns:

Type	Description
<code>ndarray</code>	The reprojected raster data.
<code>dict</code>	The updated metadata.

### Raises:

Type	Description
<code>NonMatchinCrsException</code>	Raster is already in the target CRS.

### Source code in `eis_toolkit/raster_processing/reprojecting.py`

```

55 @beartype
56 def reproject_raster(
57     raster: rasterio.io.DatasetReader, target_crs: int, resampling_method: warp.Resampling = warp.Resampling.nearest
58 ) -> Tuple[np.ndarray, dict]:
59     """Reprojects raster to match given coordinate reference system (EPSG).
60
61     Args:
62         raster: The raster to be reprojected.
63         target_crs: Target CRS as EPSG code.
64         resampling_method: Resampling method. Most suitable method depends on the dataset and context.
65             Nearest, bilinear and cubic are some common choices. This parameter defaults to nearest.
66
67     Returns:
68         The reprojected raster data.
69         The updated metadata.
70
71     Raises:
72         NonMatchinCrsException: Raster is already in the target CRS.
73     """
74     if target_crs == int(raster.crs.to_string()[5:]):
75         raise MatchingCrsException("Raster is already in the target CRS.")
76
77     out_image, out_meta = _reproject_raster(raster, target_crs, resampling_method)
78
79     return out_image, out_meta

```

## 6.6 Resampling

---

```
resample(raster, resolution, resampling_method=Resampling.bilinear)
```

Resamples raster according to given resolution.

### Parameters:

Name	Type	Description	Default
raster	DatasetReader	The raster to be resampled.	required
resolution	Number	Target resolution i.e. cell size of the output raster.	required
resampling_method	Resampling	Resampling method. Most suitable method depends on the dataset and context. Nearest, bilinear and cubic are some common choices. This parameter defaults to bilinear.	bilinear

### Returns:

Type	Description
ndarray	The resampled raster data.
dict	The updated metadata.

### Raises:

Type	Description
NumericValueSignException	Resolution is not a positive value.

#### Source code in eis\_toolkit/raster\_processing/resampling.py

```

52 @beartype
53 def resample(
54     raster: rasterio.io.DatasetReader,
55     resolution: Number,
56     resampling_method: Resampling = Resampling.bilinear,
57 ) -> Tuple[np.ndarray, dict]:
58     """Resamples raster according to given resolution.
59
60     Args:
61         raster: The raster to be resampled.
62         resolution: Target resolution i.e. cell size of the output raster.
63         resampling_method: Resampling method. Most suitable
64             method depends on the dataset and context. Nearest, bilinear and cubic are some
65             common choices. This parameter defaults to bilinear.
66
67     Returns:
68         The resampled raster data.
69         The updated metadata.
70
71     Raises:
72         NumericValueSignException: Resolution is not a positive value.
73     """
74     if resolution <= 0:
75         raise exceptions.NumericValueSignException(f"Expected a positive value for resolution: {resolution}")
76
77     out_image, out_meta = _resample(raster, resolution, resampling_method)
78     return out_image, out_meta

```

## 6.7 Snapping

---

### `snap_with_raster(raster, snap_raster)`

Snaps/aligns raster to given snap raster.

Raster is snapped from its left-bottom corner to nearest snap raster grid corner in left-bottom direction. If rasters are aligned, simply returns input raster data and metadata.

#### Parameters:

Name	Type	Description	Default
<code>raster</code>	<code>DatasetReader</code>	The raster to be clipped.	required
<code>snap_raster</code>	<code>DatasetReader</code>	The snap raster i.e. reference grid raster.	required

#### Returns:

Type	Description
<code>ndarray</code>	The snapped raster data.
<code>dict</code>	The updated metadata.

#### Raises:

Type	Description
<code>NonMatchingCredException</code>	Raster and and snap raster are not in the same CRS.
<code>MatchingRasterGridException</code>	Raster grids are already aligned.

#### Source code in `eis_toolkit/raster_processing/snapping.py`

```

66 @beartype
67 def snap_with_raster(raster: rasterio.DatasetReader, snap_raster: rasterio.DatasetReader) -> Tuple[np.ndarray, dict]:
68     """Snaps/aligns raster to given snap raster.
69
70     Raster is snapped from its left-bottom corner to nearest snap raster grid corner in left-bottom direction.
71     If rasters are aligned, simply returns input raster data and metadata.
72
73     Args:
74         raster: The raster to be clipped.
75         snap_raster: The snap raster i.e. reference grid raster.
76
77     Returns:
78         The snapped raster data.
79         The updated metadata.
80
81     Raises:
82         NonMatchingCredException: Raster and and snap raster are not in the same CRS.
83         MatchingRasterGridException: Raster grids are already aligned.
84     """
85
86     if not check_matching_crs(
87         objects=[raster, snap_raster],
88     ):
89         raise NonMatchingCredException("Raster and and snap raster have different CRS.")
90
91     if snap_raster.bounds.bottom == raster.bounds.bottom and snap_raster.bounds.left == raster.bounds.left:
92         raise MatchingRasterGridException("Raster grids are already aligned.")
93
94     out_image, out_meta = _snap(raster, snap_raster)
95     return out_image, out_meta

```

## 6.8 Unifying

---

```
unify_raster_grids(base_raster, rasters_to_unify,
                  resampling_method=Resampling.nearest, same_extent=False)
```

Unifies (reprojects, resamples, aligns and optionally clips) given rasters relative to base raster.

### Parameters:

Name	Type	Description	Default
<code>base_raster</code>	<code>DatasetReader</code>	The base raster to determine target raster grid properties.	required
<code>rasters_to_unify</code>	<code>Sequence[DatasetReader]</code>	Rasters to be unified with the base raster.	required
<code>resampling_method</code>	<code>Resampling</code>	Resampling method. Most suitable method depends on the dataset and context. Nearest, bilinear and cubic are some common choices. This parameter defaults to nearest.	<code>nearest</code>
<code>same_extent</code>	<code>bool</code>	If the unified rasters will be forced to have the same extent/bounds as the base raster. Expands smaller rasters with nodata cells. Defaults to False.	<code>False</code>

### Returns:

Type	Description
<code>List[Tuple[ndarray, dict]]</code>	List of unified rasters' data and metadata. First element is the base raster.

### Raises:

Type	Description
<code>InvalidParameterValueException</code>	Rasters to unify is empty.

Source code in `eis_toolkit/raster_processing/unifying.py` 

```
87 @beartype
88 def unify_raster_grids(
89     base_raster: rasterio.io.DatasetReader,
90     rasters_to_unify: Sequence[rasterio.io.DatasetReader],
91     resampling_method: Resampling = Resampling.nearest,
92     same_extent: bool = False,
93 ) -> List[Tuple[np.ndarray, dict]]:
94     """Unifies (reprojects, resamples, aligns and optionally clips) given rasters relative to base raster.
95
96     Args:
97         base_raster: The base raster to determine target raster grid properties.
98         rasters_to_unify: Rasters to be unified with the base raster.
99         resampling_method: Resampling method. Most suitable
100             method depends on the dataset and context. Nearest, bilinear and cubic are some
101             common choices. This parameter defaults to nearest.
102         same_extent: If the unified rasters will be forced to have the same extent/bounds
103             as the base raster. Expands smaller rasters with nodata cells. Defaults to False.
104
105     Returns:
106         List of unified rasters' data and metadata. First element is the base raster.
107
108     Raises:
109         InvalidParameterValueException: Rasters to unify is empty.
110     """
111     if len(rasters_to_unify) == 0:
112         raise InvalidParameterValueException("Rasters to unify is empty.")
113
114     out_rasters = _unify_raster_grids(base_raster, rasters_to_unify, resampling_method, same_extent)
115     return out_rasters
```

## 6.9 Windowing

---

```
extract_window(raster, center_coords, height, width)
```

Extract window from raster.

Center coordinate must be inside the raster but window can extent outside the raster in which case padding with raster nodata value is used. Args: raster: Source raster. center\_coords: Center coordinates for window in form (x, y). The coordinates should be in the raster's CRS. height: Window height in pixels. width: Window width in pixels.

### Returns:

Type	Description
ndarray	The extracted raster window.
dict	The updated metadata.

### Raises:

Type	Description
InvalidParameterValueException	Window size is too small.
CoordinatesOutOfBoundException	Window center coordinates are out of raster bounds.

#### Source code in `eis_toolkit/raster_processing/windowing.py`

```

81 @beartype
82 def extract_window(
83     raster: rasterio.io.DatasetReader,
84     center_coords: Tuple[Number, Number],
85     height: int,
86     width: int,
87 ) -> Tuple[np.ndarray, dict]:
88     """Extract window from raster.
89
90     Center coordinate must be inside the raster but window can extent outside the raster in which case padding with
91     raster nodata value is used.
92     Args:
93         raster: Source raster.
94         center_coords: Center coordinates for window in form (x, y). The coordinates should be in the raster's CRS.
95         height: Window height in pixels.
96         width: Window width in pixels.
97
98     Returns:
99         The extracted raster window.
100        The updated metadata.
101
102     Raises:
103         InvalidParameterValueException: Window size is too small.
104         CoordinatesOutOfBoundException: Window center coordinates are out of raster bounds.
105     """
106
107     if height < 1 or width < 1:
108         raise InvalidParameterValueException(f"Window size is too small: {height}, {width}.")
109
110     center_x = center_coords[0]
111     center_y = center_coords[1]
112
113     if (
114         center_x < raster.bounds.left
115         or center_x > raster.bounds.right
116         or center_y < raster.bounds.bottom
117         or center_y > raster.bounds.top
118     ):
119         raise CoordinatesOutOfBoundException("Window center coordinates are out of raster bounds.")
120
121     out_image, out_meta = _extract_window(raster, center_coords, height, width)
122
123     return out_image, out_meta

```

## 7. Training data tools

---

### 7.1 Class balancing

---

```
balance_SMOTETomek(X, y, sampling_strategy='auto', random_state=None)
```

Balances the classes of input dataset using SMOTETomek resampling method.

#### Parameters:

Name	Type	Description	Default
X	Union[DataFrame, ndarray]	The feature matrix (input data as a DataFrame).	required
y	Union[Series, ndarray]	The target labels corresponding to the feature matrix.	required
sampling_strategy	Union[float, str, dict]	Parameter controlling how to perform the resampling. If float, specifies the ratio of samples in minority class to samples of majority class, if str, specifies classes to be resampled ("minority", "not minority", "not majority", "all", "auto"), if dict, the keys should be targeted classes and values the desired number of samples for the class. Defaults to "auto", which will resample all classes except the majority class.	'auto'
random_state	Optional[int]	Parameter controlling randomization of the algorithm. Can be given a seed (number). Defaults to None, which randomizes the seed.	None

#### Returns:

Type	Description
tuple[Union[DataFrame, ndarray], Union[Series, ndarray]]	Resampled feature matrix and target labels.

#### Raises:

Type	Description
NonMatchingParameterLengthsException	If X and y have different length.



Source code in `eis_toolkit/training_data_tools/class_balancing.py` 

```

11 @beartype
12 def balance_SMOTETomek(
13     X: Union[pd.DataFrame, np.ndarray],
14     y: Union[pd.Series, np.ndarray],
15     sampling_strategy: Union[float, str, dict] = "auto",
16     random_state: Optional[int] = None,
17 ) -> tuple[Union[pd.DataFrame, np.ndarray], Union[pd.Series, np.ndarray]]:
18     """Balances the classes of input dataset using SMOTETomek resampling method.
19
20     Args:
21         X: The feature matrix (input data as a DataFrame).
22         y: The target labels corresponding to the feature matrix.
23         sampling_strategy: Parameter controlling how to perform the resampling.
24             If float, specifies the ratio of samples in minority class to samples of majority class,
25             if str, specifies classes to be resampled ("minority", "not minority", "not majority", "all", "auto"),
26             if dict, the keys should be targeted classes and values the desired number of samples for the class.
27             Defaults to "auto", which will resample all classes except the majority class.
28         random_state: Parameter controlling randomization of the algorithm. Can be given a seed (number).
29             Defaults to None, which randomizes the seed.
30
31     Returns:
32         Resampled feature matrix and target labels.
33
34     Raises:
35         NonMatchingParameterLengthsException: If X and y have different length.
36     """
37
38     if len(X) != len(y):
39         raise exceptions.NonMatchingParameterLengthsException(
40             "Feature matrix X and target labels y must have the same length."
41         )
42
43     X_res, y_res = SMOTETomek(sampling_strategy=sampling_strategy, random_state=random_state).fit_resample(X, y)
44     return X_res, y_res

```

## 8. Transformations

---

### 8.1 Binarize

---

```
binarize(raster, bands=None, thresholds=[Number], nodata=None)
```

Binarize data based on a given threshold.

Replaces values less or equal threshold with 0. Replaces values greater than the threshold with 1.

Takes one nodata value which will be re-written after transformation.

If no band/column selection specified, all bands/columns will be used. If a parameter contains only 1 entry, it will be applied for all bands. The threshold can be set for each band individually.

#### Parameters:

Name	Type	Description	Default
raster	DatasetReader	Data object to be transformed.	required
bands	Optional[Sequence[int]]	Selection of bands to be transformed.	None
thresholds	Sequence[Number]	Threshold values for transformation.	[Number]
nodata	Optional[Number]	Nodata value to be considered.	None

#### Returns:

Name	Type	Description
out_array	ndarray	The transformed data.
out_meta	dict	Updated metadata.
out_settings	dict	Log of input settings and calculated statistics if available.

#### Raises:

Type	Description
InvalidRasterBandException	The input contains invalid band numbers.
NonMatchingParameterLengthsException	The input does not match the number of selected bands.

Source code in `eis_toolkit/transformations/binarize.py` 

```

24 @beartype
25 def binarize( # type: ignore[no-any-unimported]
26     raster: rasterio.io.DatasetReader,
27     bands: Optional[Sequence[int]] = None,
28     thresholds: Sequence[Number] = [Number],
29     nodata: Optional[Number] = None,
30 ) -> Tuple[np.ndarray, dict, dict]:
31     """
32     Binarize data based on a given threshold.
33
34     Replaces values less or equal threshold with 0.
35     Replaces values greater than the threshold with 1.
36
37     Takes one nodata value which will be re-written after transformation.
38
39     If no band/column selection specified, all bands/columns will be used.
40     If a parameter contains only 1 entry, it will be applied for all bands.
41     The threshold can be set for each band individually.
42
43     Args:
44         raster: Data object to be transformed.
45         bands: Selection of bands to be transformed.
46         thresholds: Threshold values for transformation.
47         nodata: Nodata value to be considered.
48
49     Returns:
50         out_array: The transformed data.
51         out_meta: Updated metadata.
52         out_settings: Log of input settings and calculated statistics if available.
53
54     Raises:
55         InvalidRasterBandException: The input contains invalid band numbers.
56         NonMatchingParameterLengthsException: The input does not match the number of selected bands.
57     """
58     bands = list(range(1, raster.count + 1)) if bands is None else bands
59     nodata = cast_scalar_to_int(raster.nodata if nodata is None else nodata)
60
61     if check_raster_bands(raster, bands) is False:
62         raise InvalidRasterBandException("Invalid band selection.")
63
64     if check_parameter_length(bands, thresholds) is False:
65         raise NonMatchingParameterLengthsException("Invalid threshold length.")
66
67     expanded_args = expand_and_zip(bands, thresholds)
68     thresholds = [element[1] for element in expanded_args]
69
70     out_settings = {}
71
72     for i in range(0, len(bands)):
73         band_array = raster.read(bands[i])
74         initial_dtype = band_array.dtype
75
76         band_mask = np.isin(band_array, nodata)
77         band_array = binarize(band_array, threshold=thresholds[i])
78         band_array = np.where(band_mask, nodata, band_array)
79
80         if not check_dtype_for_int(nodata):
81             band_array = band_array.astype(initial_dtype)
82         else:
83             band_array = band_array.astype(np.min_scalar_type(nodata))
84
85         band_array = np.expand_dims(band_array, axis=0)
86
87         if i == 0:
88             out_array = band_array.copy()
89         else:
90             out_array = np.vstack((out_array, band_array))
91
92         current_transform = f"transformation {i + 1}"
93         current_settings = {
94             "band_origin": bands[i],
95             "threshold": thresholds[i],
96             "nodata": nodata,
97         }
98
99         out_settings[current_transform] = current_settings
100
101     out_meta = raster.meta.copy()
102     out_meta.update({"count": len(bands), "nodata": nodata, "dtype": out_array.dtype.name})
103
104     return out_array, out_meta, out_settings

```

## 8.2 Clip

---

```
clip_transform(raster, limits, bands=None, nodata=None)
```

Clips data based on specified upper and lower limits.

Takes one nodata value that will be ignored in calculations. Replaces values below the lower limit and above the upper limit with provided values, respectively. Works both one-sided and two-sided but raises error if no limits provided.

If no band/column selection specified, all bands/columns will be used. If a parameter contains only 1 entry, it will be applied for all bands. The limits can be set for each band individually.

### Parameters:

Name	Type	Description	Default
raster	DatasetReader	Data object to be transformed.	required
bands	Optional[Sequence[int]]	Selection of bands to be transformed.	None
limits	Sequence[Tuple[Optional[Number], Optional[Number]]]	Lower and upper limits (lower, upper) as real values.	required
nodata	Optional[Number]	Nodata value to be considered.	None

### Returns:

Name	Type	Description
out_array	ndarray	The transformed data.
out_meta	dict	Updated metadata.
out_settings	dict	Log of input settings and calculated statistics if available.

### Raises:

Type	Description
InvalidRasterBandException	The input contains invalid band numbers.
NonMatchingParameterLengthsException	The input does not match the number of selected bands.
InvalidParameterValueException	The input does not match the requirements (values, order of values).

Source code in `eis_toolkit/transformations/clip.py` 

```

42 @beartype
43 def clip_transform( # type: ignore[no-any-unimported]
44     raster: rasterio.io.DatasetReader,
45     limits: Sequence[Tuple[Optional[Number], Optional[Number]]],
46     bands: Optional[Sequence[int]] = None,
47     nodata: Optional[Number] = None,
48 ) -> Tuple[np.ndarray, dict, dict]:
49     """
50     Clips data based on specified upper and lower limits.
51
52     Takes one nodata value that will be ignored in calculations.
53     Replaces values below the lower limit and above the upper limit with provided values, respectively.
54     Works both one-sided and two-sided but raises error if no limits provided.
55
56     If no band/column selection specified, all bands/columns will be used.
57     If a parameter contains only 1 entry, it will be applied for all bands.
58     The limits can be set for each band individually.
59
60     Args:
61         raster: Data object to be transformed.
62         bands: Selection of bands to be transformed.
63         limits: Lower and upper limits (lower, upper) as real values.
64         nodata: Nodata value to be considered.
65
66     Returns:
67         out_array: The transformed data.
68         out_meta: Updated metadata.
69         out_settings: Log of input settings and calculated statistics if available.
70
71     Raises:
72         InvalidRasterBandException: The input contains invalid band numbers.
73         NonMatchingParameterLengthsException: The input does not match the number of selected bands.
74         InvalidParameterValueException: The input does not match the requirements (values, order of values).
75     """
76     bands = list(range(1, raster.count + 1)) if bands is None else bands
77     nodata = raster.nodata if nodata is None else nodata
78
79     if check_raster_bands(raster, bands) is False:
80         raise InvalidRasterBandException("Invalid band selection")
81
82     if check_parameter_length(bands, limits) is False:
83         raise NonMatchingParameterLengthsException("Invalid limit length.")
84
85     for item in limits:
86         if item.count(None) == len(item):
87             raise InvalidParameterValueException(f"Limit values all None: {item}.")
88
89         if not check_minmax_position(item):
90             raise InvalidParameterValueException(f"Invalid min-max values provided: {item}.")
91
92     expanded_args = expand_and_zip(bands, limits)
93     limits = [element[1] for element in expanded_args]
94
95     out_settings = {}
96
97     for i in range(0, len(bands)):
98         band_array = raster.read(bands[i])
99         initial_dtype = band_array.dtype
100
101         band_array = cast_array_to_float(band_array, cast_int=True)
102         band_array = nodata_to_nan(band_array, nodata_value=nodata)
103
104         band_array = _clip_transform(band_array, limits=limits[i])
105
106         band_array = nan_to_nodata(band_array, nodata_value=nodata)
107         band_array = cast_array_to_int(band_array, scalar=nodata, initial_dtype=initial_dtype)
108
109         band_array = np.expand_dims(band_array, axis=0)
110
111         if i == 0:
112             out_array = band_array.copy()
113         else:
114             out_array = np.vstack((out_array, band_array))
115
116         current_transform = f"transformation {i + 1}"
117         current_settings = {
118             "band_origin": bands[i],
119             "limit_lower": cast_scalar_to_int(limits[i][0]),
120             "limit_upper": cast_scalar_to_int(limits[i][1]),
121             "nodata": cast_scalar_to_int(nodata),
122         }
123
124         out_settings[current_transform] = current_settings
125
126     out_meta = raster.meta.copy()
127     out_meta.update({"count": len(bands), "nodata": nodata, "dtype": out_array.dtype.name})
128
129     return out_array, out_meta, out_settings

```

## 8.3 Linear

---

```
min_max_scaling(raster, bands=None, new_range=[(0, 1)], nodata=None)
```

Normalize data based on a specified new range.

Uses the provided new minimum and maximum to transform data into the new interval. Takes one nodata value that will be ignored in calculations.

If no band/column selection specified, all bands/columns will be used. The `new_range` can be set for each band individually. If a parameter contains only 1 entry, it will be applied for all bands.

### Parameters:

Name	Type	Description	Default
<code>raster</code>	<code>DatasetReader</code>	Data object to be transformed.	required
<code>bands</code>	<code>Optional[Sequence[int]]</code>	Selection of bands to be transformed.	None
<code>new_range</code>	<code>Sequence[Tuple[Number, Number]]</code>	The new interval data will be transformed into. First value corresponds to min, second to max.	<code>[(0, 1)]</code>
<code>nodata</code>	<code>Optional[Number]</code>	Nodata value to be considered.	None

### Returns:

Name	Type	Description
<code>out_array</code>	<code>ndarray</code>	The transformed data.
<code>out_meta</code>	<code>dict</code>	Updated metadata.
<code>out_settings</code>	<code>dict</code>	Log of input settings and calculated statistics if available.

### Raises:

Type	Description
<code>InvalidRasterBandException</code>	The input contains invalid band numbers.
<code>NonMatchingParameterLengthsException</code>	The input does not match the number of selected bands.
<code>InvalidParameterValueException</code>	The input does not match the requirements (values, order of values).

Source code in `eis_toolkit/transformations/linear.py` 

```

125 @beartype
126 def min_max_scaling( # type: ignore[no-any-unimported]
127     raster: rasterio.io.DatasetReader,
128     bands: Optional[Sequence[int]] = None,
129     new_range: Sequence[Tuple[Number, Number]] = [(0, 1)],
130     nodata: Optional[Number] = None,
131 ) -> Tuple[np.ndarray, dict, dict]:
132     """
133     Normalize data based on a specified new range.
134
135     Uses the provided new minimum and maximum to transform data into the new interval.
136     Takes one nodata value that will be ignored in calculations.
137
138     If no band/column selection specified, all bands/columns will be used.
139     The new_range can be set for each band individually.
140     If a parameter contains only 1 entry, it will be applied for all bands.
141
142     Args:
143         raster: Data object to be transformed.
144         bands: Selection of bands to be transformed.
145         new_range: The new interval data will be transformed into. First value corresponds to min, second to max.
146         nodata: Nodata value to be considered.
147
148     Returns:
149         out_array: The transformed data.
150         out_meta: Updated metadata.
151         out_settings: Log of input settings and calculated statistics if available.
152
153     Raises:
154         InvalidRasterBandException: The input contains invalid band numbers.
155         NonMatchingParameterLengthsException: The input does not match the number of selected bands.
156         InvalidParameterValueException: The input does not match the requirements (values, order of values).
157     """
158     bands = list(range(1, raster.count + 1)) if bands is None else bands
159     nodata = raster.nodata if nodata is None else nodata
160
161     if check_raster_bands(raster, bands) is False:
162         raise InvalidRasterBandException("Invalid band selection")
163
164     if check_parameter_length(bands, new_range) is False:
165         raise NonMatchingParameterLengthsException("Invalid new_range length")
166
167     for item in new_range:
168         if not check_minmax_position(item):
169             raise InvalidParameterValueException(f"Invalid min-max values provided: {item}")
170
171     expanded_args = expand_and_zip(bands, new_range)
172     new_range = [element[1] for element in expanded_args]
173
174     out_settings = {}
175     out_decimals = set_max_precision()
176
177     for i in range(0, len(bands)):
178         band_array = raster.read(bands[i])
179         band_array = cast_array_to_float(band_array, cast_int=True)
180         band_array = replace_values(band_array, values_to_replace=[nodata, np.inf], replace_value=np.nan)
181
182         band_array = _min_max_scaling(band_array.astype(np.float64), new_range=new_range[i])
183
184         band_array = truncate_decimal_places(band_array, decimal_places=out_decimals)
185         band_array = nan_to_nodata(band_array, nodata_value=nodata)
186         band_array = cast_array_to_float(band_array, scalar=nodata, cast_float=True)
187
188         band_array = np.expand_dims(band_array, axis=0)
189
190         if i == 0:
191             out_array = band_array.copy()
192         else:
193             out_array = np.vstack((out_array, band_array))
194
195         current_transform = f"transformation {i + 1}"
196         current_settings = {
197             "band_origin": bands[i],
198             "scaled_min": new_range[i][0],
199             "scaled_max": new_range[i][1],
200             "nodata": nodata,
201             "decimal_places": out_decimals,
202         }
203
204         out_settings[current_transform] = current_settings
205
206     out_meta = raster.meta.copy()
207     out_meta.update({"count": len(bands), "nodata": nodata, "dtype": out_array.dtype.name})
208
209     return out_array, out_meta, out_settings

```

`z_score_normalization(raster, bands=None, nodata=None)`

Normalize data based on mean and standard deviation.

Results will have a mean = 0 and standard deviation = 1. Takes one nodata value that will be ignored in calculations.

If no band/column selection specified, all bands/columns will be used. If a parameter contains only 1 entry, it will be applied for all bands.

**Parameters:**

Name	Type	Description	Default
<code>raster</code>	<code>DatasetReader</code>	Data object to be transformed.	required
<code>bands</code>	<code>Optional[Sequence[int]]</code>	Selection of bands to be transformed.	None
<code>nodata</code>	<code>Optional[Number]</code>	Nodata value to be considered.	None

**Returns:**

Name	Type	Description
<code>out_array</code>	<code>ndarray</code>	The transformed data.
<code>out_meta</code>	<code>dict</code>	Updated metadata.
<code>out_settings</code>	<code>dict</code>	Log of input settings and calculated statistics if available.

**Raises:**

Type	Description
<code>InvalidRasterBandException</code>	The input contains invalid band numbers.
<code>NonMatchingParameterLengthsException</code>	The input does not match the number of selected bands.



Source code in `eis_toolkit/transformations/linear.py` 

```

52 @beartype
53 def z_score_normalization( # type: ignore[no-any-unimported]
54     raster: rasterio.io.DatasetReader,
55     bands: Optional[Sequence[int]] = None,
56     nodata: Optional[Number] = None,
57 ) -> Tuple[np.ndarray, dict, dict]:
58     """
59     Normalize data based on mean and standard deviation.
60
61     Results will have a mean = 0 and standard deviation = 1.
62     Takes one nodata value that will be ignored in calculations.
63
64     If no band/column selection specified, all bands/columns will be used.
65     If a parameter contains only 1 entry, it will be applied for all bands.
66
67     Args:
68         raster: Data object to be transformed.
69         bands: Selection of bands to be transformed.
70         nodata: Nodata value to be considered.
71
72     Returns:
73         out_array: The transformed data.
74         out_meta: Updated metadata.
75         out_settings: Log of input settings and calculated statistics if available.
76
77     Raises:
78         InvalidRasterBandException: The input contains invalid band numbers.
79         NonMatchingParameterLengthsException: The input does not match the number of selected bands.
80     """
81     bands = list(range(1, raster.count + 1)) if bands is None else bands
82     nodata = raster.nodata if nodata is None else nodata
83
84     if check_raster_bands(raster, bands) is False:
85         raise InvalidRasterBandException("Invalid band selection.")
86
87     out_settings = {}
88     out_decimals = set_max_precision()
89
90     for i in range(0, len(bands)):
91         band_array = raster.read(bands[i])
92         band_array = cast_array_to_float(band_array, cast_int=True)
93         band_array = replace_values(band_array, values_to_replace=[nodata, np.inf], replace_value=np.nan)
94
95         band_array, mean_array, sd_array = _z_score_normalization(band_array.astype(np.float64))
96
97         band_array = truncate_decimal_places(band_array, decimal_places=out_decimals)
98         band_array = nan_to_nodata(band_array, nodata_value=nodata)
99         band_array = cast_array_to_float(band_array, scalar=nodata, cast_float=True)
100
101         band_array = np.expand_dims(band_array, axis=0)
102
103         if i == 0:
104             out_array = band_array.copy()
105         else:
106             out_array = np.vstack((out_array, band_array))
107
108         current_transform = f"transformation {i + 1}"
109         current_settings = {
110             "band_origin": bands[i],
111             "original_mean": truncate_decimal_places(mean_array, decimal_places=out_decimals),
112             "original_sd": truncate_decimal_places(sd_array, decimal_places=out_decimals),
113             "nodata": nodata,
114             "decimal_places": out_decimals,
115         }
116
117         out_settings[current_transform] = current_settings
118
119     out_meta = raster.meta.copy()
120     out_meta.update({"count": len(bands), "nodata": nodata, "dtype": out_array.dtype.name})
121
122     return out_array, out_meta, out_settings

```

## 8.4 Logarithmic

---

```
log_transform(raster, bands=None, log_transform=['log2'], nodata=None)
```

Perform a logarithmic transformation on the provided data.

Takes one nodata value that will be ignored in calculations. Negative values will not be considered for transformation and replaced by the specific nodata value.

If no band/column selection specified, all bands/columns will be used. If a parameter contains only 1 entry, it will be applied for all bands. The log\_transform can be set for each band individually.

### Parameters:

Name	Type	Description	Default
raster	DatasetReader	Data object to be transformed.	required
bands	Optional[Sequence[int]]	Selection of bands to be transformed.	None
log_transform	Sequence[str]	The base for logarithmic transformation. Valid values 'ln', 'log2' and 'log10'.	['log2']
nodata	Optional[Number]	Nodata value to be considered.	None

### Returns:

Name	Type	Description
out_array	ndarray	The transformed data.
out_meta	dict	Updated metadata.
out_settings	dict	Log of input settings and calculated statistics if available.

### Raises:

Type	Description
InvalidRasterBandException	The input contains invalid band numbers.
NonMatchingParameterLengthsException	The input does not match the number of selected bands
InvalidParameterValueException	The input does not match the requirements (values, order of values)

Source code in `eis_toolkit/transformations/logarithmic.py` 

```

49 @beartype
50 def log_transform( # type: ignore[no-any-unimported]
51     raster: rasterio.io.DatasetReader,
52     bands: Optional[Sequence[int]] = None,
53     log_transform: Sequence[str] = ["log2"],
54     nodata: Optional[Number] = None,
55 ) -> Tuple[np.ndarray, dict, dict]:
56     """
57     Perform a logarithmic transformation on the provided data.
58
59     Takes one nodata value that will be ignored in calculations.
60     Negative values will not be considered for transformation and replaced by the specific nodata value.
61
62     If no band/column selection specified, all bands/columns will be used.
63     If a parameter contains only 1 entry, it will be applied for all bands.
64     The log_transform can be set for each band individually.
65
66     Args:
67         raster: Data object to be transformed.
68         bands: Selection of bands to be transformed.
69         log_transform: The base for logarithmic transformation. Valid values 'ln', 'log2' and 'log10'.
70         nodata: Nodata value to be considered.
71
72     Returns:
73         out_array: The transformed data.
74         out_meta: Updated metadata.
75         out_settings: Log of input settings and calculated statistics if available.
76
77     Raises:
78         InvalidRasterBandException: The input contains invalid band numbers.
79         NonMatchingParameterLengthsException: The input does not match the number of selected bands
80         InvalidParameterValueException: The input does not match the requirements (values, order of values)
81     """
82     bands = list(range(1, raster.count + 1)) if bands is None else bands
83     nodata = raster.nodata if nodata is None else nodata
84
85     if check_raster_bands(raster, bands) is False:
86         raise InvalidRasterBandException("Invalid band selection")
87
88     if check_parameter_length(bands, log_transform) is False:
89         raise NonMatchingParameterLengthsException("Invalid length for log-base values.")
90
91     for item in log_transform:
92         if not (item == "ln" or item == "log2" or item == "log10"):
93             raise InvalidParameterValueException(f"Invalid method: {item}.")
94
95     expanded_args = expand_and_zip(bands, log_transform)
96     log_transform = [element[1] for element in expanded_args]
97
98     out_settings = {}
99     out_decimals = set_max_precision()
100
101     for i in range(0, len(bands)):
102         band_array = raster.read(bands[i])
103         band_array = cast_array_to_float(band_array, cast_int=True)
104         band_array = replace_values(band_array, values_to_replace=[nodata, np.inf], replace_value=np.nan)
105         band_array[band_array <= 0] = np.nan
106
107         if log_transform[i] == "ln":
108             band_array = _log_transform_ln(band_array.astype(np.float64))
109         elif log_transform[i] == "log2":
110             band_array = _log_transform_log2(band_array.astype(np.float64))
111         elif log_transform[i] == "log10":
112             band_array = _log_transform_log10(band_array.astype(np.float64))
113
114         band_array = truncate_decimal_places(band_array, decimal_places=out_decimals)
115         band_array = nan_to_nodata(band_array, nodata_value=nodata)
116         band_array = cast_array_to_float(band_array, scalar=nodata, cast_float=True)
117
118         band_array = np.expand_dims(band_array, axis=0)
119
120         if i == 0:
121             out_array = band_array.copy()
122         else:
123             out_array = np.vstack((out_array, band_array))
124
125         current_transform = f"transformation {i + 1}"
126         current_settings = {
127             "band_origin": bands[i],
128             "log_transform": log_transform[i],
129             "nodata": nodata,
130             "decimal_places": out_decimals,
131         }
132
133         out_settings[current_transform] = current_settings
134
135     out_meta = raster.meta.copy()
136     out_meta.update({"count": len(bands), "nodata": nodata, "dtype": out_array.dtype.name})
137
138     return out_array, out_meta, out_settings

```

## 8.5 Sigmoid

---

```
sigmoid_transform(raster, bands=None, bounds=[(0, 1)], slope=[1], center=True,
nodata=None)
```

Transform data into a sigmoid-shape based on a specified new range.

Uses the provided new minimum and maximum, shift and slope parameters to transform the data. Takes one nodata value that will be ignored in calculations.

If no band/column selection specified, all bands/columns will be used. If a parameter contains only 1 entry, it will be applied for all bands. The bounds and slope values can be set for each band individually.

### Parameters:

Name	Type	Description	Default
raster	DatasetReader	Data object to be transformed.	required
bands	Optional[Sequence[int]]	Selection of bands to be transformed.	None
bounds	Sequence[Tuple[Number, Number]]	Boundaries for the calculation of the sigmoid function (lower, upper).	[(0, 1)]
slope	Sequence[Number]	Value which modifies the slope of the resulting sigmoid-curve.	[1]
center	bool	Center array values around mean = 0 before sigmoid transformation.	True
nodata	Optional[Number]	Nodata value to be considered.	None

### Returns:

Name	Type	Description
out_array	ndarray	The transformed data.
out_meta	dict	Updated metadata.
out_settings	dict	Log of input settings and calculated statistics if available.

### Raises:

Type	Description
InvalidRasterBandException	The input contains invalid band numbers.
NonMatchingParameterLengthsException	The input does not match the number of selected bands.
InvalidParameterValueException	The input does not match the requirements (values, order of values)

Source code in `eis_toolkit/transformations/sigmoid.py` 

```

42 @beartype
43 def sigmoid_transform( # type: ignore[no-any-unimported]
44     raster: rasterio.io.DatasetReader,
45     bands: Optional[Sequence[int]] = None,
46     bounds: Sequence[Tuple[Number, Number]] = [(0, 1)],
47     slope: Sequence[Number] = [1],
48     center: bool = True,
49     nodata: Optional[Number] = None,
50 ) -> Tuple[np.ndarray, dict, dict]:
51     """
52     Transform data into a sigmoid-shape based on a specified new range.
53
54     Uses the provided new minimum and maximum, shift and slope parameters to transform the data.
55     Takes one nodata value that will be ignored in calculations.
56
57     If no band/column selection specified, all bands/columns will be used.
58     If a parameter contains only 1 entry, it will be applied for all bands.
59     The bounds and slope values can be set for each band individually.
60
61     Args:
62         raster: Data object to be transformed.
63         bands: Selection of bands to be transformed.
64         bounds: Boundaries for the calculation of the sigmoid function (lower, upper).
65         slope: Value which modifies the slope of the resulting sigmoid-curve.
66         center: Center array values around mean = 0 before sigmoid transformation.
67         nodata: Nodata value to be considered.
68
69     Returns:
70         out_array: The transformed data.
71         out_meta: Updated metadata.
72         out_settings: Log of input settings and calculated statistics if available.
73
74     Raises:
75         InvalidRasterBandException: The input contains invalid band numbers.
76         NonMatchingParameterLengthsException: The input does not match the number of selected bands.
77         InvalidParameterValueException: The input does not match the requirements (values, order of values)
78     """
79     bands = list(range(1, raster.count + 1)) if bands is None else bands
80     nodata = raster.nodata if nodata is None else nodata
81
82     if check_raster_bands(raster, bands) is False:
83         raise InvalidRasterBandException("Invalid band selection")
84
85     for parameter_name, parameter in [("bounds", bounds), ("slope", slope)]:
86         if check_parameter_length(bands, parameter) is False:
87             raise NonMatchingParameterLengthsException(f"Invalid length for {parameter_name}.")
88
89     for item in bounds:
90         if check_minmax_position(item) is False:
91             raise InvalidParameterValueException(f"Invalid min-max values provided: {item}.")
92
93     expanded_args = expand_and_zip(bands, bounds, slope)
94     bounds = [element[1] for element in expanded_args]
95     slope = [element[2] for element in expanded_args]
96
97     out_settings = {}
98     out_decimals = set_max_precision()
99
100    for i in range(0, len(bands)):
101        band_array = raster.read(bands[i])
102        band_array = cast_array_to_float(band_array, cast_int=True)
103        band_array = replace_values(band_array, values_to_replace=[nodata, np.inf], replace_value=np.nan)
104
105        band_array = _sigmoid_transform(band_array.astype(np.float64), bounds=bounds[i], slope=slope[i], center=center)
106
107        band_array = truncate_decimal_places(band_array, decimal_places=out_decimals)
108        band_array = nan_to_nodata(band_array, nodata_value=nodata)
109        band_array = cast_array_to_float(band_array, scalar=nodata, cast_float=True)
110
111        band_array = np.expand_dims(band_array, axis=0)
112
113        if i == 0:
114            out_array = band_array.copy()
115        else:
116            out_array = np.vstack((out_array, band_array))
117
118        current_transform = f"transformation {i + 1}"
119        current_settings = {
120            "band_origin": bands[i],
121            "bound_lower": truncate_decimal_places(bounds[i][0], decimal_places=out_decimals),
122            "bound_upper": truncate_decimal_places(bounds[i][1], decimal_places=out_decimals),
123            "slope": slope[i],
124            "center": center,
125            "nodata": nodata,
126            "decimal_places": out_decimals,
127        }
128
129        out_settings[current_transform] = current_settings
130
131    out_meta = raster.meta.copy()
132    out_meta.update({"count": len(bands), "nodata": nodata, "dtype": out_array.dtype.name})
133
134    return out_array, out_meta, out_settings

```

## 8.6 Winsorize

---

```
winsorize(raster, percentiles, bands=None, inside=False, nodata=None)
```

Winsorize data based on specified percentile values.

Takes one nodata value that will be ignored in calculations. Replaces values between [minimum, lower percentile] and [upper percentile, maximum] if provided. Works both one-sided and two-sided but raises error if no percentile values provided.

Percentiles are symmetrical, i.e. percentile\_lower = 10 corresponds to the interval [min, 10%]. And percentile\_upper = 10 corresponds to the interval [90%, max]. I.e. percentile\_lower = 0 refers to the minimum and percentile\_upper = 0 to the data maximum.

Calculation of percentiles is ambiguous. Users can choose whether to use the value for replacement from inside or outside of the respective interval. Example: Given the np.array[5 10 12 15 20 24 27 30 35] and percentiles(10, 10), the calculated percentiles are (5, 35) for inside and (10, 30) for outside. This results in [5 10 12 15 20 24 27 30 35] and [10 10 12 15 20 24 27 30 30], respectively.

If no band/column selection specified, all bands/columns will be used. If a parameter contains only 1 entry, it will be applied for all bands. The percentiles can be set for each band individually, but inside parameter is same for all bands.

### Parameters:

Name	Type	Description	Default
raster	DatasetReader	Data object to be transformed.	required
bands	Optional[Sequence[int]]	Selection of bands to be transformed.	None
percentiles	Sequence[Tuple[Optional[Number], Optional[Number]]]	Lower and upper percentile values (lower, upper) between [0, 100].	required
inside	bool	Whether to use the value for replacement from the left or right of the calculated percentile.	False
nodata	Optional[Number]	Nodata value to be considered.	None

### Returns:

Name	Type	Description
out_array	ndarray	The transformed data.
out_meta	dict	Updated metadata.
out_settings	dict	Log of input settings and calculated statistics if available.

**Raises:**

<b>Type</b>	<b>Description</b>
InvalidRasterBandException	The input contains invalid band numbers.
NonMatchingParameterLengthsException	The input does not match the number of selected bands.
InvalidParameterValueException	The input does not match the requirements (values, order of values)

Source code in [eis\\_toolkit/transformations/winsorize.py](#) 



54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
68  
69  
70  
71  
72  
73  
74  
75  
76  
77  
78  
79  
80  
81  
82  
83  
84  
85  
86  
87  
88  
89  
90  
91  
92  
93  
94  
95  
96  
97  
98  
99  
100  
101  
102  
103  
104  
105  
106  
107  
108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161

```

162 @beartype
163 def winsorize( # type: ignore[no-any-unimported]
    raster: rasterio.io.DatasetReader,
    percentiles: Sequence[Tuple[Optional[Number], Optional[Number]]],
    bands: Optional[Sequence[int]] = None,
    inside: bool = False,
    nodata: Optional[Number] = None,
) -> Tuple[np.ndarray, dict, dict]:
    """
    Winsorize data based on specified percentile values.

    Takes one nodata value that will be ignored in calculations.
    Replaces values between [minimum, lower percentile] and [upper percentile, maximum] if provided.
    Works both one-sided and two-sided but raises error if no percentile values provided.

    Percentiles are symmetrical, i.e. percentile_lower = 10 corresponds to the interval [min, 10%].
    And percentile_upper = 10 corresponds to the interval [90%, max].
    I.e. percentile_lower = 0 refers to the minimum and percentile_upper = 0 to the data maximum.

    Calculation of percentiles is ambiguous. Users can choose whether to use the value
    for replacement from inside or outside of the respective interval. Example:
    Given the np.array([5 10 12 15 20 24 27 30 35]) and percentiles([10, 10]), the calculated
    percentiles are (5, 35) for inside and (10, 30) for outside.
    This results in [5 10 12 15 20 24 27 30 35] and [10 10 12 15 20 24 27 30 30], respectively.

    If no band/column selection specified, all bands/columns will be used.
    If a parameter contains only 1 entry, it will be applied for all bands.
    The percentiles can be set for each band individually, but inside parameter is same for all bands.

    Args:
        raster: Data object to be transformed.
        bands: Selection of bands to be transformed.
        percentiles: Lower and upper percentile values (lower, upper) between [0, 100].
        inside: Whether to use the value for replacement from the left or right of the calculated percentile.
        nodata: Nodata value to be considered.

    Returns:
        out_array: The transformed data.
        out_meta: Updated metadata.
        out_settings: Log of input settings and calculated statistics if available.

    Raises:
        InvalidRasterBandException: The input contains invalid band numbers.
        NonMatchingParameterLengthsException: The input does not match the number of selected bands.
        InvalidParameterValueException: The input does not match the requirements (values, order of values)
    """
    bands = list(range(1, raster.count + 1)) if bands is None else bands
    nodata = raster.nodata if nodata is None else nodata

    if check_raster_bands(raster, bands) is False:
        raise InvalidRasterBandException("Invalid band selection")

    if check_parameter_length(bands, percentiles) is False:
        raise NonMatchingParameterLengthsException("Invalid length for percentiles.")

    for item in percentiles:
        if item.count(None) == len(item):
            raise InvalidParameterValueException(f"Percentile values all None: {item}.")

        if None not in item and sum(item) >= 100:
            raise InvalidParameterValueException(f"Sum >= 100: {item}.")

        if item[0] is not None and not (0 < item[0] < 100):
            raise InvalidParameterValueException(f"Invalid lower percentile value: {item}.")

        if item[1] is not None and not (0 < item[1] < 100):
            raise InvalidParameterValueException(f"Invalid upper percentile value: {item}.")

    expanded_args = expand_and_zip(bands, percentiles)
    percentiles = [element[1] for element in expanded_args]

    out_settings = {}

    for i in range(0, len(bands)):
        band_array = raster.read(bands[i])
        initial_dtype = band_array.dtype

        band_array = cast_array_to_float(band_array, cast_int=True)
        band_array = nodata_to_nan(band_array, nodata_value=nodata)

        band_array, calculated_lower, calculated_upper = _winsorize(
            band_array, percentiles=percentiles[i], inside=inside
        )

        band_array = nan_to_nodata(band_array, nodata_value=nodata)
        band_array = cast_array_to_int(band_array, scalar=nodata, initial_dtype=initial_dtype)

        band_array = np.expand_dims(band_array, axis=0)

        if i == 0:
            out_array = band_array.copy()
        else:
            out_array = np.vstack((out_array, band_array))

        current_transform = f"transformation {i + 1}"
        current_settings = {
            "band_origin": bands[i],
            "percentile_lower": cast_scalar_to_int(percentiles[i][0]),
            "percentile_upper": cast_scalar_to_int(percentiles[i][1]),
            "calculated_lower": cast_scalar_to_int(calculated_lower),
            "calculated_upper": cast_scalar_to_int(calculated_upper),
            "nodata": cast_scalar_to_int(nodata),
        }

        out_settings[current_transform] = current_settings

    out_meta = raster.meta.copy()
    out_meta.update({"count": len(bands), "nodata": nodata, "dtype": out_array.dtype.name})

```

```
return out_array, out_meta, out_settings
```

## 9. Validation

---

### 9.1 Calculate AUC

---

#### `calculate_auc(x_values, y_values)`

Calculate area under curve (AUC).

Calculates AUC for curve. X-axis should be either proportion of area or false positive rate. Y-axis should be always true positive rate. AUC is calculated with `sklearn.metrics.auc` which uses trapezoidal rule for calculation.

#### Parameters:

Name	Type	Description	Default
<code>x_values</code>	<code>ndarray</code>	Either proportion of area or false positive rate values.	required
<code>y_values</code>	<code>ndarray</code>	True positive rate values.	required

#### Returns:

Type	Description
<code>float</code>	The area under curve.

#### Raises:

Type	Description
<code>InvalidParameterValueException</code>	<code>x_values</code> or <code>y_values</code> are out of bounds.

#### Source code in `eis_toolkit/validation/calculate_auc.py`

```

13 @beartype
14 def calculate_auc(x_values: np.ndarray, y_values: np.ndarray) -> float:
15     """Calculate area under curve (AUC).
16
17     Calculates AUC for curve. X-axis should be either proportion of area or false positive rate. Y-axis should be
18     always true positive rate. AUC is calculated with sklearn.metrics.auc which uses trapezoidal rule for calculation.
19
20     Args:
21         x_values: Either proportion of area or false positive rate values.
22         y_values: True positive rate values.
23
24     Returns:
25         The area under curve.
26
27     Raises:
28         InvalidParameterValueException: x_values or y_values are out of bounds.
29     """
30     if x_values.max() > 1 or x_values.min() < 0:
31         raise InvalidParameterValueException("x_values should be within range 0-1")
32
33     if y_values.max() > 1 or y_values.min() < 0:
34         raise InvalidParameterValueException("y_values should be within range 0-1")
35
36     auc_value = _calculate_auc(x_values=x_values, y_values=y_values)
37     return auc_value

```

## 9.2 Calculate base metrics

---

```
calculate_base_metrics(raster, deposits, band=1, negatives=None)
```

Calculate true positive rate, proportion of area and false positive rate values for different thresholds.

Function calculates true positive rate, proportion of area and false positive rate values for different thresholds which are determined from inputted deposit locations and mineral prospectivity map. Note that calculation of false positive rate is optional and is only done if negative point locations are provided.

### Parameters:

Name	Type	Description	Default
raster	DatasetReader	Mineral prospectivity map or evidence layer.	required
deposits	GeoDataFrame	Mineral deposit locations as points.	required
band	int	Band index of the mineral prospectivity map. Defaults to 1.	1
negatives	Optional[GeoDataFrame]	Negative locations as points.	None

### Returns:

Type	Description
DataFrame	DataFrame containing true positive rate, proportion of area, threshold values and false positive rate (optional) values.

### Raises:

Type	Description
NonMatchingCrsException	The raster and point data are not in the same CRS.
NotApplicableGeometryTypeException	The input geometries contain non-point features.

Source code in `eis_toolkit/validation/calculate_base_metrics.py` 

```

73 @beartype
74 def calculate_base_metrics(
75     raster: rasterio.io.DatasetReader,
76     deposits: geopandas.GeoDataFrame,
77     band: int = 1,
78     negatives: Optional[geopandas.GeoDataFrame] = None,
79 ) -> pd.DataFrame:
80     """Calculate true positive rate, proportion of area and false positive rate values for different thresholds.
81
82     Function calculates true positive rate, proportion of area and false positive rate values for different thresholds
83     which are determined from inputted deposit locations and mineral prospectivity map. Note that calculation of false
84     positive rate is optional and is only done if negative point locations are provided.
85
86     Args:
87         raster: Mineral prospectivity map or evidence layer.
88         deposits: Mineral deposit locations as points.
89         band: Band index of the mineral prospectivity map. Defaults to 1.
90         negatives: Negative locations as points.
91
92     Returns:
93         DataFrame containing true positive rate, proportion of area, threshold values and false positive
94         rate (optional) values.
95
96     Raises:
97         NonMatchingCrsException: The raster and point data are not in the same CRS.
98         NotApplicableGeometryTypeException: The input geometries contain non-point features.
99     """
100     if negatives is not None:
101         geometries = pd.concat([deposits, negatives]).geometry
102     else:
103         geometries = deposits["geometry"]
104
105     if not check_matching_crs(
106         objects=[raster, geometries],
107     ):
108         raise NonMatchingCrsException("The raster and deposits are not in the same CRS.")
109
110     if not check_geometry_types(
111         geometries=geometries,
112         allowed_types=["Point"],
113     ):
114         raise NotApplicableGeometryTypeException("The input geometries contain non-point features.")
115
116     base_metrics = _calculate_base_metrics(raster=raster, deposits=deposits, band=band, negatives=negatives)
117
118     return base_metrics

```

## 9.3 Get P-A plot intersection point

---

```
get_pa_intersection(true_positive_rate_values, proportion_of_area_values,
threshold_values)
```

Calculate the intersection point for prediction rate and area curves in (P-A plot).

Threshold\_values values act as x-axis for both curves. Prediction rate curve uses true positive rate for y-axis. Area curve uses inverted proportion of area as y-axis.

#### Parameters:

Name	Type	Description	Default
true_positive_rate_values	ndarray	True positive rate values, values should be within range 0-1.	required
proportion_of_area_values	ndarray	Proportion of area values, values should be within range 0-1.	required
threshold_values	ndarray	Threshold values that were used to calculate true positive rate and proportion of area.	required

#### Returns:

Type	Description
Tuple[float, float]	X and y coordinates of the intersection point.

#### Raises:

Type	Description
InvalidParameterValueException	true_positive_rate_values or proportion_of_area_values values are out of bounds.

#### Source code in eis\_toolkit/validation/get\_pa\_intersection.py

```

20 @beartype
21 def get_pa_intersection(
22     true_positive_rate_values: np.ndarray, proportion_of_area_values: np.ndarray, threshold_values: np.ndarray
23 ) -> Tuple[float, float]:
24     """Calculate the intersection point for prediction rate and area curves in (P-A plot).
25
26     Threshold_values values act as x-axis for both curves. Prediction rate curve uses true positive rate for y-axis.
27     Area curve uses inverted proportion of area as y-axis.
28
29     Args:
30         true_positive_rate_values: True positive rate values, values should be within range 0-1.
31         proportion_of_area_values: Proportion of area values, values should be within range 0-1.
32         threshold_values: Threshold values that were used to calculate true positive rate and proportion of area.
33
34     Returns:
35         X and y coordinates of the intersection point.
36
37     Raises:
38         InvalidParameterValueException: true_positive_rate_values or proportion_of_area_values values are out of bounds.
39     """
40     if true_positive_rate_values.max() > 1 or true_positive_rate_values.min() < 0:
41         raise InvalidParameterValueException("true_positive_rate_values values should be within range 0-1")
42
43     if proportion_of_area_values.max() > 1 or proportion_of_area_values.min() < 0:
44         raise InvalidParameterValueException("proportion_of_area_values values should be within range 0-1")
45
46     intersection = _get_pa_intersection(
47         true_positive_rate_values=true_positive_rate_values,
48         proportion_of_area_values=proportion_of_area_values,
49         threshold_values=threshold_values,
50     )
51
52     return intersection.x, intersection.y

```



## 9.4 Plot correlation matrix

---

```
plot_correlation_matrix(matrix, annotate=True, cmap=None, plot_title=None,
                        **kwargs)
```

Create a Seaborn heatmap to visualize correlation matrix.

### Parameters:

Name	Type	Description	Default
<code>matrix</code>	<code>DataFrame</code>	Correlation matrix as a DataFrame.	required
<code>annotate</code>	<code>bool</code>	If plot squares should display the correlation values. Defaults to True.	<code>True</code>
<code>cmap</code>	<code>Optional[ListedColormap]</code>	Colormap for plotting. Optional parameter. Defaults to None, in which case a default colormap is used.	<code>None</code>
<code>plot_title</code>	<code>Optional[str]</code>	Title of the plot. Optional parameter, defaults to none (no title).	<code>None</code>
<code>**kwargs</code>	<code>dict</code>	Additional parameters to pass to Seaborn and matplotlib.	<code>{}</code>

### Returns:

Type	Description
<code>Axes</code>	Matplotlib axes object with the produced plot.

### Raises:

Type	Description
<code>EmptyDataFrameException</code>	Input matrix is empty.

Source code in `eis_toolkit/validation/plot_correlation_matrix.py` 

```

11 def plot_correlation_matrix(
12     matrix: pd.DataFrame,
13     annotate: bool = True,
14     cmap: Optional[matplotlib.colors.ListedColormap] = None,
15     plot_title: Optional[str] = None,
16     **kwargs: dict
17 ) -> matplotlib.axes.Axes:
18     """
19     Create a Seaborn heatmap to visualize correlation matrix.
20
21     Args:
22         matrix: Correlation matrix as a DataFrame.
23         annotate: If plot squares should display the correlation values. Defaults to True.
24         cmap: Colormap for plotting. Optional parameter. Defaults to None, in which
25             case a default colormap is used.
26         plot_title: Title of the plot. Optional parameter, defaults to none (no title).
27         **kwargs: Additional parameters to pass to Seaborn and matplotlib.
28
29     Returns:
30         Matplotlib axes object with the produced plot.
31
32     Raises:
33         EmptyDataFrameException: Input matrix is empty.
34     """
35     if matrix.empty:
36         raise exceptions.EmptyDataFrameException("Input matrix DataFrame is empty.")
37
38     # Mask for the upper triangle of the heatmap
39     mask = np.triu(np.ones_like(matrix, dtype=bool))
40
41     if cmap is None:
42         # Generate a default diverging colormap
43         cmap = sns.diverging_palette(230, 20, as_cmap=True)
44
45     ax = sns.heatmap(
46         matrix,
47         mask=mask,
48         cmap=cmap,
49         vmax=0.3,
50         center=0,
51         square=True,
52         linewidths=0.5,
53         annot=annotate,
54         cbar_kws={"shrink": 0.5},
55         **kwargs
56     )
57     if plot_title is not None:
58         ax.set_title(plot_title)
59
60     return ax

```

## 9.5 Plot prediction-area (P-A) curves

---

```
plot_prediction_area_curves(true_positive_rate_values,
                             proportion_of_area_values, threshold_values)
```

Plot prediction-area (P-A) plot.

Plots prediction area plot that can be used to evaluate mineral prospectivity maps and evidential layers. See e.g., Yousefi and Carranza (2015).

### Parameters:

Name	Type	Description	Default
true_positive_rate_values	ndarray	True positive rate values.	required
proportion_of_area_values	ndarray	Proportion of area values.	required
threshold_values	ndarray	Threshold values.	required

### Returns:

Type	Description
Figure	P-A plot figure object.

### Raises:

Type	Description
InvalidParameterValueException	true_positive_rate_values or proportion_of_area_values values are out of bounds.

### References ▼

Yousefi, Mahyar, and Emmanuel John M. Carranza. "Fuzzification of continuous-value spatial evidence for mineral prospectivity mapping." *Computers & Geosciences* 74 (2015): 97-109.

Source code in `eis_toolkit/validation/plot_prediction_area_curves.py` 

```

41 @beartype
42 def plot_prediction_area_curves(
43     true_positive_rate_values: np.ndarray, proportion_of_area_values: np.ndarray, threshold_values: np.ndarray
44 ) -> matplotlib.figure.Figure:
45     """Plot prediction-area (P-A) plot.
46
47     Plots prediction area plot that can be used to evaluate mineral prospectivity maps and evidential layers. See e.g.,
48     Yousefi and Carranza (2015).
49
50     Args:
51         true_positive_rate_values: True positive rate values.
52         proportion_of_area_values: Proportion of area values.
53         threshold_values: Threshold values.
54
55     Returns:
56         P-A plot figure object.
57
58     Raises:
59         InvalidParameterValueException: true_positive_rate_values or proportion_of_area_values values are out of bounds.
60
61     References:
62         Yousefi, Mahyar, and Emmanuel John M. Carranza. "Fuzzification of continuous-value spatial evidence for mineral
63         prospectivity mapping." Computers & Geosciences 74 (2015): 97-109.
64     """
65     if true_positive_rate_values.max() > 1 or true_positive_rate_values.min() < 0:
66         raise InvalidParameterValueException("true_positive_rate values should be within range 0-1")
67
68     if proportion_of_area_values.max() > 1 or proportion_of_area_values.min() < 0:
69         raise InvalidParameterValueException("proportion_of_area values should be within range 0-1")
70
71     fig = _plot_prediction_area_curves(
72         true_positive_rate_values=true_positive_rate_values,
73         proportion_of_area_values=proportion_of_area_values,
74         threshold_values=threshold_values,
75     )
76     return fig

```

## 9.6 Plot rate curve

---

```
plot_rate_curve(x_values, y_values, plot_type='success_rate')
```

Plot success rate, prediction rate or ROC curve.

Plot type depends on `plot_type` argument. Y-axis is always true positive rate, while x-axis can be either false positive rate (roc) or proportion of area (success and prediction rate) depending on plot type.

### Parameters:

Name	Type	Description	Default
<code>x_values</code>	<code>ndarray</code>	False positive rate values or proportion of area values.	required
<code>y_values</code>	<code>ndarray</code>	True positive rate values.	required
<code>plot_type</code>	<code>str</code>	Plot type. Can be either: "success_rate", "prediction_rate" or "roc".	'success_rate'

### Returns:

Type	Description
<code>Figure</code>	Success rate, prediction rate or ROC plot figure object.

### Raises:

Type	Description
<code>InvalidParameterValueException</code>	Invalid plot type.
<code>InvalidParameterValueException</code>	<code>x_values</code> or <code>y_values</code> are out of bounds.

Source code in `eis_toolkit/validation/plot_rate_curve.py` 

```

27 @beartype
28 def plot_rate_curve(
29     x_values: np.ndarray,
30     y_values: np.ndarray,
31     plot_type: str = "success_rate",
32 ) -> matplotlib.figure.Figure:
33     """Plot success rate, prediction rate or ROC curve.
34
35     Plot type depends on plot_type argument. Y-axis is always true positive rate, while x-axis can be either false
36     positive rate (roc) or proportion of area (success and prediction rate) depending on plot type.
37
38     Args:
39         x_values: False positive rate values or proportion of area values.
40         y_values: True positive rate values.
41         plot_type: Plot type. Can be either: "success_rate", "prediction_rate" or "roc".
42
43     Returns:
44         Success rate, prediction rate or ROC plot figure object.
45
46     Raises:
47         InvalidParameterValueException: Invalid plot type.
48         InvalidParameterValueException: x_values or y_values are out of bounds.
49     """
50     if plot_type == "success_rate":
51         label = "Success rate"
52         xlab = "Proportion of area"
53     elif plot_type == "prediction_rate":
54         label = "Prediction rate"
55         xlab = "Proportion of area"
56     elif plot_type == "roc":
57         label = "ROC"
58         xlab = "False positive rate"
59     else:
60         raise InvalidParameterValueException("Invalid plot type")
61
62     if x_values.max() > 1 or x_values.min() < 0:
63         raise InvalidParameterValueException("x_values should be within range 0-1")
64
65     if y_values.max() > 1 or y_values.min() < 0:
66         raise InvalidParameterValueException("y_values should be within range 0-1")
67
68     fig = _plot_rate_curve(x_values=x_values, y_values=y_values, label=label, xlab=xlab)
69
70     return fig

```

## 10. Vector processing

---

### 10.1 Cell-Based Association

---

```
cell_based_association(cell_size, geodata, output_path, column=None,  
subset_target_attribute_values=None, add_name=None, add_buffer=None)
```

Creation of CBA matrix.

Initializes a CBA matrix from a vector file. The mesh is calculated according to the geometries contained in this file and the size of cells. Allows to add multiple vector data to the matrix, based on targeted shapes and/or attributes.

**Parameters:**



<b>Name</b>	<b>Type</b>	<b>Description</b>	<b>Default</b>
cell_size	int	Size of the cells.	required
geodata	List[GeoDataFrame]	GeoDataFrame to create the CBA matrix. Additional GeoDataFrame(s) can be imputed to add to the CBA matrix.	required
output_path	str	Name of the saved .tif file.	required
column	Optional[List[str]]	Name of the column of interest. If no attribute is specified, then an artificial attribute is created representing the presence or absence of the geometries of this file for each cell of the CBA grid. A categorical attribute will generate as many columns (binary) in the CBA matrix than values considered of interest (dummification). See parameter . Additional column(s) can be imputed for each added GeoDataFrame(s).	None
subset_target_attribute_values	Optional[List[Union[None, list, str]]]	List of values of interest of the target attribute, in case a categorical target attribute has been specified. Allows to filter a subset of relevant values. Additional values can be imputed for each added GeoDataFrame(s).	None
add_name	Optional[List[Union[str, None]]]	Name of the column(s) to add to the matrix.	None
add_buffer	Optional[List[Union[Number, bool]]]	Allow the use of a buffer around shapes before the intersection with CBA cells for the added	None


<b>Name</b>	<b>Type</b>	<b>Description</b>	<b>Default</b>
		GeoDataFrame(s). Minimize border effects or allow increasing positive samples (i.e. cells with mineralization). The size of the buffer is computed using the CRS (if projected CRS in meters: value in meters).	

---

**Returns:**

<b>Type</b>	<b>Description</b>
GeoDataFrame	CBA matrix is created.

---

Source code in `eis_toolkit/vector_processing/cell_based_association.py` 

```

18 @beartype
19 def cell_based_association(
20     cell_size: int,
21     geodata: List[gpd.GeoDataFrame],
22     output_path: str,
23     column: Optional[List[str]] = None,
24     subset_target_attribute_values: Optional[List[Union[None, list, str]]] = None,
25     add_name: Optional[List[Union[str, None]]] = None,
26     add_buffer: Optional[List[Union[Number, bool]]] = None,
27 ) -> gpd.GeoDataFrame:
28     """Creation of CBA matrix.
29
30     Initializes a CBA matrix from a vector file. The mesh is calculated
31     according to the geometries contained in this file and the size of cells.
32     Allows to add multiple vector data to the matrix, based on targeted shapes
33     and/or attributes.
34
35     Args:
36         cell_size: Size of the cells.
37         geodata: GeoDataFrame to create the CBA matrix. Additional
38             GeoDataFrame(s) can be imputed to add to the CBA matrix.
39         output_path: Name of the saved .tif file.
40         column: Name of the column of interest. If no attribute is specified,
41             then an artificial attribute is created representing the presence
42             or absence of the geometries of this file for each cell of the CBA
43             grid. A categorical attribute will generate as many columns (binary)
44             in the CBA matrix than values considered of interest (dummification).
45             See parameter <subset_target_attribute_values>. Additional
46             column(s) can be imputed for each added GeoDataFrame(s).
47         subset_target_attribute_values: List of values of interest of the
48             target attribute, in case a categorical target attribute has been
49             specified. Allows to filter a subset of relevant values. Additional
50             values can be imputed for each added GeoDataFrame(s).
51         add_name: Name of the column(s) to add to the matrix.
52         add_buffer: Allow the use of a buffer around shapes before the
53             intersection with CBA cells for the added GeoDataFrame(s). Minimize
54             border effects or allow increasing positive samples (i.e. cells
55             with mineralization). The size of the buffer is computed using the
56             CRS (if projected CRS in meters: value in meters).
57
58     Returns:
59         CBA matrix is created.
60     """
61
62     # Swapping None to list values
63     if column is None:
64         column = [""]
65     if add_buffer is None:
66         add_buffer = [False]
67
68     # Consistency checks on input data
69     for frame in geodata:
70         if frame.empty:
71             raise exceptions.EmptyDataFrameException("The input GeoDataFrame is empty.")
72
73     if cell_size <= 0:
74         raise exceptions.InvalidParameterValueException("Expected cell size to be positive and non-zero.")
75
76     add_buffer = [False if x == 0 else x for x in add_buffer]
77     if any(num < 0 for num in add_buffer):
78         raise exceptions.InvalidParameterValueException("Expected buffer value to be positive, null or False.")
79
80     for i, name in enumerate(column):
81         if column[i] == "":
82             if subset_target_attribute_values[i] is not None:
83                 raise exceptions.InvalidParameterValueException("Can't use subset of values if no column is targeted.")
84             elif column[i] not in geodata[i]:
85                 raise exceptions.InvalidColumnException("Targeted column not found in the GeoDataFrame.")
86
87     for i, subset in enumerate(subset_target_attribute_values):
88         if subset is not None:
89             for value in subset:
90                 if value not in geodata[i][column[i]].unique():
91                     raise exceptions.InvalidParameterValueException(
92                         "Subset of value(s) not found in the targeted column."
93                 )
94
95     # Computation
96     for i, data in enumerate(geodata):
97         if i == 0:
98             # Initialization of the CBA matrix
99             grid, cba = _init_from_vector_data(cell_size, geodata[0], column[0], subset_target_attribute_values[0])
100         else:
101             # If necessary, adding data to matrix
102             cba = _add_layer(
103                 cba,
104                 grid,
105                 geodata[i],
106                 column[i],
107                 subset_target_attribute_values[i],
108                 add_name[i - 1],
109                 add_buffer[i - 1],
110             )
111
112     # Export
113     _to_raster(cba, output_path)
114
115     return cba

```

## 10.2 Distance computation

---

### distance\_computation(raster\_profile, geometries)

Calculate distance from raster cell to nearest geometry.

#### Parameters:

Name	Type	Description	Default
raster_profile	Union[Profile, dict]	The raster profile of the raster in which the distances to the nearest geometry are determined.	required
geometries	GeoDataFrame	The geometries to determine distance to.	required

#### Returns:

Type	Description
ndarray	A 2D numpy array with the distances computed.

#### Source code in eis\_toolkit/vector\_processing/distance\_computation.py

```

21 @beartype
22 def distance_computation(raster_profile: Union[profiles.Profile, dict], geometries: gpd.GeoDataFrame) -> np.ndarray:
23     """Calculate distance from raster cell to nearest geometry.
24
25     Args:
26         raster_profile: The raster profile of the raster in which the distances
27             to the nearest geometry are determined.
28         geometries: The geometries to determine distance to.
29
30     Returns:
31         A 2D numpy array with the distances computed.
32
33     """
34     if raster_profile.get("crs") != geometries.crs:
35         raise exceptions.NonMatchingCrsException("Expected coordinate systems to match between raster and geometries. ")
36     if geometries.shape[0] == 0:
37         raise exceptions.EmptyDataFrameException("Expected GeoDataFrame to not be empty.")
38
39     raster_width = raster_profile.get("width")
40     raster_height = raster_profile.get("height")
41
42     if not isinstance(raster_width, int) or not isinstance(raster_height, int):
43         raise exceptions.InvalidParameterValueException(
44             f"Expected raster_profile to contain integer width and height. {raster_profile}"
45         )
46
47     raster_transform = raster_profile.get("transform")
48
49     if not isinstance(raster_transform, transform.Affine):
50         raise exceptions.InvalidParameterValueException(
51             f"Expected raster_profile to contain an affine transformation. {raster_profile}"
52         )
53
54     return _distance_computation(
55         raster_width=raster_width, raster_height=raster_height, raster_transform=raster_transform, geometries=geometries
56     )

```

## 10.3 IDW

---

```
idw(geodataframe, target_column, resolution, extent=None, power=2)
```

Calculate inverse distance weighted (IDW) interpolation.

### Parameters:

Name	Type	Description	Default
<code>geodataframe</code>	<code>GeoDataFrame</code>	The vector dataframe to be interpolated.	required
<code>target_column</code>	<code>str</code>	The column name with values for each geometry.	required
<code>resolution</code>	<code>Tuple[Number, Number]</code>	The resolution i.e. cell size of the output raster as (pixel_size_x, pixel_size_y).	required
<code>extent</code>	<code>Optional[Tuple[Number, Number, Number, Number]]</code>	The extent of the output raster as (x_min, x_max, y_min, y_max). If None, calculate extent from the input vector data.	None
<code>power</code>	<code>Number</code>	The value for determining the rate at which the weights decrease. As power increases, the weights for distant points decrease rapidly. Defaults to 2.	2

### Returns:

Type	Description
<code>Tuple[ndarray, dict]</code>	Rasterized vector data and metadata.

### Raises:

Type	Description
<code>EmptyDataFrameException</code>	The input GeoDataFrame is empty.
<code>InvalidParameterValueException</code>	Invalid resolution or target_column.

Source code in `eis_toolkit/vector_processing/idw_interpolation.py` 

```

79 @beartype
80 def idw(
81     geodataframe: gpd.GeoDataFrame,
82     target_column: str,
83     resolution: Tuple[Number, Number],
84     extent: Optional[Tuple[Number, Number, Number, Number]] = None,
85     power: Number = 2,
86 ) -> Tuple[np.ndarray, dict]:
87     """Calculate inverse distance weighted (IDW) interpolation.
88
89     Args:
90         geodataframe: The vector dataframe to be interpolated.
91         target_column: The column name with values for each geometry.
92         resolution: The resolution i.e. cell size of the output raster as (pixel_size_x, pixel_size_y).
93         extent: The extent of the output raster as (x_min, x_max, y_min, y_max).
94             If None, calculate extent from the input vector data.
95         power: The value for determining the rate at which the weights decrease.
96             As power increases, the weights for distant points decrease rapidly.
97             Defaults to 2.
98
99     Returns:
100         Rasterized vector data and metadata.
101
102     Raises:
103         EmptyDataFrameException: The input GeoDataFrame is empty.
104         InvalidParameterValueException: Invalid resolution or target_column.
105     """
106
107     if geodataframe.shape[0] == 0:
108         raise EmptyDataFrameException("Expected geodataframe to contain geometries.")
109
110     if target_column not in geodataframe.columns:
111         raise InvalidParameterValueException(
112             f"Expected target_column ({target_column}) to be contained in geodataframe columns."
113         )
114
115     if resolution[0] <= 0 or resolution[1] <= 0:
116         raise InvalidParameterValueException("Expected height and width greater than zero.")
117
118     interpolated_values, out_meta = _idw_interpolation(geodataframe, target_column, resolution, power, extent)
119
120     return interpolated_values, out_meta

```

## 10.4 Kriging interpolation

---

```
kriging(data, target_column, resolution, extent=None, variogram_model='linear',  
coordinates_type='geographic', method='ordinary')
```

Perform Kriging interpolation on the input data.

**Parameters:**

<b>Name</b>	<b>Type</b>	<b>Description</b>	<b>Default</b>
<code>data</code>	<code>GeoDataFrame</code>	GeoDataFrame containing the input data.	required
<code>target_column</code>	<code>str</code>	The column name with values for each geometry.	required
<code>resolution</code>	<code>Tuple[Number, Number]</code>	The resolution i.e. cell size of the output raster as ( <code>pixel_size_x</code> , <code>pixel_size_y</code> ).	required
<code>extent</code>	<code>Optional[Tuple[Number, Number, Number, Number]]</code>	The extent of the output raster as ( <code>x_min</code> , <code>x_max</code> , <code>y_min</code> , <code>y_max</code> ). If None, calculate extent from the input vector data.	None
<code>variogram_model</code>	<code>Literal[linear, power, gaussian, spherical, exponential]</code>	Variogram model to be used. Either 'linear', 'power', 'gaussian', 'spherical' or 'exponential'. Defaults to 'linear'.	'linear'
<code>coordinates_type</code>	<code>Literal[euclidean, geographic]</code>	Determines are coordinates on a plane ('euclidean') or a sphere ('geographic'). Used only in ordinary kriging. Defaults to 'geographic'.	'geographic'
<code>method</code>	<code>Literal[ordinary, universal]</code>	Ordinary or universal kriging. Defaults to 'ordinary'.	'ordinary'

**Returns:**

<b>Type</b>	<b>Description</b>
<code>Tuple[ndarray, dict]</code>	Grid containing the interpolated values and metadata.



**Raises:**

Type	Description
EmptyDataFrameException	The input GeoDataFrame is empty.
InvalidParameterValueException	Target column name is invalid or resolution is not greater than zero.

**Source code in `eis_toolkit/vector_processing/kriging_interpolation.py`**

```

58 @beartype
59 def kriging(
60     data: gpd.GeoDataFrame,
61     target_column: str,
62     resolution: Tuple[Number, Number],
63     extent: Optional[Tuple[Number, Number, Number, Number]] = None,
64     variogram_model: Literal["linear", "power", "gaussian", "spherical", "exponential"] = "linear",
65     coordinates_type: Literal["euclidean", "geographic"] = "geographic",
66     method: Literal["ordinary", "universal"] = "ordinary",
67 ) -> Tuple[np.ndarray, dict]:
68     """
69     Perform Kriging interpolation on the input data.
70
71     Args:
72     data: GeoDataFrame containing the input data.
73     target_column: The column name with values for each geometry.
74     resolution: The resolution i.e. cell size of the output raster as (pixel_size_x, pixel_size_y).
75     extent: The extent of the output raster as (x_min, x_max, y_min, y_max).
76         If None, calculate extent from the input vector data.
77     variogram_model: Variogram model to be used.
78         Either 'linear', 'power', 'gaussian', 'spherical' or 'exponential'. Defaults to 'linear'.
79     coordinates_type: Determines are coordinates on a plane ('euclidean') or a sphere ('geographic').
80         Used only in ordinary kriging. Defaults to 'geographic'.
81     method: Ordinary or universal kriging. Defaults to 'ordinary'.
82
83     Returns:
84     Grid containing the interpolated values and metadata.
85
86     Raises:
87     EmptyDataFrameException: The input GeoDataFrame is empty.
88     InvalidParameterValueException: Target column name is invalid or resolution is not greater than zero.
89     """
90
91     if data.empty:
92         raise EmptyDataFrameException("The input GeoDataFrame is empty.")
93
94     if target_column not in data.columns:
95         raise InvalidParameterValueException(
96             f"Expected target_column {{target_column}} to be contained in geodataframe columns."
97         )
98
99     if resolution[0] <= 0 or resolution[1] <= 0:
100         raise InvalidParameterValueException("The resolution must be greater than zero.")
101
102     data_interpolated, out_meta = kriging(
103         data, target_column, resolution, extent, variogram_model, coordinates_type, method
104     )
105
106     return data_interpolated, out_meta

```

## 10.5 Rasterize vector

---

```
rasterize_vector(geodataframe, resolution=None, value_column=None,
default_value=1.0, fill_value=0.0, base_raster_profile=None, buffer_value=None,
merge_strategy='replace')
```

Transform vector data into raster data.

### Parameters:

Name	Type	Description	Default
geodataframe	GeoDataFrame	The vector dataframe to be rasterized.	required
resolution	Optional[float]	The resolution i.e. cell size of the output raster. Optional if base_raster_profile is given.	None
value_column	Optional[str]	The column name with values for each geometry. If None, then default_value is used for all geometries.	None
default_value	float	Default value burned into raster cells based on geometries.	1.0
base_raster_profile	Optional[Union[Profile, dict]]	Base raster profile to be used for determining the grid on which vectors are burned in. If None, the geometries and provided resolution value are used to compute grid.	None
fill_value	float	Value used outside the burned/ rasterized geometry cells.	0.0
buffer_value	Optional[float]	For adding a buffer around passed geometries before rasterization.	None
merge_strategy	Literal[replace, add]	How to handle overlapping geometries. "add" causes overlapping geometries to add together the values while "replace" does not. Adding them together is the basis for density computations where the density can be calculated by using a default value of 1.0 and the sum in each cell is the count of intersecting geometries.	'replace'

### Returns:

Type	Description
Tuple[ndarray, dict]	Rasterized vector data and metadata.

Source code in `eis_toolkit/vector_processing/rasterize_vector.py` 

```

11 @beartype
12 def rasterize_vector(
13     geodataframe: gpd.GeoDataFrame,
14     resolution: Optional[float] = None,
15     value_column: Optional[str] = None,
16     default_value: float = 1.0,
17     fill_value: float = 0.0,
18     base_raster_profile: Optional[Union[profiles.Profile, dict]] = None,
19     buffer_value: Optional[float] = None,
20     merge_strategy: Literal["replace", "add"] = "replace",
21 ) -> Tuple[np.ndarray, dict]:
22     """Transform vector data into raster data.
23
24     Args:
25         geodataframe: The vector dataframe to be rasterized.
26         resolution: The resolution i.e. cell size of the output raster.
27             Optional if base_raster_profile is given.
28         value_column: The column name with values for each geometry.
29             If None, then default_value is used for all geometries.
30         default_value: Default value burned into raster cells based on geometries.
31         base_raster_profile: Base raster profile
32             to be used for determining the grid on which vectors are
33             burned in. If None, the geometries and provided resolution
34             value are used to compute grid.
35         fill_value: Value used outside the burned/rasterized geometry cells.
36         buffer_value: For adding a buffer around passed
37             geometries before rasterization.
38         merge_strategy: How to handle overlapping geometries.
39             "add" causes overlapping geometries to add together the
40             values while "replace" does not. Adding them together is the
41             basis for density computations where the density can be
42             calculated by using a default value of 1.0 and the sum in
43             each cell is the count of intersecting geometries.
44
45     Returns:
46         Rasterized vector data and metadata.
47     """
48
49     if geodataframe.shape[0] == 0:
50         # Empty GeoDataFrame
51         raise exceptions.EmptyDataFrameException("Expected geodataframe to contain geometries.")
52
53     if resolution is None and base_raster_profile is None:
54         raise exceptions.InvalidParameterValueException(
55             "Expected either resolution or base_raster_profile to be given."
56         )
57     if resolution is not None and resolution <= 0:
58         raise exceptions.NumericValueSignException(
59             f"Expected a positive value resolution ({dict(resolution=resolution)})"
60         )
61     if value_column is not None and value_column not in geodataframe.columns:
62         raise exceptions.InvalidParameterValueException(
63             f"Expected value_column ({value_column}) to be contained in geodataframe columns."
64         )
65     if buffer_value is not None and buffer_value < 0:
66         raise exceptions.NumericValueSignException(
67             f"Expected a positive buffer_value ({dict(buffer_value=buffer_value)})"
68         )
69
70     if base_raster_profile is not None and not isinstance(base_raster_profile, (profiles.Profile, dict)):
71         raise exceptions.InvalidParameterValueException(
72             f"Expected base_raster_profile ({type(base_raster_profile)}) to be dict or rasterio.profiles.Profile."
73         )
74
75     if buffer_value is not None:
76         geodataframe = geodataframe.copy()
77         geodataframe["geometry"] = geodataframe["geometry"].apply(lambda geom: geom.buffer(buffer_value))
78
79     return rasterize_vector(
80         geodataframe=geodataframe,
81         value_column=value_column,
82         default_value=default_value,
83         fill_value=fill_value,
84         base_raster_profile=base_raster_profile,
85         resolution=resolution,
86         merge_alg=getattr(MergeAlg, merge_strategy),
87     )

```

## 10.6 Reproject vector

---

`reproject_vector(geodataframe, target_crs)`

Reprojects vector data to match given coordinate reference system (EPSG).

### Parameters:

Name	Type	Description	Default
<code>geodataframe</code>	<code>GeoDataFrame</code>	The vector dataframe to be reprojected.	required
<code>target_crs</code>	<code>int</code>	Target CRS as an EPSG code.	required

### Returns:

Type	Description
<code>GeoDataFrame</code>	Reprojected vector data.

### Source code in `eis_toolkit/vector_processing/reproject_vector.py`

```

7 @beartype
8 def reproject_vector(geodataframe: geopandas.GeoDataFrame, target_crs: int) -> geopandas.GeoDataFrame:
9     """Reprojects vector data to match given coordinate reference system (EPSG).
10
11     Args:
12         geodataframe: The vector dataframe to be reprojected.
13         target_crs: Target CRS as an EPSG code.
14
15     Returns:
16         Reprojected vector data.
17     """
18
19     if geodataframe.crs.to_epsg() == target_crs:
20         raise MatchingCredException("Vector data is already in the target CRS.")
21
22     reprojected_gdf = geodataframe.to_crs("epsg:" + str(target_crs))
23     return reprojected_gdf

```

## 10.7 Vector density

---

```
vector_density(geodataframe, resolution=None, base_raster_profile=None,
               buffer_value=None, statistic='density')
```

Compute density of geometries within raster.

### Parameters:

Name	Type	Description	Default
<code>geodataframe</code>	<code>GeoDataFrame</code>	The dataframe with vectors of which density is computed.	required
<code>resolution</code>	<code>Optional[float]</code>	The resolution i.e. cell size of the output raster. Optional if <code>base_raster_profile</code> is given.	None
<code>base_raster_profile</code>	<code>Optional[Union[Profile, dict]]</code>	Base raster profile to be used for determining the grid on which vectors are burned in. If None, the geometries and provided resolution value are used to compute grid.	None
<code>buffer_value</code>	<code>Optional[float]</code>	For adding a buffer around passed geometries before computing density.	None

### Returns:

Type	Description
<code>Tuple[ndarray, dict]</code>	Computed density of vector data and metadata.

Source code in `eis_toolkit/vector_processing/vector_density.py` 

```

10 @beartype
11 def vector_density(
12     geodataframe: gpd.GeoDataFrame,
13     resolution: Optional[float] = None,
14     base_raster_profile: Optional[Union[profiles.Profile, dict]] = None,
15     buffer_value: Optional[float] = None,
16     statistic: Literal["density", "count"] = "density",
17 ) -> Tuple[np.ndarray, dict]:
18     """Compute density of geometries within raster.
19
20     Args:
21         geodataframe: The dataframe with vectors
22             of which density is computed.
23         resolution: The resolution i.e. cell size of the output raster.
24             Optional if base_raster_profile is given.
25         base_raster_profile: Base raster profile
26             to be used for determining the grid on which vectors are
27             burned in. If None, the geometries and provided resolution
28             value are used to compute grid.
29         buffer_value: For adding a buffer around passed
30             geometries before computing density.
31
32     Returns:
33         Computed density of vector data and metadata.
34     """
35     out_raster_array, out_metadata = rasterize_vector(
36         geodataframe=geodataframe,
37         resolution=resolution,
38         base_raster_profile=base_raster_profile,
39         buffer_value=buffer_value,
40         value_column=None,
41         default_value=1.0,
42         fill_value=0.0,
43         merge_strategy="add",
44     )
45     max_count = np.max(out_raster_array)
46     if statistic == "count" or np.isclose(max_count, 0.0):
47         return out_raster_array, out_metadata
48     else:
49         return (out_raster_array / max_count), out_metadata

```