

# Shapefile Analysis and Similarity Calculation Report

Your Name

May 25, 2024

## Abstract

This report details the process and methodologies applied in exploring and analyzing shapefiles, as well as calculating similarities between different data points using various similarity metrics. The implementation is discussed with reference to the provided Python scripts.

## 1 Introduction

The purpose of this project is to develop a tool for exploring and analyzing shapefiles, with a particular focus on calculating similarity between different data points. The tool is implemented using Python and leverages several libraries including pandas, numpy, scikit-learn, and Tkinter for the user interface. The scripts provided perform tasks such as loading shapefiles, determining attribute types, and calculating similarities using various metrics.

## 2 Methodology

The project is divided into several scripts, each responsible for specific tasks. The main functionalities are detailed below.

### 2.1 Shapefile Explorer

The `ShpExplorer.py` script is responsible for loading shapefiles and determining the types of attributes contained within them. It includes functions to check if a value is a number, date, or string, and uses these checks to classify the type of each attribute.

```
import os
from osgeo import ogr
import datetime
from dateutil.parser import parse
import dateparser
```

```

import numpy as np
from collections import Counter
from concurrent.futures import ProcessPoolExecutor, as_completed
import pandas as pd

class ShpExplorer:
    ...

```

## 2.2 Similarity Calculation

The `Similarity.py` script provides a class `SimilarityCalculator` that computes similarity metrics between data points. The metrics implemented include cosine similarity, Euclidean distance, and Jaccard similarity.

```

import pandas as pd
import numpy as np
from sklearn.metrics.pairwise import cosine_similarity

class SimilarityCalculator:
    ...

```

## 2.3 User Interface

The `ui.py` script creates a graphical user interface using Tkinter, allowing users to upload shapefiles, select attributes, and compute similarity metrics between selected data points.

```

import sys
import os
import tkinter as tk
from tkinter import filedialog, messagebox
from tkinter import ttk
import pandas as pd
import geopandas as gpd
from ShpExplorer import ShpExplorer
from similarity.Similarity import SimilarityCalculator
...

```

# 3 Implementation

## 3.1 Loading Shapefiles

The shapefiles are loaded using the `ogr` library, and the attributes are extracted and classified. The classification is done in parallel using `ProcessPoolExecutor` to improve performance.

### 3.2 Calculating Similarity

The similarity between data points is calculated using the `SimilarityCalculator` class. This class provides methods to calculate cosine similarity, Euclidean distance, and Jaccard similarity, allowing users to select the most appropriate metric for their needs.

### 3.3 User Interface

The user interface provides a seamless experience for users to upload shapefiles, select the attributes they are interested in, and compute similarities. The results are displayed in a user-friendly manner, making it easy to interpret the similarities between data points.

## 4 Conclusions

This project successfully developed a tool for exploring and analyzing shapefiles, as well as calculating similarities between different data points. The use of parallel processing and a graphical user interface enhances the usability and performance of the tool. Future work could include the addition of more similarity metrics and the ability to handle larger datasets more efficiently.