

Task1: Data orchestration tools

Data orchestration tools are software platforms or frameworks designed to manage and streamline the process of collecting, processing, transforming, and moving data across various systems and applications within an organization. These tools play a crucial role in modern data pipelines, ensuring data flows efficiently and reliably from source to destination while often incorporating data quality and governance features. Here are some popular data orchestration tools as of my last knowledge update in September 2021:

1. **Apache NiFi:** An open-source data integration tool that provides an intuitive interface for designing data flows and automating data movement and transformation.
2. **Apache Airflow:** While primarily known as a workflow automation tool, Airflow is commonly used for data orchestration. It allows you to schedule, monitor, and manage complex data workflows.
3. **Talend:** Talend is an open-source data integration platform that offers data orchestration capabilities for ETL (Extract, Transform, Load) and data integration tasks.
4. **Apache Oozie:** A workflow scheduler system for Apache Hadoop that can be used to define and manage data processing jobs.
5. **Microsoft Azure Data Factory:** A cloud-based data integration service that allows you to create, schedule, and manage data pipelines.
6. **Google Cloud Dataflow:** A fully managed stream and batch data processing service that enables data transformation and orchestration on Google Cloud Platform.
7. **AWS Step Functions:** Part of Amazon Web Services, AWS Step Functions allows you to coordinate multiple AWS services into serverless workflows, including data processing tasks.
8. **Prefect:** An open-source data workflow management system that helps you schedule, orchestrate, and monitor data workflows with Python.
9. **Luigi:** An open-source Python module that helps you build complex and dynamic data pipelines. It was developed by Spotify.
10. **StreamSets:** A data integration platform designed for efficiently ingesting, processing, and moving data across various sources and destinations.
11. **Conductor by Netflix:** An open-source workflow orchestration platform that can be used for data pipeline management.
12. **Apache Beam:** Although primarily a data processing framework, Apache Beam includes tools for defining and executing data pipelines in various data processing engines, making it a versatile option for data orchestration.
13. **Dagster:** An open-source data orchestrator that provides a programming model for building data pipelines and includes features for testing, scheduling, and monitoring.

Task 2: SMOTE Analysis

SMOTE (Synthetic Minority Over-sampling Technique):

SMOTE is a technique used to address class imbalance in supervised machine learning. Class imbalance occurs when one class of the target variable has significantly fewer instances than the other class(es). This can lead to models that perform poorly, as they tend to be biased towards the majority class.

SMOTE works by creating synthetic examples of the minority class to balance the class distribution. Here's how it works:

1. For each instance in the minority class, SMOTE selects its k-nearest neighbors (typically, k is specified by the user).
2. It then generates synthetic examples by interpolating between the selected instance and one or more of its nearest neighbors. This creates new instances that belong to the minority class.
3. These synthetic instances are added to the original dataset, effectively balancing the class distribution.

SMOTE helps prevent the model from being biased towards the majority class and can lead to improved classification performance, especially when dealing with imbalanced datasets.

In summary, SMOTE is a technique used to address class imbalance in machine learning by generating synthetic examples of the minority class. It is not typically referred to as "SMOTE analysis" but rather as a data preprocessing step in the machine learning pipeline.

1. Class Imbalance Problem:

Class imbalance is a common issue in many machine learning datasets. It occurs when one class (the minority class) has significantly fewer examples than another class (the majority class). For example, in a binary classification problem where you're trying to predict whether a transaction is fraudulent or not, most transactions might be non-fraudulent, while only a small percentage are fraudulent.

2. The Problem with Class Imbalance:

Class imbalance can cause machine learning models to perform poorly, especially for the minority class. This is because the models tend to be biased towards the majority class, as they have more examples to learn from. As a result, they may have difficulty correctly identifying and classifying instances of the minority class.

3. SMOTE's Solution:

SMOTE is a technique designed to alleviate the class imbalance problem by oversampling the minority class. Here's how it works step by step:

a. **Identify Minority Class Instances:** SMOTE starts by identifying instances belonging to the minority class. These are the instances that need to be oversampled.

b. **Select Nearest Neighbors:** For each instance in the minority class, SMOTE selects a specified number of its nearest neighbors. The number of neighbors to select (often denoted as "k") is determined by the user.

c. **Create Synthetic Examples:** SMOTE generates synthetic examples by interpolating between the feature vectors of the selected instance and its nearest neighbours. It creates new instances that are combinations of the existing minority class instances.

d. **Balancing the Dataset:** These synthetic instances are then added to the original dataset, effectively increasing the number of minority class instances. This balances the class distribution.

4. Effect on Machine Learning Models:

By using SMOTE to oversample the minority class, machine learning models have a more balanced dataset to train on. This helps prevent them from being biased towards the majority class and allows them to better learn the characteristics of the minority class. As a result, the models are often more capable of correctly classifying minority class instances, leading to improved overall performance, especially in terms of metrics like precision, recall, and F1-score.

Task3: Comparison between famous web scrapping tools

1. Octoparse

Octoparse is an easy-to-use web scraping tool developed to accommodate complicated web scraping for non-coders. As an intelligent web scraper on both Windows and Mac OS, it automatically “guesses” the desired data fields for users, which saves a large amount of time and energy as you don’t need to manually select the data. It is powerful enough to deal with dynamic websites and interact with any sites in various ways, such as authentication, text input, selecting from drop-down menus, hovering over dynamic menus, infinite scrolling and many more. Octoparse offers cloud-based extraction (paid features) as well as local extraction (free). For precise scraping, Octoparse also has built-in XPath and Regular Expression tools to help users scrape data with high accuracy.

2. **Parsehub**

Parsehub is another non-programmer friendly software. Being a desktop application, Parsehub is supported in various systems such as Windows, Mac OS X, and Linux. Like Octoparse, Parsehub can deal with complicated web scraping scenarios mentioned earlier. Although Parsehub intends to offer an easy web scraping experience, a typical user will still need to be a bit technical to fully grasp many of its advanced functionalities.

3. **Dexi.io**

Dexi.io is a cloud-based web scraper providing development, hosting and scheduling services. Dexi.io can be very powerful but does require more advanced programming skills compared to Octoparse and Parsehub. With Dexi, three kinds of robots are available: extractors, crawlers, pipes. Dexi supports integration with many third-party services such as captcha solvers, cloud storage and many more.

4. **Mozenda**

Mozenda offers a cloud-based web scraping service, similar to that of Octoparse cloud extraction. Being one of the “oldest” web scraping software in the market, Mozenda performs with a high-level of consistency, has nice looking UI and everything else anyone may need to start on a web scraping project. There are two parts to Mozenda: the Mozenda Web Console and Agent Builder. The Mozenda agent builder is a Windows application used for building scraping projects and the web console is a web application allowing users to set schedules to run the projects or access the extracted data. Similar to Octoparse, Mozenda also relies on a Windows system and can be a bit tricky for Mac users.

5. **Import.io**

Famous for its “Magic” – automatically turning any website into structured data, Import.io has gained in popularity. However, many users found out that it was not really “magical” enough to handle various kinds of websites. Besides that, Import.io does have a nice well-guided interface,

supports real-time data retrieval through JSON REST-based and streaming APIs and it is a web application that can be run in various systems

Task4:RPA Tool in AWS

Amazon Web Services (AWS) offers a service called AWS RoboMaker, which is not a traditional Robotic Process Automation (RPA) tool but is designed for robotics development and simulation. AWS RoboMaker allows you to develop, test, and deploy robotic applications in the cloud. It provides tools and services for building and managing robotic workflows.

Robotic Process Automation (RPA) tools have gained popularity for automating repetitive, rule-based tasks in business processes. There are numerous RPA tools available, both open-source and commercial. Here are some well-known RPA tools as of my last knowledge update in September 2021:

1. **UiPath:** UiPath is one of the leading commercial RPA platforms known for its user-friendly interface and extensive capabilities. It provides a range of automation tools, including Studio for building automation workflows, Orchestrator for managing robots, and Robots for executing tasks.
2. **Automation Anywhere:** Automation Anywhere offers a comprehensive RPA platform that includes tools for building, deploying, and managing bots. It supports both attended and unattended automation and provides features for task scheduling and analytics.
3. **Blue Prism:** Blue Prism is an RPA platform known for its enterprise-grade capabilities. It offers a Digital Workforce that can automate various business processes. Blue Prism provides a robust control room for managing and monitoring bots.
4. **OpenSpan (now part of Pegasystems):** OpenSpan, now integrated into Pegasystems, provides RPA capabilities as part of a broader automation suite. It's known for its process analytics and workforce analytics features.
5. **AutomationEdge:** AutomationEdge is an RPA and IT automation platform that also includes cognitive automation capabilities. It offers automation for IT, business, and customer service processes.
6. **Kofax RPA (formerly Kapow):** Kofax RPA is designed for automating data-intensive tasks and processes. It includes features for web data extraction, automation of repetitive tasks, and integration with other systems.
7. **WorkFusion:** WorkFusion combines RPA with AI capabilities to automate data-driven tasks and processes. It includes features for data extraction, document processing, and cognitive automation.
8. **NICE Robotic Automation:** NICE offers a suite of automation solutions, including RPA, attended automation, and unattended automation. It's known for its capabilities in automating contact center operations.
9. **AntWorks:** AntWorks offers an integrated intelligent automation platform that combines RPA with machine learning and artificial intelligence. It focuses on automating complex, unstructured data processes.
10. **Microsoft Power Automate (formerly Microsoft Flow):** Part of the Microsoft Power Platform, Power Automate provides low-code automation for a wide range of Microsoft and

third-party applications. While not a traditional RPA tool, it can be used for workflow automation.

11. **ElectroNeek:** ElectroNeek is an RPA platform designed for small and medium-sized businesses. It offers features for building and managing bots, as well as automating various tasks.
12. **TagUI:** TagUI is an open-source RPA tool that uses a simple scripting language for automating web interactions. It's known for its simplicity and ease of use.
13. **Robocop:** Robocop offers an open-source RPA platform called Robcloud, which includes development tools, a robot orchestrator, and support for both attended and unattended automation.

Task5: Libraries used in webscrabbing

1. **Beautiful Soup:** Beautiful Soup is a Python library that is widely used for parsing HTML and XML documents. It provides a convenient way to extract data from web pages and navigate the HTML or XML structure.
2. **Requests:** While not a scraping library per se, the Python **requests** library is often used in conjunction with Beautiful Soup or other libraries to make HTTP requests to websites and retrieve web pages' HTML content.
3. **Scrapy:** Scrapy is a Python framework for web scraping. It provides a more structured approach to building web scrapers and offers features like automatic crawling and following links.
4. **Selenium:** Selenium is a versatile tool used for web scraping and automation. It allows you to automate interactions with web pages, including filling out forms, clicking buttons, and scrolling. It's particularly useful for websites with dynamic content loaded via JavaScript.
5. **Puppeteer:** Puppeteer is a Node.js library developed by Google for controlling headless Chrome or Chromium browsers. It's often used for web scraping and automating tasks that require browser interaction.
6. **XPath:** XPath is a query language for navigating XML and HTML documents. While not a library itself, XPath expressions can be used with libraries like lxml (for Python) or in browser developer tools to locate and extract specific elements from web pages.
7. **lxml:** lxml is a Python library that provides a fast and efficient way to parse and manipulate HTML and XML documents. It supports both parsing and element extraction using XPath.
8. **Goutte:** Goutte is a web scraping library for PHP. It provides a simple API for making HTTP requests and extracting data from web pages.
9. **Nokogiri:** Nokogiri is a popular web scraping library for Ruby. It allows you to parse and manipulate HTML and XML documents.
10. **HtmlAgilityPack:** HtmlAgilityPack is a .NET library for parsing and manipulating HTML documents. It's commonly used for web scraping in C#.

11. **Jsoup:** Jsoup is a Java library for working with HTML documents. It provides a convenient API for parsing and manipulating HTML content in Java applications.
12. **Cheerio:** Cheerio is a fast, flexible, and lean implementation of jQuery for Node.js. It's often used for web scraping and parsing HTML documents in JavaScript.
13. **PyQuery:** PyQuery is a Python library that allows you to make jQuery queries on XML documents. It's useful for parsing and extracting data from HTML and XML files.
14. **MechanicalSoup:** MechanicalSoup is a Python library that simplifies the process of interacting with websites and filling out forms. It's built on top of BeautifulSoup and provides session management capabilities.

Task6: DOM in JS

In JavaScript, the Document Object Model (DOM) is a programming interface for web documents. It represents the structure of an HTML or XML document as a tree-like structure, where each element in the document is a node in the tree. The DOM allows developers to interact with and manipulate web pages dynamically using JavaScript.

Here are some key concepts related to the DOM in JavaScript:

1. **Document Object:** The **document** object is the root of the DOM tree and represents the entire web page. You can access it using **document** in JavaScript. For example, **document.getElementById('myElement')** would find an element with the specified ID in the document.
2. **DOM Nodes:** DOM nodes are the building blocks of the DOM tree. There are various types of nodes, including elements, text, attributes, and more. Elements represent HTML tags, while text nodes represent the text within elements.
3. **DOM Tree Structure:** The DOM tree structure represents the hierarchy of elements in a web page. Each element can have child nodes (other elements or text) and parent nodes (the element that contains it).
4. **Accessing Elements:** You can access elements in the DOM using methods like **getElementById**, **getElementsByTagName**, **getElementsByClassName**, and **querySelector**. These methods allow you to retrieve specific elements based on their attributes or CSS selectors.
5. **Modifying Elements:** You can change the content, attributes, and style of elements in the DOM dynamically using JavaScript. Common methods for manipulation include **innerHTML**, **textContent**, **setAttribute**, and **style**.
6. **Creating Elements:** You can create new elements and add them to the DOM using methods like **createElement**, **appendChild**, **insertBefore**, and **insertAdjacentHTML**.
7. **Event Handling:** You can attach event listeners to DOM elements to respond to user interactions like clicks, key presses, and mouse movements. Common event handling methods include **addEventListener** and inline event attributes like **onclick**.

8. **Traversing the DOM:** You can navigate the DOM tree by moving between nodes. Common methods for traversal include **parentNode**, **childNodes**, **nextSibling**, **previousSibling**, and **querySelectorAll**.
9. **Document Manipulation:** You can manipulate the entire document, including its structure, content, and attributes, using the **document** object.