

Student Information:

- Name: Anshu Kumari
- University: IMS
- Degree Program: BCA
- Email: kumarianshu301130@gmail.com
- GitHub/Portfolio: anshu-codes

Project Title:

Analyzing the Impact of Biodiversity Genome Datasets in Scientific Publications Using Python, Machine Learning, and Text Processing

Synopsis:

The use of genomic datasets in biodiversity research has surged in recent years, providing vital insights into the genetic diversity, conservation, and evolution of species. However, the full impact of these datasets on scientific publications has yet to be systematically assessed. Given the vast amount of literature on biodiversity genomics, extracting meaningful insights from these scientific publications requires advanced data processing, text mining, and machine learning techniques.

This project aims to develop a Python-based tool that utilizes machine learning and natural language processing (NLP) techniques to analyze the impact of biodiversity genome datasets in scientific publications. The tool will identify patterns in the usage, citation, and influence of genomic datasets, and generate metrics on their contribution to scientific discoveries and conservation efforts.

By applying machine learning models for text classification, citation analysis, and dataset mention recognition, this project will provide an analytical framework to quantify how biodiversity genome datasets are being used in scientific research. It will also offer a visualization interface to explore trends in dataset utilization across various species, research areas, and geographic locations.

Benefits to the Community:

1. Improved Understanding of Genomic Dataset Usage:

- This project will provide insights into how genomic data is used in biodiversity research, offering a better understanding of its role in advancing knowledge in conservation genetics, ecology, and species conservation.

2. Open-Source Tool Development:

- The project will deliver an open-source tool with machine learning models and text mining capabilities that will be available for researchers, conservationists, and policymakers to analyze genomic data in their respective domains.

3. Automated Citation and Impact Analysis:

- By leveraging Python and machine learning, the tool will automate the analysis of citation patterns and dataset mentions, allowing for large-scale and real-time analysis of scientific literature without manual intervention.

4. Identification of Research Gaps:

- The tool will help identify underexplored areas of biodiversity genomics and highlight the gaps in dataset usage, providing guidance on areas that need more attention from the scientific community.

5. Influencing Future Research and Conservation:

- The findings from this analysis will guide future research agendas, especially in biodiversity conservation efforts, by showing which datasets are most influential and what species or ecosystems require more genomic attention.

Deliverables:

1. Data Collection and Dataset Integration:

- Collect genomic datasets from major biodiversity databases (e.g., GenBank, ENA, DDBJ).
- Gather metadata on scientific publications related to biodiversity genomics (e.g., from PubMed, Scopus, or Web of Science).

2. Machine Learning Models:

- Train machine learning models to recognize mentions of genomic datasets in scientific literature.
- Develop text classification models to identify the relevance of genomic datasets in the context of different research topics (e.g., species conservation, evolutionary genetics).

3. Citation and Impact Analysis:

- Create algorithms for analyzing the citation frequency and impact of genomic datasets in scientific publications.
- Measure the influence of datasets based on citation metrics, co-authorship patterns, and references in policy documents or conservation reports.

4. Visualization and Reporting:

- Develop interactive visualizations (using tools like Plotly or D3.js) to display trends in dataset usage, citation patterns, and geographic distribution.
- Generate comprehensive reports summarizing the impact of biodiversity genome datasets, including insights on dataset contribution to key research areas and conservation efforts.

5. Documentation and Open-Source Release:

- Provide comprehensive documentation on the tools, machine learning models, and analysis techniques developed.
- Release the code and tools as open-source software, enabling the scientific community to adapt and extend the tool for their research.

Technical Approach:

1. Data Collection and Preprocessing:

- Genomic Datasets: Retrieve biodiversity genome datasets from online repositories such as GenBank, ENA, and DDBJ. Use their APIs to automate the process of data collection.
- Publications Metadata: Gather scientific publications from platforms like PubMed and Scopus using APIs or web scraping (where allowed).
- Data Cleaning: Clean and preprocess the raw data to ensure consistency and format compatibility for downstream analysis (e.g., normalizing

citation formats, removing duplicates).

2. Text Mining with Natural Language Processing (NLP):

- Dataset Mention Recognition: Use NLP techniques to extract mentions of genomic datasets from scientific articles. This will involve training a named entity recognition (NER) model to detect dataset names, using annotated datasets for training.
- Text Classification: Develop a text classification model (using approaches like TF-IDF or BERT) to categorize publications based on their focus area (e.g., species conservation, population genetics, evolution).
- Citation Context Analysis: Use NLP techniques to identify the context in which genomic datasets are cited—whether they are used for primary research, hypothesis testing, or conservation policy recommendations.

3. Machine Learning Models for Citation and Impact Analysis:

- Citation Analysis: Implement machine learning models to predict the impact of publications that use genomic datasets. Use citation counts, co-authorship networks, and other bibliometric measures.
- Co-occurrence Networks: Build models to visualize how often datasets are cited together with other research topics, helping to identify research trends and collaborations within the scientific community.

4. Visualization:

- Create interactive visualizations to display the relationships between genomic datasets and scientific publications.
- Use tools like Plotly and D3.js to enable users to interact with citation data, explore trends, and analyze dataset usage across species, regions, and time.

5. Testing and Evaluation:

- Test the machine learning models and NLP techniques using a set of annotated data to ensure high accuracy in dataset recognition and citation analysis.
- Perform end-to-end testing of the visualization tool to ensure its functionality and usability.

Timeline:

Community Bonding Period (May 18 - May 27):

- Set up communication channels with mentors and contributors.
- Familiarize with tools, repositories, and libraries to be used for the project.
- Review the literature on biodiversity genomics, dataset citation, and machine learning techniques.

Phase 1 (May 27 - June 30): Data Collection & Preprocessing

- Retrieve genomic datasets and scientific publications metadata from APIs.
- Clean and preprocess the collected data to ensure consistency and compatibility for analysis.
- Set up the initial framework for text mining and dataset recognition.

Phase 2 (July 1 - July 31): NLP and Machine Learning Model Development

- Develop and train NLP models to extract genomic dataset mentions from publications.
- Begin building text classification models to categorize research papers based on genomic dataset usage.
- Implement citation analysis algorithms and co-authorship network generation.

Phase 3 (August 1 - August 15): Visualization & Reporting

- Develop interactive visualizations for dataset usage, citation trends, and research impact.
- Generate reports summarizing key insights from the dataset analysis.
- Test the visualization tools with sample data to ensure accuracy and user-friendliness

Phase 4 (August 16 - August 24): Final Refinements & Documentation

- Refine models and tools based on feedback from mentors and testing.
- Finalize documentation and prepare the project for open-source release.
- Conduct a final evaluation of the tools and visualization

Final Evaluation (August 24 - August 31):

- Submit the final deliverables, including source code, documentation, and reports.

- Participate in a final review with mentors and project reviewers.

Skills Required:

- Python Programming:
- Experience with Python libraries such as Pandas, NumPy, and Scikit-learn for data analysis and machine learning.
- Natural Language Processing (NLP): Familiarity with NLP techniques and libraries like SpaCy, NLTK, or Hugging Face's Transformers for text mining and dataset recognition.
- Machine Learning: Knowledge of supervised learning, text classification, and citation analysis using models such as BERT or SVMs.
- Data Visualization: Proficiency with data visualization tools like Plotly or D3.js to create interactive graphs and charts.
- Scientific Research Methods: Understanding of citation metrics and bibliometric analysis, with a focus on genomic datasets.

Why Me?

I am passionate about both biodiversity conservation and machine learning, with a strong background in Python programming and data science. My experience with natural language processing and machine learning will enable me to effectively tackle the challenges of this project. Furthermore, I am eager to contribute to open-source tools that can have a significant impact on biodiversity research, helping to streamline the analysis of genomic data and its influence on scientific publications.

Future Directions:

After the completion of this project, there is potential to expand the tool to support more types of genomic data (e.g., microbiomes, environmental genomics), integrate additional data sources (e.g., environmental or ecological datasets), and enhance the machine learning models for improved prediction accuracy. Additionally, future work could focus on creating a community-driven platform for sharing insights and datasets related to biodiversity genomics.

Mentorship:

I will be guided by mentors who are experts in biodiversity genomics, machine learning, and data science. Their feedback will ensure the project remains scientifically sound and technically robust.