# Experiment No.: 08

**Aim:** Explore different dataset finders e.g. Google Dataset Search, Kaggle, mendeley etc.

**Course Outcome:** Analyze different forms of data with respect to different phases of Machine Learning.

**Theory:**

A dataset finder is a specialized online platform or search engine designed to help users locate datasets from a variety of fields—such as education, healthcare, business, environment, and technology. These platforms index datasets stored across the web, allowing users to search, preview, and download data in multiple formats like CSV, JSON, XML, and Excel.
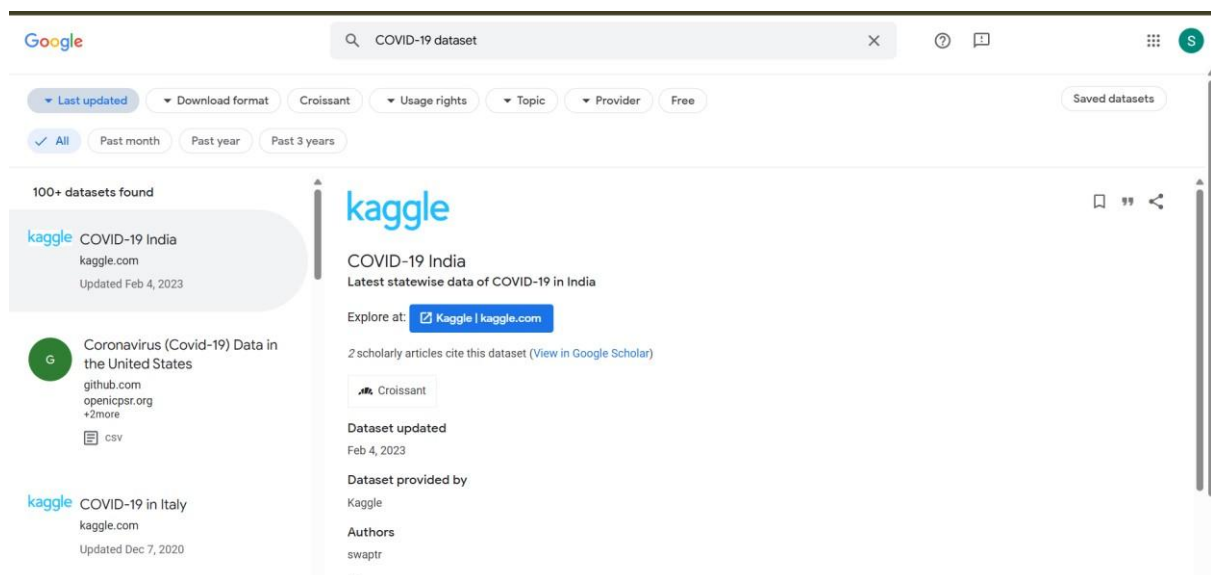
Dataset finders play a vital role in data science and research because they:
   • Save time by aggregating datasets from various sources.
   • Ensure datasets are well-documented and accessible.
   • Promote open data sharing and collaboration.
   • Encourage reproducibility in research and data-driven decision-making.

Below are some commonly used dataset finders:

## 1. Google Dataset Search
   • A search engine developed by Google for datasets.
   • Works similarly to Google Search but specifically indexes datasets.
   • Displays dataset title, publisher, description, and file format.
   • Offers filters for file type, date, and license.
   • Widely used by researchers for academic and professional data discovery.
   • URL: https://datasetsearch.research.google.com



## 2. Kaggle
   • One of the largest online communities for data scientists.
   • Offers thousands of public datasets across multiple domains.

- Datasets are available in clean, structured formats.
- Provides cloud-based notebooks for analysis and visualization.
- Encourages collaboration and learning through competitions.
- URL: https://www.kaggle.com/datasets



## 3. Mendeley Data

- A research data repository developed by Elsevier.
- Allows researchers to publish, store, and share datasets securely.
- Every dataset is assigned a **DOI (Digital Object Identifier)** for citation.
- Used mostly in academic and scientific research.
- Supports multiple file formats and ensures long-term data preservation.
- U RL: https://data.mendeley.com



Platform Summary

| Platform Name | Website Link | Data Formats | Area of Use | Key Feature |
|---|---|---|---|---|
| Google Dataset Search | datasetsearch.research.google.com | CSV, XLSX, JSON | General research | Aggregates datasets from various websites |
| Kaggle | kaggle.com/datasets | CSV, JSON, ZIP | Data science, ML | Community-based data sharing |
| Mendeley Data | data.mendeley.com | ZIP, CSV, PDF | Academic & research | DOI and citation support |
| Platform Name | Website Link | Data Formats | Area of Use | Key Feature |
| UCI ML Repository | archive.ics.uci.edu | CSV, ARFF | Machine learning | Ready-to-use research datasets |

Feature Comparison

| Feature | Google Dataset Search | Kaggle | Mendeley Data |
|---|---|---|---|
| Login Required | No | Yes | Yes |
| Community Support | Moderate | High | Low |
| Dataset Size Range | Small to Large | Small to Very Large | Moderate |
| Ideal For | Quick search and discovery | Data science projects | Academic research |
| Dataset Citation (DOI) | No | Optional | Yes |

**Conclusion:** I have successfully explored different dataset finders such as Google Dataset Search, Kaggle, and Mendeley. These platforms help users easily find, access, and use datasets for research, analysis, and learning purposes.