



## International Journal of Research and Technology (IJRT)

International Open-Access, Peer-Reviewed, Refereed, Online Journal

ISSN (Print): 2321-7510 | ISSN (Online): 2321-7529

Conference “Innovation and Intelligence: A Multidisciplinary Research on Artificial Intelligence and its Contribution to Commerce and Beyond”

Organized by the IQAC, KHMW College of Commerce (December 2025)

### General AI & Machine Learning

Miss Shaista Shaikh

Miss Rahila Sayyed

Bachelor of Commerce

K.H.M.W Degree College

#### Abstract

General AI and Machine Learning transform how systems learn, reason, and make decisions. By analyzing vast data, they enable automation, predictive accuracy, and intelligent problem-solving across industries. Their advances drive innovation in healthcare, business, robotics, and everyday applications, shaping a more efficient and adaptive digital future.

**Keywords:** Artificial Intelligence, Machine Learning, Supervised and Unsupervised Learning, Neural Networks, Data-Driven Decision Making

#### Introduction

Artificial Intelligence (AI) and Machine Learning (ML) have emerged as pivotal forces driving innovation in scientific, industrial, and socio-economic sectors. As these technologies evolve, the significance of data—its availability, quality, structure, and context—has become increasingly critical in influencing the accuracy, fairness, and generalizability of ML models. Given that the collection of high-quality primary data is often expensive, time-intensive, and logistically complex, researchers commonly turn to secondary data: datasets that were initially compiled for different purposes but are subsequently repurposed for new analytical or predictive aims. The utilization of secondary data has long been a recognized practice in disciplines such as economics, health sciences, and demography, yet its importance has surged dramatically with the advent of data-driven AI research. Secondary datasets, which include government census records, institutional surveys, benchmark repositories, and extensive scraped databases, enable ML researchers to swiftly prototype, validate, and compare algorithms. These datasets have played a crucial role in facilitating numerous methodological advancements over the past decade. For example, benchmark datasets like ImageNet for computer vision or the UCI repository datasets for tabular modelling have allowed researchers to assess algorithms under standardized conditions, promoting reproducibility and transparent comparisons across various studies. The existence of well-documented secondary datasets also democratizes research by granting students, small institutions, and emerging researcher's immediate access to data that would otherwise necessitate considerable resources to gather.

Despite these benefits, utilizing secondary data in machine learning (ML) poses several significant challenges. Secondary datasets are not created with the current goals of researchers in mind; therefore, they may include measurement biases, missing values, inconsistent sampling methods, historical inaccuracies, or insufficient documentation. If these datasets are integrated into ML pipelines without thorough examination, models may inherit or exacerbate existing biases, misrepresent demographic groups, or show diminished predictive validity



when applied to new situations. Additionally, ethical and legal issues—such as consent, privacy, data origin, and responsible governance—have become increasingly significant as secondary datasets expand in size and sensitivity. These challenges necessitate that researchers find a balance between technical efficiency and accountability, transparency, and fairness. In light of these complexities, academics from various fields have stressed the importance of structured frameworks for secondary data analysis (SDA) within ML. These frameworks underscore critical steps such as evaluating dataset provenance, clarifying the original intent behind data collection, documenting preprocessing choices, and recognizing limitations when making conclusions or deploying models. When applied rigorously, SDA can provide insights that are comparable to those obtained from primary data collection while remaining cost-effective and methodologically robust.

## **Literature review**

### **1. Secondary data analysis**

Secondary data analysis (SDA) refers to the re-examination of data that has been meeting by other researchers. SDA has a rich history in the fields of social science and health research, serving as a cost-efficient method to produce new insights without the need for gathering new primary data. The existing literature outlines both the technical benefits (such as scale and longitudinal depth) and the methodological drawbacks (including limited control over variables, issues of missing data, and measurement discrepancies). Secondary data analysis is the process of using data originally collected for a different purpose to answer a new research question

### **2.Secondary data reuse in ML: opportunities and risks**

Recent surveys focused on machine learning indicate that the creation and reuse of datasets are pivotal for advancement. Furthermore, the selection of datasets influences model presence and fairness outcomes. Paullada et al. (2021) examine data-generation methodologies and underscore issues such as inadequate documentation and biases in dataset construction. Thylstrup (2022) and others have criticized the 'politics' surrounding reuse, which refers to how reuse practices incorporate underlying assumptions, governance decisions, and regulatory challenges. More contemporary research emphasizes the importance of domain- and sensitivity-aware frameworks for the reuse of datasets in production machine learning systems.

### **3. Ethical and practical considerations**

The reuse of secondary data brings forth concerns regarding consent, privacy, and provenance. It is advised to document the provenance, evaluate the representativeness, manage any missing or biased measurements, and consider legal and ethical limitations prior to reuse. Numerous guides in the fields of health and social science offer checklists for ethical secondary data analysis that can be adapted for machine learning environments.



### **Research question & rationale**

Research question: What straightforward and reproducible steps are required to prepare a real public secondary dataset (UCI Adult) for a fundamental supervised task (predicting an income greater than 50K), what descriptive patterns emerge, and what significant reuse risks should a practitioner consider? Rationale: The UCI Adult dataset is a recognized benchmark for income prediction and is extensively reused; thus, it provides a practical and realistic example for illustrating best practices in Software Development and Analysis (SDA) within Machine Learning

### **Data and methods**

#### **1. Dataset (secondary source)**

Dataset: UCI Adult, also known as Census Income, was initially derived from the 1994 US Census and made available in the UCI Machine Learning Repository (Becker & Kohavi, 1996). The unprocessed file comprises 48,842 instances along with 14 attributes and a binary target ( $\leq 50K$  /  $> 50K$ ). There are missing values present, represented as which necessitate appropriate handling. This dataset is publicly accessible and frequently utilized as a standard benchmark. UCI Machine Learning Repository Citation (dataset): Becker, B., & Kohavi, R. (1996). Adult [Dataset]. UCI Machine Learning Repository. UCI Machine Learning Repository

#### **2. Reproducible analysis plan (described so reader can repeat)**

The dataset is publicly available, allowing any reader to replicate the following steps. The analysis proceeds as follows:

1. Load the data from UCI (CSV).
2. Examine the rows and columns, identify any missing values (?), and determine the total number of instances.
3. Clean the data: Replace? With NaN. For this demonstration, drop rows that contain missing values (note: alternative strategies such as imputation will be discussed later).
4. Encode categorical variables using either one-hot or label encoding as appropriate.
5. Calculate descriptive statistics, including counts and proportions for the target class, as well as distributions based on sex and education.
6. Fit a basic logistic regression model (baseline) to predict whether income exceeds 50K, using a limited set of covariates: age, education-num, hours-per-week, and sex (encoded), to demonstrate the direction of effects.
7. Assess the model's performance using accuracy, precision, recall, and a confusion matrix on a held-out test split (for example, 80/20).
8. Record all preprocessing decisions and maintain both raw and processed copies along with metadata.

#### **3. Justification of methods**



Logistic regression serves as a clear baseline to demonstrate how secondary data can uncover associations, as well as the sensitivity of these associations to pre-processing and selection bias. Descriptive statistics highlight the importance of transparency prior to modelling. These practices are standard and recommended in SDA/ML method guides.

### **Simple data analysis (results — summary & illustrative numbers)**

Important note: the figures provided below are derived from documented summaries of the adult dataset of the standard pre-processing decisions outlined at published analyses; they were included here for the sake of clarity and to demonstrate of the workflow. Individuals who replicate the analysis locally may achieve slightly varying counts based on their specific preprocessing choices (for instance, whether the test/train split is pre-supplied and how missing rows are addressed).

#### **1. Basic counts and missingness**

Total instances (raw): 48,842. Rows with missing values: typically around ~3,620 rows, resulting in approximately ~45,222 complete rows after the straightforward removal of missing values. (Various authors utilize different methods; this is the frequently referenced preprocessing outcome.)

#### **2. Target distribution**

Proportion exceeding 50K: around 24% of the total rows (that is, about one in every four occurrences). While the class exhibits some imbalance, it is not excessively pronounced; numerous published studies categorize these as a mild to moderate level of imbalance.

#### **3. Example descriptive patterns (illustrative)**

Sex differences: The dataset reveals that men are disproportionately represented among individuals earning over \$50,000, a trend that is often noted in analyses of this dataset and emphasized in fairness studies that utilize Adult as a case study. [yanhan.github.io+1](https://github.com/yanhan-io/1) Education / income: The data indicates that higher levels of education, both in terms of numerical education and categorical education classifications, correspond to steadily increasing proportions of individuals earning over \$50,000 (for instance, those with a college education or higher tend to represent the larger share of this income bracket), which is consistent with conventional findings in income-prediction research.

#### **4. Illustrative logistic regression (conceptual summary)**

When fitting a logistic regression model using predictors such as age, education level, hours worked per week, and sex (with males coded as 1), a common outcome observed in similar published studies is as follows:

Positive and statistically significant coefficients for age, education level, and hours worked per week, indicating that higher values increase the odds of earning more than \$50,000.

A positive coefficient for males, suggesting that males have higher odds when other variables are held constant.



The baseline accuracy in published models generally falls within the range of 80% to 88%, which can vary based on preprocessing techniques and the type of model used (for instance, logistic regression versus gradient boosting). For instance, certain studies have reported validation accuracies approaching 88% for optimized gradient-boosted models, whereas simpler models tend to produce lower yet still acceptable accuracy rates. arXiv+1

### **Discussion: what the example shows about secondary data reuse in ML**

#### **1. Benefits illustrated**

Speed and cost: Publicly available secondary datasets, such as the UCI Adult dataset, facilitate quick prototyping and comparison of methods. The Open University

Reproducibility and benchmarking: Datasets shared within the community provide the means for establishing baselines and conducting replicable comparisons among various algorithms. Science Direct

#### **2. Risks & limitations highlighted**

Provenance uncertainty: The Adult dataset is the derived subset from the 1994 Census; discrepancies in timing or sampling choices may restrict its generalizability. UCI Machine Learning Repository

Bias & fairness: Recognized demographic imbalances (gender, ethnicity, and nationality) can result in models that perpetuate social inequalities unless they are specifically tackled. Science Direct

Documentation gaps: The absence of metadata (the rationale behind variable selection, encoding decisions) hinders reuse and replication; both dataset creators and users should provide a comprehensive README and provenance documentation.

#### **3. Practical recommendations for responsible reuse**

1. Transformations / Missingness plan: Document Provenance checklist: Always record the original source, collection date, and the methods for handling missing values; consider multiple imputation when suitable.
2. Bias assessment: Conduct subgroup analyses and fairness audits prior to deployment.
3. Ethics & legal check: Confirm consent, licensing, and privacy implications; for sensitive areas, consult ethics boards.
4. Share artifacts: Publish preprocessing scripts, random seeds, and small metadata files to ensure transparency and reproducibility in reuse.

#### **Conclusion**

Secondary data play a crucial role in machine learning research by providing quick access to real-world scale and the facilitating reproducible benchmarking. Nevertheless, it is essential to approach reuse with caution: one must document provenance, identify and assess potential biases, manage missing data transparently, and consider a ethical limitations. The brief, illustrative analysis utilizing the UCI Adult dataset highlights both the immediate benefits of reuse and the common challenges that researchers need to address.





## **References**

1. (Additional resource pointers used in the review are listed inline above; readers should consult the cited items for deeper detail.)
2. Becker, B., & Kohavi, R. (1996). Adult [Dataset]. UCI Machine Learning Repository
- Cheng, H. G., & Phillips, M. R. (2014). Secondary analysis of existing data: opportunities and implementation. [Article summarizing SDA rationale and sources]
3. Bhagat, P. H., & Shaikh, S. A. (2025). Managing health care in the digital world: A comparative analysis on customers using health care services in Mumbai suburbs and Pune city. IJCRT. Registration ID: IJCRT\_216557.
4. Chakrabarty, N., & Biswas, S. (2018). A Statistical Approach to Adult Census Income Level Prediction. arXiv preprint (example of model performance reporting on Adult dataset).
5. Chougale, Z. S., & Shaikh, S. (2022). To understand the impact of Ayurvedic health-care business & its importance during COVID-19 with special reference to “Patanjali Products”. In Proceedings of the National Conference on Sustainability of Business during COVID-19, IJCRT, 10(1),
6. Machine Learning Mastery. (2020). Imbalanced Classification with the Adult Income Dataset (descriptive preprocessing summary).
7. Parikh, V. (2023). Whistleblowing in B-Schools, Education and Society, Vol-47, Issue – 1, Pg. 183-1
8. Parikh, V. C. (2022) Strategic talent management in education sector around organizational life cycle stages! JOURNAL OF THE ASIATIC SOCIETY OF MUMBAI, SSN: 0972-0766, Vol. XCV, No.11.
9. Paullada, A., et al. (2021). Data and its (dis)contents: A survey of dataset development and use in machine learning. [Survey article on dataset practices]
10. Shaikh, S. A. (2024). Empowering Gen Z and Gen Alpha: A comprehensive approach to cultivating future leaders. In Futuristic Trends in Management (IIP Series, Vol. 3, Book 9, Part 2, Chapter 2). IIP Series. <https://doi.org/10.58532/V3BHMA9P2CH2>
11. Shaikh, S. A., & Jagirdar, A. H. (2026). Beyond AI dependence: Pedagogical approaches to strengthen student reasoning and analytical skills. In S. Khan & P. Pringuet (Eds.), Empowering learners with AI: Strategies, ethics, and frameworks (Chapter 8, pp. 1–16). IGI Global. <https://doi.org/10.4018/979-8-3373-7386-7.ch008>
12. Thylstrup, N. B. (2022). Politics of data reuse in machine learning systems: Theorizing data re-use entanglements. [Article on political/ethical aspects of data reuse].
13. Tripathy, J. P. (2013). Secondary Data Analysis: Ethical Issues and Challenges. [Discussion of ethical issues in SDA].
14. Wickham, R. J. (2019). Secondary Analysis Research — review of methods and pitfalls. [Overview of SDA challenges].