

Database Summary: Simulated Loan Risk Assessment Data

1. Origin and Construction

This repository consists of synthetic mortgage approval and borrower risk data, modeled after prominent public sources such as Kaggle's established loan evaluation datasets. By leveraging advanced sampling techniques—specifically, Synthetic Minority Oversampling for Mixed Data (SMOTENC)—the cohort has been engineered to counteract sample imbalance and reflect diverse attribute distributions frequently exhibited in real-world lending scenarios.

With **25,000 entries** and **14 distinct data fields**, the database is meticulously shaped to mirror the nuanced relationships between applicant profiles, debt applications, and repayment risk criteria, making it especially suitable for exploratory, predictive, and benchmarking purposes.

2. Field Definitions

Every row equates to a single application event, incorporating a blend of demographic background, financial situation, and lending conditions. Below are the column definitions:

- **applicant_age (Float)**: Represents the age in years, signaling the applicant's position in the lifecycle and its effect on eligibility and risk.
- **applicant_gender (Categorical)**: Identifies the sex of the individual, supplying demographic details relevant for compliance and analysis.
- **education_level (Categorical)**: Catalogs the highest degree or diploma held, affecting earning capacity and stability.
- **annual_earnings (Float)**: States yearly income, considered key for debt service evaluation.
- **employment_years (Integer)**: Denotes period of job tenure, useful for quantifying career consistency.
- **residence_type (Categorical)**: Records housing status—such as owning, renting, or mortgaging—informing financial health assessments.
- **requested_loan (Float)**: The monetary figure sought by the borrower, capturing aggregate demand for credit.
- **loan_purpose (Categorical)**: Specifies the loan's intended use, e.g., consuming, investing, tuition payments, or business investment.
- **interest_rate (Float)**: Sets the percentage rate for repayment, a key variable in cost and risk calculations.
- **loan_to_income_ratio (Float)**: The proportion of the loan request against annual income, gauging borrower leverage.

- **history_duration (Float):** Tracks years of recorded credit history, indicative of long-term financial patterns.
 - **bureau_score (Integer):** Numeric rating summarizing the applicant's credibility, often sourced or modeled after national agencies.
 - **has_prior_defaults (Categorical):** Encodes past incidents of loan non-payment or adverse credit events.
 - **application_outcome (Integer):** The central label showing whether the request resulted in funding (1) or denial (0).
-

3. Intended Analysis and Usage Scenarios

The **Simulated Loan Risk Assessment Data** has broad applicability across numerous research and industry requirements:

- **Descriptive and Exploratory Analysis:**
Investigators may profile applicants, review income distribution, identify historical default trends, and dissect the impacts of education or housing status on repayment likelihood.
 - **Machine Learning for Approval Prediction:**
The dataset is ready-made for supervised learning, enabling users to train algorithms like logistic regression, decision trees, or ensemble classifiers to forecast loan approval outcomes directly from an array of applicant and loan features.
 - **Modeling Risk and Segmentation:**
Institutions or scholars can construct tailored risk models, isolate high-risk subgroups, and validate the effectiveness of credit policies by simulating potential shifts in acceptance criteria or underwriting rules.
 - **Synthetic Data Testing and Benchmarking:**
Designed with class balance in mind, thanks to SMOTENC oversampling, the database is ideal for evaluating the robustness and fairness of algorithms under synthetic but realistic class distributions.
-

4. Data Strengths and Practical Considerations

- **Comprehensive Feature Design:**
The database includes both quantitative (numeric) and qualitative (categorical) fields, supporting complex feature engineering and comparative modeling workflows.
- **Balanced Sampling:**
SMOTENC's integration aids in counteracting natural class skew, enhancing minority representation and reducing bias encountered in traditional financial datasets.

- **Synthetic Realism:**
While the records are artificial, strong underlying relationships and attribute statistics are maintained, ensuring meaningful training and testing—without risking sensitive data exposure.
 - **Quality Control:**
Users are encouraged to vet and preprocess fields for outliers, especially improbable age values or uncommon income entries, to ensure analytical integrity during development and deployment.
-

5. Value for Lenders and Data Practitioners

By furnishing an array of demographic, financial, and risk attributes in a manageable and realistic structure, the dataset empowers:

- **Credit Model Development:**
Users can build new risk scoring techniques or strengthen decision support systems that evaluate client applications more transparently and robustly.
- **Business Strategy Simulation:**
The dataset facilitates scenario planning, including the prospective impact of policy or regulatory modifications on client approval rates or default probabilities.
- **Compliance and Fairness Testing:**
Synthetic data offers a safe environment for examining discrimination or bias in automated decision frameworks, critical for regulatory readiness.