
OpenL2D Fraud Detection Dataset

Anonymous Author(s)

Affiliation

Address

email

1 Motivation

2 Q1: For what purpose was the dataset created?

3 A1: The target of this dataset, comprised by a *learning to defer* (L2D) training scenario with
4 limited expert predictions and a set of expert predictions, is to contribute to the development
5 and evaluation of L2D algorithms. We focus particularly in testing fairness and performance
6 in dynamic conditions, in order to *stress-test* L2D methods.

7 Q2: Who created the dataset (e.g., which team, research group) and on behalf of which
8 entity (e.g., company, institution, organization)?

9 A2: Currently anonymous for double-blind review process.

10 Q3: Who funded the creation of the dataset?

11 A3: The dataset is synthetically generated using the OpenL2D framework. There was no specific
12 funding for the creation of this dataset.

13 Composition

14 Q4: What do the instances that comprise the dataset represent (e.g., documents, photos,
15 people, countries)?

16 A4: The OpenL2D Fraud Detection Dataset is comprised of:

- 17 • Input dataset - we utilize the bank-account-fraud (BAF) dataset's base variant.
18 Each instance in the bank-account-fraud dataset represents a synthetic, feature-
19 engineered bank account opening application in tabular format. For more information
20 on this dataset, please consult [https://www.kaggle.com/datasets/sgpjesus/](https://www.kaggle.com/datasets/sgpjesus/bank-account-fraud-dataset-neurips-2022?datasetId=2673949)
21 [bank-account-fraud-dataset-neurips-2022?datasetId=2673949](https://www.kaggle.com/datasets/sgpjesus/bank-account-fraud-dataset-neurips-2022?datasetId=2673949).
- 22 • Limited expert decision dataset - Subset of the input dataset, where each instance is
23 associated with a decision from either the model or an expert. It was used to develop
24 our L2D methods with limited human decision data.
- 25 • Expert decision table - This table contains every expert's decision for every instance in
26 the BAF dataset.
- 27 • Capacity constraint tables - For any given subset of the data (e.g. test split), each
28 capacity constraint is defined by a pair of tables:
 - 29 – Batch tables - define which batch every instance of the input dataset belongs to.
 - 30 – Capacity tables - define the maximum number of cases that can be deferred to each
31 expert for a given batch.

32 Q5: How many instances are there in total (of each type, if appropriate)?

33 A5: • The input dataset has 1M instances
 34 • The limited expert decision dataset has 506118 instances
 35 • The expert decision table contains 1M instances
 36 • The batch tables contain the same number of instances as the data for which they are
 37 generated (e.g batch table generated for the limited expert decision dataset contains
 38 506118 instances). The capacity tables contain one instance for each batch.

39 Q6: **Does the dataset contain all possible instances or is it a sample (not necessarily random)**
 40 **of instances from a larger set?**

41 A6: • The input dataset was used in its totality
 42 • The limited expert decision dataset includes instances from a subset of the BAF dataset,
 43 corresponding to instances from months 4 to 7.
 44 • The expert prediction table contains all instances

45 Q7: **What data does each instance consist of?**

46 A7: In the limited expert decision dataset, the generated fields are:

- 47 • **model_score** (numeric) : Model score obtained with our ML Model.
- 48 • **batch** (categorical): Batch to which the instance belongs, defined by capacity con-
 49 straints.
- 50 • **assignment** (categorical): Entity (expert or model) to which instance's decision was
 51 deferred.
- 52 • **decision** (numeric): Decision on said instance. In case of deferral to an expert, it is
 53 either 0 (accept), or 1 (reject). In case of deferral to the ML model, it is equal to that
 54 model's score for said instance.

55 In our expert decision table, every instance is matched to each entity's decision. Each
 56 column corresponds to one of the 51 decision making entities: (50 experts) and one
 57 ML Model. Experts are identified by their group and a numeric identifier, in the for-
 58 mat "{group_name}#{id}". The ML model is identified as "model#0". Expert decisions
 59 are binary, while the model's decision is comprised of the model score for each instance,
 60 allowing for posterior thresholding.

61 In the batch tables the columns are:

- 62 • **case_id** - identifier for each instance in the input dataset
- 63 • **batch_id** - identifier of the batch that the instance belongs to

64 In the capacity tables the columns are the same as the expert decision table, and each instance
 65 is identified by the "batch_id".

66 The fields of the input dataset, related to information on the bank account application, are
 67 the same as the base BAF variant. These are also present in the limited expert decision
 68 dataset:

- 69 • **income** (numeric): Annual income of the applicant (in decile form). Ranges between
 70 [0.1, 0.9].
- 71 • **name_email_similarity** (numeric): Metric of similarity between email and applicant's
 72 name. Higher values represent higher similarity. Ranges between [0, 1].
- 73 • **prev_address_months_count** (numeric): Number of months in previous registered
 74 address of the applicant, *i.e.* the applicant's previous residence, if applicable. Ranges
 75 between [-1, 380] months (-1 is a missing value).
- 76 • **current_address_months_count** (numeric): Months in currently registered address of
 77 the applicant. Ranges between [-1, 429] months (-1 is a missing value).
- 78 • **customer_age** (numeric): Applicant's age in years, rounded to the decade. Ranges
 79 between [10, 90] years.

- 80 • **days_since_request** (numeric): Number of days passed since application was done.
81 Ranges between [0, 79] days.
- 82 • **intended_balcon_amount** (numeric): Initial transferred amount for application.
83 Ranges between [-16, 114] (negatives are missing values).
- 84 • **payment_type** (categorical): Credit payment plan type. 5 possible (anonymized)
85 values.
- 86 • **zip_count_4w** (numeric): Number of applications within same zip code in last 4 weeks.
87 Ranges between [1, 6830].
- 88 • **velocity_6h** (numeric): Velocity of total applications made in last 6 hours *i.e.*, average
89 number of applications per hour in the last 6 hours. Ranges between [-175, 16818].
- 90 • **velocity_24h** (numeric): Velocity of total applications made in last 24 hours *i.e.*, average
91 number of applications per hour in the last 24 hours. Ranges between [1297, 9586]
- 92 • **velocity_4w** (numeric): Velocity of total applications made in last 4 weeks, *i.e.*, average
93 number of applications per hour in the last 4 weeks. Ranges between [2825, 7020].
- 94 • **bank_branch_count_8w** (numeric): Number of total applications in the selected bank
95 branch in last 8 weeks. Ranges between [0, 2404].
- 96 • **date_of_birth_distinct_emails_4w** (numeric): Number of emails for applicants with
97 same date of birth in last 4 weeks. Ranges between [0, 39].
- 98 • **employment_status** (categorical): Employment status of the applicant. 7 possible
99 (anonymized) values.
- 100 • **credit_risk_score** (numeric): Internal score of application risk. Ranges between
101 [-191, 389].
- 102 • **email_is_free** (binary): Domain of application email (either free or paid).
- 103 • **housing_status** (categorical): Current residential status for applicant. 7 possible
104 (anonymized) values.
- 105 • **phone_home_valid** (binary): Validity of provided home phone.
- 106 • **phone_mobile_valid** (binary): Validity of provided mobile phone.
- 107 • **bank_months_count** (numeric): How old is previous account (if held) in months.
108 Ranges between [-1, 32] months (-1 is a missing value).
- 109 • **has_other_cards** (binary): If applicant has other cards from the same banking com-
110 pany.
- 111 • **proposed_credit_limit** (numeric): Applicant's proposed credit limit. Ranges between
112 [200, 2000].
- 113 • **foreign_request** (binary): If origin country of request is different from bank's country.
- 114 • **source** (categorical): Online source of application. Either browser (INTERNET) or
115 app (TELEAPP).
- 116 • **session_length_in_minutes** (numeric): Length of user session in banking website in
117 minutes. Ranges between [-1, 107] minutes (-1 is a missing value).
- 118 • **device_os** (categorical): Operative system of device that made request. Possible values
119 are: Windows, macOS, Linux, X11, or other.
- 120 • **keep_alive_session** (binary): User option on session logout.
- 121 • **device_distinct_emails** (numeric): Number of distinct emails in banking website from
122 the used device in last 8 weeks. Ranges between [-1, 2] emails (-1 is a missing value).
- 123 • **device_fraud_count** (numeric): Number of fraudulent applications with used device.
124 Ranges between [0, 1].
- 125 • **month** (numeric): Month where the application was made. Ranges between [0, 7].
- 126 • **fraud_bool** (binary): If the application is fraudulent or not.

127 Q8: Is there a label or target associated with each instance?

128 A8: Yes, the label is contained in the **fraud_bool** field. A positive value (fraud_bool=1) repre-
 129 sents a fraudulent bank account application. A negative value (fraud_bool=0) represents a
 130 legitimate bank account application.
 131 When accepted, all accounts are opened with access to credit.
 132 For additional information on how the labels were obtained, consult the BAF datasheet,
 133 available at

134 **Q9: Is any information missing from individual instances?**
 135 A9: There is no missing information from individual instances.

136 **Q10: Are relationships between individual instances made explicit (e.g., users' movie ratings, social network links)?**
 137
 138 A10: There are no relationships between individual instances.

139 **Q11: Are there recommended data splits (e.g., training, development/validation, testing)?**
 140 A11: The performed data splits are based on the temporal information of the dataset. To this
 141 end, we use the column **month**. Practitioners can test different temporal L2D development
 142 strategies (e.g. Training an ML Model, training an assignment system, deploying the
 143 system).

144 **Q12: Are there any errors, sources of noise, or redundancies in the dataset?**
 145 A12: Not applicable to the synthetically generated decisions. There may be sources of error associ-
 146 ated with the BAF dataset used as input to generate the expert decisions. Please refer to BAF
 147 datasheet for more details [https://github.com/feedzai/bank-account-fraud/](https://github.com/feedzai/bank-account-fraud/blob/main/documents/datasheet.pdf)
 148 [blob/main/documents/datasheet.pdf](https://github.com/feedzai/bank-account-fraud/blob/main/documents/datasheet.pdf).

149 **Q13: Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g., websites, tweets, other datasets)?**
 150 The generated expert predictions are self-contained. The input dataset is the
 151 publicly available BAF dataset [https://www.kaggle.com/datasets/sgpjesus/](https://www.kaggle.com/datasets/sgpjesus/bank-account-fraud-dataset-neurips-2022)
 152 [bank-account-fraud-dataset-neurips-2022](https://www.kaggle.com/datasets/sgpjesus/bank-account-fraud-dataset-neurips-2022)
 153

154 **Q14: Does the dataset contain data that might be considered confidential (e.g., data that is protected by legal privilege or by doctor-patient confidentiality, data that includes the content of individuals' non-public communications)?**
 155
 156
 157 A14: There is no confidential data in this dataset.

158 **Q15: Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, or might otherwise cause anxiety?**
 159
 160 A15: No.

161 **Q16: Does the dataset relate to people?**
 162 A16: This dataset relates to synthetically generated human experts, it does not relate to real people.

163 **Q17: Does the dataset identify any subpopulations (e.g., by age, gender)?**
 164 A17: The synthetically generated dataset using the OpenL2D framework does not identify
 165 any subpopulations. The BAF dataset used as input dataset identifies age groups on
 166 the synthetically generated bank account applications. Please refer to BAF datasheet
 167 for more details [https://github.com/feedzai/bank-account-fraud/blob/main/](https://github.com/feedzai/bank-account-fraud/blob/main/documents/datasheet.pdf)
 168 [documents/datasheet.pdf](https://github.com/feedzai/bank-account-fraud/blob/main/documents/datasheet.pdf).

169 **Q18: Is it possible to identify individuals (i.e., one or more natural persons), either directly or indirectly (i.e., in combination with other data) from the dataset?**
 170
 171 A18: No, there is no information that allows the identification of individuals.

172 Q19: **Does the dataset contain data that might be considered sensitive in any way (e.g., data**
173 **that reveals racial or ethnic origins, sexual orientations, religious beliefs, political**
174 **opinions or union memberships, or locations; financial or health data; biometric or**
175 **genetic data; forms of government identification, such as social security numbers;**
176 **criminal history)?**
177 Our synthetically generated expert decisions do not relate to characteristics of specific
178 individuals.

179 Collection Process

180 Q21: **How was the data associated with each instance acquired?**

181 A21: Our training dataset was obtained by associating an entity's decision with a given instance
182 of the BAF dataset, respecting the generated work capacity constraints for each entity (See
183 Section 3.3 of the paper). Our expert predictions were obtained by utilizing our synthetic
184 expert decision generation method. The ML model's decision is represented by its score
185 relating to the positive class (rejecting the application).

186 Q22: **What mechanisms or procedures were used to collect the data (e.g., hardware apparatus**
187 **or sensor, manual human curation, software program, software API)?**

188 A22: The details of our synthetic expert generation method are available in Section 3.1 of the paper.
189 These are generated by applying noise to the BAF dataset's label, with an instance-dependent
190 noise approach.

191 The L2D training dataset was generated according to the previous answer.

192 Q23: **If the dataset is a sample from a larger set, what was the sampling strategy (e.g.,**
193 **deterministic, probabilistic with specific sampling probabilities)?**

194 A23: The training dataset corresponds to instances from months 4 to 7 in the BAF dataset. This
195 split is deterministic, and is done across months due to the temporal nature of the dataset.

196 Expert predictions are available for every instance in this dataset.

197 Q24: **Who was involved in the data collection process (e.g., students, crowdworkers, contrac-**
198 **tors) and how were they compensated (e.g., how much were crowdworkers paid)?**

199 A24: There was no data collection process. All the new data provided was synthetically generated.

200 Q25: **Over what timeframe was the data collected?**

201 A25: The generation of our synthetic expert decisions and generation of training scenarios is not
202 time dependant.

203 Q26: **Were any ethical review processes conducted (e.g., by an institutional review board)?**

204 A26: No.

205 Q27: **Does the dataset relate to people?**

206 A27: This dataset relates to synthetic expert decisions and work capacity constraints.

207 Q28: **Did you collect the data from the individuals in question directly, or obtain it via third**
208 **parties or other sources (e.g., websites)?**

209 A28: Not applicable.

210 Q29: **Were the individuals in question notified about the data collection?**

211 A29: Not applicable..

212 Q30: **Did the individuals in question consent to the collection and use of their data?**

213 A30: Not applicable..

214 Q31: **If consent was obtained, were the consenting individuals provided with a mechanism**
215 **to revoke their consent in the future or for certain uses?**

216 A31: Not applicable..

217 Q32: **Has an analysis of the potential impact of the dataset and its use on data subjects (e.g.,**
218 **a data protection impact analysis) been conducted?**

219 A32: No. The datasets are synthetic and should not be used to train fraud detection models or
220 human-AI collaboration systems to be used in real-world fraud applications. The use of
221 these datasets should be self-contained for L2D experimentation. Our synthetic expert
222 predictions should not be a replacement for real human behaviour data.

223 **Preprocessing/cleaning/labeling**

224 Q33: **Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucket-**
225 **ing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances,**
226 **processing of missing values)?**

227 A33: No.

228 Q34: **Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to**
229 **support unanticipated future uses)?**

230 A34: Not applicable.

231 Q35: **Is the software used to preprocess/clean/label the instances available?**

232 A35: Not applicable.

233 **Uses**

234 Q36: **Has the dataset been used for any tasks already?**

235 A36: Our generated synthetic experts and training dataset have only been used for the experiments
236 detailed in the paper.

237 Q37: **Is there a repository that links to any or all papers or systems that use the dataset?**

238 A37: There are still no applications of the presented datasets. We intend to keep track of its uses
239 in the project GitHub repo ¹.

240 Q38: **What (other) tasks could the dataset be used for?**

241 A38: These datasets should be used for the context of benchmarking L2D methods.

242 Q39: **Is there anything about the composition of the dataset or the way it was collected and**
243 **preprocessed/cleaned/labeled that might impact future uses?**

244 A39: Assuming the dataset is used exclusively for the evaluation or development of L2D algo-
245 rithms, the composition of the dataset should not impact future uses.

246 Q40: **Are there tasks for which the dataset should not be used?**

247 A40: Using models trained in these datasets for real-world bank account opening fraud detection
248 (or any other related application) directly should be avoided. The same applies to assignment
249 systems trained in this dataset. The patterns and behaviours observed in these applications
250 are highly dynamic and context-dependant, and using these models can result in unexpected
251 low performances and biased decisions.

¹<https://anonymous.4open.science/r/openl2d-7BD3>

Distribution

Q41: Will the dataset be distributed to third parties outside of the entity (e.g., company, institution, organization) on behalf of which the dataset was created?

A41: Yes, the training dataset, expert prediction data, and the BAF base variant are all publicly accessible.

Q42: How will the dataset will be distributed (e.g., tarball on website, API, GitHub)?

A42: It will be distributed on GitHub.

Q43: When will the dataset be distributed?

A43: The suite is publicly available as of today on GitHub. It will be updated and code will be organized over the following days. There are no plans in removing the datasets from public usage.

Q44: Will the dataset be distributed under a copyright or other intellectual property (IP) license, and/or under applicable terms of use (ToU)?

A44: The suite is licensed under the Creative Commons CC BY-NC-ND 4.0 license.

Q45: Have any third parties imposed IP-based or other restrictions on the data associated with the instances?

A45: No.

Q46: Do any export controls or other regulatory restrictions apply to the dataset or to individual instances?

A46: No.

Maintenance

Q47: Who is supporting/hosting/maintaining the dataset?

A47: The dataset is supported and maintained by currently anonymous personnel.

Q48: How can the owner/curator/manager of the dataset be contacted (e.g., email address)?

A48: The authors are currently anonymous.

Q49: Is there an erratum?

A49: No, there is no erratum as of yet. If necessary in the future, an erratum will be developed for the suite, as well as for this document.

Q50: Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete instances)?

A50: There are no current plans on updating the current version of the datasets. This can change in the future, to correct any undetected bug in the generated datasets.

Q51: If the dataset relates to people, are there applicable limits on the retention of the data associated with the instances (e.g., were individuals in question told that their data would be retained for a fixed period of time and then deleted)?

A51: There are no applicable retention limits of the data.

Q52: Will older versions of the dataset continue to be supported/hosted/maintained?

A52: Currently, there is only the initial version. If any updates are published, previous versions will be available.

Q53: If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so?

A53: There are no current mechanisms to contribute to the suite of datasets. Novel ideas and variants of the dataset should be submitted via email to the authors or as an issue on GitHub.

295 **Author Statement**

296 The authors confirm the data in the BAF suite is under the Creative Commons CC BY-NC-ND 4.0
297 license. The authors bear responsibility in case of violation of copyrights.